

Insurance, biases, discrimination & fairness

Arthur Charpentier

2024

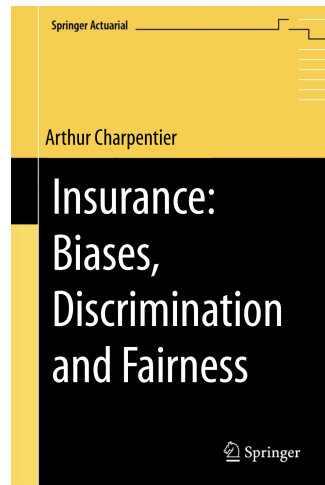


Reference book

Insurance, Biases, Discrimination and Fairness

ISBN : 978-3-031-49782-7

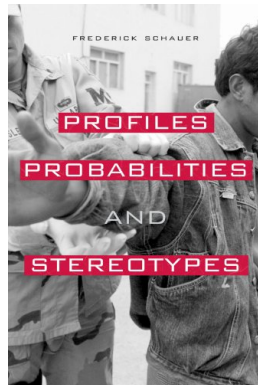
Pitch: **Discrimination** and **fairness** of **predictive models**, in **insurance**, in the context of **data enrichment** ("big data") and **opaque models** ("machine learning", not to say "artificial intelligence").



Definition 1.1: Actuaries, Schauer (2006)

To be an **actuary** is to be a specialist in generalization, and actuaries engage in a form of decision making that is sometimes called actuarial. Actuaries guide insurance companies in making decisions about large categories that have the effect of attributing to the entire category certain characteristics that are probabilistically indicated by membership in the **category**, but that still may not be possessed by a particular member of the category.

See **Barry and Charpentier (2020)** on personalization of insurance prices.



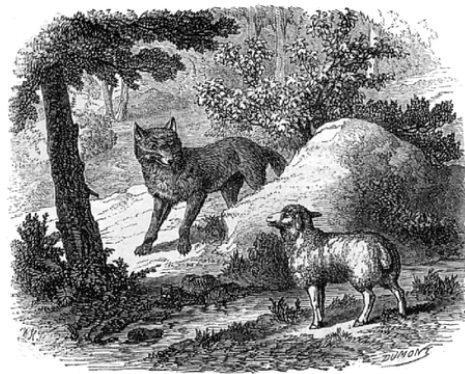
Preliminaries

...

- *Tu la troubles, reprit cette bête cruelle,
Et je sais que de moi tu médis l'an passé.*
- *Comment l'aurais-je fait si je n'étais pas né ?*
- Reprit l'Agneau, je tette encor ma mère.*
- *Si ce n'est toi, c'est donc ton frère.*
- *Je n'en ai point.*
- *C'est donc quelqu'un des tiens.*

...

de La Fontaine (1668), *Le Loup et l'Agneau*.



Definition 1.2: Discrimination, [Merriam-Webster \(2022\)](#)

Discrimination is the act, practice, or an instance of separating or distinguishing categorically rather than individually.

Definition 1.3: Prejudice, [Merriam-Webster \(2022\)](#)

Prejudice is (1) preconceived judgment or opinion, or an adverse opinion or leaning formed without just grounds or before sufficient knowledge; (2) an instance of such judgment or opinion; (3) an irrational attitude of hostility directed against an individual, a group, a race, or their supposed characteristics.

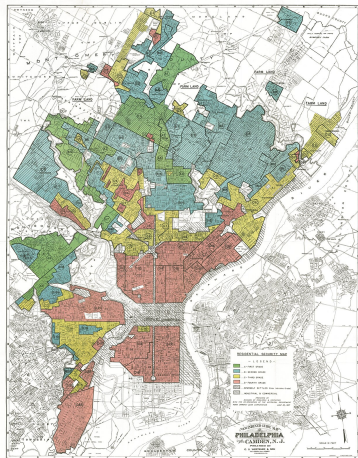
Definition 1.4: Disparate treatment, [Merriam-Webster \(2022\)](#)

Disparate treatment corresponds to the treatment of an individual (as an employee or prospective juror) that is less favorable than treatment of others for discriminatory reasons (as race, religion, national origin, sex, or disability).

Definition 1.5: Disparate impact, [Merriam-Webster \(2022\)](#)

Disparate impact corresponds to an unnecessary discriminatory effect on a protected class caused by a practice or policy (as in employment or housing) that appears to be nondiscriminatory.

Motivation (1. Redlining)



1937 HOLC (Home Owners' Loan Corporation) "residential security" map of Philadelphia

RESIDENTIAL SECURITY MAP

- L E G E N D -

- A - FIRST GRADE
- B - SECOND GRADE
- C - THIRD GRADE
- D - FOURTH GRADE
- SPARSELY SETTLED (Color Indicates Grade)
- INDUSTRIAL & COMMERCIAL

PREPARED BY
 DIVISION OF RESEARCH & STATISTICS
 WITH THE CO-OPERATION OF THE APPRAISAL DEPARTMENT
 HOME OWNERS' LOAN CORPORATION JUNE 25, 1937

MS FORM-6
2-3-37

AREA DESCRIPTION
(For Instructions see Reverse Side)

1. NAME OF CITY Philadelphia, Pa. SECURITY GRADE C AREA NO. 6

2. DESCRIPTION OF TERRAIN. Level

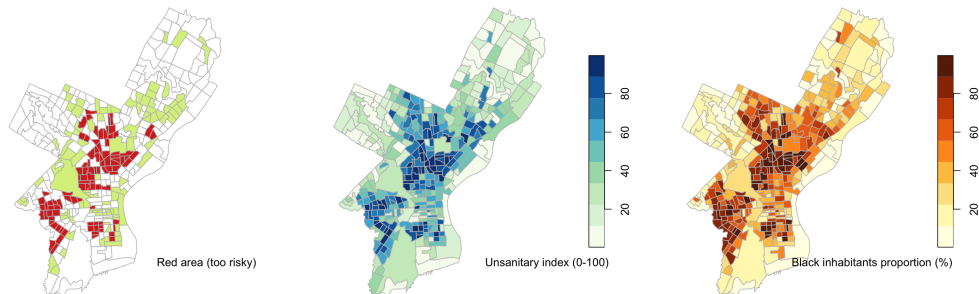
3. FAVORABLE INFLUENCES. Good transportation, particularly in eastern part, -Near to industrial plants of major consequence to entire Philadelphia area.

4. DETRIMENTAL INFLUENCES. Minimal

5. INHABITANTS:
 a. Type Skilled labor ; b. Estimated annual family income \$1,500 - \$1,800.
 c. Foreign-born nominal ; \$; d. Negro no ; \$;
(Nationality) *(Yes or No)*
 e. Infiltration of no ; f. Relief families moderate ;
 g. Population is increasing ; decreasing ; static.

6. BUILDINGS:
 a. Type or types 2 story rows ; b. Type of construction brick ;
 c. Average age 20 - 40 ; d. Repair fair

Motivation (1. Redlining)



(Fictitious maps, inspired by a Home Owners' Loan Corporation map from 1937)

- ▶ Federal Home Loan Bank Board (FHLBB) "*residential security maps*" (for real-estate investments), [Crossney \(2016\)](#) and [Rhyhart \(2020\)](#)
- ▶ Unsanitary index and proportion of Black inhabitants

Motivation (1. Redlining)

Definition 2.1: Redline, [Merriam-Webster \(2022\)](#)

To [redline](#) is (1) to withhold home-loan funds or insurance from neighborhoods considered poor economic risks; (2) to discriminate against in housing or insurance.

See <https://evolutionofraceandinsurance.org/> for some historical perspective, [Squires and Velez \(1988\)](#), or more recently [Squires \(2003\)](#)

... but still a concern see, e.g., [Li \(1996\)](#) about homosexuals.

Motivation (2. “Gender directive”, 2004/113/EC)

Treaty on European Union (26.10.2012, [C326](#))

– Article 2 –

The Union is founded on the values of respect for human dignity, freedom, democracy, equality, the rule of law and respect for human rights, including the rights of persons belonging to minorities. These values are common to the Member States in a society in which pluralism, non-discrimination, tolerance, justice, solidarity and equality between women and men prevail.

– Article 3 –

(...) It shall combat social exclusion and discrimination, and shall promote social justice and protection, equality between women and men, solidarity between generations and protection of the rights of the child.



Motivation (2. “Gender directive”, 2004/113/EC)

Charter of Fundamental Rights of the European Union (18.12.2000 , [C364](#))

– Article 21 (Non discrimination) –

Any discrimination based on any ground such as sex, race, colour, ethnic or social origin, genetic features, language, religion or belief, political or any other opinion, membership of a national minority, property, birth, disability, age or sexual orientation shall be prohibited.

– Article 23 (Equality between men and women) –

Equality between men and women must be ensured in all areas, including employment, work and pay.

The principle of equality shall not prevent the maintenance or adoption of measures providing for specific advantages in favour of the under-represented sex.

Motivation (2. “Gender directive”, 2004/113/EC)

EU Directive ([2004/113/EC](#)), 2004 version

– Article 5 (Actuarial factors) –

1. Member States shall ensure that in all new contracts concluded after 21 December 2007 at the latest, the use of sex as a factor in the calculation of premiums and benefits for the purposes of insurance and related financial services shall not result in differences in individuals’ premiums and benefits.

2. Notwithstanding paragraph 1, Member States may decide before 21 December 2007 to permit proportionate differences in individuals’ premiums and benefits where the use of sex is a determining factor in the assessment of risk based on relevant and accurate actuarial and statistical data. The Member States concerned shall inform the Commission and ensure that accurate data relevant to the use of sex as a determining actuarial factor are compiled, published and regularly updated.

Motivation (2. “Gender directive”, 2004/113/EC)

- › There was initially (2004) an **opt-out clause** (Article 5(2)).
- › Where gender is a determining factor in the assessment of risk based on relevant and accurate actuarial and statistical data then proportionate differences in individual premiums or benefits are allowed.
- › March 2011, the European Court of Justice issued its judgement into the “Test-Achats case”. The ECJ ruled Article 5(2) was invalid.
- › Insurers were no longer able to use gender as a risk factor when pricing policies, “**unisex pricing**”.

”Machine learning won’t give you anything like gender neutrality ‘for free’ that you didn’t explicitly ask for ”, Kearns and Roth (2019)

Motivation (2. “Gender directive”, 2004/113/EC)

“Ten Oever” judgement (*Gerardus Cornelis Ten Oever v Stichting Bedrijfspensioenfonds voor het Glazenwassers – en Schoonmaakbedrijf*, in April 1993), the Advocate General Van Gerven argued that *“the fact that women generally live longer than men has no significance at all for the life expectancy of a specific individual and it is not acceptable for an individual to be penalized on account of assumptions which are not certain to be true in his specific case,”* as mentioned in [De Baere and Goessens \(2011\)](#).



[Schanze \(2013\)](#) used the term *“injustice by generalization,”* from [Britz \(2008\)](#) (“[Generalisierungsunrecht](#)”)

Motivation (2. “Gender directive”, 2004/113/EC)

The Telegraph News Sport Money Business Opinion

Men are still charged more than women for car insurance, despite EU rule change

Car insurers are dodging European equality laws by making gender judgements based on people's jobs, an economist has found

By Kate Palmer
10 April 2015 • 12:33pm



Insurers will price by occupation, and female-dominated jobs tend to attract cheaper premiums | CREDIT: Photo: Rex Features

CAR COSTS: Insurance according to job

Job	Proportion of men	Approximate average premium for a Fiat 500 driver
Dental Nurse	Less than 1pc male	£840
Solicitor	59pc male	£848
Sports and leisure assistants	56pc male	£880
Civil engineer	92pc male	£910
Social worker	21pc male	£920
Plasterer	98pc male	£950

McDonald, 'Indirect Gender Discrimination' (2015); ONS occupation data (2008)

(data source: [Mcdonald \(2015\)](#))

Motivation (3. Québec)

Au Québec, Charte des droits et libertés de la personne (C-12)

– Article 10 –

Toute personne a droit à la reconnaissance et à l'exercice, en pleine égalité, des droits et libertés de la personne, sans distinction, exclusion ou préférence fondée sur la race, la couleur, le sexe, l'identité ou l'expression de genre, la grossesse, l'orientation sexuelle, l'état civil, l'âge sauf dans la mesure prévue par la loi, la religion, les convictions politiques, la langue, l'origine ethnique ou nationale, la condition sociale, le handicap ou l'utilisation d'un moyen pour pallier ce handicap.

Il y a **discrimination** lorsqu'une telle distinction, exclusion ou préférence a pour effet de détruire ou de compromettre ce droit.



Motivation (3. Québec)

Au Québec, Charte des droits et libertés de la personne (C-12)

– Article 20.1 –

Dans un **contrat d'assurance** ou de rente, un régime d'avantages sociaux, de retraite, de rentes ou d'assurance ou un régime universel de rentes ou d'assurance, une distinction, exclusion ou préférence fondée sur l'âge, le sexe ou l'état civil est **réputée non discriminatoire lorsque son utilisation est légitime et que le motif qui la fonde constitue un facteur de détermination de risque, basé sur des données actuarielles.**



Motivation (4. Colorado)

Andrus et al. (2021), "*What we can't measure, we can't understand*"



First Regular Session | 74th General Assembly

Colorado General Assembly

September 27, 2023, the Colorado Division of Insurance exposed a new proposed regulation entitled **Concerning Quantitative Testing of External Consumer Data and Information Sources, Algorithms, and Predictive Models Used for Life Insurance Underwriting for Unfairly Discriminatory Outcomes**



Motivation (4. Colorado)

– Section 4 (Definitions) –

Bayesian Improved First Name Surname Geocoding, or “BIFSG” means, for the purposes of this regulation, the statistical methodology developed by the RAND corporation for estimating race and ethnicity.

External Consumer Data and Information Source, or “ECDIS” means, for the purposes of this regulation, a data source or an information source that is used by a life insurer to supplement or supplant traditional underwriting factors. This term includes credit scores, credit history, social media habits, purchasing habits, home ownership, educational attainment, licensures, civil judgments, court records, occupation that does not have a direct relationship to mortality, morbidity or longevity risk, consumer-generated Internet of Things data, biometric data, and any insurance risk scores derived by the insurer or third-party from the above listed or similar data and/or information source.

Motivation (4. Colorado)

– Section 5 (Estimating Race and Ethnicity) –

Insurers shall estimate the race or ethnicity of all proposed insureds that have applied for coverage on or after the insurer's initial adoption of the use of ECDIS, or algorithms and predictive models that use ECDIS, including a third party acting on behalf of the insurer that used ECDIS, or algorithms and predictive models that used ECDIS, in the underwriting decision-making process, by utilizing:

1. BIFSG and the insureds' or proposed insureds' name and geolocation (information included in the applications) for life insurance shall be used to estimate the race and ethnicity of each insured or proposed insured.
2. For the purposes of BIFSG, the following racial and ethnic categories shall be used: Hispanic, Black, Asian Pacific Islander (API), and White.

Motivation (4. Colorado)

– Section 6 (Application Approval Decision Testing Requirements) –

Using the BIFSG estimated race and ethnicity of proposed insureds and the following methodology, insurers shall calculate whether Hispanic, Black, and API proposed insureds are disapproved at a statistically significant different rate relative to White applicants for whom the insurer, or a third party acting on behalf of the insurer, used ECDIS, or an algorithm or predictive model that used ECDIS, in the underwriting decision-making process.

1. Logistic regression shall be used to model the binary underwriting outcome of either approved or denied.
2. The following factors may be accounted for as control variables in the regression model: policy type, face amount, age, gender, and tobacco use.
3. The estimated race or ethnicity of the proposed insureds shall be accounted for by including Hispanic, Black, and Asian Pacific Islander (API) as separate dummy variables in the regression model.

Motivation (4. Colorado)

4. Determine if there is a statistically significant difference in approval rates for each BIFSG estimated race or ethnicity variable as indicated by a p -value of less than .05.

a. If there is not a statistically significant difference in approval rates, no further testing is required.

b. If there is a statistically significant difference in approval rates, the insurer shall determine whether the difference in approval rates is five (5) percentage points or greater as indicated by the marginal effects value of each BIFSG estimated race or ethnicity variable. (...)

Motivation (4. Colorado)

– Section 7 (Premium Rate Testing Requirements) –

Using the insureds' BIFSG estimated race and ethnicity, insurers shall determine if there is a statistically significant difference in the premium rate per \$1,000 of face amount for policies issued to Hispanic, Black, and API insureds relative to White insureds for whom the insurer, or a third party acting on behalf of the insurer, used ECDIS, or an algorithm or predictive model that used ECDIS, in the underwriting decision-making process.

1. Linear regression shall be used to model the continuous numerical outcome of premium rate per \$1,000 of face amount.
2. The following factors may be accounted for as control variables in the regression model: policy type, face amount, age, gender, and tobacco use.
3. The estimated race or ethnicity of the proposed insureds shall be accounted for by including Hispanic, Black, and Asian Pacific Islander (API) as separate dummy variables in the regression model.

Motivation (4. Colorado)

4. Determine if there is a statistically significant difference in the premium rate per \$1,000 of face amount for each BIFSG estimated race or ethnicity variable as indicated by a p-value of less than .05.


a. If there is not a statistically significant difference in premium rate per \$1,000 of face amount, no further testing is required.

b. If there is a statistically significant difference in premium rate per \$1,000 of face amount, determine whether the premium rate per \$1,000 of face amount is at least 5% more than the average premium rate per \$1,000 for all policies.

i. If the difference in premium rate per \$1,000 of face amount is less than 5%, no further testing is required.

ii. If the difference in premium rate per \$1,000 of face amount is 5% or greater, further testing is required as described in Section 8.

Motivation (4. Colorado)

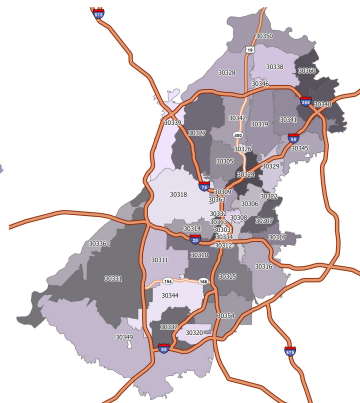
In Elliott et al. (2009), BIFSG¹, `library(eiCompare)`. , consider 12 people living near Atlanta, GA (Fulton & Gwinnett counties), and `eiCompare::wru_predict_race_wrapper`

	last	first	county	city	zipcode	whi	bla	his	asi
1	LOCKLER	GABRIELLA	Fulton	Atlanta	30318	0	0	0	0
2	RADLEY	OLIVIA	Fulton	Fairburn	30213	14	83	1	0
3	BOORSE	KEISHA	Fulton	Atlanta	30331	97	0	3	0
4	MAZ	SAVANNAH	Gwinnett	Norcross	30093	5	6	76	13
5	GAULE	NATASHIA	Gwinnett	Snellville	30078	67	19	14	0
6	MCMELLEN	ISMAEL	Gwinnett	Lilburn	30047	73	15	6	3
7	RIDEOUT	LUQMAN	Gwinnett	Snellville	30078	77	18	2	0
8	WASHINGTON	BRYN	Gwinnett	Norcross	30093	0	95	3	0
9	KULENOVIC	EVELYN	Gwinnett	Buford	30518	100	0	0	0
10	HERNANDEZ	SAMANTHA	Gwinnett	Duluth	30096	3	1	94	1
11	LONG	BESSIE	Gwinnett	Duluth	30096	53	39	1	1
12	HE	JOSE	Gwinnett	Lawrenceville	30045	2	3	4	89

¹Bayesian Improved First Name Surname Geocoding

Motivation (3. Colorado)

We have 12 people,
in two counties near Atlanta
(about 10 zip-codes)



Motivation (3. Colorado)

- Use `eiCompare::wru_predict_race_wrapper` on a revised dataset with the same name “Savannah Maz”

```
1      last      first      county      city      zipcode      whi      bla      his      asi
2  1      MAZ      SAVANNAH      Fulton      Atlanta      30318      0      0      0      100
3  2      MAZ      SAVANNAH      Fulton      Fairburn      30213      13      61      22      3
4  3      MAZ      SAVANNAH      Fulton      Atlanta      30331      3      77      19      1
5  4      MAZ      SAVANNAH      Gwinnett      Norcross      30093      5      6      76      13
6  5      MAZ      SAVANNAH      Gwinnett      Snellville      30078      13      18      69      0
7  6      MAZ      SAVANNAH      Gwinnett      Lilburn      30047      28      22      34      16
8  7      MAZ      SAVANNAH      Gwinnett      Snellville      30078      53      3      40      3
9  8      MAZ      SAVANNAH      Gwinnett      Norcross      30093      5      6      76      13
10 9      MAZ      SAVANNAH      Gwinnett      Buford      30518      79      4      14      2
11 10     MAZ      SAVANNAH      Gwinnett      Duluth      30096      32      8      38      22
12 11     MAZ      SAVANNAH      Gwinnett      Duluth      30096      55      19      22      5
13 12     MAZ      SAVANNAH      Gwinnett      Lawrenceville      30045      15      19      62      4
```

Motivation (3. Colorado)

- Use `eiCompare::wru_predict_race_wrapper` on a revised dataset with the same name “Bryn Washington”

```
1      last first   county      city zipcode whi bla his asi
2 1 WASHINGTON BRYN   Fulton    Atlanta 30318  0  0  0 100
3 2 WASHINGTON BRYN   Fulton    Fairburn 30213  0 99  0  0
4 3 WASHINGTON BRYN   Fulton    Atlanta 30331  0 99  0  0
5 4 WASHINGTON BRYN Gwinnett   Norcross 30093  0 95  3  0
6 5 WASHINGTON BRYN Gwinnett   Snellville 30078  0 96  1  0
7 6 WASHINGTON BRYN Gwinnett   Lilburn 30047  1 98  0  0
8 7 WASHINGTON BRYN Gwinnett   Snellville 30078  6 87  2  0
9 8 WASHINGTON BRYN Gwinnett   Norcross 30093  0 95  3  0
10 9 WASHINGTON BRYN Gwinnett   Buford 30518  7 92  1  0
11 10 WASHINGTON BRYN Gwinnett   Duluth 30096  2 96  1  0
12 11 WASHINGTON BRYN Gwinnett   Duluth 30096  1 96  0  0
13 12 WASHINGTON BRYN Gwinnett Lawrenceville 30045  0 98  1  0
```

Motivation (3. Colorado)

- Use `eiCompare::wru_predict_race_wrapper` on a revised dataset with the same name “Samantha Hernandez”

	last	first	county	city	zipcode	whi	bla	his	asi
1	HERNANDEZ	SAMANTHA	Fulton	Atlanta	30318	0	0	0	100
2	HERNANDEZ	SAMANTHA	Fulton	Fairburn	30213	2	12	85	0
3	HERNANDEZ	SAMANTHA	Fulton	Atlanta	30331	0	16	81	0
4	HERNANDEZ	SAMANTHA	Gwinnett	Norcross	30093	0	0	99	0
5	HERNANDEZ	SAMANTHA	Gwinnett	Snellville	30078	1	1	97	0
6	HERNANDEZ	SAMANTHA	Gwinnett	Lilburn	30047	3	3	92	1
7	HERNANDEZ	SAMANTHA	Gwinnett	Snellville	30078	5	0	94	0
8	HERNANDEZ	SAMANTHA	Gwinnett	Norcross	30093	0	0	99	0
9	HERNANDEZ	SAMANTHA	Gwinnett	Buford	30518	17	1	81	0
10	HERNANDEZ	SAMANTHA	Gwinnett	Duluth	30096	3	1	94	1
11	HERNANDEZ	SAMANTHA	Gwinnett	Duluth	30096	8	4	86	0
12	HERNANDEZ	SAMANTHA	Gwinnett	Lawrenceville	30045	1	2	97	0

Motivation (3. Colorado)

- Use `eiCompare::wru_predict_race_wrapper` on a revised dataset with the same name “Jose He”

	last	first	county	city	zipcode	whi	bla	his	asi
1	HE	JOSE	Fulton	Atlanta	30318	0	0	0	100
2	HE	JOSE	Fulton	Fairburn	30213	2	9	2	84
3	HE	JOSE	Fulton	Atlanta	30331	1	27	3	55
4	HE	JOSE	Gwinnett	Norcross	30093	0	0	2	98
5	HE	JOSE	Gwinnett	Snellville	30078	13	18	30	0
6	HE	JOSE	Gwinnett	Lilburn	30047	1	1	1	97
7	HE	JOSE	Gwinnett	Snellville	30078	8	1	3	86
8	HE	JOSE	Gwinnett	Norcross	30093	0	0	2	98
9	HE	JOSE	Gwinnett	Buford	30518	19	1	2	78
10	HE	JOSE	Gwinnett	Duluth	30096	1	0	0	98
11	HE	JOSE	Gwinnett	Duluth	30096	6	2	1	85
12	HE	JOSE	Gwinnett	Lawrenceville	30045	2	3	4	89

Motivation (5. Motor Insurance in the U.S.)

via [The Zebra \(2022\)](#),

California

Allowed (with applicable limitations): driving experience, marital status, address/zip code

Prohibited (or effectively prohibited): gender, age, credit history, education, occupation, employment status, residential status, insurance history

Notes & Clarifications: California's insurance commissioner banned gender as of January 2019. Occupation and education are permitted for use in group plans (i.e. for alumni associations and other membership programs).

Georgia

Allowed (with applicable limitations): gender, age, years of driving experience, credit history, marital status, residential status, address/zip code, insurance history

Prohibited (or effectively prohibited): occupation, education, and employment status

Notes & Clarifications: none

Hawaii

Allowed (with applicable limitations): address/zip code, insurance history

Prohibited (or effectively prohibited): gender, age, years of driving experience, credit history, education, occupation, employment status, marital status, residential status

Notes & Clarifications: none

Illinois

Allowed (with applicable limitations): gender, age, years of driving experience, credit history, education, occupation, employment status, marital status, residential status, address/zip code, insurance history

Prohibited (or effectively prohibited): none

Notes & Clarifications: none

Massachusetts

Allowed (with applicable limitations): years of driving experience, address/zip code, insurance history

Prohibited (or effectively prohibited): gender, age, credit history, education, occupation, employment status, marital status, residential status

Notes & Clarifications: none

Michigan

Allowed (with applicable limitations): gender (group-rated policies), age, years of driving experience, credit history, education, occupation, employment status, marital status (group-rated policies), residential status, address/zip code, insurance history

Prohibited (or effectively prohibited): gender (non-group policies), marital status (non-group policies)

Notes & Clarifications: Gender and marital status are permitted only in rate-making for group plans (i.e. for alumni associations and other membership programs). **UPDATE:** [Michigan lawmakers approved a major insurance reform bill](#) in May 2019 that will ban insurers in the state from using gender, marital status, address/zipcode, residential status, education and occupation in rate setting. The ban will be enforced starting in July 2020. Insurers will be permitted to use "territory" as approved by the state regulators instead of zip code.

New York

Allowed (with applicable limitations): gender, age, years of driving experience, credit history, marital status, residential status, address/zip code, insurance history

Prohibited (or effectively prohibited): occupation, education, employment status

Notes & Clarifications: none

see also



[Avraham et al. \(2013\)](#)

Motivation (6. Admission in Graduate Program, UC Berkeley)

	Total	Men	Women	Proportions
Total	5233/12763 ~ 41%	3714/8442 ~ 44%	1512/4321 ~ 35%	66%-34%
Top 6	1745/4526 ~ 39%	1198/2691 ~ 45%	557/1835 ~ 30%	59%-41%
A	597/933 ~ 64%	512/825 ~ 62%	89/108 ~ 82%	88%-12%
B	369/585 ~ 63%	353/560 ~ 63%	17/ 25 ~ 68%	96%- 4%
C	321/918 ~ 35%	120/325 ~ 37%	202/593 ~ 34%	35%-65%
D	269/792 ~ 34%	138/417 ~ 33%	131/375 ~ 35%	53%-47%
E	146/584 ~ 25%	53/191 ~ 28%	94/393 ~ 24%	33%-67%
F	43/714 ~ 6%	22/373 ~ 6%	24/341 ~ 7%	52%-48%

Data from [Bickel et al. \(1975\)](#) (discussed as an illustration of "[Simpson's paradox](#)")

Formalize the later, S is the (binary) genre, Y the admission and X the program (category),

Motivation (6. Admission in Graduate Program, UC Berkeley)

$$\begin{aligned} & \mathbb{P}[Y = \text{yes} \mid S = \text{men}] \geq \mathbb{P}[Y = \text{yes} \mid S = \text{women}] \\ & \quad \text{overall admission} \\ & \mathbb{P}[Y = \text{yes} \mid X = x, S = \text{men}] \leq \mathbb{P}[Y = \text{yes} \mid X = x, S = \text{women}], \forall x. \\ & \quad \text{conditional on program} \end{aligned}$$

“the bias in the aggregated data stems not from any pattern of discrimination on the part of admissions committees, which seems quite fair on the whole, but apparently from prior screening at earlier levels of the educational system. Women are shunted by their socialization and education toward fields of graduate study that are generally more crowded, less productive of completed degrees, and less well funded, and that frequently offer poorer professional employment prospects,” Bickel et al. (1975)

Motivation (6'. Admission in hospitals)

Consider the following mortality rates in two hospitals (fake data)

	Total	Healthy	Pre-condition	Proportions
Hospital A	800/1000 = 80%	590/600 ~ 98%	210/400 ~ 53%	60%-40%
Hospital B	900/1000 = 90%	870/900 ~ 97%	30/100 ~ 30%	90%-10%

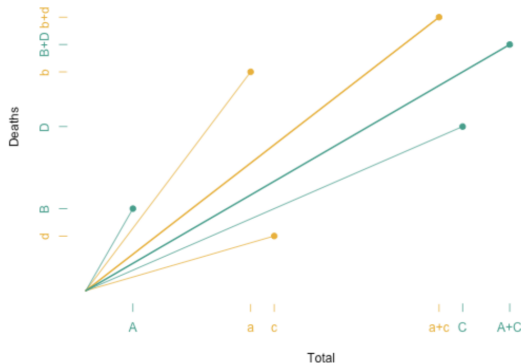
There is no mathematical "paradox", *per se*.

We could have

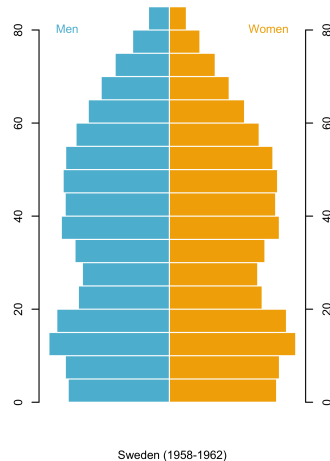
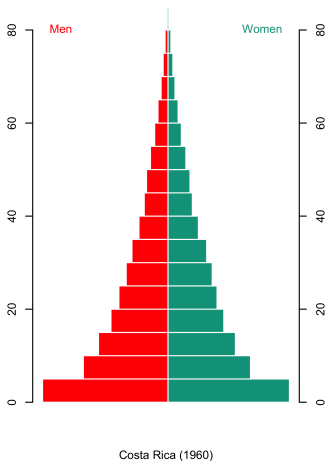
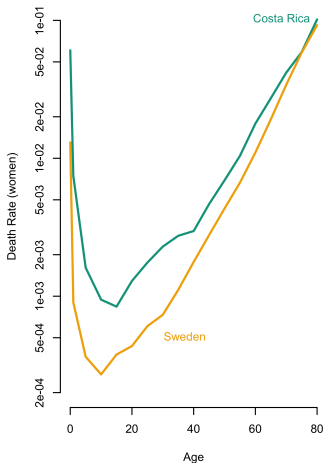
$$\frac{A}{B} \geq \frac{a}{b} \text{ and } \frac{C}{D} \geq \frac{c}{d}$$

and at the same time

$$\frac{A+C}{B+D} \leq \frac{a+c}{b+d}$$



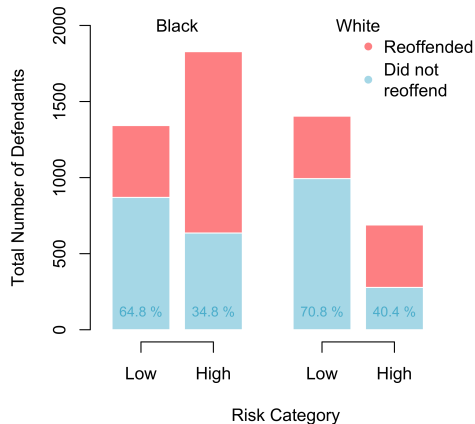
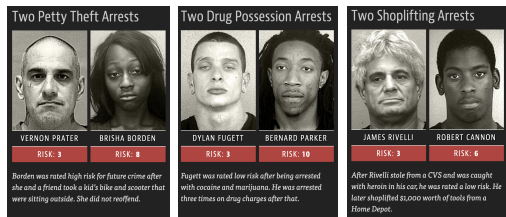
Motivation (6". Mortality in Costa Rica and Sweden)



Overall mortality rate for women, **8.12‰** in Costa Rica, against **9.29‰** in Sweden.

Motivation (7. Propublica, Actuarial Justice)

- Concept of "actuarial justice" as coined in **Feeley and Simon (1994)**
- Correctional **Offender Management Profiling for Alternative Sanctions (COMPAS)**, **Perry (2013)**

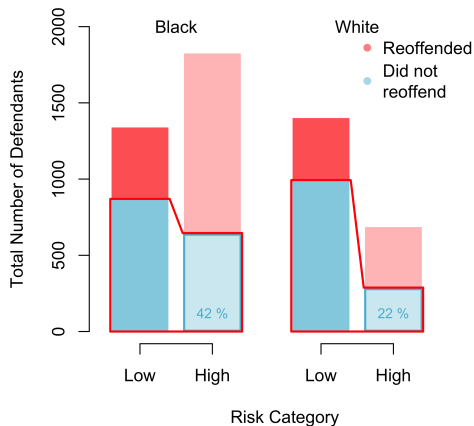


<https://github.com/propublica/compas-analysis>

- **Angwin et al. (2016)** Machine Bias
- Dressel and Farid (2018)**

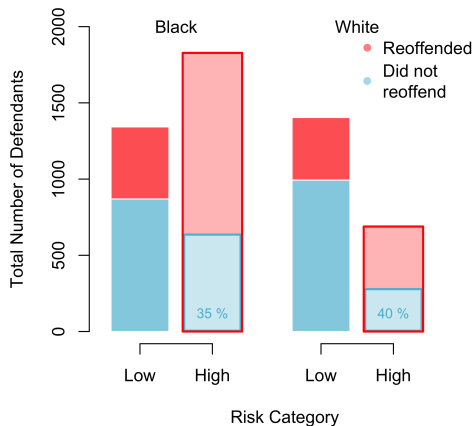
Motivation (7. Propublica, Actuarial Justice)

- From Feller et al. (2016),
 - ▶ for White people, among those who did not re-offend, 22% were wrongly classified,
 - ▶ for Black people, among those who did not re-offend, 42% were wrongly classified,
 - ▶ problem, since $42\% \gg 22\%$



Motivation (7. Propublica, Actuarial Justice)

- From Dieterich et al. (2016),
 - ▶ for White people, among those who were classified as high risk, 40% did not re-offend,
 - ▶ for Black people, among those who were classified as high risk, 35% did not re-offend,
 - ▶ no problem, since $40\% \approx 35\%$



Motivation (7. Propublica, Actuarial Justice)

Formalize the later,

$\left\{ \begin{array}{l} S : \text{race (binary), black \& white} \\ Y : \text{re-offense (binary), no \& yes} \\ \hat{Y} : \text{classifier (risk category), low \& high} \end{array} \right.$

$\mathbb{P}[\hat{Y} = \text{high} | Y = \text{no}, S = \text{black}] = 42\% \stackrel{?}{=} \mathbb{P}[\hat{Y} = \text{high} | Y = \text{no}, S = \text{white}] = 22\%$

sensitive (blue arrow from 'black' to 'white')

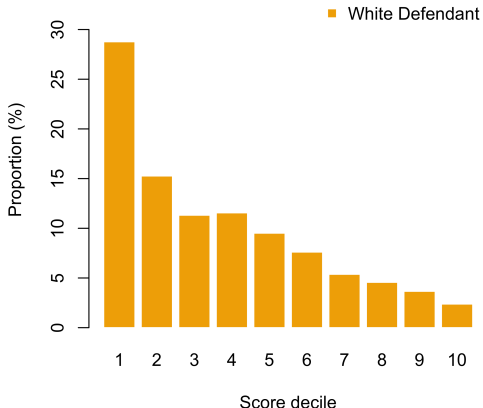
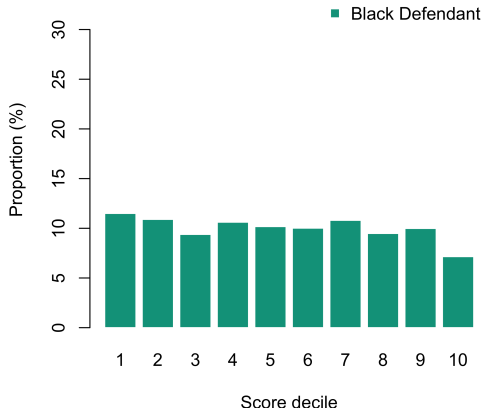
sensitive (orange arrow from 'black' to 'white')


↑ false positive rate

$\mathbb{P}[Y = \text{no} | \hat{Y} = \text{high}, S = \text{black}] = 35\% \stackrel{?}{=} \mathbb{P}[Y = \text{no} | \hat{Y} = \text{high}, S = \text{white}] = 40\%$

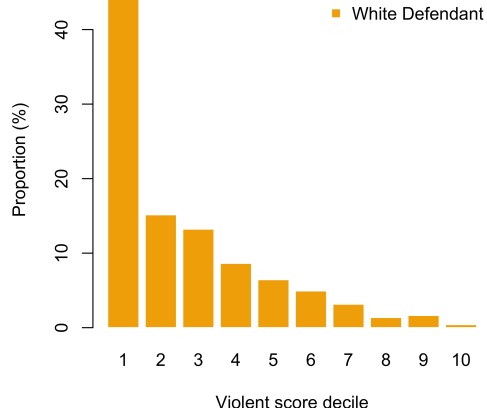
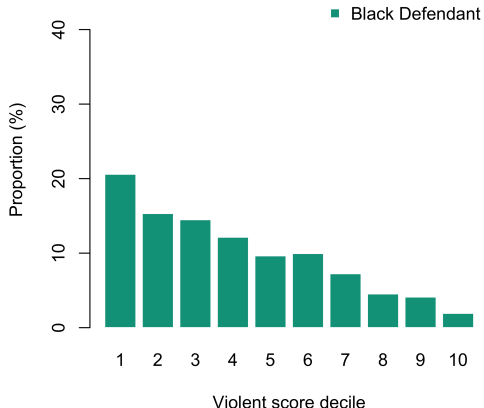
↑ false discovery rate

Motivation (7. Propublica, Actuarial Justice)



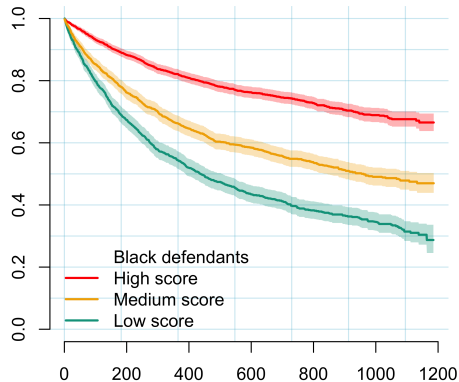
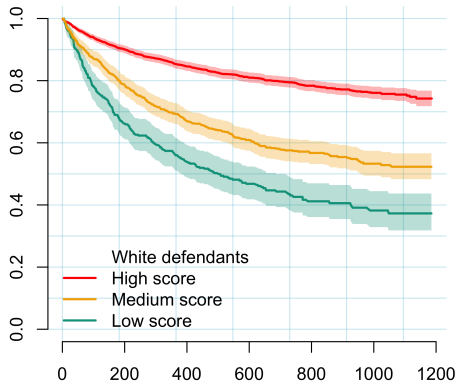
Look at score distributions, **black** and **white** defendant, [Larson et al. \(2016\)](#) .


Motivation (7. Propublica, Actuarial Justice)



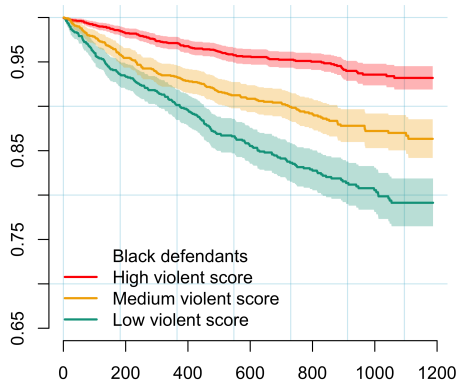
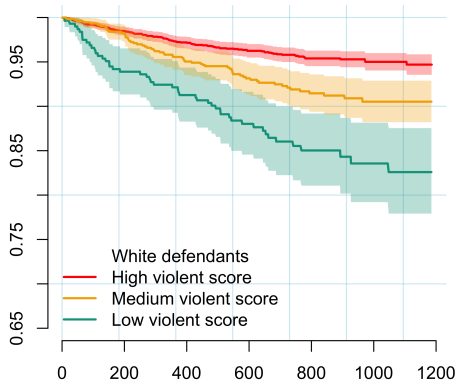
Look at score distributions, **black** and **white** defendant, [Larson et al. \(2016\)](#) .


Motivation (7. Propublica, Actuarial Justice)



Cox Proportional Hazards model, **black** and **white** defendant, [Larson et al. \(2016\)](#) 

Motivation (7. Propublica, Actuarial Justice)



Cox Proportional Hazards model, black and white defendant, Larson et al. (2016) 

Motivation (8. Intention)

En France, Loi n° 2008-496 du 27 mai 2008

– Article 1 –

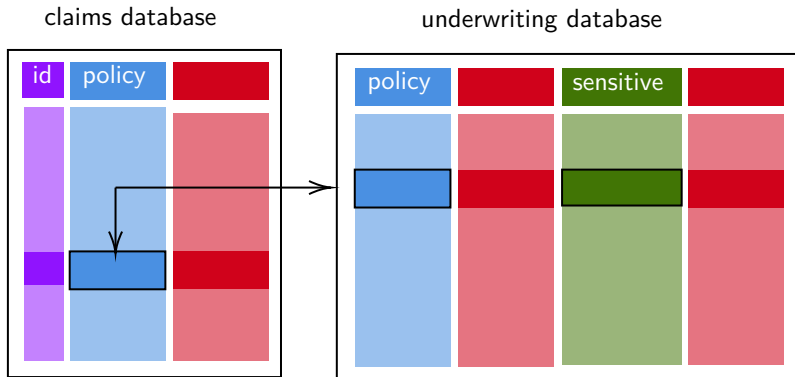
Constitue une **discrimination indirecte** une disposition, un critère ou une pratique neutre en apparence, mais susceptible d'entraîner, pour l'un des motifs mentionnés au premier alinéa, un désavantage particulier pour des personnes par rapport à d'autres personnes, à moins que cette disposition, ce critère ou cette pratique ne soit objectivement justifié par un but légitime et que les moyens pour réaliser ce but ne soient nécessaires et appropriés.

Extention de la "Loi n° 72-546 du 1 juillet 1972", qui supprima l'exigence de l'intention spécifique.

" *Technology is neither good nor bad; nor is it neutral* " , **Kranzberg (1986)**

Motivation (9. Biases, biases everywhere...)

- underwriters biases
 - commercial discounts
 - inferred data
 - multiple decisions
- claims biases
 - fraud detection
 - sexist mechanic
 - ageist manager

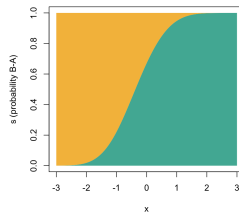
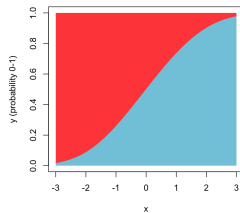
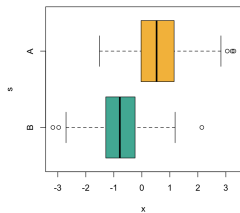
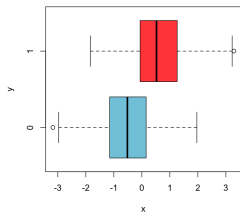


Datasets

> toydata1

Consider a confounding Gaussian variable X_0 , $X_0 \sim \mathcal{N}(0, 1)$, and

$$\begin{cases} X = X_0 + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1/2^2), \\ S = \mathbf{1}(X_0 + \eta > 0), \quad \eta \sim \mathcal{N}(0, 1/2^2), \quad s \in \{A, B\}, \\ Y = \mathbf{1}(X_0 + \nu > 0), \quad \nu \sim \mathcal{N}(0, 1/2^2), \quad y \in \{0, 1\}. \end{cases}$$

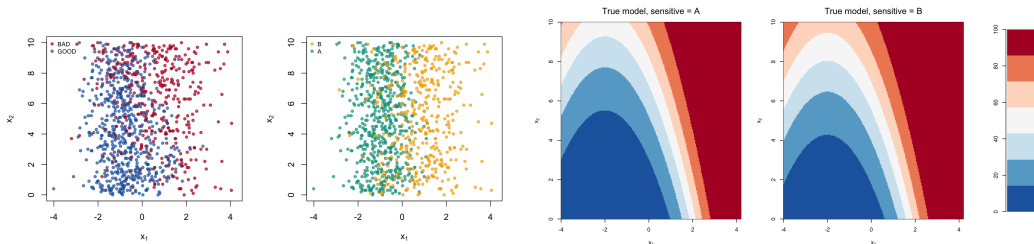


$x \mapsto \mathbb{P}[Y = 0|X = x]$ (left-hand side) and $x \mapsto \mathbb{P}[S = A|X = x]$ (right-hand side)

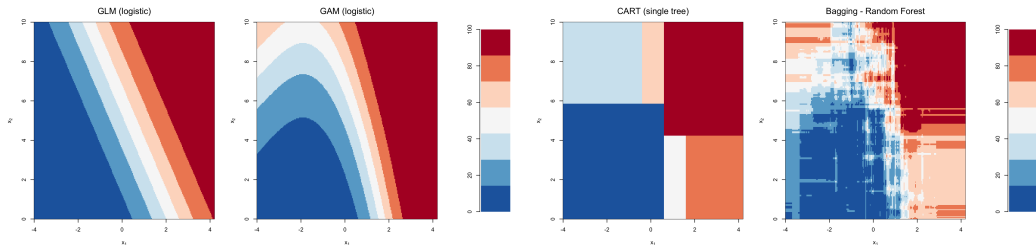
Datasets

toydata2

- ▶ binary sensitive attribute, $s \in \{A, B\}$, (60% and 40%)
- ▶ $(x_1, x_3) \sim \mathcal{N}(\boldsymbol{\mu}_s, \boldsymbol{\Sigma}_s)$, $r_{s=A} = 0.4$ and $r_{s=B} = 0.7$
- ▶ $x_2 \sim \mathcal{U}([0, 10])$, independent of x_1 and x_3
- ▶ $\eta = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 \mathbf{1}_B(s)$, that does not depend on x_3
- ▶ $y \sim \mathcal{B}(p)$ where $p = \exp(\eta) / [1 + \exp(\eta)] = \mu(x_1, x_2, s)$.

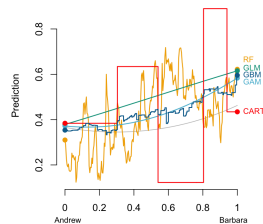
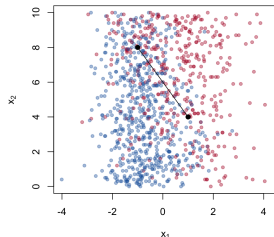


Datasets



Five models are considered

- ▶ plain GLM (logistic)
- ▶ GAM (cubic splines)
- ▶ CART (classification tree)
- ▶ RF (random forest)
- ▶ GBM (gradient boosting)



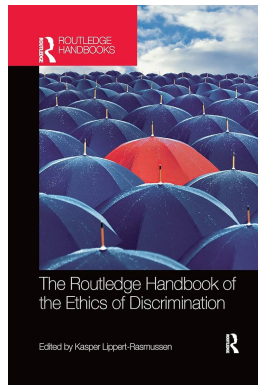
Datasets

- **GermanCredit**, $m = 1,000$
 - ▶ binary sensitive attribute, $s \in \{A, B\}$, (64% and 36%) corresponding to gender
 - ▶ y denotes a default (30%)
 - ▶ x_1, \dots, x_k denote legitimate credit variables (Duration, Purpose, Credit_amount, Age, Housing, Existing_credits, Foreign_worker, Resident_since, etc)
- **FrenchMotor** (policy observe over one year), $n = 12,437$
 - ▶ binary sensitive attribute, $s \in \{A, B\}$, (31% and 69%) corresponding to gender
 - ▶ y denotes the occurrence of a car accident (8.67%, unbalanced data)
 - ▶ x_1, \dots, x_k denote legitimate credit variables (MariStat, VehAge, SocioCateg, DrivAge, VehBody, VehEnergy, VehMaxSpeed, Garage, VehUsage, etc)

– Part 1 –
Insurance

Discrimination and Insurance

"What is unique about insurance is that even statistical discrimination which by definition is absent of any malicious intentions, poses significant moral and legal challenges. Why? Because on the one hand, policy makers would like insurers to treat their insureds equally, without discriminating based on race, gender, age, or other characteristics, even if it makes statistical sense to discriminate (...) On the other hand, at the core of insurance business lies discrimination between risky and non-risky insureds. But riskiness often statistically correlates with the same characteristics policy makers would like to prohibit insurers from taking into account." Avraham (2017)



Discrimination and Insurance

Definition 2.2: Mutuality, Wilkie (1997)

Mutuality is considered as the normal form of commercial private insurance, where participants contribute to the risk pool through a premium that relates to their particular risk at the time of the application, i.e., the higher the risk that they bring to the pool, the higher the premium required.

Definition 2.3: Solidarity, Wilkie (1997)

Solidarity is the basis of most national or social insurance schemes. Participation in such state-run schemes is generally compulsory and individuals have no discretion over their level of cover. All participants normally have the same level of cover. In solidarity schemes the contributions are not based on the expected risk of each participant.

Insurance Pricing and Predictive Modeling

“*Humans think in stories rather than facts, numbers or equations - and the simpler the story, the better,*” Harari (2018). For insurers, it is often a mixture of both.

For Glenn (2000), insurer’s risk selection process has two sides:

- the one presented to regulators and policyholders (numbers, statistics and objectivity),
- the other presented to underwriters (stories, character and subjective judgment).

The rhetoric of insurance exclusion – numbers, objectivity and statistics – forms what Brian Glenn calls “*the myth of the actuary,*” “*a powerful rhetorical situation in which decisions appear to be based on objectively determined criteria when they are also largely based on subjective ones*” or “*the subjective nature of a seemingly objective process*”.

Glenn (2003) claimed that there are many ways to rate accurately. Insurers can rate risks in many different ways depending on the stories they tell on which characteristics

Insurance Pricing and Predictive Modeling

are important and which are not. “*The fact that the selection of risk factors is subjective and contingent upon narratives of risk and responsibility has in the past played a far larger role than whether or not someone with a wood stove is charged higher premiums.*” Going further, “*virtually every aspect of the insurance industry is predicated on stories first and then numbers.*”

“*all models are wrong but some models are useful,*” Box et al. (2011) (in other words, any model is at best a useful fable).

Insurance Pricing and Predictive Modeling

Definition 3.1: Pure premium (homogeneous risks)

Let Y be the non-negative random variable corresponding to the total annual loss associated with a given policy, then the **pure premium** is $\mathbb{E}[Y]$.

Proposition 3.1: Law of Large Numbers (2)

Consider an infinite collection of i.i.d. random variables $Y, Y_1, Y_2, \dots, Y_n, \dots$ in a probabilistic space $(\Omega, \mathcal{F}, \mathbb{P})$, with finite expected value, then

$$\underbrace{\frac{1}{n} \sum_{i=1}^n Y_i}_{\text{(empirical) average}} \xrightarrow{\text{a.s.}} \underbrace{\mathbb{E}(Y)}_{\text{expected value}}, \text{ as } n \rightarrow \infty.$$

Insurance Pricing and Predictive Modeling

More realistically, population is heterogeneous (with respect to risks), with some covariates \mathbf{x} (legitimate, or not).

Definition 3.2: Pure premium (heterogeneous risks)

Let Y be the non-negative random variable corresponding to the total annual loss associated with a given policy, with covariates \mathbf{x} , then the **pure premium** is $\mu(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$.

In this general setting, \mathbf{x} consist in numeric or categorical variables.

Proposition 3.2: Law of Large Numbers (2')

Consider an infinite collection of i.i.d. random pairs (\mathbf{X}, Y) , (\mathbf{X}_1, Y_1) , $(\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_n, Y_n), \dots$ in a probabilistic space $(\Omega, \mathcal{F}, \mathbb{P})$, with finite expected value, then for any $\mathcal{A} \subset \mathcal{X}$ such that $\mathbb{P}[\mathbf{X} \in \mathcal{A}] > 0$,

$$\frac{\sum_{i=1}^n Y_i \mathbf{1}(\mathbf{X}_i \in \mathcal{A})}{\sum_{i=1}^n \mathbf{1}(\mathbf{X}_i \in \mathcal{A})} = \frac{1}{n_{\mathcal{A}}} \underbrace{\sum_{i \in \mathcal{I}_n(\mathcal{A})} Y_i}_{\text{conditional average}} \xrightarrow{\text{a.s.}} \underbrace{\mathbb{E}(Y | \mathbf{X} \in \mathcal{A})}_{\text{conditional expected value}}, \text{ as } n \rightarrow \infty,$$

where $\mathcal{I}_n(\mathcal{A}) = \{i : \mathbf{X}_i \in \mathcal{A}\} \subset \{1, 2, \dots, n\}$ and $n_{\mathcal{A}} = \text{Card}(\mathcal{I}_n(\mathcal{A}))$.

Insurance Pricing and Predictive Modeling

➤ Excerpt from the Men and Women life tables in 1720 (source: [Struyck \(1912\)](#)). Mortality, as a function of the **age** and the **gender** of the individual.



Table des Hommes.

Années	Per-sonnes	Années	Per-sonnes	Années	Per-sonnes	Années	Per-sonnes	Années	Per-sonnes	Années	Per-sonnes
5	710	20	607	35	474	50	313	65	142	80	33
6	697	21	599	36	464	51	301	66	132	81	29
7	688	22	591	37	454	52	289	67	123	82	25
8	681	23	583	38	444	53	277	68	114	83	22
9	675	24	575	39	434	54	265	69	105	84	19
10	670	25	567	40	424	55	253	70	97	85	16
11	665	26	558	41	414	56	241	71	89	86	13
12	660	27	549	42	404	57	229	72	82	87	10
13	654	28	540	43	393	58	217	73	75	88	8
14	648	29	531	44	382	59	206	74	68	89	6
15	642	30	522	45	371	60	195	75	61	90	4
16	635	31	513	46	360	61	184	76	54	91	3
17	628	32	504	47	349	62	173	77	48	92	2
18	621	33	494	48	337	63	162	78	43	93	1
19	614	34	484	49	325	64	152	79	38	94	

Table des femmes.

Années	Per-sonnes	Années	Per-sonnes	Années	Per-sonnes	Années	Per-sonnes	Années	Per-sonnes	Années	Per-sonnes
5	711	20	624	35	508	50	373	65	205	80	55
6	700	21	617	36	500	51	362	66	194	81	47
7	692	22	610	37	492	52	351	67	183	82	40
8	685	23	603	38	484	53	340	68	172	83	34
9	679	24	596	39	476	54	329	69	161	84	29
10	674	25	588	40	468	55	318	70	150	85	24
11	669	26	580	41	459	56	306	71	140	86	20
12	664	27	572	42	450	57	294	72	130	87	17
13	660	28	564	43	441	58	282	73	120	88	14
14	656	29	556	44	432	59	271	74	110	89	11
15	652	30	548	45	423	60	260	75	100	90	8
16	647	31	540	46	414	61	249	76	90	91	6
17	642	32	532	47	404	62	238	77	81	92	4
18	636	33	524	48	394	63	227	78	72	93	2
19	630	34	516	49	384	64	216	79	63	94	1

Insurance Pricing and Predictive Modeling

- Excerpt from the Men and Women life tables in 1720 (source: [Struyck \(1912\)](#))
Mortality, as a function of the **age** and the **gender** of the individual.

men			women								
x	L_x	${}_5p_x$	x	L_x	${}_5p_x$	x	L_x	${}_5p_x$	x	L_x	${}_5p_x$
0	1000	29.0%	45	371	16.6%	0	1000	28.9%	45	423	11.8%
5	710	5.6%	50	313	19.2%	5	711	5.2%	50	373	14.7%
10	670	4.2%	55	253	22.9%	10	674	3.3%	55	318	18.2%
15	642	5.5%	60	195	27.2%	15	652	4.3%	60	260	21.2%
20	607	6.6%	65	142	31.7%	20	624	5.8%	65	205	26.8%
25	567	7.9%	70	97	37.1%	25	588	6.8%	70	150	33.3%
30	522	9.2%	75	61	45.9%	30	548	7.3%	75	100	45.0%
35	474	10.5%	80	33	51.5%	35	508	7.9%	80	55	56.4%
40	424	12.5%	85	16		40	468	9.6%	85	24	

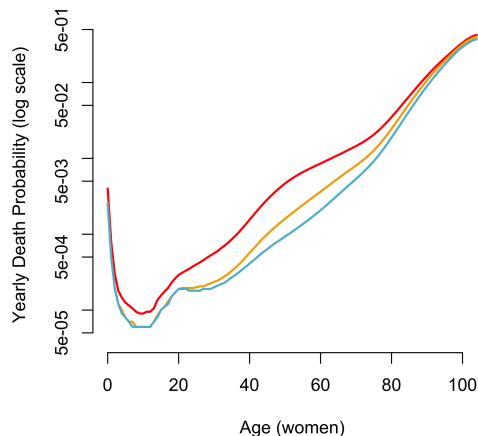
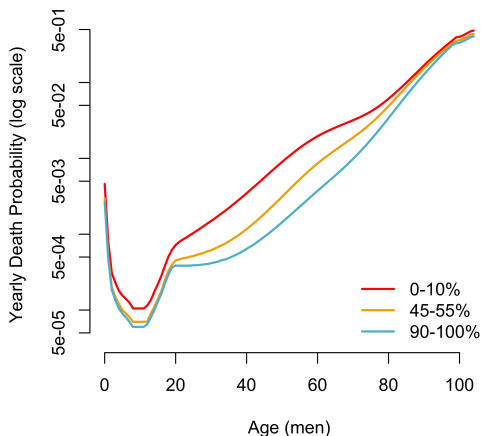
Insurance Pricing and Predictive Modeling

- › Excerpt from the Men and Women life tables in 2016 (source: [Blanpain \(2018\)](#))
Mortality, as a function of the **age**, the **gender** and the **wealth** of the individual.

men			
x	0-5%	45-50%	95-100%
0	100000	100000	100000
10	99299	99566	99619
20	99024	99396	99469
30	97930	98878	99094
40	95595	98058	98627
50	90031	96172	97757
60	77943	91050	95649
70	59824	79805	90399
80	38548	59103	76115
90	13337	23526	38837
100	530	1308	3231

women			
x	0-5%	45-50%	95-100%
0	100000	100000	100000
10	99385	99608	99623
20	99227	99506	99526
30	98814	99302	99340
40	97893	98960	99074
50	95021	97959	98472
60	88786	95543	97192
70	79037	90408	94146
80	63224	79117	85825
90	31190	45750	55918
100	2935	5433	8717

Insurance Pricing and Predictive Modeling



Force of mortality (log scale) for various income quantile, in France, [Blanpain \(2018\)](#).

Insurance Pricing and Predictive Modeling

U.S. DECENNIAL LIFE TABLES FOR 1969-71

Volume I, Number 1



United States Life Tables: 1969-71

1. Life table for the total population: United States, 1969-71----- 6
2. Life table for males: United States, 1969-71----- 8
3. Life table for females: United States, 1969-71----- 10
4. Life table for the white population: United States, 1969-71----- 12
5. Life table for white males: United States, 1969-71----- 14
6. Life table for white females: United States, 1969-71----- 16
7. Life table for the population other than white: United States, 1969-71--- 18
8. Life table for males other than white: United States, 1969-71----- 20
9. Life table for females other than white: United States, 1969-71----- 22
10. Life table for the Negro population: United States, 1969-71----- 24

TABLE 10. LIFE TABLE FOR THE NEGRO POPULATION: UNITED STATES, 1969-71

AGE INTERVAL	PROPORTION DYING	OF 100,000 BORN ALIVE		STATIONARY POPULATION		AVERAGE REMAINING LIFETIME
	PERIOD OF LIFE BETWEEN TWO AGES	PROPORTION OF PERSONS ALIVE AT BEGINNING OF AGE INTERVAL DURING INTERVAL	NUMBER LIVING AT BEGINNING OF AGE INTERVAL	NUMBER DYING DURING AGE INTERVAL	IN THE AGE INTERVAL	IN THIS AND ALL SUBSEQUENT AGE INTERVALS
(1)	(2)	(3)	(4)	(5)	(6)	(7)
x To x + f	q_x	l_x	d_x	l'_x	T_x	e_x
DAYS						
0-1.....	0.01348	100,000	1,348	272	6,431,264	64.11
1-7.....	0.0668	98,652	6,59	1,610	6,715,992	64.99
7-28.....	0.0256	91,993	2,39	5,631	6,409,376	65.61
28-365.....	0.1027	97,744	1,013	89,778	6,405,745	65.52

Mortality, gender and “race”

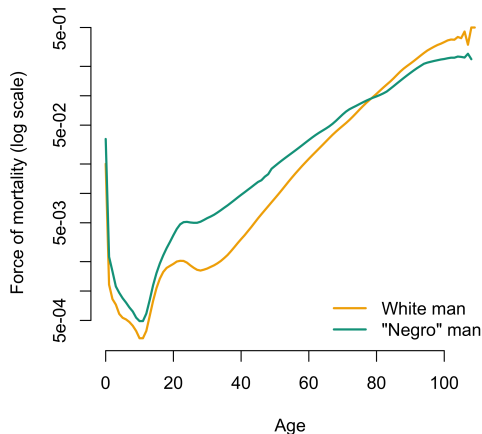
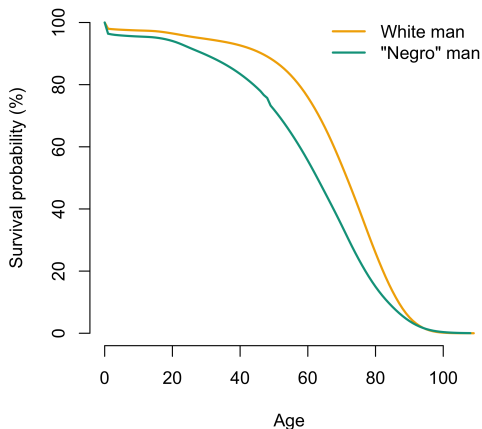


Frederick L. Hoffman
Hoffman (1896, 1918, 1931)

Insurance Pricing and Predictive Modeling

White, men			"Negro", men		
x	L_x	${}_5p_x$	x	L_x	${}_5p_x$
0	100000	2.3%	55	83001	8.5%
5	97671	0.2%	60	75969	12.7%
10	97441	0.2%	65	66343	18.4%
15	97208	0.7%	70	54138	25.5%
20	96480	1.0%	75	40324	35.8%
25	95524	0.8%	80	25885	47.7%
30	94716	0.9%	85	13527	62.1%
35	93843	1.3%	90	5125	75.1%
40	92631	2.1%	95	1274	85.2%
45	90725	3.3%	100	189	90.5%
50	87690	5.3%	105	18	100.0%

Insurance Pricing and Predictive Modeling



Force of mortality (log scale) white men and "Negro" men, 1968-71, U.S.

Insurance Pricing and Predictive Modeling

Definition 3.3: Balance Property

A pricing function m satisfies the **balance property** if $\mathbb{E}_{\mathbf{X}}[m(\mathbf{X})] = \mathbb{E}_Y[Y]$.

Proposition 3.3: Law of total expectations

$$\mathbb{E}_Y[Y] = \mathbb{E}_{\mathbf{X}}[\mathbb{E}_{Y|\mathbf{X}}[Y|\mathbf{X}]] = \mathbb{E}_{\mathbf{X}}[\mu(\mathbf{X})].$$

Proof Since $\mathbb{E}(Y) = \int yf_y(y)dy$ and $\mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = \int yf_{y|\mathbf{x}}(y|\mathbf{x})dy$,

$$\begin{aligned}\mathbb{E}(\mathbb{E}(X|Y)) &= \int \left(\int x\mathbb{P}[X = x|Y = y]dx \right) \mathbb{P}[Y = y]dy = \int \int x\mathbb{P}[X = x, Y = y]dxdy \\ &= \int x \left(\int \mathbb{P}[X = x, Y = y]dy \right) dx = \int x\mathbb{P}[X = x]dx = \mathbb{E}(X).\end{aligned}$$

Insurance Pricing and Predictive Modeling

Homogeneous risk sharing

	Policyholder	Insurer
Loss	$\mathbb{E}[Y]$	$Y - \mathbb{E}[Y]$
Average loss	$\mathbb{E}[Y]$	0
Variance	0	$\text{Var}[Y]$

$\mathbb{E}[Y]$ is the premium paid, and Y the total loss,
from [De Wit and Van Eeghen \(1984\)](#) and [Denuit and Charpentier \(2004\)](#)

Insurance Pricing and Predictive Modeling

Heterogeneous risk sharing, with perfect information

	Policyholder	Insurer
Loss	$\mathbb{E}[Y \Theta]$	$Y - \mathbb{E}[Y \Theta]$
Average loss	$\mathbb{E}[Y]$	0
Variance	$\text{Var}[\mathbb{E}[Y \Theta]]$	$\text{Var}[Y - \mathbb{E}[Y \Theta]]$

where Θ denotes the heterogeneous risk factor.

The term on the bottom right is $\mathbb{E}[\text{Var}[Y|\Theta]]$, corresponding to the standard [variance decomposition](#) (or Pythagoras theorem)

$$\text{Var}[Y] = \text{Var}[\mathbb{E}[Y|\Theta]] + \mathbb{E}[\text{Var}[Y|\Theta]].$$

to go further  (for more details on Lebesgue spaces, and L^2)



Proposition 3.4: Variance decomposition (1)

For any measurable random variable Y with finite variance

$$\text{Var}[Y] = \underbrace{\mathbb{E}[\text{Var}[Y|\Theta]]}_{\rightarrow \text{insurer}} + \underbrace{\text{Var}[\mathbb{E}[Y|\Theta]]}_{\rightarrow \text{policyholder}}.$$

Proof:

$$\begin{aligned}\text{Var}[Y] &= \mathbb{E}[Y^2] - \mathbb{E}[Y]^2 = \mathbb{E}[\text{Var}[Y|\Theta] + \mathbb{E}[Y|\Theta]^2] - \mathbb{E}[\mathbb{E}[Y|\Theta]]^2 \\ &= (\mathbb{E}[\text{Var}[Y|\Theta]]) + (\mathbb{E}[\mathbb{E}[Y|\Theta]^2] - \mathbb{E}[\mathbb{E}[Y|\Theta]]^2) = \mathbb{E}[\text{Var}[Y|\Theta]] + \text{Var}[\mathbb{E}[Y|\Theta]].\end{aligned}$$

Insurance Pricing and Predictive Modeling

Heterogeneous risk sharing, with imperfect information

	Policyholder	Insurer
Loss	$\mathbb{E}[Y \mathbf{X}]$	$Y - \mathbb{E}[Y \mathbf{X}]$
Average loss	$\mathbb{E}[Y]$	0
Variance	$\text{Var}[\mathbb{E}[Y \mathbf{X}]]$	$\mathbb{E}[\text{Var}[Y \mathbf{X}]]$

$$\mathbb{E}[\text{Var}[Y|\mathbf{X}]] = \underbrace{\mathbb{E}[\text{Var}[Y|\Theta]]}_{\text{perfect ratemaking}} + \underbrace{\mathbb{E}\{\text{Var}[\mathbb{E}[Y|\Theta]|\mathbf{X}]\}}_{\text{misclassification}}$$

This “misclassification” term (on the right) is called “*subsidierende solidariteit*” in [De Pril and Dhaene \(1996\)](#), or “*subsidiary solidarity*”, as opposed to “*kanssolidariteit*” or “*random solidarity*” term (on the left).

Proposition 3.5: Variance decomposition (2)

For any measurable random variable Y with finite variance

$$\text{Var}[Y] = \underbrace{\mathbb{E}[\text{Var}[Y|\mathbf{X}]]}_{\rightarrow \text{insurer}} + \underbrace{\text{Var}[\mathbb{E}[Y|\mathbf{X}]]}_{\rightarrow \text{policyholder}},$$

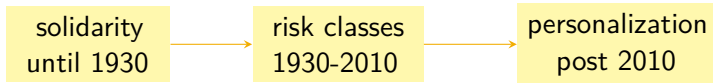
where

$$\begin{aligned} \mathbb{E}[\text{Var}[Y|\mathbf{X}]] &= \mathbb{E}[\mathbb{E}[\text{Var}[Y|\Theta]|\mathbf{X}]] + \mathbb{E}[\text{Var}[\mathbb{E}[Y|\Theta]|\mathbf{X}]] \\ &= \underbrace{\mathbb{E}[\text{Var}[Y|\Theta]]}_{\text{perfect ratemaking}} + \underbrace{\mathbb{E}\{\text{Var}[\mathbb{E}[Y|\Theta]|\mathbf{X}]\}}_{\text{misclassification}}. \end{aligned}$$

Clubs, Group and Categories

- Groups, or risk classes, are built on the basis of available data, and exist primarily as the product of actuarial models.
- For example, as mentioned in [Bailey and Simon \(1959\)](#), in motor insurance five risk classes can be considered, with rate surcharges relative to the first class (used here as a reference)
 - ▶ “*pleasure, no male operator under 25,*” (reference),
 - ▶ “*pleasure, non-principal male operator under 25,*” +65%,
 - ▶ “*business use,*” +65%,
 - ▶ “*married owner or principal operator under 25,*” +65%,
 - ▶ “*unmarried owner or principal operator under 25,*” +140%.
- There is no “physical basis” for group members to identify other members of *their* group, in the sense that they usually don’t share anything, except some common characteristics, [Gandy \(2016\)](#).

Clubs, Group and Categories



- In ancient Rome, a *collegium* (plural *collegia*) was an association, such as military *collegia*, [Verboven \(2011\)](#).
- As explained in [Ginsburg \(1940\)](#), upon the completion of his service a veteran had the right to join one of the many *collegia veteranorum* in each legion.
- In case of retirement, upon the completion of his term of service, the soldier would received a a lump sum which helped him somewhat to arrange the rest of his life. The membership in a *collegium* gave him a mutual insurance against “*unforeseen risks*.” These *collegia*, besides being cooperative insurance companies, had other functions.



Clubs, Group and Categories

- › In the early 1660th, the [Pirate's Code](#) was supposedly written by Portuguese buccaneer Bartolomeu Portuguê.
- › A section is explicitly dedicated to insurance and benefits: “*a standard compensation is provided for maimed and mutilated buccaneers. Thus they order for the loss of a right arm six hundred pieces of eight, or six slaves; for the loss of a left arm five hundred pieces of eight, or five slaves; for a right leg five hundred pieces of eight, or five slaves; for the left leg four hundred pieces of eight, or four slaves; for an eye one hundred pieces of eight, or one slave; for a finger of the hand the same reward as for the eye,*” see [Barbour \(1911\)](#) (or more recently [Leeson \(2009\)](#) and [Fox \(2013\)](#) about this piratical schemes).



Clubs, Group and Categories

- In the XIX-th century, in Europe, **mutual aid societies** involved a group of individuals who made regular payments into a common fund in order to provide for themselves in later, unforeseeable moments of financial hardship or of old age. As mentioned by **Garrioch (2011)**, in 1848, there were in Paris 280 mutual aid societies with well over 20,000 members.
- For example, the *Société des Arts Graphiques*, was created in 1808. It admitted only men over twenty and under fifty, and it charged much higher admission and annual fees for those who joined at a more advanced age. In return, they received benefits if they were unable to work, reducing over a period of time, but in case of serious illness the Society would pay the admission fee for a hospice. In England, there were “**friendly societies**,” as described in **Ismay (2018)**.



Clubs, Group and Categories

- › The money collected through contributions came to the rescue of unfortunate workers, who would no longer have any reason to radicalize. It was proposed that insurance should become compulsory (Bismark proposed this in Germany in 1883), but the idea was rejected in favor of giving workers the freedom to contribute, as the only way to moralize the working classes, as [Da Silva \(2023\)](#) explains.
- › In 1852, of the 236 mutual funds created, 21 were on a professional basis, while the other 215 were on a territorial basis. And from 1870 onwards, mutual funds diversified the professional profile of contributors beyond blue-collar workers, and expanded to include employees, civil servants, the self-employed and artists.
- › The amount of the premium is not linked to the risk.



Clubs, Group and Categories

- As [Da Silva \(2023\)](#) puts it, “*mutual insurers see in the actuarial figure the programmed end of solidarity.*” For mutual funds, solidarity is essential, with everyone contributing according to their means and receiving according to their needs. Around the same time, in France, the first insurance companies appeared, based on risk selection, and the first mathematical approaches to calculating premiums.
- [Hubbard \(1852\)](#) advocates the introduction of an “*English-style scientific organization*” in their management. For its members, they had to be able to know “*the probable average of the claims*” that they should cover, like insurance companies. The development of tables should lead insurers to adopt the principle of contributions varying according to the age of entry and the specialization of contributions and funds (health/retirement).
- For [Stone \(1993\)](#) and [Gowri \(2014\)](#) the defining feature of “modern insurance” is its reliance on [segmenting the risk pool into distinct categories](#), each receiving a price

Clubs, Group and Categories

corresponding to the particular risk that the individuals assigned to that category are expected to represent (as accurately as can be estimated by actuaries).

- Once heterogeneity with respect to the risk was observed in portfolios, insurers have operated by categorizing individuals into **risk classes** and assigning corresponding tariffs. This ongoing process of categorization ensures that the sums collected, on average, are sufficient to address the realized risks within specific groups.
- The aim of **risk classification**, as explained in **Wortham (1986)**, is to identify the specific characteristics that are supposed to determine an individual's propensity to suffer an adverse event, forming groups within which the risk is (approximately) equally shared. The problem, of course, is that the characteristics associated with various types of risk are almost infinite; as they cannot all be identified and priced in every risk classification system, there will necessarily be unpriced sources of heterogeneity between individuals in a given risk class.

Clubs, Group and Categories

- In 1915, as mentioned in [Rothstein \(2003\)](#), the president of the Association of Life Insurance Medical Directors of America noted that the question asked almost universally of the Medical Examiner was “*What is your opinion of the risk? Good, bad, first-class, second-class, or not acceptable?*” Historically, insurance prices were a (finite) collection of prices (maybe more than than the two classes mentioned, “*first-class*” and “*second-class*”).
- In the early 1920’s, Albert Henry Mowbray, who worked for New York Life Insurance Company and later Liberty Mutual (and was also an actuary for state-level insurance commissions in New Carolina and California, and the National Council on Workmen’s Insurance) gives his perspective on insurance rate making. See [Mowbray \(1921\)](#).



Clubs, Group and Categories

“Classification of risks in some manner forms the basis of rate making in practically all branches of insurance. It would appear therefore that there should be some fundamental principle to which a correct system of classification in any branch of insurance should conform (...) As long ago as the days of ancient Greece and Rome the gradual transition of natural phenomena was observed and set down in the Latin maxim, ‘natura non agit per altum’. If each risk, therefore is to be precisely rated, it would be necessary to recognize very minute differences and precisely measure them. (...) Since we are not capable of covering a large field fully and at the same time recognizing small differences in all parts of the field, it is natural that we resort to subdivision of the field by means of classification, thereby concentrating our attention on a smaller interval which may again be subdivided by further classification, and the system so carried on to the limit to which we find it necessary or desirable to go. But however far we may go in any system of classification, whether in the field of pure or applied science including the business or insurance, we shall always find difficulties presented by the borderline case, difficulties which arise from the continuous character

Clubs, Group and Categories

of natural phenomena which we are attempting to place in more or less arbitrary divisions. While thus acknowledging that classification will never completely solve the problem of recognizing differences between individuals, nevertheless classification seems to be necessary at least as a preliminary step toward such recognition in any field of study. The fact that a complete and final solution cannot be made is, therefore, no justification for completely discarding classification as a method of approach. Since it is insurance hazards that we undertake to measure and classify, the preliminary step in studying classification theory may well be to ask what is an insurance hazard and how it may be determined. It must be evident to the members of this Society that an insurance hazard is what is termed "a mathematical expectation," that is a product of a sum at risk and the probability of loss from the conditions insured against, e.g., the destruction of a piece of property by fire, the death of an individual, etc. If the net premiums collected are so determined on the basis of the true natural probability and there is a sufficient spread then the sums collected will just cover the losses and this is what should be," Mowbray (1921).

Clubs, Group and Categories

- “1. The classification should bring together risks which have inherent in their operation the same causes of loss.*
- 2. The variation from risk to risk in the strength of each cause or at least of the more important should not be greater than can be handled by the formula by which the classification is subdivided, i.e., the Schedule and / or Experience Rating Plan used.*
- 3. The classification should not cover risks which include, as important elements of their hazard, causes which are not common to all.*
- 4. The classification system and the formula for its extension (Schedule and / or Experience Rating Plans) should be harmonious.*
- 5. The basis throughout should be the outward, recognizable indicia of the presence and potency of the several inherent causes of loss including extent as well as occurrence of loss,” Mowbray (1921).*

Clubs, Group and Categories

- Several articles and textbooks in sociology tried to understand how classification mechanisms establish symbolic boundaries that reinforce group identities, such as Bourdieu (2018), Massey (2007), Fourcade and Healy (2013).
- But here, those “groups” or “classes” do not share any identity, and Simon (1988) or Harcourt (2015) use the term “actuarial classification” (where “actuarial” designates any decision-making technique that relies on predictive statistical methods, replacing more holistic or subjective forms of judgment). In those class-based systems, based on insurance rating table (or grid), results are determined by assigning individuals to a group in which each person is positioned as “average” or “typical”.
- [Most] “*actuaries cannot think of individuals except as members of groups*” claimed Brilmayer et al. (1979). Each individual is assigned the same value as all other members of the group to which it is assigned.



Clubs, Group and Categories

- › Simon (1987, 1988), and then Feeley and Simon (1992), defined “actuarialism,” that designate the use of statistics to guide “*class-based decision-making*,” used to price pensions and insurance. As explained in Harcourt (2015), this “actuarial classification” is the constitution of groups with no experienced social significance for the participants. A person classified as a particular risk by an insurance company shares nothing with the other people so classified, apart from a series of formal characteristics (e.g. age, sex, marital status, etc.).
- › For Austin (1983) and Simon (1988), categories used by the insurance company when grouping risks are “*singularly sterile*,” resulting in inert, immobile and deactivated communities, corresponding to “*artificial*” groups. These are not groups organized around a shared history, common experiences or active commitment, forming some “*aggregates*” – living only in the imagination of the actuary who calculates and tabulates, not in any lived form of human association.



Clubs, Group and Categories

- › If [Hacking \(1990\)](#) observed that standard classes creates coherent group identities (causing possible stereotypes and discrimination, [Simon \(1988\)](#)), provocatively suggests that actuarial classifications can in turn “*undo people's identity*.”
- › As mentioned in [Abraham \(1986\)](#), the goal for actuaries is to create groups, or “*classes*” made up of individuals who share a series of common characteristics and are therefore presumed to represent the same risk. Following [François \(2022\)](#), we could claim that actuarial techniques reduce individuals to a series of formal roles that have no “*moral density*” and therefore do not grant an “*identity*” that organizes a coherent sense of self. And the inclusion of nominally “*demoralized categories*,” such as gender, in class-based rating systems makes their total demoralization difficult to achieve – and is in itself an issue of struggle. [Heimer \(1985\)](#) used the term “*community of fate*.”
- › [Rouvroy et al. \(2013\)](#) and [Cheney-Lippold \(2017\)](#) point out that scoring technologies are continually swapping predictors, “*shuffling the cards*,” so that there is no stable basis for constructing group memberships, or a coherent sense.

Clubs, Group and Categories

“The price which a person pays for automobile insurance depends on age, sex, marital status, place of residence and other factors. This risk classification system produces widely differing prices for the same coverage for different people. Questions have been raised about the fairness of this system, and especially about its reliability as a predictor of risk for a particular individual. While we have not tried to judge the propriety of these groupings, and the resulting price differences, we believe that the questions about them warrant careful consideration by the State insurance departments. In most States the authority to examine classification plans is based on the requirement that insurance rates are neither inadequate, excessive, nor unfairly discriminatory. The only criterion for approving classifications in most States is that the classifications be statistically justified – that is, that they reasonably reflect loss experience. Relative rates with respect to age, sex, and marital status are based on the analysis of national data. A youthful male driver, for example, is charged twice as much as an older driver all over the country (...). It has also been claimed that insurance companies engage in redlining – the arbitrary denial of insurance to everyone

Clubs, Group and Categories

living in a particular neighborhood. Community groups and others have complained that State regulators have not been diligent in preventing redlining and other forms of improper discrimination that make insurance unavailable in certain areas. In addition to outright refusals to insure, geographic discrimination can include such practices as: selective placement of agents to reduce business in some areas, terminating agents and not renewing their book of business, pricing insurance at un-affordable levels, and instructing agents to avoid certain areas. We reviewed what the State insurance departments were doing in response to these problem. To determine if redlining exists, it is necessary to collect data on a geographic oasis. Such data should include current insurance policies, new policies being written, cancellations, and non-renewals. It is also important to examine data on losses by neighborhoods within existing rating territories because marked discrepancies within territories would cast doubt on the validity of territorial boundaries. Yet, not even a fifth of the States collect anything other than loss data, and that data is gathered on a territory-wide basis," Havens (1979)

Clubs, Group and Categories

“On the other hand, the opinion that distinctions based on sex, or any other group variable, necessarily violate individual rights reflects ignorance of the basic rules of logical inference in that it would arbitrarily forbid the use of relevant information. It would be equally fallacious to reject a classification system based on socially acceptable variables because the results appear discriminatory. For example, a classification system may be built on use of car, mileage, merit rating, and other variables, excluding sex. However, when verifying the average rates according to sex one may discover significant differences between males and females. Refusing to allow such differences would be attempting to distort reality by choosing to be selectively blind. The use of rating territories is a case in point. Geographical divisions, however designed, are often correlated with socio-demographic factors such as income level and race because of natural aggregation or forced segregation according to these factors. Again we conclude that insurance companies should be free to delineate territories and assess territorial differences as well as they can. At the same time, insurance companies should recognize that it is in their best interest to be objective and use clearly relevant

Clubs, Group and Categories

factors to define territories lest they be accused of invidious discrimination by the public. (...) ” Casey et al. (1976)

“One possible standard does exist for exception to the counsel that particular rating variables should not be proscribed. What we have called ‘equal treatment’ standard of fairness may precipitate a societal decision that the process of differentiating among individuals on the basis of certain variables is discriminatory and intolerable. This type of decision should be made on a specific, statutory basis. Once taken, it must be adhered to in private and public transactions alike and enforced by the insurance regulator. This is, in effect, a standard for conduct that by design transcends and preempts economic considerations. Because it is not applied without economic cost, however, insurance regulators and the industry should participate in and inform legislative deliberations that would ban the, use of particular rating variables as discriminatory.” Casey et al. (1976)

Price Optimization

- › Decision theory under uncertainty (see [Charpentier \(2014\)](#)),

$$X \preceq Y \iff \mathcal{R}(X) \leq \mathcal{R}(Y),$$

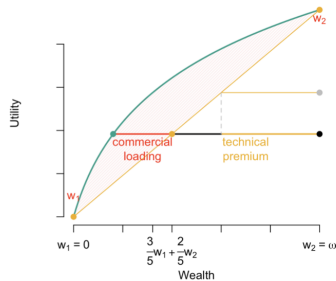
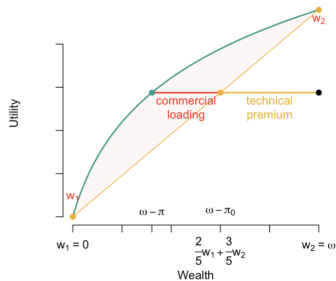
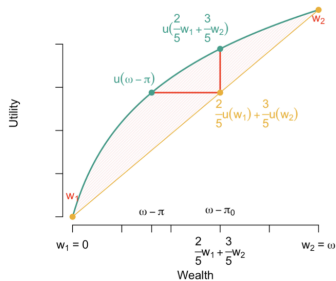
- › A classical representation is $\mathcal{R}(Y) = \mathbb{E}[u(\omega - Y)]$, as in [Neumann and Morgenstern \(1947\)](#), where ω is the initial wealth.
- › u denotes the utility of the agent
- › Let π denote the premium asked to transfer risk (loss) Y ,

$$\begin{cases} u(\omega - \pi) > \mathbb{E}[u(\omega - Y)] : & \text{purchases insurance} \\ u(\omega - \pi) < \mathbb{E}[u(\omega - Y)] : & \text{does not purchase insurance} \end{cases}$$

Price Optimization

Definition 3.4: Indifference utility principle

Let Y be the non-negative random variable corresponding to the total annual loss associated with a given policy, for a policyholder with utility u and wealth w , the **indifference premium** is $\pi = \omega - u^{-1}(\mathbb{E}[u(\omega - Y)])$.



– Part 2 –

Machine Learning

Proposition 4.1: Law of Large Numbers (1)

Consider an infinite collection of i.i.d. random variables $Y, Y_1, Y_2, \dots, Y_n, \dots$ in a probabilistic space $(\Omega, \mathcal{F}, \mathbb{P})$, then

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \mathbf{1}(Y_i \in \mathcal{A})}_{\text{(empirical) frequency}} \xrightarrow{\text{a.s.}} \underbrace{\mathbb{P}(\{Y \in \mathcal{A}\})}_{\text{probability}} = \mathbb{P}[Y \in \mathcal{A}], \text{ as } n \rightarrow \infty.$$

“law of the unconscious statistician” (Ross (2014) and Casella and Berger (1990)),
“statisticians make liberal use of conditioning arguments to shorten what would otherwise be long proofs,” Proschan and Presnell (1998)

$$\mathbb{P}(Y \in \mathcal{A} | X = x) = \lim_{\epsilon \rightarrow 0} \frac{\mathbb{P}(\{Y \in \mathcal{A}\} \cap \{|X - x| \leq \epsilon\})}{\mathbb{P}(\{|X - x| \leq \epsilon\})} = \lim_{\epsilon \rightarrow 0} \mathbb{P}(Y \in \mathcal{A} | |X - x| \leq \epsilon).$$

Statistical Learning

This frequentist approach is unable to make sense of the probability of a "single singular event", as noted by von Mises (1928, 1939).

"When we speak of the 'probability of death', the exact meaning of this expression can be defined in the following way only. We must not think of an individual, but of a certain class as a whole, e.g., 'all insured men forty-one years old living in a given country and not engaged in certain dangerous occupations'. A probability of death is attached to the class of men or to another class that can be defined in a similar way. We can say nothing about the probability of death of an individual even if we know his condition of life and health in detail. The phrase 'probability of death', when it refers to a single person, has no meaning for us at all."



Definition 4.1: Loss ℓ

A **loss function** ℓ is a function defined on $\mathcal{Y} \times \mathcal{Y}$ such that $\ell(y, y') \geq 0$ and $\ell(y, y) = 0$.

Definition 4.2: Risk \mathcal{R}

For a fitted model \hat{m} , its **risk** is

$$\mathcal{R}(\hat{m}) = \mathbb{E}_{\mathbb{P}} \left[\ell(Y, \hat{m}(\mathbf{X})) \right] = \int \ell(y, \hat{m}(\mathbf{x})) d\mathbb{P}(y, \mathbf{x}).$$

Definition 4.3: Empirical risk $\hat{\mathcal{R}}_n$

Given a sample $\{(y_i, \mathbf{x}_i), i = 1, \dots, n\}$, define the **empirical risk**

$$\hat{\mathcal{R}}_n(\hat{m}) = \frac{1}{n} \sum_{i=1}^n \ell(\hat{m}(\mathbf{x}_i), y_i).$$

➤ Following **Vapnik (1991)**, the "**empirical risk minimization principle**" states that the learning algorithm \hat{m}^* is

$$\hat{m}^* = \underset{\hat{m} \in \mathcal{M}}{\operatorname{argmin}} \{ \hat{\mathcal{R}}_n(\hat{m}) \}.$$

Proposition 4.2: Optimal Decision, "Bayes decision rule"

For each \mathbf{x} choose the prediction $m_{\mathbf{x}}^*$ that minimizes the conditional expected loss,

$$m_{\mathbf{x}}^* \in \operatorname{argmin}_{z \in \mathcal{Y}} \left\{ \int \ell(y, z) d\mathbb{P}_{Y|\mathbf{X}}(y|\mathbf{x}) \right\}$$

➤ It is straightforward since $d\mathbb{P}_{Y, \mathbf{X}}(y, \mathbf{x}) = d\mathbb{P}_{Y|\mathbf{X}}(y|\mathbf{x}) \cdot d\mathbb{P}_{\mathbf{X}}(\mathbf{x})$,

$$\mathcal{R}(\hat{m}) = \int \left[\int \ell(y, \hat{m}(\mathbf{x})) d\mathbb{P}_{Y|\mathbf{X}}(y|\mathbf{x}) \right] d\mathbb{P}_{\mathbf{X}}(\mathbf{x}).$$

by definition, $m_{\mathbf{x}}^*$ minimizes the term in blue, i.e., for any \hat{m}

$$\mathcal{R}(\hat{m}) \geq \int \left[\int \ell(y, m_{\mathbf{x}}^*) d\mathbb{P}_{Y|\mathbf{X}}(y|\mathbf{x}) \right] d\mathbb{P}_{\mathbf{X}}(\mathbf{x}) = \mathcal{R}(m^*).$$



Statistical Learning

- It is coined "Bayes decision rule" because the conditional distribution $Y|\mathbf{X}$ is sometimes be referred to as the "posterior" distribution of Y given data \mathbf{X} .

Definition 4.4: Misclassification loss, $\ell_{0/1}$

$$\ell_{0/1}(y, \hat{y}) = \mathbf{1}(y \neq \hat{y}).$$

In the case of a binary classifier, observe that

$$\begin{aligned}\mathcal{R}(\hat{m}) &= \mathbb{E}[\ell(\hat{m}(\mathbf{X}), Y)] = \mathbb{E}[\mathbb{E}[\ell(\hat{m}(\mathbf{X}), Y) | \mathbf{X}]] \\ &= \mathbb{E}[\ell(\hat{m}(\mathbf{X}), 1) \cdot \mathbb{P}(Y = 1 | \mathbf{X}) + \ell(\hat{m}(\mathbf{X}), 0) \cdot \mathbb{P}(Y = 0 | \mathbf{X})] \\ &= \mathbb{E}[\mathbf{1}[\hat{m}(\mathbf{X}) \neq 1] \cdot \mu(\mathbf{X}) + \mathbf{1}[\hat{m}(\mathbf{X}) \neq 0] \cdot (1 - \mu(\mathbf{X}))] \\ &= \mathbb{E}[\mathbf{1}[\hat{m}(\mathbf{X}) \neq 1] \cdot \mu(\mathbf{X}) + (1 - \mathbf{1}[\hat{m}(\mathbf{X}) \neq 1]) \cdot (1 - \mu(\mathbf{X}))] \\ &= \mathbb{E}[\mathbf{1}[\hat{m}(\mathbf{X}) \neq 1] \cdot (2\mu(\mathbf{X}) - 1) + 1 - \mu(\mathbf{X})].\end{aligned}$$

Statistical Learning

Since $\hat{m} : \mathcal{X} \rightarrow \{0, 1\}$, this expectation is minimized by choosing $\hat{m} = m^*$, where

$$m^*(\mathbf{x}) = \mathbf{1}(\mu(\mathbf{x}) > 1/2) = \begin{cases} 1 & \text{if } \mu(\mathbf{x}) > 1/2 \\ 0 & \text{if } \mu(\mathbf{x}) \leq 1/2 \end{cases}$$

The optimal risk ("Bayes risk") is $\mathcal{R}(m^*) = \inf_m \{\mathcal{R}(m)\}$.

Definition 4.5: Excess of risk of \hat{m}

For any model \hat{m} , the excess of risk is $\mathcal{R}(\hat{m}) - \mathcal{R}(m^*)$.

For a classifier

$$\mathcal{R}(\hat{m}) - \mathcal{R}(m^*) = \mathbb{E}[|2\mu(\mathbf{X}) - 1| \cdot \mathbf{1}(\hat{m}(\mathbf{X}) \neq m^*(\mathbf{X}))].$$

Since we do not know μ consider a classifier based on \hat{m}

Definition 4.6: Plug-in Estimator

Estimate $\hat{\mu}$ and use, as a classifier, $\mathbf{1}(\hat{\mu}(\mathbf{x}) > 1/2)$.

Proposition 4.3

For any model $\hat{\mu}$, the risk of the plug-in classifier $\hat{m}(\mathbf{x}) = \mathbf{1}(\hat{\mu}(\mathbf{x}) > 1/2)$ satisfies

$$\mathcal{R}(\hat{m}) - \mathcal{R}(m^*) \leq 2\mathbb{E}|\mu(\mathbf{X}) - \hat{\mu}(\mathbf{X})|.$$

Proof We have seen that

$$\mathcal{R}(\hat{m}) - \mathcal{R}(m^*) = \mathbb{E}(1[\hat{m}(\mathbf{X}) \neq 1] - 1[m^*(\mathbf{X}) \neq 1]) \cdot (2\mu(\mathbf{X}) - 1).$$

Statistical Learning

But

$$\begin{aligned} & (\mathbf{1}[\hat{m}(\mathbf{X}) \neq 1] - \mathbf{1}[m^*(\mathbf{X}) \neq 1]) (2\mu(\mathbf{X}) - 1) \\ &= \mathbf{1}[\hat{m}(\mathbf{X}) \neq m^*(\mathbf{X})] (\mathbf{1}[\hat{m}(\mathbf{X}) \neq 1] - \mathbf{1}[m^*(\mathbf{X}) \neq 1]) (2\mu(\mathbf{X}) - 1) \\ &= \begin{cases} \mathbf{1}[\hat{m}(\mathbf{X}) \neq m^*(\mathbf{X})] (2\mu(\mathbf{X}) - 1) & \text{if } 2\mu(\mathbf{X}) - 1 > 0, \\ \mathbf{1}[\hat{m}(\mathbf{X}) \neq m^*(\mathbf{X})] (-1)(2\mu(\mathbf{X}) - 1) & \text{if } 2\mu(\mathbf{X}) - 1 \leq 0. \end{cases} \end{aligned}$$

(from the definition of m^*)

$$= \mathbf{1}[\hat{m}(\mathbf{X}) \neq m^*(\mathbf{X})] \cdot |2\mu(\mathbf{X}) - 1|,$$

$$\mathcal{R}(\hat{m}) - \mathcal{R}(m^*) = \mathbb{E}(\mathbf{1}[\hat{m}(\mathbf{X}) \neq m^*(\mathbf{X})]) \cdot 2|\mu(\mathbf{X}) - 1/2|.$$

If $\hat{m}(\mathbf{x}) \neq m^*(\mathbf{x})$, it means that $\hat{\mu}(\mathbf{x})$ and $\mu(\mathbf{x})$ lie on opposite sides of $1/2$,

$$|\hat{\mu}(\mathbf{x}) - \mu(\mathbf{x})| = |\hat{\mu}(\mathbf{x}) - 1/2| + \underbrace{|1/2 - \mu(\mathbf{x})|}_{\geq 0} \geq |\hat{\mu}(\mathbf{x}) - 1/2|$$

i.e.

$$|\hat{\mu}(\mathbf{x}) - \mu(\mathbf{x})| \geq |\hat{\mu}(\mathbf{x}) - 1/2| \cdot \mathbf{1}[\hat{m}(\mathbf{X}) \neq m^*(\mathbf{X})]$$

which is also valid when $\hat{m}(\mathbf{x}) = m^*(\mathbf{x})$, thus

$$\mathcal{R}(\hat{m}) - \mathcal{R}(m^*) = 2\mathbb{E}(\mathbf{1}[\hat{m}(\mathbf{X}) \neq m^*(\mathbf{X})]) \cdot |\mu(\mathbf{X}) - 1/2| \leq 2\mathbb{E}[|\hat{\mu}(\mathbf{X}) - \mu(\mathbf{X})|].$$

- This $\ell_{0/1}$ loss function may be difficult to directly optimize, as shown in [Bartlett et al. \(2006\)](#). One could consider some [surrogate loss](#) $\tilde{\ell}$ which is easier to optimize.

Definition 4.7: Elicitation, [Brier \(1950\)](#), [Good \(1952\)](#)

A statistical functional $\mathcal{I}(Y)$ is said to be elicitable if it minimizes expected loss for some loss function s , in the sense that

$$\mathcal{I}(Y) = \operatorname{argmin}_{y \in \mathbb{R}} \{ \mathbb{E}[s(Y, y)] \}$$

- Important properties for risk measures and backtesting. "*The elicibility of a risk measure means that the risk measure can be obtained by minimizing the expectation of a forecasting objective function. Elicitability is closely related to backtesting, whose objective is to evaluate the performance of a risk forecasting model. If a risk measure is elicitable, then the sample average forecasting error based on the objective function can be used for backtesting the risk measure,*" [He et al. \(2022\)](#)

Loss Functions

- In a regression problem, a quadratic loss function ℓ_2 is used

Definition 4.8: Quadratic loss, ℓ_2

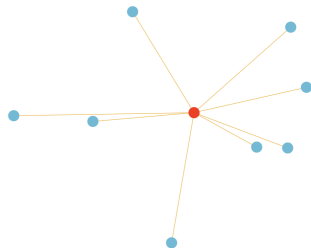
$$\ell_2(y, \hat{y}) = (y - \hat{y})^2, \text{ and the risk is then } \mathcal{R}_2(\hat{m}) = \mathbb{E}[(Y - \hat{m}(\mathbf{X}))^2].$$

- Observe that

$$\mathbb{E}[Y] = \operatorname{argmin}_{m \in \mathbb{R}} \{ \mathcal{R}_2(m) \} = \operatorname{argmin}_{m \in \mathbb{R}} \{ \mathbb{E}[\ell_2(Y, m)] \}.$$

The expected value is “ellicitable” (for the $s = \ell_2$ loss).

The empirical risk minimizer is the “**least-square**” estimate.



Loss Functions

- See [Huttegger \(2013\)](#), explaining why the expected value is also called “best estimate”.
- Up to a monotonic transformation (the square root function), the distance here is the expectation of the quadratic loss function. With the terminology of [Angrist and Pischke \(2009\)](#), the regression function μ is the function of \mathbf{x} that serves as “*the best predictor of y , in the mean-squared error sense.*”

Proposition 4.4: Optimal Decision, “*Bayes decision rule*”

For the quadratic loss ℓ_2 , Bayes decision rule is the (conditional) expected value,
 $m_{\mathbf{x}}^* = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}] = \mu(\mathbf{x})$.

Loss Functions

Definition 4.9: Inner product

An **inner product** on \mathcal{H} is the application $(f, g) \mapsto \langle f, g \rangle_{\mathcal{H}}$ (taking value in \mathbb{R}) bilinear, symmetric, definite positive:

- ▶ $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$
- ▶ $\langle \alpha f + \beta g, h \rangle_{\mathcal{H}} = \alpha \langle f, h \rangle_{\mathcal{H}} + \beta \langle g, h \rangle_{\mathcal{H}}$
- ▶ $\langle f, f \rangle_{\mathcal{H}} \geq 0$ and $\langle f, f \rangle_{\mathcal{H}} = 0$ if and only if $f = 0$.

Example : $\mathcal{H} = \mathbb{R}^n$, $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^{\top} \mathbf{y}$

Example : $\mathcal{H} = \mathbb{R}^n$, let Σ denote some symmetric $n \times n$ positive definite matrix. Then

$\langle \mathbf{x}, \mathbf{y} \rangle_{\Sigma} = \mathbf{x}^{\top} \Sigma^{-1} \mathbf{y}$ is an inner product on \mathbb{R}^n .

Example : $\mathcal{H} = \ell^2 = \left\{ u : \sum_{i=1}^{\infty} u_i^2 < \infty \right\}$, $\langle u, v \rangle = \sum_{i=1}^{\infty} u_i v_i$

Loss Functions

Example : $\mathcal{H} = L^2(\mu) = \left\{ f : \int f(x)^2 d\mu(x) < \infty \right\}$, $\langle f, g \rangle = \int f(x)g(x) d\mu(x)$

Example : Consider the vector space \mathcal{V} that consists of all real-valued random variables defined on $(\Omega, \mathcal{F}, \mathbb{P})$. Given $k \in [1, \infty)$, define

$$\|X\|_k = \left[\mathbb{E} \left(|X|^k \right) \right]^{1/k}.$$

to go further  (for more details on Lebesgue spaces, and L^2)

Loss Functions

A **norm** $\|\cdot\|$, in \mathbb{R}^n , satisfies

- ▶ homogeneity, $\|a\vec{u}\| = |a| \cdot \|\vec{u}\|$, $\forall a$
- ▶ triangle inequality, $\|\vec{u} + \vec{v}\| \leq \|\vec{u}\| + \|\vec{v}\|$
- ▶ positivity, $\|\vec{u}\| \geq 0$
- ▶ definiteness, $\|\vec{u}\| = 0 \iff \vec{u} = \vec{0}$

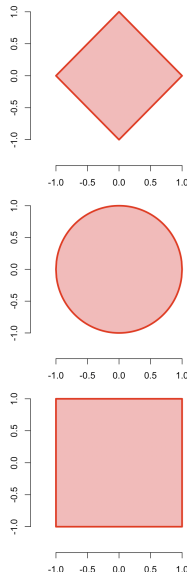
ℓ_1 norm: $\|\mathbf{x}\|_{\ell_1} = |x_1| + \dots + |x_n|$,

ℓ_2 norm: $\|\mathbf{x}\|_{\ell_2} = \sqrt{x_1^2 + \dots + x_n^2}$,

ℓ_p norm: with $p \geq 1$, $\|\mathbf{x}\|_{\ell_p} = (|x_1|^p + \dots + |x_n|^p)^{1/p}$

ℓ_∞ norm: $\|\mathbf{x}\|_{\ell_\infty} = \max\{x_i\}$

Unit balls ($\|\mathbf{x}\|_{\ell_p} \leq 1$) are convex sets



Proposition 4.5: Gradient of ℓ_p norms

$$\frac{\partial}{\partial x_j} \|\mathbf{x}\|_{\ell_p} = \frac{1}{p} \left(\sum_i |x_i|^p \right)^{\frac{1}{p}-1} \cdot p|x_j|^{p-1} \text{sign}(x_j) = \left(\frac{|x_j|}{\|\mathbf{x}\|_{\ell_p}} \right)^{p-1} \text{sign}(x_j).$$

$$\begin{aligned} \frac{\partial}{\partial x_j} \|\mathbf{x}\|_{\ell_p} &= \frac{\partial}{\partial x_j} \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} = \frac{1}{p} \left(\sum_{i=1}^n |x_i|^p \right)^{(1/p)-1} \frac{\partial}{\partial x_j} \left(\sum_{i=1}^n |x_i|^p \right) \\ &= \left[\left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} \right]^{1-p} \sum_{i=1}^n |x_i|^{p-1} \delta_{ij} \frac{x_i}{|x_i|} = \left(\frac{|x_j|}{\|\mathbf{x}\|_{\ell_p}} \right)^{p-1} \text{sign}(x_j). \end{aligned}$$

Loss Functions

Definition 4.10: Quantile loss, $\ell_{q,\alpha}$

The **quantile loss** $\ell_{q,\alpha}$ for some $\alpha \in (0, 1)$ is

$$\ell_{q,\alpha}(y, \hat{y}) = \max \{ \alpha(y - \hat{y}), (1 - \alpha)(\hat{y} - y) \} = (y - \hat{y})(\alpha - \mathbf{1}_{(y < \hat{y})}).$$

- This loss is not symmetric $\ell_{q,\alpha}(y, \hat{y}) \neq \ell_{q,\alpha}(\hat{y}, y)$ (if $\alpha \neq 1/2$).
- It is called “quantile” loss since

$$Q(\alpha) = F^{-1}(\alpha) \in \underset{q \in \mathbb{R}}{\operatorname{argmin}} \left\{ \mathbb{E} \left[\ell_{q,\alpha}(Y, q) \right] \right\},$$

(quantiles are also “elicitable” functionals, elicited by

$$s(y, \hat{y}) = \alpha(y - \hat{y})_+ + (1 - \alpha)(y - \hat{y})_-)$$

Loss Functions

➤ Indeed, the first order condition of

$$\min_{q \in \mathbb{R}} \left\{ (\alpha - 1) \int_{-\infty}^q (y - q) dF_Y(y) + \alpha \int_q^{\infty} (y - q) dF_Y(y) \right\},$$

can be written, using Leibniz integral rule,

$$(1 - \alpha) \int_{-\infty}^{q^*} dF_Y(y) - \alpha \int_{q^*}^{\infty} dF_Y(y) = 0$$

i.e. $F_Y(q^*) - \alpha = 0$.

Loss Functions

Definition 4.11: Expectile loss, $\ell_{e,\alpha}$

The **expectile loss** $\ell_{e,\alpha}$, for some $\alpha \in (0, 1)$ is

$$\ell_{e,\alpha}(y, \hat{y}) = (y - \hat{y})^2 \cdot (\alpha - \mathbf{1}_{(y < \hat{y})})$$

$$E(\alpha) = \operatorname{argmin}_{e \in \mathbb{R}} \left\{ \mathbb{E} \left[\ell_{e,\alpha}(Y, e) \right] \right\},$$

(expectiles are elicited by $s(x, y) = \alpha(y - x)_+^2 + (1 - \alpha)(y - x)_-^2$).

“Expectiles have properties that are similar to quantiles” Newey and Powell (1987)

Loss Functions

Portnoy and Koenker (1997), "*The Gaussian Hare and the Laplacian Tortoise*"



to go further → (for more details on optimization issues)

Loss and Generalized Linear Models

- › In GLM, the scaled deviance ($-2 \times$ the log-likelihood) of the exponential model is

$$D^* = \sum_{i=1}^n d^*(y_i, \hat{y}_i), \text{ where } d^*(y_i, \hat{y}_i) = 2(\log \mathcal{L}_i(y_i) - \log \mathcal{L}_i(\hat{y}_i)).$$

that can be related to in-sample empirical risk

$$\hat{\mathcal{R}}_n(\hat{m}) = \sum_{i=1}^n \ell(y_i, \hat{m}(\mathbf{x}_i)),$$

- › For the Poisson distribution (with a log-link), the loss would be

$$\ell(y_i, \hat{y}_i) = \begin{cases} 2(y_i \log y_i - y_i \log \hat{y}_i - y_i + \hat{y}_i) & y_i > 0 \\ 2\hat{y}_i & y_i = 0, \end{cases}$$

while for a logistic regression, we have the standard binary cross-entropy loss

$$\ell(y_i, \hat{y}_i) = -(y_i \log[\hat{y}_i] + (1 - y_i) \log[1 - \hat{y}_i]).$$

Distance Between Distributions

Definition 4.12: Distance (or metric)

A distance d on a set E is a function $E \times E \rightarrow \mathbb{R}_+$ such that

- ▶ d is symmetric, $\forall (a, b) \in E^2, d(a, b) = d(b, a),$
- ▶ d is separable, $\forall (a, b) \in E^2, d(a, b) = 0 \Leftrightarrow a = b,$
- ▶ d satisfies $\forall (a, b, c) \in E^3, d(a, c) \leq d(a, b) + d(b, c)$

In a vector space, with norm $\| \cdot \|$ the induced distance is $d(x, y) = \|y - x\|.$

Conversely, if

- ▶ d invariant by translation, $d(x, y) = d(x + a, y + a)$
- ▶ d is homogeneous, $d(\alpha x, \alpha y) = |\alpha|d(x, y)$

Distance Between Distributions

then $\|x\| = d(x, 0)$ is a norm.

Proposition 4.6

If d is a distance on E , and if $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is an increasing function such that $\psi(0) = 0$ and $\psi(t) > 0$ for all $t > 0$. If ψ is subadditive ($\psi(s+t) \leq \psi(s) + \psi(t)$), then $\delta(a, b) = \psi(d(a, b))$ is also a distance on E .

Proposition 4.7:

If d is a distance on E , then d^2 is not necessarily a distance.

Distance Between Distributions

› Consider the Euclidean distance in $E = \mathbb{R}^2$, i.e.

$d(\mathbf{z}_1, \mathbf{z}_2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$. d^2 is not a distance, see

$$\begin{cases} d^2(-\mathbf{1}, +\mathbf{1}) = 2^2 + 2^2 = 8 \\ d^2(-\mathbf{1}, \mathbf{0}) = 1^2 + 1^2 = 2 \\ d^2(\mathbf{0}, +\mathbf{1}) = 1^2 + 1^2 = 2 \end{cases}$$

i.e. d^2 does not satisfy the triangular inequality

$$d^2(-\mathbf{1}, +\mathbf{1}) > d^2(-\mathbf{1}, \mathbf{0}) + d^2(\mathbf{0}, +\mathbf{1}),$$

while

$$d(-\mathbf{1}, +\mathbf{1}) \leq d(-\mathbf{1}, \mathbf{0}) + d(\mathbf{0}, +\mathbf{1}).$$

(functions that generalize squared distance are sometimes referred to as divergences)

Distance Between Distributions

- › In addition to "distance", similar terms are used, including "dissimilarity", "deviance", "deviation", "discrepancy", "discrimination", and "divergence"

(... all denoted " d ", or " D ")

- › A fundamental problem in statistics and machine learning is to come up with useful measures of "distance" between pairs of probability distributions. Two desirable properties of a distance function are symmetry and the triangle inequality.
- › Unfortunately, many notions of "distance" between probability distributions do not satisfy these properties. Weaker notions of distance are often used, such as dissimilarity measures and divergences.
- › See [Cha \(2007\)](#) for a comprehensive list of distances...

Distance Between Distributions

Definition 4.13: Dissimilarity measure

A dissimilarity measure D on a set E is a function $E \times E \rightarrow \mathbb{R}_+$ such that D is positive and separable, i.e., $\forall (a, b) \in E^2$, $D(a, b) = 0 \Leftrightarrow a = b$,

Definition 4.14: Divergence on \mathbb{R}^n

A divergence D on a set $E \subset \mathbb{R}^n$ is a function $E \times E \rightarrow \mathbb{R}_+$ such that

- ▶ D is separable, $\forall (\mathbf{x}, \mathbf{y}) \in E^2$, $D(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} = \mathbf{y}$,
- ▶ D admits development

$$\forall (\mathbf{x}, \mathbf{x} + \boldsymbol{\epsilon}) \in E^2, D(\mathbf{x}, \mathbf{x} + \boldsymbol{\epsilon}) = \frac{1}{2} \sum A_{i,j}(\boldsymbol{\epsilon}) \epsilon_i \epsilon_j + O(|\boldsymbol{\epsilon}|^3),$$

where $A(\boldsymbol{\epsilon})$ is definite positive.

Distance Between Distributions

Definition 4.15: Scale sensitive divergence, Zolotarev (1976)

A divergence D is **scale sensitive** (of order $\beta > 0$) if $D(c\mathbf{x}, c\mathbf{y}) \leq |c|^\beta D(\mathbf{x}, \mathbf{y})$

Definition 4.16: Bregman Divergence, Bregman (1967)

Let $\psi : \mathcal{X} \rightarrow \mathbb{R}$ be a strictly convex function that is continuously differentiable. Then the **Bregman divergence** $D_\psi(\mathbf{x}, \mathbf{y})$ is defined as

$$D_\psi(\mathbf{x}, \mathbf{y}) = \psi(\mathbf{x}) - \psi(\mathbf{y}) - \langle \nabla \psi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle.$$

➤ If $\psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2$ (strictly convex), then $D_\psi(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2$.

(recall that $\nabla \|\mathbf{x}\|^2 = 2\mathbf{x}$)

Proposition 4.8: Bregman Divergence

Let $\psi : \mathcal{X} \rightarrow \mathbb{R}$ be a strictly convex function that is continuously differentiable. Then **Bregman divergence** $D_\psi(\mathbf{x}, \mathbf{y})$ is

- ▶ strictly convex in \mathbf{x} ,
- ▶ (generally) non-convex in \mathbf{y} ,
- ▶ non-negative $D_\psi(\mathbf{x}, \mathbf{y}) \geq 0$,
- ▶ separable, $D_\psi(\mathbf{x}, \mathbf{y}) = 0$ if and only if $\mathbf{x} = \mathbf{y}$,
- ▶ (generally) asymmetric.

Distance Between Distributions

- › If $\mathcal{X} = \mathbb{R}^n$, and $\psi(\mathbf{x}) = \frac{1}{2} \sum_{ij} A_{ij} x_i x_j = \frac{1}{2} \mathbf{x}^\top A \mathbf{x}$ for some $n \times n$ matrix A definite positive, then

$$D_\psi(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \sum_{ij} A_{ij} (x_i - y_i)(x_j - y_j) = (\mathbf{x} - \mathbf{y})^\top A (\mathbf{x} - \mathbf{y})$$

(see Mahalanobis distance).

- › If $\mathcal{X} = \mathbb{R}^n$, and $\psi(\mathbf{x}) = - \sum_i \log(x_i)$ then

$$D_\psi(\mathbf{x}, \mathbf{y}) = \sum_i \frac{x_i}{y_i} - \log \frac{x_i}{y_i} - 1$$

See [Banerjee et al. \(2005\)](#) for more examples.

Distance Between Distributions

We have defined norms

› on \mathbb{R}^n , e.g.,

$$\|\mathbf{x}\|_{\ell_2} = \left(|x_1|^2 + \dots + |x_n|^2\right)^{1/2} = \left(\sum_{i=1}^n |x_i|^2\right)^{1/2}$$

that could be extended

› on \mathbb{R} -valued random variables, e.g.,

$$\|X\|_2 = \left(\mathbb{E} \left[|X|^2\right]\right)^{1/2} = \left(\sum |x|^2 p(x)\right)^{1/2} = \left(\int |x|^2 f(x) dx\right)^{1/2}$$

We can also define "distances", "dissimilarity" measures, and "divergences"

› on \mathbb{R}^n , e.g.,

$$\ell_2(\mathbf{x}, \mathbf{y}) = d(\mathbf{x}, \mathbf{y}) = \left(|x_1 - y_1|^2 + \dots + |x_n - y_n|^2\right)^{1/2} = \left(\sum_{i=1}^n |x_i - y_i|^2\right)^{1/2}$$

Distance Between Distributions

that could be extended

› on \mathbb{R} -valued random variables as components of a random vector, e.g.,

$$D(X, Y) = \left(\mathbb{E} \left[|X - Y|^2 \right] \right)^{1/2} = \left(\sum |x - y|^2 p(x, y) \right)^{1/2} = \left(\int |x - y|^2 f(x, y) dx dy \right)^{1/2}$$

where p or f is the joint distribution of (X, Y) , e.g., for a Gaussian vector

$$D(X, Y) = (\mu_x - \mu_y)^2 + (\sigma_x - \sigma_y)^2 + 2\sigma_x\sigma_y(1 - \rho).$$

› on \mathbb{R} -valued random variables assuming that random variables are independent, e.g.,

$$D_{\perp}(X, Y) = \left(\sum |x - y|^2 p_x(x) p_y(y) \right)^{1/2} = \left(\int |x - y|^2 f_x(x) f_y(y) dx dy \right)^{1/2}$$

e.g., for two Gaussian distributions

$$D_{\perp}(X, Y) = (\mu_x - \mu_y)^2 + \sigma_x^2 + \sigma_y^2.$$

Distance Between Distributions

and one can consider some distance

› on \mathbb{R} -valued distributions, e.g.,

$$D(\mathcal{N}(\mu_x, \sigma_x^2), \mathcal{N}(\mu_y, \sigma_y^2)) = (\mu_x - \mu_y)^2 + (\sigma_x - \sigma_y)^2.$$

In the context of "probabilistic forecasts" (as in [Gneiting et al. \(2007\)](#)), a "distance"

› on pairs $\mathbb{R} \times \mathbb{R}$ -valued distributions, e.g.,

$$D(x, \mathcal{N}(\mu_y, \sigma_y^2)) = (x - \mu_y)^2 + \sigma_y^2.$$

Distance Between Distributions

Definition 4.17: Sum invariant divergence, [Zolotarev \(1976\)](#)

A divergence D is **sum invariant** if $D(X + Z, Y + Z) \leq D(X, Y)$ whenever $Z \perp\!\!\!\perp X, Y$

Example: if D is 1-scale sensitive, $D(\mathbf{1}_0, \mathbf{1}_1) \leq \frac{1}{2}D(\mathbf{1}_0, \mathbf{1}_2)$

Example: if D is sum invariant, $D(\mathbf{1}_0, \mathbf{1}_1) = D(\mathbf{1}_1, \mathbf{1}_2)$

See [Bellemare et al. \(2017a\)](#).

Distance Between Distributions

Consider sample $\{x_1, \dots, x_n\}$ an i.i.d. sample, with empirical measure $\hat{p}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i}$

Definition 4.18: Divergence based inference

Consider some parametric family $\mathcal{Q} = \{q_\theta, \theta \in \Theta\}$. Given a divergence D , we want to find

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} \{D(\underbrace{p}_{\text{unknown}}, \underbrace{q_\theta}_{q_\theta \in \mathcal{Q}})\}$$

or its empirical version

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \{D(\underbrace{\hat{p}_n}_{\text{estimated}}, q_\theta)\}$$

Distance Between Distributions

Definition 4.19: Unbiased sample gradients, [Bellemare et al. \(2017a\)](#)

A divergence D has **unbiased sample gradients** when the expected gradient of the sample loss equals the gradient of the true loss for all p and n ,

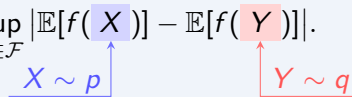
$$\mathbb{E}(\nabla_{\theta} D(\hat{p}_n, q_{\theta})) = \nabla_{\theta} D(p, q_{\theta}).$$

- Then D is a **proper scoring rule** (see [Gneiting and Raftery \(2007\)](#)).
- If this is not satisfied, stochastic gradient descent may not converge...

Distance Between Distributions

Definition 4.20: Integral probability metric, Müller (1997)

Integral probability metrics (IPMs) are distances on the space of distributions over a set \mathcal{X} , defined by a class \mathcal{F} of real-valued functions on \mathcal{X} as

$$D_{\mathcal{F}}(p, q) = \sup_{f \in \mathcal{F}} |\mathbb{E}[f(\mathbf{X})]| - \mathbb{E}[f(\mathbf{Y})]|.$$


Discussed also in [Dedecker and Merlevède \(2007\)](#)

Distance Between Distributions

- Note that it is still possible to define projections with deviance (that will not be "orthogonal" projections since divergence are not related to inner products)

Definition 4.21: Projection, [Bregman \(1967\)](#), [Bauschke et al. \(1997\)](#)

Given a strictly convex function continuously differentiable ψ and the associated Bregman divergence D_ψ , a closed convex $K \subset \mathcal{X}$ and a point $\mathbf{x} \in \mathcal{X}$. The Bregman projection of \mathbf{x} onto K is

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{y} \in K} \{D_\psi(\mathbf{x}, \mathbf{y})\}$$

- If $\psi(\mathbf{x}) = \|\mathbf{x}\|_{\ell_2}^2$, Bregman projection is the standard orthogonal projection onto a convex set,

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{y} \in K} \{\|\mathbf{x} - \mathbf{y}\|_{\ell_2}^2\}$$

Distance Between Distributions

► With Bregman divergence D_ψ , we have a generalized version of the Pythagorean theorem

$$D_\psi(\mathbf{x}, \mathbf{y}) \geq D_\psi(\mathbf{x}, \mathbf{x}^*) + D_\psi(\mathbf{x}^*, \mathbf{y})$$

Numerically, one can use (cyclical) Dykstra algorithm with Bregman projections (see [Censor and Reich \(1998\)](#), [Bauschke and Lewis \(2000\)](#)) to compute \mathbf{x}^* : suppose that

$K = \bigcap_{i=1}^m K_i$ where K_i 's are convex sets (e.g. half-planes - when K is some polyhedral).

Let P_i^ψ denote the orthogonal on K_i based on D_ψ ,

$$P_i^\psi : \mathbf{x} \mapsto \operatorname{argmin}_{\mathbf{y} \in K_i} \{ \|\mathbf{x} - \mathbf{y}\|_{\ell_2}^2 \}$$

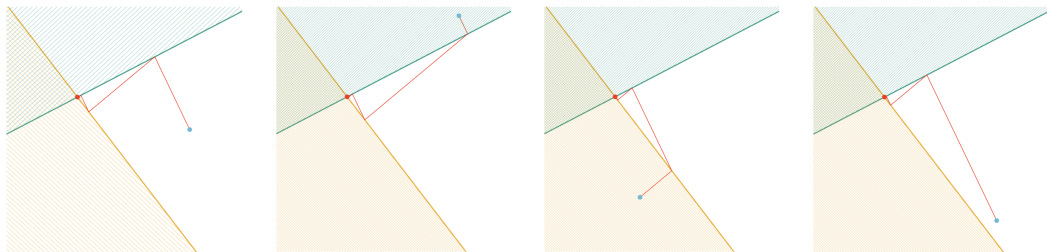
and consider the following iterative sequence of projections, until some $\mathbf{x}_j \in K$,

$$\mathbf{x}_0 \xrightarrow{P_1^\psi} \mathbf{x}_1 \xrightarrow{P_2^\psi} \mathbf{x}_2 \xrightarrow{P_3^\psi} \cdots \xrightarrow{P_{m-1}^\psi} \mathbf{x}_{m-1} \xrightarrow{P_m^\psi} \mathbf{x}_m \xrightarrow{P_1^\psi} \mathbf{x}_{m+1} \xrightarrow{P_2^\psi} \mathbf{x}_{m+2} \cdots$$

Distance Between Distributions

$$\mathbf{x}_0 \xrightarrow{P_1^\psi} \mathbf{x}_1 \xrightarrow{P_2^\psi} \mathbf{x}_2 \xrightarrow{P_3^\psi} \cdots \xrightarrow{P_{m-1}^\psi} \mathbf{x}_{m-1} \xrightarrow{P_m^\psi} \mathbf{x}_m \xrightarrow{P_1^\psi} \mathbf{x}_{m+1} \xrightarrow{P_2^\psi} \mathbf{x}_{m+2} \cdots$$

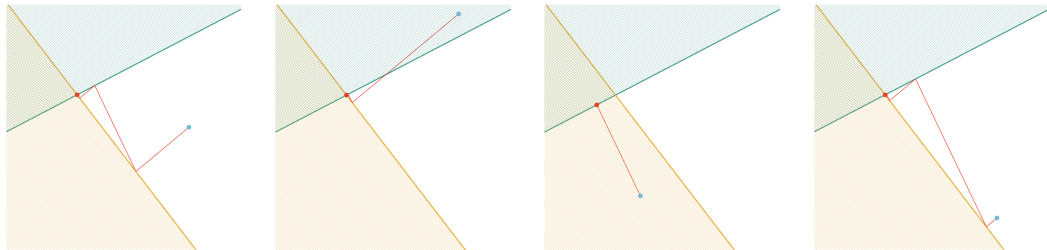
until some $\mathbf{x}_j \in K$, see [Boyle and Dykstra \(1986\)](#) for the original idea



Distance Between Distributions

$$\mathbf{x}_0 \xrightarrow{P_1^\psi} \mathbf{x}_1 \xrightarrow{P_2^\psi} \mathbf{x}_2 \xrightarrow{P_3^\psi} \dots \xrightarrow{P_{m-1}^\psi} \mathbf{x}_{m-1} \xrightarrow{P_m^\psi} \mathbf{x}_m \xrightarrow{P_1^\psi} \mathbf{x}_{m+1} \xrightarrow{P_2^\psi} \mathbf{x}_{m+2} \dots$$

until some $\mathbf{x}_j \in K$, see [Boyle and Dykstra \(1986\)](#) for the original idea



with half spaces, $\|\mathbf{x}_j - \mathbf{x}^*\|_{\ell_2} \leq cr^j \|\mathbf{x} - \mathbf{x}^*\|_{\ell_2}$ for some $r \in (0, 1)$ and $c > 0$ (or "linear convergence" since $\|\mathbf{x}_j - \mathbf{x}^*\|_{\ell_2} \leq r \|\mathbf{x}_{j-1} - \mathbf{x}^*\|_{\ell_2}$).

Distance Between Distributions

Definition 4.22: Hellinger distance, [Hellinger \(1909\)](#)

For two discrete distributions p and q , [Hellinger distance](#) is

$$d_H(p, q)^2 = \frac{1}{2} \sum_i \left(\sqrt{p(i)} - \sqrt{q(i)} \right)^2 = 1 - \sum_i \sqrt{p(i)q(i)} \in [0, 1],$$

and for absolutely continuous distributions, if p and q are densities,

$$d_H(p, q)^2 = \frac{1}{2} \int_{\mathbb{R}} \left(\sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx \text{ or } \frac{1}{2} \int_{\mathbb{R}^k} \left(\sqrt{p(\mathbf{x})} - \sqrt{q(\mathbf{x})} \right)^2 d\mathbf{x}.$$

See [Pardo \(2018\)](#).

Distance Between Distributions

Proposition 4.9: Distance between Beta variables

Consider two Beta distribution, then $d_H^2(\mathcal{B}(a_1, b_1), \mathcal{B}(a_2, b_2))$ is

$$1 - \frac{1}{\sqrt{B(a_1, b_1)B(a_2, b_2)}} B\left(\frac{a_1 + a_2}{2}, \frac{b_1 + b_2}{2}\right)$$

Proof

$$1 - \int_0^1 \sqrt{f_1(t)f_2(t)} dt = 1 - \frac{1}{\sqrt{B(a_1, b_1)B(a_2, b_2)}} \int_0^1 t^{(a_1+a_2)/2-1} (1-t)^{(b_1+b_2)/2-1} dt,$$

$$\text{then use } B(a, b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}.$$

Distance Between Distributions

Proposition 4.10: Distance between Gaussian vectors

Consider two Gaussian distributions, then $d_{\text{H}}^2(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2))$ is

$$2 - 2 \frac{|\boldsymbol{\Sigma}_1|^{\frac{1}{4}} |\boldsymbol{\Sigma}_2|^{\frac{1}{4}}}{|\bar{\boldsymbol{\Sigma}}|^{\frac{1}{2}}} \exp\left(-\frac{1}{8} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^{\top} \bar{\boldsymbol{\Sigma}}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)\right)$$

where $\bar{\boldsymbol{\Sigma}} = \frac{1}{2}(\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)$.

Note that it is a Bregman divergence D_{ψ} with $\psi(\mathbf{x}) = \sum_{i=1}^2 x_i^2$

Distance Between Distributions

Definition 4.23: Pearson/Neyman χ -square divergences Nielsen and Nock (2013)

For two discrete distributions p and q , Pearson chi-square divergence is

$$d_{P\chi}(p\|q)^2 = \sum_i \frac{[p(i) - q(i)]^2}{q(i)},$$

while Neyman chi-square divergence is

$$d_{N\chi}(p\|q)^2 = \sum_i \frac{[p(i) - q(i)]^2}{p(i)} = d_{P\chi}(q\|p),$$

Distance Between Distributions

- Note that both are Bregman divergences D_ψ with $\psi_P(\mathbf{x}) = -2 \sum_{i=1} \sqrt{x_i}$ and

$$\psi_N(\mathbf{x}) = \sum_{i=1} x_i^{-1}.$$

- d_χ can be extended to the case of continuous distributions, e.g.,

$$d_{P_\chi}(p\|q)^2 = \int \left(\frac{p(x)}{q(x)} - 1 \right)^2 p(x) dx$$

Distance Between Distributions

Definition 4.24: Total Variation, [Jordan \(1881\)](#); [Rudin \(1966\)](#)

For two distributions p and q , the **total variation distance** between p and q is

$$d_{\text{TV}}(p, q) = \sup_{\mathcal{A}} \{ |p(\mathcal{A}) - q(\mathcal{A})| \}.$$

Proposition 4.11: Total Variation

For two univariate distributions p and q , the **total variation distance** between p and q is

$$d_{\text{TV}}(p, q) = \frac{1}{2} \sum_i |p(i) - q(i)| = \frac{1}{2} \|p - q\|_{\ell_1} = \sum_{i:p(i) \geq q(i)} (p(i) - q(i))$$

See Proposition 4.2 in [Levin and Peres \(2017\)](#).

Distance Between Distributions

- › Equivalently,

$$d_{\text{TV}}(p, q) = \frac{1}{2} \sup_{f: \mathbb{R}^k \rightarrow \{0,1\}} \left\{ \int f d p - \int f d q \right\}$$

(see e.g. <https://djalil.chafai.net/blog/>, with $f: \mathbb{R}^k \rightarrow \{-1, 1\}$, $f = \mathbf{1}_{\mathcal{A}} - \mathbf{1}_{\mathcal{A}^c}$)

- › It is an IPM with $\mathcal{F} = \{f: \mathcal{X} \rightarrow \{0, 1\}\}$, so that \mathcal{F} is a set of indicator functions for any event.
- › For Gaussian distributions, the distance has no explicit formula, see, e.g., [Devroye et al. \(2018\)](#).

Distance Between Distributions

Proposition 4.12: Total Variation, Scheffé theorem, Billingsley (2017)

For two distributions p and q on \mathbb{R}^k ,

$$d_{\text{TV}}(p, q) = \frac{1}{2} \int_{\mathbb{R}^d} |p(\mathbf{x}) - q(\mathbf{x})| d\mathbf{x},$$

$$d_{\text{TV}}(p, q) = 1 - \int_{\mathbb{R}^d} \min \{p(\mathbf{x}) - q(\mathbf{x})\} d\mathbf{x},$$

$$d_{\text{TV}}(p, q) = p(\mathcal{A}) - q(\mathcal{A}) \text{ where } \mathcal{A} = \{\mathbf{x} : p(\mathbf{x}) \geq q(\mathbf{x})\}.$$

Distance Between Distributions

In the univariate case, we can restrict \mathcal{A} to half-lines $(-\infty, t]$

Definition 4.25: Kolmorov-Smirnov, Kolmogorov (1933); Smirnov (1948)

For two distributions p and q , **Kolmorov-Smirnov distance** between p and q is

$$d_{\text{KS}}(p, q) = \sup_{t \in \mathbb{R}} \{ |p((-\infty, t]) - q((-\infty, t])| \} = \sup_{t \in \mathbb{R}} \{ |F_p(t) - F_q(t)| \} = \|F_p - F_q\|_{\infty},$$

where F_p and F_q are the respective cumulative distribution functions.

Distance Between Distributions

Definition 4.26: Entropy, Shannon (1948)

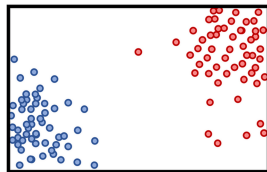
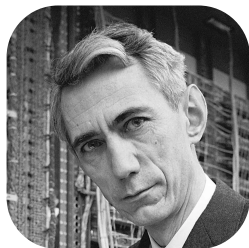
The entropy associated with distribution p is

$$\mathcal{E}_p(p) = - \sum_i p(i) \log p(i) = \mathbb{E}_p[-\log p(X)].$$

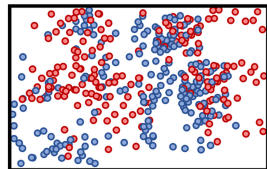
and define cross-entropy (of q relative to p) as

$$\mathcal{E}_q(p) = - \sum_i p(i) \log q(i) = \mathbb{E}_p[-\log q(X)].$$

See [Amari \(2016\)](#) or [Chambert-Loir \(2023\)](#) for more details.



Low Entropy



High Entropy

Distance Between Distributions

Definition 4.27: Kullback–Leibler, [Kullback and Leibler \(1951\)](#)

For two discrete distributions p and q , [Kullback–Leibler divergence](#) of p , with respect to q is

$$D_{\text{KL}}(p\|q) = \sum_i p(i) \log \frac{p(i)}{q(i)},$$

and for absolutely continuous distributions,

$$D_{\text{KL}}(p\|q) = \int_{\mathbb{R}} p(x) \log \frac{p(x)}{q(x)} dx \text{ or } \int_{\mathbb{R}^k} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x},$$

in higher dimension.

Also called relative entropy, since $D_{\text{KL}}(p\|q) = \mathcal{E}_q(p) - \mathcal{E}_p(p)$.

Distance Between Distributions

Proposition 4.13: Divergence for Gaussian vectors

Consider two Gaussian distributions, then $D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \parallel \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2))$ is

$$\frac{1}{2} \left[(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) - \log \frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_2|} - k \right]$$

where k is the dimension, see [Polyanskiy and Wu \(2022\)](#).

Distance Between Distributions

The entropy of X according to p is smaller than or equal to the cross-entropy of p and q , or equivalently

Proposition 4.14: Gibbs' inequality

$D_{\text{KL}}(p||q)$ is positive and separable, i.e. $D_{\text{KL}}(p||q) \geq 0$ and $D_{\text{KL}}(p||q) = 0$ if and only if $p = q$.

Proof: $\sum_{x \in I} p(x) \log \frac{p(x)}{q(x)} \geq 0$ where I is the set of all x for which $p(x) > 0$. Recall that $\log x \leq x - 1$ (with equality only when $x = 1$), thus $\log(1/x) \geq 1 - x$, and

$$\sum_{x \in I} p(x) \ln \frac{p(x)}{q(x)} \geq \sum_{x \in I} p(x) \left(1 - \frac{q(x)}{p(x)}\right) = \underbrace{\sum_{x \in I} p(x)}_{=1} - \underbrace{\sum_{x \in I} q(x)}_{\leq 1} \geq 0.$$

Distance Between Distributions

Proposition 4.15: Additivity for independence distributions

$D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) = D_{\text{KL}}(p_x \parallel q_x) + D_{\text{KL}}(p_y \parallel q_y)$ if $\mathbf{p}(x, y) = p_x(x)p_y(y)$ and $\mathbf{q}(x, y) = q_x(x)q_y(y)$.

Proof By definition

$$D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) = \int_{\mathcal{X}} \int_{\mathcal{Y}} p(x, y) \cdot \log \frac{p(x, y)}{q(x, y)} dy dx .$$

and since $\mathbf{p}(x, y) = p_x(x)p_y(y)$ and $\mathbf{q}(x, y) = q_x(x)q_y(y)$,

$$D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) = \int_{\mathcal{X}} \int_{\mathcal{Y}} p_1(x) p_2(y) \cdot \log \frac{p_1(x) p_2(y)}{q_1(x) q_2(y)} dy dx .$$

Distance Between Distributions

$$\begin{aligned}D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q}) &= \int_{\mathcal{X}} \int_{\mathcal{Y}} p_x(x) p_y(y) \cdot \left(\log \frac{p_x(x)}{q_x(x)} + \log \frac{p_y(y)}{q_y(y)} \right) dy dx \\&= \int_{\mathcal{X}} \int_{\mathcal{Y}} p_x(x) p_y(y) \cdot \log \frac{p_x(x)}{q_x(x)} dy dx + \int_{\mathcal{X}} \int_{\mathcal{Y}} p_x(x) p_y(y) \cdot \log \frac{p_y(y)}{q_y(y)} dy dx \\&= \int_{\mathcal{X}} p_x(x) \cdot \log \frac{p_x(x)}{q_x(x)} \int_{\mathcal{Y}} p_y(y) dy dx + \int_{\mathcal{Y}} p_y(y) \cdot \log \frac{p_y(y)}{q_y(y)} \int_{\mathcal{X}} p_x(x) dx dy \\&= \int_{\mathcal{X}} p_x(x) \cdot \log \frac{p_x(x)}{q_x(x)} dx + \int_{\mathcal{Y}} p_y(y) \cdot \log \frac{p_y(y)}{q_y(y)} dy \\&= D_{\text{KL}}(p_x \parallel q_x) + D_{\text{KL}}(p_y \parallel q_y).\end{aligned}$$

But for other distances,

$$\begin{cases}d_{\text{H}}(\mathbf{p}, \mathbf{q})^2 \leq d_{\text{H}}(p_x, q_x)^2 + d_{\text{H}}(p_y, q_y)^2 \\d_{\text{TV}}(\mathbf{p}, \mathbf{q}) \leq d_{\text{TV}}(p_x, q_x) + d_{\text{TV}}(p_y, q_y).\end{cases}$$

Distance Between Distributions

- It is only defined in this way if, for all x , $q(x) = 0$ implies $p(x) = 0$ (“absolute continuity” with respect to p).

Proposition 4.16

The KL divergence has unbiased sample gradients, but is not scale sensitive.

Proof [Bellemare et al. \(2017b\)](#).

- In a Bayesian setting, $D_{\text{KL}}(p||q)$ is a measure of the information gained by revising one’s beliefs from the prior probability distribution q to the posterior probability distribution p (it is the amount of information lost when q is used to approximate p).
- If $\psi(\mathbf{x}) = \sum x_i \log(x_i)$ (strictly convex), then Bregman divergence is

$$D_{\psi}(\mathbf{x}, \mathbf{y}) = \sum x_i \log \frac{x_i}{y_i} = D_{\text{KL}}(\mathbf{x}||\mathbf{y})$$

Distance Between Distributions

$$D_{\text{KL}}(\mathcal{B}(p) \parallel \mathcal{B}(q)) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$$

$$D_{\text{KL}}(\mathcal{B}(n, p) \parallel \mathcal{B}(n, q)) = np \log \frac{p}{q} + n(1-p) \log \frac{1-p}{1-q} = n D_{\text{KL}}(\mathcal{B}(p) \parallel \mathcal{B}(q))$$

$$D_{\text{KL}}(\mathcal{U}([a_1, b_1]) \parallel \mathcal{U}([a_2, b_2])) = \log \frac{b_2 - a_2}{b_1 - a_1}$$

$$D_{\text{KL}}(\mathcal{N}(\mu_1, \sigma_1^2) \parallel \mathcal{N}(\mu_2, \sigma_2^2)) = \frac{1}{2} \left[\frac{(\mu_1 - \mu_2)^2}{\sigma_2^2} + \frac{\sigma_1^2}{\sigma_2^2} - \log \frac{\sigma_1^2}{\sigma_2^2} - 1 \right]$$

$$D_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \parallel \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) = \frac{1}{2} \left[(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) - \log \frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_2|} - n \right]$$

Distance Between Distributions

- Consider some distribution p_θ , as in [Nielsen \(2022\)](#). Using Taylor expansion,

$$D_{\text{KL}}(p_\theta \| p_{\theta+d\theta}) = \frac{1}{2} d\theta^\top I(\theta) d\theta \approx \frac{1}{2} ds_\theta^2.$$

Definition 4.28: Jeffreys (symmetric) divergence [Jeffreys \(1946\)](#)

The [Jeffrey divergence](#) is a symmetric divergence induced by Kullback-Liebler divergence,

$$D_J(p_1, p_2) = \frac{1}{2} D_{\text{KL}}(p_1 \| p_2) + \frac{1}{2} D_{\text{KL}}(p_2 \| p_1).$$

Distance Between Distributions

Definition 4.29: Jensen-Shannon, Lin (1991)

The **Jensen-Shannon divergence** is a symmetric divergence induced by Kullback-Liebler divergence,

$$D_{\text{JS}}(p_1, p_2) = \frac{1}{2}D_{\text{KL}}(p_1 \| q) + \frac{1}{2}D_{\text{KL}}(p_2 \| q),$$

where $q = \frac{1}{2}(p_1 + p_2)$.

Endres and Schindelin (2003) proved that $\sqrt{D_{\text{JS}}(p_1, p_2)}$ is a proper distance.

› See `philentropy` package.

Distance Between Distributions

Definition 4.30: f -divergence, Rényi (1961), Ali and Silvey (1966)

Given a continuous convex function $f : [0, \infty) \rightarrow \overline{\mathbb{R}}$, define

$$D_f(p\|q) = \sum_i q(i) \cdot f\left(\frac{p(i)}{q(i)}\right)$$

and for absolutely continuous function

$$D_f(p\|q) = \int_{\mathbb{R}} q(x) f\left(\log \frac{p(x)}{q(x)}\right) dx \text{ or } \int_{\mathbb{R}^k} q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x},$$

➤ $D_f(p\|q)$ is properly defined when $p \ll q$, see also Csiszár (1964, 1967).

If $f(u) = u \log u$, $D_f(p\|q) = D_{\text{KL}}(p, q)$

If $f(u) = |u - 1|$, $D_f(p\|q) = d_{\text{TV}}(p, q)$

Distance Between Distributions

$$\text{If } f(u) = \frac{1}{2}(\sqrt{u} - 1)^2, D_f(p\|q) = d_H(p, q)^2$$

$$\text{If } f(u) = \frac{1}{2} \left(u \log u - (u + 1) \log \left(\frac{u + 1}{2} \right) \right), D_f(p\|q) = d_{JS}(p, q)$$

► One can define $D_f(p\|q)$ when $p \ll q$: Since f is convex, and $f(1) = 0$, the function $\frac{f(x)}{x-1}$ must nondecrease, so there exists $f'(\infty) := \lim_{x \rightarrow \infty} f(x)/x$, taking value in $(-\infty, +\infty]$. And since for any $p(x) > 0$, we have $\lim_{q(x) \rightarrow 0} q(x) f\left(\frac{p(x)}{q(x)}\right) = p(x) f'(\infty)$.

Proposition 4.17

$$D_f(p\|q) \text{ is linear in } f, D_{af+bg}(p\|q) = aD_f(p, q) + bD_g(p\|q).$$

Distance Between Distributions

Proposition 4.18

$D_f = D_g$ if and only if $f(x) = g(x) + c(x - 1)$ for some $c \in \mathbb{R}$.

- The only f -divergence that is also a Bregman ψ -divergence is the KL divergence
- The only f -divergence that is also an integral probability metric is the total variation.
- There is a [variational representation](#) of D_f , in [Polyanskiy and Wu \(2022\)](#).

Distance Between Distributions

- Since f is convex, let f^* be the convex conjugate of f . Let $\text{effdom}(f^*)$ be the effective domain of f^* (i.e., $\text{effdom}(f^*) = \{y : f^*(y) < \infty\}$)

$$D_f(p; q) = \sup_{g: \Omega \rightarrow \text{effdom}(f^*)} \mathbb{E}_p[g] - \mathbb{E}_q[f^* \circ g]$$

- For example, with the total variation, $f(x) = \frac{1}{2}|x - 1|$, its convex conjugate is

$$f^*(x^*) = \begin{cases} x^* & \text{on } [-1/2, 1/2], \\ +\infty & \text{else.} \end{cases}, \text{ and we obtain}$$

$$d_{\text{TV}}(p, q) = \sup_{|g| \leq 1/2} \mathbb{E}_p[g(X)] - \mathbb{E}_q[g(X)].$$

Distance Between Distributions

- Extending Rényi entropy of order α , $H_\alpha(X) = \frac{1}{1-\alpha} \log \left(\sum_i p(i)^\alpha \right)$, define

Definition 4.31: Rényi α -divergence, Rényi (1961)

Given $\alpha \in (0, \infty)$, define

$$D_\alpha(p\|q) = \frac{1}{\alpha - 1} \log \left(\sum_i \frac{p(i)^\alpha}{q(i)^{\alpha-1}} \right)$$

and for absolutely continuous function

$$D_\alpha(p\|q) = \frac{1}{\alpha - 1} \log \left(\int_{\mathbb{R}} \frac{p(x)^\alpha}{q(x)^{\alpha-1}} dx \right) \text{ or } \frac{1}{\alpha - 1} \log \left(\int_{\mathbb{R}^k} \frac{p(\mathbf{x})^\alpha}{q(\mathbf{x})^{\alpha-1}} d\mathbf{x} \right).$$

Distance Between Distributions

- › Recall that

$$D_\alpha(p\|q) = \frac{1}{\alpha - 1} \log \left(\sum_i \frac{p(i)^\alpha}{q(i)^{\alpha-1}} \right) \text{ when } \alpha \in (0, \infty).$$

- › One can define limiting cases, $D_0(P\|Q) = -\log Q(\{i : p_i > 0\})$ and $D_\infty(P\|Q) = \log \sup_i \frac{p_i}{q_i}$
- › Observe also that $D_1(p\|q) = D_{\text{KL}}(p\|q)$

Distance Between Distributions

Definition 4.32: Cramér, [Cramér \(1928a,b\)](#) and [Székely \(2003\)](#)

Consider two measures on p and q on \mathbb{R} . Then define **Cramér distance**

$$C_k(p, q) = \left(\int_{-\infty}^{\infty} |F_p(x) - F_q(x)|^k dx \right)^{1/k}, \text{ for } k \geq 1$$

- C_2 is named "energy-distance" in [Székely \(2003\)](#) and [Rizzo and Székely \(2016\)](#), and "continuous ranked probability score" in [Gneiting et al. \(2007\)](#).
- It is an Integral Probability Metrics (IPM), since

$$C_k(p, q) = \sup_{f \in \mathcal{F}_{k'}} |\mathbb{E}[f(X)] - \mathbb{E}[f(Y)]|.$$

(Note: In the original image, 'k' is in a green box, 'X' is in a blue box, 'Y' is in a red box, and 'f ∈ F_{k'}' is in a green box. Arrows point from 'k' to 'k^{-1} + k'^{-1} = 1', from 'X' to 'X ~ p', and from 'Y' to 'Y ~ q'.)

where $\mathcal{F}_{k'}$ is the set of absolutely continuous functions such that $\|\nabla f\|_{k'} \leq 1$.

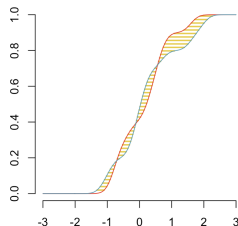
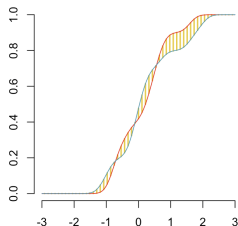
- For example, if $k = 1$, $\|\nabla f\|_{\infty} \leq 1$ (corresponding to 1-Lipschitz functions).

Distance Between Distributions

Definition 4.33: Wasserstein, Wasserstein (1969)

Consider two measures on p and q on \mathbb{R} . Then define Wasserstein distance

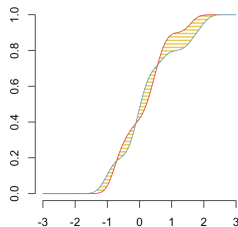
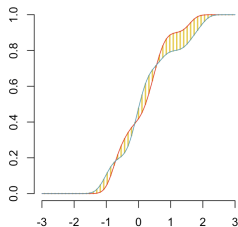
$$W_k(p, q) = \left(\int_0^1 |F_p^{-1}(u) - F_q^{-1}(u)|^k du \right)^{1/k}, \text{ for } k \geq 1$$



```
1 > c2 = function(x) (pnorm(x,0,1)-pnorm(x,1,2))^2
2 > w2 = function(u) (qnorm(u,0,1)-qnorm(u,1,2))^2
3 > sqrt(integrate(c2,-Inf,Inf)$value)
4 [1] 0.5167714
5 > sqrt(integrate(w2,0,1)$value)
6 [1] 1.414214
```

where F^{-1} denotes the generalized inverse of F , $F^{-1}(u) = \inf_{x \in \mathbb{R}} \{F(x) \geq u\}$.

Distance Between Distributions



```
1 > c1 = function(x) abs(pnorm(x,0,1) -  
2   pnorm(x,1,2))  
3 > w1 = function(x) abs(qnorm(x,0,1) -  
4   qnorm(x,1,2))  
5 > integrate(c1,-Inf,Inf)$value  
6 [1] 1.166631  
7 > integrate(w1,0,1)$value  
8 [1] 1.166636
```

Proposition 4.19: C_1 and W_1

Consider two measures on p and q on \mathbb{R} .

$$W_1(p, q) = \int_0^1 |F_p^{-1}(u) - F_q^{-1}(u)| du = \int_{-\infty}^{\infty} |F_p(x) - F_q(x)| dx = C_1(p, q).$$

Proof See [Prokhorov \(1956\)](#), [Dall'Aglio \(1956\)](#) and [Vallender \(1974\)](#).

Distance Between Distributions

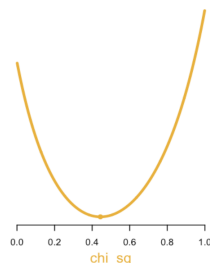
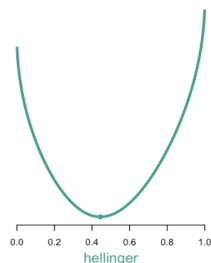
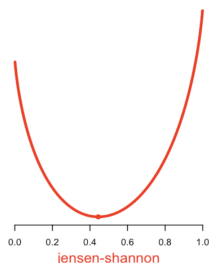
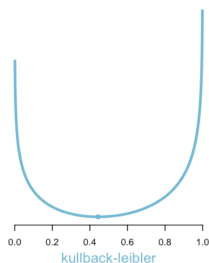
Instead of the geometric proof (see plot above), observe that

$$\begin{aligned}\int_0^1 |F_p^{-1}(u) - F_q^{-1}(u)| du &= \int_0^1 \int_{-\infty}^{\infty} g(u, x) dx du, \quad g(u, x) = 1 \text{ if } \begin{cases} x \in [F_p^{-1}(u), F_q^{-1}(u)] \\ x \in [F_q^{-1}(u), F_p^{-1}(u)] \end{cases} \\ &= \int_{-\infty}^{\infty} \int_0^1 h(u, x) du dx, \quad h(u, x) = 1 \text{ if } \begin{cases} u \in [F_p(x), F_q(x)] \\ u \in [F_q(x), F_p(x)] \end{cases} \\ &= \int_{-\infty}^{\infty} |F_p(x) - F_q(x)| dx\end{aligned}$$

(see Proposition 2.17 in [Santambrogio \(2015\)](#) for a proper justification)

Distance Between Distributions

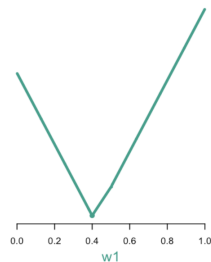
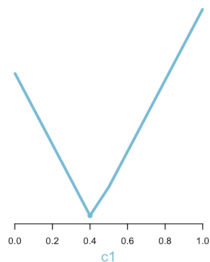
- μ : multinomial distribution on $\{0, 1, 10\}$, with $\mathbf{p} = (.5, .1, .4)$
- ν_θ : binomial type distribution on $\{0, 10\}$, with $\mathbf{q}_\theta = (1 - \theta, \theta)$
- Let $\theta^* = \operatorname{argmin}\{d(p, q_\theta)\}$ or $\theta^* = \operatorname{argmin}\{d(p\|q_\theta)\}$



- with $d_{\text{KL}}(p\|q_\theta)$, $d_{\text{JS}}(p, q_\theta)$, $d_{\text{H}}(p, q_\theta)$ and $d_{\text{H}_{\chi^2}}(p\|q_\theta)$.

Distance Between Distributions

- μ : multinomial distribution on $\{0, 1, 10\}$, with $\mathbf{p} = (.5, .1, .4)$
- ν_θ : binomial type distribution on $\{0, 10\}$, with $\mathbf{q}_\theta = (1 - \theta, \theta)$
- Let $\theta^* = \operatorname{argmin}\{d(\mathbf{p}, \mathbf{q}_\theta)\}$



- with $C_1(p, q_\theta)$, $C_2(p, q_\theta)$, $W_1(p, q_\theta)$ and $W_2(p, q_\theta)$.

Distance Between Distributions

Proposition 4.20

The Wasserstein metric is scale and sum invariant, but does not have unbiased sample gradients.

Proof [Bellemare et al. \(2017b\)](#)

Example If x_i are drawn from a Bernoulli distribution

› Non-vanishing minimax bias:

$$\forall n, \exists p, q_\theta, |\mathbb{E}(\nabla_\theta W_k^k(\hat{p}_n, q_\theta)) - \nabla_\theta W_k^k(p, q_\theta)| \geq 2e^{-2}$$

› Wrong minimum: in general,

$$\hat{\theta}_n = \operatorname{argmin} \left\{ \mathbb{E}((W_k^k(\hat{p}_n, q_\theta))) \right\} \neq \operatorname{argmin} \left\{ W_k^k(\mathbb{P}, \mathbb{Q}_\theta) \right\} = \theta$$

Distance Between Distributions

Proposition 4.21

The Cramér metric is scale and sum invariant.

- › $C_k(X + Z, Y + Z) \leq C_k(X, Y)$ whenever $Z \perp\!\!\!\perp X, Y$ and $k \geq 1$, and $C_k(cX, cY) \leq |c|^{1/k} C_k(X, Y)$.

Proposition 4.22

C_2 has unbiased sample gradients (only $k = 2$),

$$\mathbb{E}(\nabla_{\theta} C_2(\hat{p}_n, q_{\theta})) = \nabla_{\theta} C_2(p, q_{\theta}).$$

Distance Between Distributions

- Consider first W_1 (earth mover's distance), which was the only distance discussed in [Wasserstein \(1969\)](#). See also [Vallender \(1974\)](#) for an extensive review.
- W_1 is an IPM where \mathcal{F} the set of 1-Lipschitz functions, [Kantorovich and Rubinstein \(1958\)](#), i.e., if p and q have bounded support,

$$W_1(p, q) = \sup_{f \in \mathcal{F}} \left\{ \int_{-\infty}^{+\infty} f(x) d(p - q)(x) \right\},$$

\mathcal{F} being the class of 1-Lipschitz functions

Proposition 4.23: W_1 and First Order Dominance

Suppose that $X_1 \preceq X_2$ (first order dominance, $F_2^{-1}(u) \geq F_1^{-1}(u), \forall u \in (0, 1)$),

$$W_1(p_1, p_2) = \mathbb{E}[X_2] - \mathbb{E}[X_1].$$

Distance Between Distributions

Proof

$$W_1(p_1, p_2) = \int_0^1 |F_2^{-1}(u) - F_1^{-1}(u)| du = \int_0^1 F_2^{-1}(u) du - \int_0^1 F_1^{-1}(u) du$$

(Note: In the original image, a red arrow labeled ≥ 0 points to the absolute value term. A blue arrow labeled $\mathbb{E}[X_2]$ points to the first integral, and a green arrow labeled $\mathbb{E}[X_1]$ points to the second integral.)

then (property discussed later)

$$W_1(p_1, p_2) = \inf_C \int \int |x_2 - x_1| dC(F_1(x_1), F_2(x_2)) = \inf_C \int \int |F_2^{-1}(v) - F_1^{-1}(u)| dC(u, v)$$

(Note: In the original image, a red arrow labeled $\mathbb{E}[|X_1 - X_2|]$ points to the absolute value term in the second integral.)

As discussed in [Vallender \(1974\)](#),

$$\begin{aligned} \mathbb{E}[|X_1 - X_2|] &= \int [\mathbb{P}[X_1 < t, X_2 \geq t] + \mathbb{P}[X_1 \geq t, X_2 < t]] dt \\ &= \int [\mathbb{P}[X_1 < t] + \mathbb{P}[X_2 < t] - 2\mathbb{P}[X_1 < t, X_2 < t]] dt \end{aligned}$$

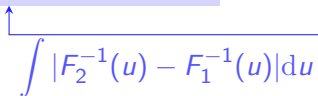
Distance Between Distributions

$$\mathbb{E}[|X_1 - X_2|] = \int [F_1(t) + F_2(t) - 2C(F_1(t), F_2(t))] dt$$

From Fréchet-Hoeffding bounds, $C(u, v) \leq M(u, v) = \min\{u, v\}$ and

$$F_1(t) + F_2(t) - 2C(F_1(t), F_2(t)) \geq F_1(t) + F_2(t) - 2M(F_1(t), F_2(t))$$

$$\mathbb{E}[|X_1 - X_2|] \geq \int \int |F_2^{-1}(v) - F_1^{-1}(u)| dM(u, v)$$


$$\int |F_2^{-1}(u) - F_1^{-1}(u)| du$$

Example let $p_1 \leq p_2$

$$W_1(\mathcal{B}(p_1), \mathcal{B}(p_2)) = p_2 - p_1.$$

Distance Between Distributions

- › We can also consider W_2

Proposition 4.24: C_2 and W_2

Consider two measures on p and q on \mathbb{R} .

$$W_2(p, q)^2 = \int_0^1 |F_p^{-1}(u) - F_q^{-1}(u)|^2 du \text{ while } C_2(p, q) = \int_{-\infty}^{\infty} |F_p(x) - F_q(x)|^2 dx.$$

Distance Between Distributions

Proposition 4.25: W_2 for Gaussian / Bernoulli distributions

Consider two Gaussian distributions, then

$$W_2(\mathcal{N}(\mu_1, \sigma_1^2), \mathcal{N}(\mu_2, \sigma_2^2))^2 = (\mu_1 - \mu_2)^2 + (\sigma_1 - \sigma_2)^2,$$

and for two Bernoulli distributions, if $p_1 \leq p_2$

$$W_2(\mathcal{B}(p_1), \mathcal{B}(p_2)) = \sqrt{p_2 - p_1}.$$

Distance Between Distributions

Proposition 4.26: Representation for W_2

Consider two measures on p and q on \mathbb{R} .

$$W_2(p, q)^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (F_p(\min\{x, y\}) - F_q(\max\{x, y\}))_+ + (F_q(\min\{x, y\}) - F_p(\max\{x, y\}))_+ dx dy$$

or

$$W_2(p, q)^2 = 2 \int_{-\infty}^{\infty} \int_x^{\infty} [(F_p(x) - F_q(y))_+ + (F_q(x) - F_p(y))_+] dx dy$$

Proof Since $W_2(p, q)^2 = \int_0^1 |F_p^{-1}(u) - F_q^{-1}(u)|^2 du$ observe that

$$F_p^{-1}(u) - F_q^{-1}(u) = F_p^{-1}(u) - F_p^{-1}(F_p(F_q^{-1}(u))) = F_p^{-1}(u) - F_p^{-1}(G(u)) \text{ where } G = F_p \circ F_q^{-1}.$$

Distance Between Distributions

Suppose that F_q is continuously differentiable, so that $H = F'_p \circ F_p^{-1}$, then

$$F_p^{-1}(u) - F_q^{-1}(u) = \int_{G(u)}^u \frac{dt}{H(t)}$$

and write

$$(F_p^{-1}(u) - F_q^{-1}(u))^2 = \int_{G(u)}^u \frac{dt}{H(t)} \frac{dv}{H(v)}$$

and depending on whether $G(u) \leq u$ or $u \leq G(u)$, we can write

$$\int_0^1 (F_p^{-1}(u) - F_q^{-1}(u))^2 du = \int_0^1 \int_0^1 (G^{-1}(\min\{t, v\}) - \max\{r, v\})_+ + (\min\{t, v\} - G^{-1}(\max\{t, v\}))_+ dt dv.$$

And finally, let $t = F_p(x)$ and $v = F_q(y)$, so that $G^{-1}(t) = F_q(x)$ and $G^{-1}(v) = F_p(y)$, and we get the desired expression.

Distance Between Distributions

- › We can finally consider W_∞

Proposition 4.27: W_∞

Consider two measures on p and q on \mathbb{R} .

$$W_\infty(p, q) = \sup_{u \in (0,1)} |F_p^{-1}(u) - F_q^{-1}(u)|.$$

Furthermore, $W_\infty(p, q)$ is the infimum over all $h \geq 0$ such that

$$F_q(x - h) \leq F_p(x) \leq F_q(x + h), \text{ for all } x \in \mathbb{R}.$$

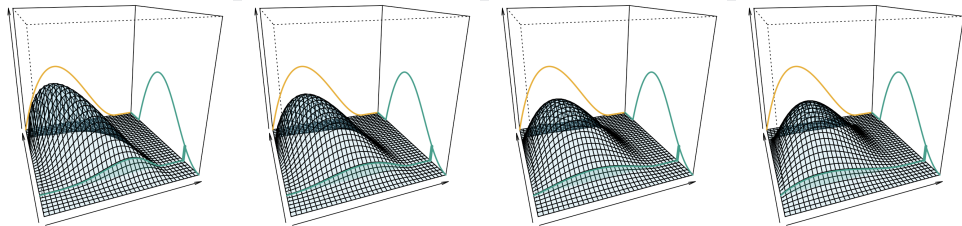
Optimal transport and Wasserstein distance

Definition 4.34: Wasserstein, **Wasserstein (1969)**

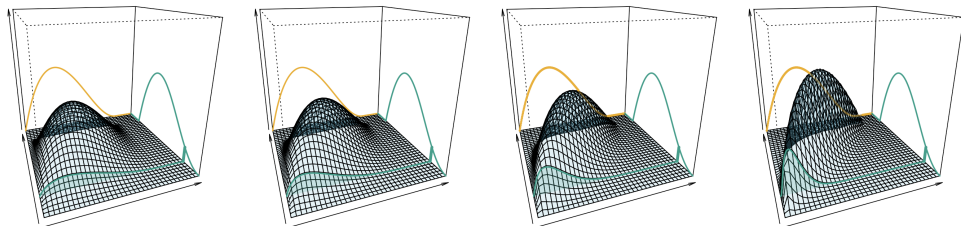
Consider two measures on p and q on \mathbb{R}^k , with a norm $\|\cdot\|$ (on \mathbb{R}^k). Then define

$$W_k(p, q) = \left(\inf_{\pi \in \Pi(p, q)} \int_{\mathbb{R}^k \times \mathbb{R}^k} \|\mathbf{x} - \mathbf{y}\|^k d\pi(\mathbf{x}, \mathbf{y}) \right)^{1/k},$$

where $\Pi(p, q)$ is the set of all couplings of p and q .



Optimal transport and Wasserstein distance



Definition 4.35: Kantorovich Problem

Kantorovich Problem is defined as

$$W_c(p, q) = \inf_{\pi \in \Pi(p, q)} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) d\pi(\mathbf{x}, \mathbf{y}),$$

for cost function c (or loss function).



Optimal transport and Monge mapping

Definition 4.36: Push-Forward and Transport Map

Given two metric spaces \mathcal{X} and \mathcal{Y} , a measurable map $T : \mathcal{X} \rightarrow \mathcal{Y}$ and a measure μ on \mathcal{X} . The **push-forward** of μ by T is the measure $\nu = T_{\#}\mu$ on \mathcal{Y} defined by

$$\forall \mathcal{B} \subset \mathcal{Y}, T_{\#}\mu(\mathcal{B}) = \mu(T^{-1}(\mathcal{B})).$$

› By the change-of-variable formula

Proposition 4.28: Push-Forward and Transport Map

For all measurable and bounded $\varphi : \mathcal{Y} \rightarrow \mathbb{R}$,

$$\int_{\mathcal{Y}} \varphi(y) dT_{\#}\mu(y) = \int_{\mathcal{X}} \varphi(T(x)) d\mu(x).$$

Optimal transport and Monge mapping

- › If \mathcal{Y} is a finite set $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$,

$$T_{\#}\mu = \sum_{i=1}^n \mu(T^{-1}(\{\mathbf{y}_i\})) \cdot \delta_{\{\mathbf{y}_i\}}$$

- › If \mathcal{X} is a single atom, $\{\mathbf{x}\}$, $\mu = \delta_{\mathbf{x}}$ and $T_{\#}\mu(\mathcal{B}) = \mu(T^{-1}(\mathcal{B})) = \delta_{T(\mathbf{x})}$. If $\text{Card}(\text{support}(\nu)) > 1$, there is no transport map.
- › One solution is to allow mass to split, leading to Kantorovich's relaxation of Monge's problem

Proposition 4.29: Existence of a map

If $\mathcal{X} = \mathcal{Y}$ is a compact subset of \mathbb{R}^k , if μ and ν are two measures, and if μ is atomless, then there exists T such that $\nu = T_{\#}\mu$.

see [Santambrogio \(2015\)](#).

Optimal transport and Monge mapping

- If \mathcal{X} and \mathcal{Y} are two sets of \mathbb{R}^k , and if measures μ and ν are absolutely continuous, with densities f and g (w.r.t. Lebesgue measure),

$$\int_{\mathcal{Y}} \varphi(\mathbf{y})g(\mathbf{y})d\mathbf{y} = \int_{\mathcal{X}} \varphi(T(\mathbf{x})) \cdot \underbrace{g(T(\mathbf{x})) \det \nabla T(\mathbf{x})}_{=f(\mathbf{x})} \cdot d\mathbf{x}.$$

Definition 4.37: Monge Problem

Monge problem

$$\inf_{T_{\#P=A}} \int_{\mathcal{X}} c(\mathbf{x}, T(\mathbf{x}))dP_{\mathbf{A}}(\mathbf{x}),$$

for cost function c .



- Note that the constraint and the objective function are non-convex.

Optimal transport and Monge mapping

Theorem 4.1: Optimal map for continuous univariate distributions

The optimal Monge map T^* for some strictly convex cost c such that $T_{\#}^* \mathbb{P}_A = \mathbb{P}_B$ is $T^* = F_B^{-1} \circ F_A$.

› T^* is an **increasing mapping**.

Example Univariate Gaussian

$$x_B = T^*(x_A) = \mu_B + \sigma_B \sigma_A^{-1} (x_A - \mu_A).$$

Optimal transport and Monge mapping

Theorem 4.2: Optimal map for continuous multivariate distributions, Brenier (1991)

With a quadratic cost, the optimal Monge map T^* is unique, and it is the gradient of a convex function, $T^* = \nabla\varphi$.

Example Multidimensional Gaussian

$$\mathbf{x}_B = \mathcal{T}^*(\mathbf{x}_A) = \boldsymbol{\mu}_B + \mathbf{A}(\mathbf{x}_A - \boldsymbol{\mu}_A),$$

where \mathbf{A} is a symmetric positive matrix that satisfies $\mathbf{A}\boldsymbol{\Sigma}_A\mathbf{A} = \boldsymbol{\Sigma}_B$, which has a unique solution given by $\mathbf{A} = \boldsymbol{\Sigma}_A^{-1/2}(\boldsymbol{\Sigma}_A^{1/2}\boldsymbol{\Sigma}_B\boldsymbol{\Sigma}_A^{1/2})^{1/2}\boldsymbol{\Sigma}_A^{-1/2}$, where $\mathbf{M}^{1/2}$ is the square root of the square (symmetric) positive matrix \mathbf{M} based on the Schur decomposition ($\mathbf{M}^{1/2}$ is a positive symmetric matrix), as described in Higham (2008).

Optimal transport and Monge mapping

Gangbo (1999) proved, when $\mathcal{X} = \mathcal{Y}$ is a compact subset of \mathbb{R} , the infimum in Monge problem and the minimum in Kantorovich problem coincide, if μ is atomless,

Proposition 4.30: Monge/Kantorovich Problems

$\mathcal{X} = \mathcal{Y}$ is a compact subset of \mathbb{R}^k and if μ is atomless,

$\min\{\text{Monge problem, see Def. 4.37}\} = \min\{\text{Kantorovich problem, see Def. 8.16}\}.$

Optimal transport (discrete)

- One can consider optimal transport for empirical measures, $\mathbb{P} = \sum_{i=1}^n \omega_i \delta_{\mathbf{x}_i}$.
- With uniform weights and n points for \mathbb{P}_A and \mathbb{P}_B , W_k^k is the **optimal matching cost** (Hungarian algorithm, [Kuhn \(1955, 1956\)](#)), cast as a linear program

$$W_k(\mathbb{P}_A, \mathbb{P}_B) = \left(\min_{s \in \mathcal{S}_n} \frac{1}{n} \sum_{i=1}^n d(\mathbf{x}_i, \mathbf{y}_{s(i)})^k \right)^{1/k},$$

where \mathcal{S}_n is the set of permutations on $\{1, 2, \dots, n\}$.

Optimal transport (discrete)

- › Consider the set of $n \times n$ doubly-stochastic matrices,

$$D_n = \{M \in \mathbb{R}_+^{n \times n} : M\mathbf{1}_n = \mathbf{1}_n \text{ and } M^\top \mathbf{1}_n = \mathbf{1}_n\},$$

and the subset of permutation matrices,

$$U_n = \{M \in \{0, 1\}^{n \times n} : M\mathbf{1}_n = \mathbf{1}_n \text{ and } M^\top \mathbf{1}_n = \mathbf{1}_n\}.$$

- › Let C denote the cost matrix, $C_{i,j} = d(x_i, y_j)^k$, then

$$W_k(\mathbf{x}, \mathbf{y})^k = \operatorname{argmin}_{P \in U_n} \{ \langle P, C \rangle \}, \text{ where } \langle P, C \rangle = \sum_{i=1}^n \sum_{j=1}^n P_{i,j} C_{i,j} \quad (1)$$

and “optimal transport” permutation matrix

$$P^* \in \operatorname{argmin}_{P \in U_n} \{ \langle P, C \rangle \} \quad (2)$$

Optimal transport (discrete)

	7	8	9	10	11	12
1	0.41	0.55	0.22	0.64	0.04	0.25
2	0.28	0.24	0.73	0.22	0.64	0.80
3	0.28	0.47	0.32	0.52	0.16	0.37
4	0.28	0.62	0.81	0.25	0.64	0.85
5	0.41	0.37	0.89	0.25	0.81	0.97
6	0.66	0.76	0.21	0.89	0.22	0.14

	7	8	9	10	11	12
1	1	.
2	.	1
3	.	.	1	.	.	.
4	1
5	.	.	.	1	.	.
6	1

1 ↔ 11

2 ↔ 8

3 ↔ 9

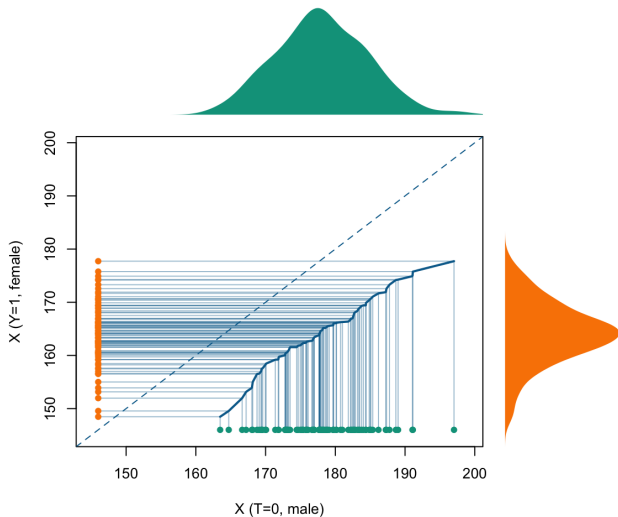
4 ↔ 7

5 ↔ 10

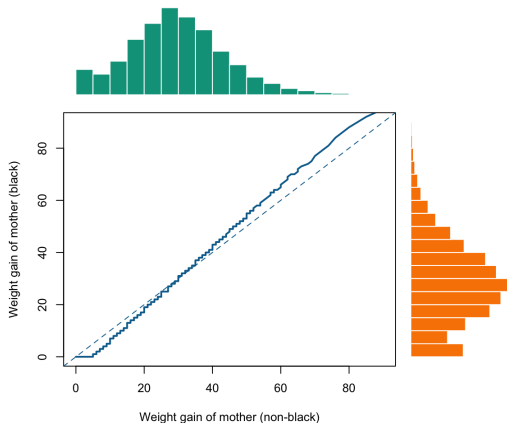
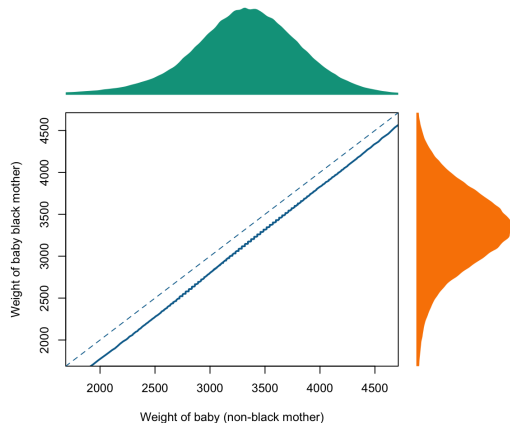
6 ↔ 12

Optimal transport (discrete)

- Consider two samples, with the height of **men** and **women** (both groups of size n).
- On the following graph, we can visualize the optimal matching of individuals in the two groups.
- It is a **monotone mapping**.



Optimal transport (discrete)



- Two groups, with **black** and **non-black** mothers, delivering babies (in the U.S.)
- $x_1 \leftrightarrow x_1$ (newborn weight) and $x_2 \leftrightarrow x_2$ (weight gain of the mother)

Optimal transport (discrete)

Proposition 4.31: Hardy–Littlewood–Pólya inequality, Hardy et al. (1952)

Given $x_1 \leq \dots \leq x_n$ and $y_1 \leq \dots \leq y_n$ n pairs of ordered real numbers, for every permutation σ of $\{1, 2, \dots, n\}$,

$$\sum_{i=1}^n x_i y_{n+1-i} \leq \sum_{i=1}^n x_i y_{\sigma(i)} \leq \sum_{i=1}^n x_i y_i.$$

- various implications, e.g. bounds on the covariance, and the correlation, see Proposition 4.35.
- This can be extended to more general function $\Phi(x_i, y_j)$.



Optimal transport (discrete)

Definition 4.38: Supermodular, Topkis (1998)

Function $\Phi : \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}$ is **supermodular** if for any $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^k$,

$$\Phi(\mathbf{z} \wedge \mathbf{z}') + \Phi(\mathbf{z} \vee \mathbf{z}') \geq \Phi(\mathbf{z}) + \Phi(\mathbf{z}'),$$

where $\mathbf{z} \wedge \mathbf{z}'$ and $\mathbf{z} \vee \mathbf{z}'$ denote respectively the maximum and the minimum componentwise. If $-\Phi$ is supermodular, Φ is said to be submodular.

Optimal transport (discrete)

Proposition 4.32: Hardy–Littlewood–Pólya inequality, [Hardy et al. \(1952\)](#)

Given $x_1 \leq \dots \leq x_n$ and $y_1 \leq \dots \leq y_n$ n pairs of ordered real numbers, and some supermodular function $\Phi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, for every permutation σ of $\{1, 2, \dots, n\}$,

$$\sum_{i=1}^n \Phi(x_i, y_{n+1-i}) \leq \sum_{i=1}^n \Phi(x_i, y_{\sigma(i)}) \leq \sum_{i=1}^n \Phi(x_i, y_i),$$

while if $\Phi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is submodular,

$$\sum_{i=1}^n \Phi(x_i, y_i) \leq \sum_{i=1}^n \Phi(x_i, y_{\sigma(i)}) \leq \sum_{i=1}^n \Phi(x_i, y_{n+1-i}).$$

➤ Functions $\Phi(x, y) = \gamma(x - y)$ for some concave function $\gamma : \mathbb{R} \rightarrow \mathbb{R}$, such as $\Phi(x, y) = -|x - y|^k$ with $k \geq 1$, are supermodular.

Optimal transport (discrete)

```
1 > permutations = function(n){
2 +   if(n==1){
3 +     return(matrix(1))
4 +   } else {
5 +     sp = permutations(n-1)
6 +     p = nrow(sp)
7 +     A = matrix(nrow=n*p, ncol=n)
8 +     for(i in 1:n){
9 +       A[(i-1)*p+1:p,] =
10 +         cbind(i, sp+(sp>=i))
11 +     }
12 +     return(A)
13 +   }
14 + }
```

$\Phi(x, y) = (x - y)^2$, submodular function,

➤ Consider $x_1 \leq \dots \leq x_n$

```
1 > Phi = function(x,y) sum((x-y)^2)
2 > set.seed(1)
3 > x = sort(x)
4 > y = y[1:6]
5 > vect = permutations(6)
6 > MY = matrix(vect, ncol=6)
7 > MPhi = function(i) Phi(x, y[MY[i,]])
8 > S = Vectorize(MPhi)(1:nrow(MY))
9 > y[MY[which.min(S),]]
10 [1] 0.046 0.288 0.409 0.788 0.883 0.940
```

Optimal transport (discrete)

➤ In a very general setting (with $n_A \neq n_B$), if $\mathbf{a}_A \in \mathbb{R}_+^{n_A}$ and $\mathbf{a}_B \in \mathbb{R}_+^{n_B}$ satisfy $\mathbf{a}_A^\top \mathbf{1}_{n_A} = \mathbf{a}_B^\top \mathbf{1}_{n_B}$ (identical sums), define

$$U(\mathbf{a}_A, \mathbf{a}_B) = \{M \in \mathbb{R}_+^{n_A \times n_B} : M\mathbf{1}_{n_B} = \mathbf{a}_A \text{ and } M^\top \mathbf{1}_{n_A} = \mathbf{a}_B\}.$$

This set of matrices is a **convex transportation polytope** (see [Brualdi \(2006\)](#)).

➤ In our case, let U_{n_A, n_B} denote $U\left(\mathbf{1}_{n_A}, \frac{n_A}{n_B} \mathbf{1}_{n_B}\right)$ ($U_{n,n}$ is the set of permutation matrices associated with \mathcal{S}_n). Let C denote the cost matrix, $C_{i,j} = d(x_i, y_j)^k$.

$$W_k(\mathbf{x}, \mathbf{y})^k = \operatorname{argmin}_{P \in U_{n_A, n_B}} \{\langle P, C \rangle\}, \text{ where } \langle P, C \rangle = \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} P_{i,j} C_{i,j} \quad (3)$$

and “optimal transport”

$$P^* \in \operatorname{argmin}_{P \in U_{n_A, n_B}} \{\langle P, C \rangle\} \quad (4)$$

Optimal transport (discrete)

	7	8	9	10	11	12	13	14	15	16
1	0.41	0.55	0.22	0.64	0.04	0.25	0.24	0.77	0.74	0.55
2	0.28	0.24	0.73	0.22	0.64	0.80	0.76	0.76	0.12	0.10
3	0.28	0.47	0.32	0.52	0.16	0.37	0.27	0.68	0.63	0.45
4	0.28	0.62	0.81	0.25	0.64	0.85	0.58	0.32	0.51	0.48
5	0.41	0.37	0.89	0.25	0.81	0.97	0.91	0.81	0.05	0.25
6	0.66	0.76	0.21	0.89	0.22	0.14	0.33	0.96	0.99	0.79

	7	8	9	10	11	12	13	14	15	16
1	.	.	1/5	.	3/5	.	1/5	.	.	.
2	.	2/5	3/5
3	3/5	2/5	.	.	.
4	.	.	.	2/5	.	.	.	3/5	.	.
5	.	1/5	.	1/5	3/5	.
6	.	.	2/5	.	.	3/5

Optimal transport (discrete)

- From Kantorovich (1942), one can use the dual linear programming problem

$$W_k(\mathbf{a}, \mathbf{b})^k = \begin{cases} \text{primal}(\mathbf{a}, \mathbf{b}, \mathbf{C}) = \min_{P \in U_{\mathbf{a}, \mathbf{b}}} \{ \langle P, \mathbf{C} \rangle \} \\ \text{or} \\ \text{dual}(\mathbf{a}, \mathbf{b}, \mathbf{C}) = \max_{(\mathbf{u}, \mathbf{v}) \in M_{\mathbf{C}}} \{ \mathbf{u}^\top \mathbf{a} + \mathbf{v}^\top \mathbf{b} \} \end{cases}$$

where $M_{\mathbf{C}} = \{(\mathbf{u}, \mathbf{v}) \in \mathbb{R}^{n_A + n_B} \mid u_i + v_j \leq C_{ij}\}$.

- If $n_A \sim n_B \sim n$, $O(n^3 \log(n))$ problem.

Set $\psi_{\mathbf{b}}(\mathbf{a}, \mathbf{C}) = \max_{(\mathbf{u}, \mathbf{v}) \in M_{\mathbf{C}}} \{ \mathbf{u}^\top \mathbf{a} + \mathbf{v}^\top \mathbf{b} \}$, $\mathbf{a} \mapsto \psi_{\mathbf{b}}(\mathbf{a}, \mathbf{C})$ is a convex non-smooth map.

- The dual optimum \mathbf{u}^* is subgradient of $\mathbf{a} \mapsto \psi_{\mathbf{b}}(\mathbf{a}, \mathbf{C})$.
- If $k = 2$ (Euclidean distance), convex quadratic problem.

Optimal transport (discrete)

- › Given $P \in U_{n_A, n_B}$, define the **entropy** as

$$\mathcal{E}(P) = - \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} P_{i,j} \log P_{i,j} \text{ or } \mathcal{E}'(P) = - \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} P_{i,j} [\log P_{i,j} - 1]$$

and consider the γ -regularized optimal transport problem

$$P_{\gamma}^* = \underset{P \in U_{n_A, n_B}}{\operatorname{argmin}} \left\{ \langle P, C \rangle - \gamma \mathcal{E}(P) \right\} \quad (5)$$

since the problem is strictly convex.

- › The Lagrangian is here

$$\mathcal{L}(P, \lambda_A, \lambda_B) = \langle P, C \rangle - \gamma \mathcal{E}(P) - \langle \lambda_A, P \mathbf{1}_{n_B} - \mathbf{1}_{n_A} \rangle - \langle \lambda_B, P^{\top} \mathbf{1}_{n_A} - \mathbf{1}_{n_B} \rangle$$

Optimal transport (discrete)

and the first order conditions are

$$C_{i,j} + \gamma \log(P_{i,j}) - \lambda_{A,i} - \lambda_{B,j} = 0,$$

i.e.

$$P_{i,j} = \exp[\lambda_{A,i} - C_{i,j} + \lambda_{B,j}] \text{ or } P = D_A \exp[-C] D_B$$

where D_A and D_B are diagonal matrices.

- This can be related to the **Doubly Stochastic Scaling Problem**: let A be some $n \times n$ matrix with positive coefficients, we want to find D_A and D_B two positive diagonal matrices ($n \times n$) such that $D_A A D_B$ is doubly stochastic (see [Parlett and Landis \(1982\)](#))
- More generally, this corresponds to the **Matrix Scaling Problem**: Let A be some $n_A \times n_B$ matrix with positive coefficients, we want to find D_A and D_B two positive diagonal matrices (respectively $n_A \times n_A$ and $n_B \times n_B$) such that $D_A A D_B$ is in $U(\mathbf{a}_A, \mathbf{a}_B)$.

Optimal transport (discrete)

Theorem 4.3: Sinkhorn - Matrix Scaling, Sinkhorn (1962)

For any matrix \mathbf{A} $n \times m$ with positive entries, for any \mathbf{a} and \mathbf{b} in the simplex, there exist unique $\mathbf{u} \in \mathbb{R}_+^n$ and $\mathbf{v} \in \mathbb{R}_+^m$ such that

$$\text{diag}[\mathbf{u}] \mathbf{A} \text{diag}[\mathbf{v}] \in U_{\mathbf{a}, \mathbf{b}}.$$

➤ Sinkhorn and Knopp (1967) (extending Sinkhorn (1962, 1964, 1966)) suggested the following algorithm (updating alternatively D_A and D_B)

$$\begin{cases} D_A^{(t)} = \text{diag}(\mathbf{a}_A / (AD_B)^{(t-1)}) \\ D_B^{(t)} = \text{diag}(\mathbf{a}_B / (AD_A)^{(t)}) \end{cases}$$

(where the division here is element-wise).

Optimal transport (discrete)

- › An alternative way to write the entropic optimization problem is

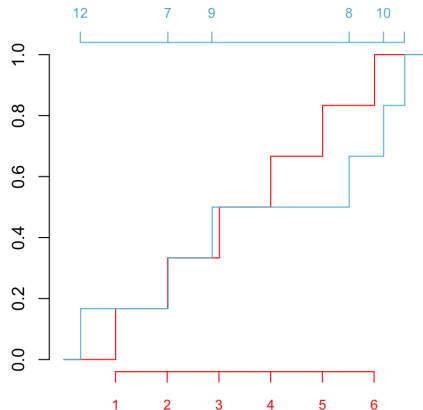
$$P_{\gamma}^* = \underset{P \in U_{\mathbf{a}_A, \mathbf{a}_B}}{\operatorname{argmin}} \left\{ \langle P, C \rangle + \gamma \cdot d_{\text{KL}}(P \| \mathbf{a}_A \otimes \mathbf{a}_B) \right\} \quad (6)$$

Using mutual information here makes it easier to extend to the continuous case...

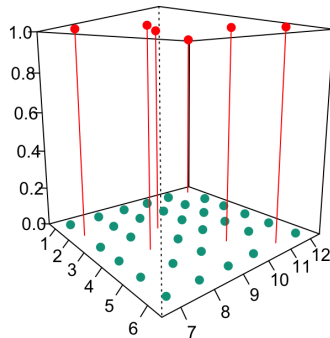
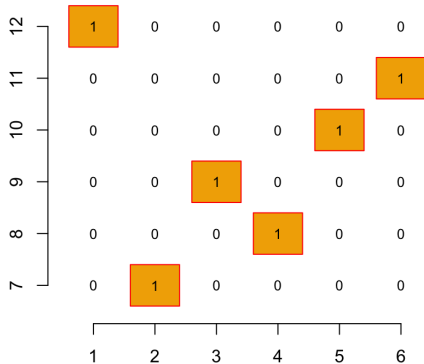
- › The extension of Sinkhorn algorithm is the coordinate descent/ascent algorithm.

Optimal transport (discrete)

```
1 > set.seed(123)
2 > x = (1:6)/7
3 > y = runif(9)
4 > x
5 [1] 0.14 0.29 0.43 0.57 0.71 0.86
6 > y[1:6]
7 [1] 0.29 0.79 0.41 0.88 0.94 0.05
8 > library(T4transport)
9 > Wxy = wasserstein(x,y[1:6])
10 > Wxy$plan
```

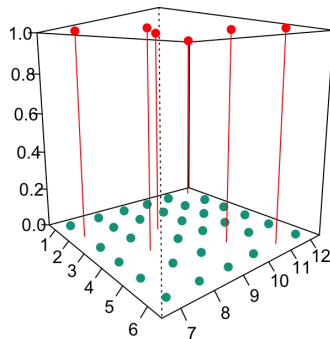
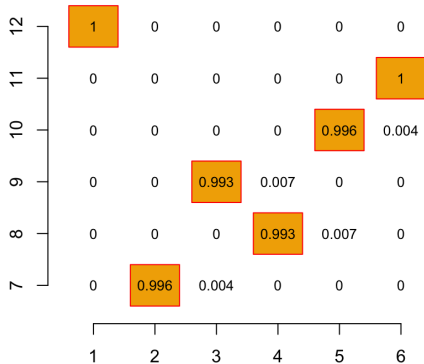


Optimal transport (discrete)



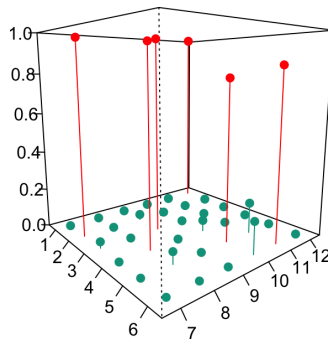
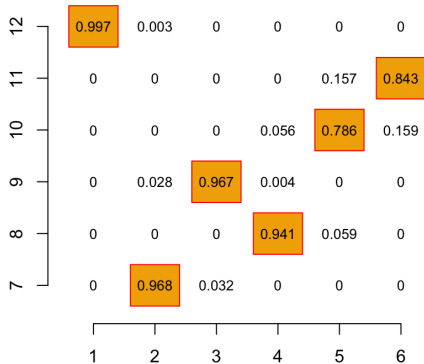
```
1 > Wxy = wasserstein(x,y[1:6])  
2 > Wxy$plan
```

Optimal transport (discrete)



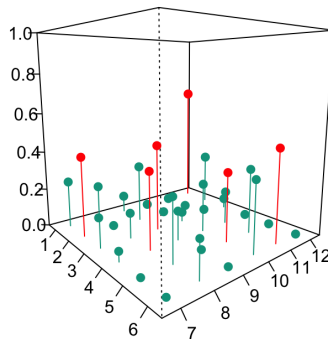
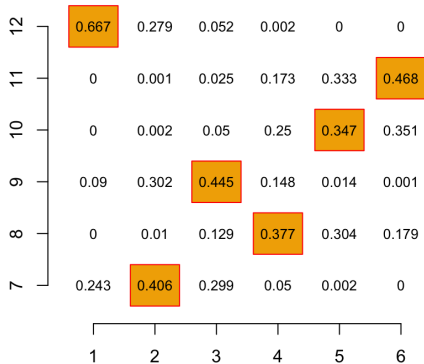
```
1 > Sxy = sinkhorn(x, y[1:6], p = 2, lambda = 0.001)
2 > Sxy$plan
```

Optimal transport (discrete)



```
1 > Sxy = sinkhorn(x, y[1:6], p = 2, lambda = 0.005)
2 > Sxy$plan
```

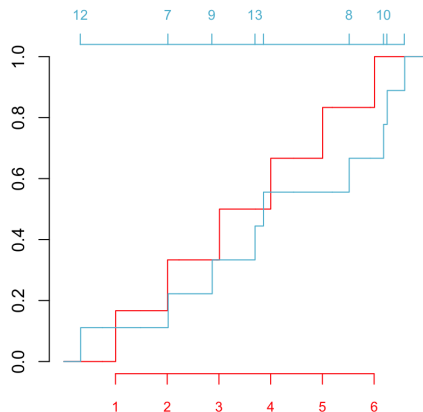

Optimal transport (discrete)



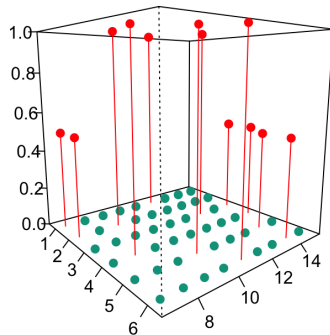
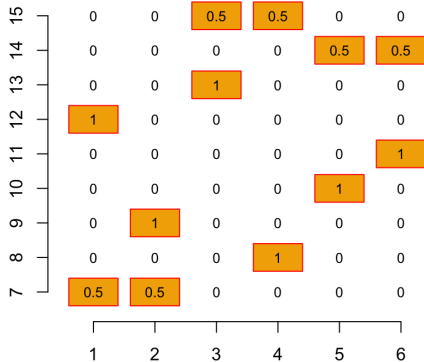
```
1 > Sxy = sinkhorn(x, y[1:6], p = 2, lambda = 0.05)
2 > Sxy$plan
```

Optimal transport (discrete)

```
1 > y
2 [1] 0.29 0.79 0.41 0.88 0.94 0.05
3 [7] 0.53 0.89 0.55
4 > library(T4transport)
5 > Wxy = wasserstein(x,y)
6      [,1] [,2] [,3] [,4] [,5] [,6]
7 [1,] 0.5 0.5 0.0 0.0 0.0 0.0
8 [2,] 0.0 0.0 0.0 1.0 0.0 0.0
9 [3,] 0.0 1.0 0.0 0.0 0.0 0.0
10 [4,] 0.0 0.0 0.0 0.0 1.0 0.0
11 [5,] 0.0 0.0 0.0 0.0 0.0 1.0
12 [6,] 1.0 0.0 0.0 0.0 0.0 0.0
13 [7,] 0.0 0.0 1.0 0.0 0.0 0.0
14 [8,] 0.0 0.0 0.0 0.0 0.5 0.5
15 [9,] 0.0 0.0 0.5 0.5 0.0 0.0
```

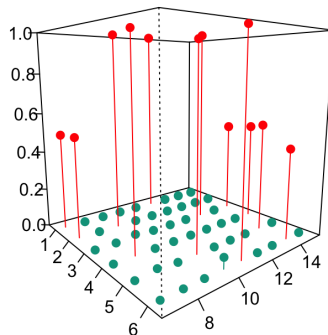
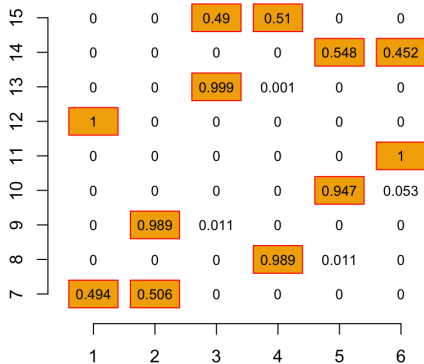


Optimal transport (discrete)



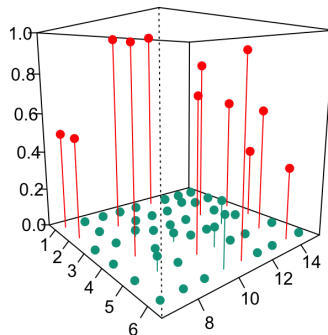
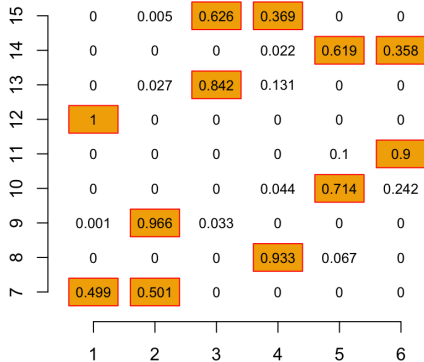
```
1 > Wxy = wasserstein(x,y)
2 > Wxy$plan
```

Optimal transport (discrete)



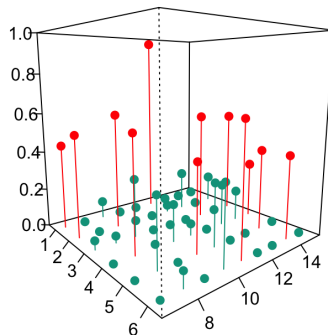
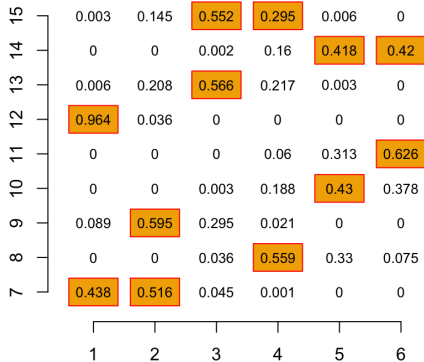
```
1 > Sxy = sinkhorn(x, y, p = 2, lambda = 0.001)
2 > Sxy$plan
```

Optimal transport (discrete)



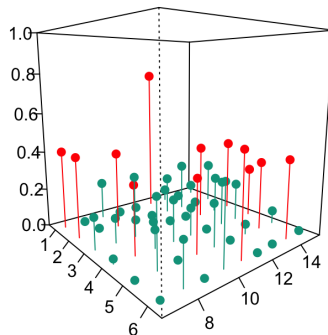
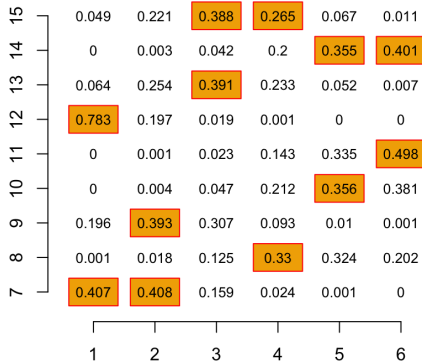
```
1 > Sxy = sinkhorn(x, y, p = 2, lambda = 0.005)
2 > Sxy$plan
```

Optimal transport (discrete)



```
1 > Sxy = sinkhorn(x, y, p = 2, lambda = 0.02)
2 > Sxy$plan
```

Optimal transport (discrete)



```
1 > Sxy = sinkhorn(x, y, p = 2, lambda = 0.05)
2 > Sxy$plan
```

Optimal transport (discrete)

Theorem 4.4: Optimal transport for discrete univariate distributions

Consider n points each group, on \mathbb{R} , $\{x_1, \dots, x_n\}$ and $\{y_1, \dots, y_n\}$, ordered in the senses that $x_1 \leq x_2 \leq \dots \leq x_n$ and $y_1 \leq y_2 \leq \dots \leq y_n$, for any $k \geq 1$,

$$W_k = \left(\frac{1}{n} \sum_{i=1}^n |x_i - y_i|^k \right)^{1/k}$$

Theorem 4.5: Optimal transport for continuous univariate distributions

$$W_k = \left(\int_0^1 |F_x^{-1}(u) - F_y^{-1}(u)|^k du \right)^{1/k}$$

Optimal transport (discrete)

Theorem 4.6: Optimal transport for continuous univariate distributions

Let \mathbb{P}_A and \mathbb{P}_B be two probability measures on \mathbb{R} , and suppose that $c(x, y) = h(x - y)$ for some strictly convex function h . Then there exists a unique $\pi \in \Pi(\mathbb{P}_A, \mathbb{P}_B)$ such that

- ▶ π is optimal to Kantorovich problem (8.16)
- ▶ π is the comonotone joint distribution with marginals \mathbb{P}_A and \mathbb{P}_B .

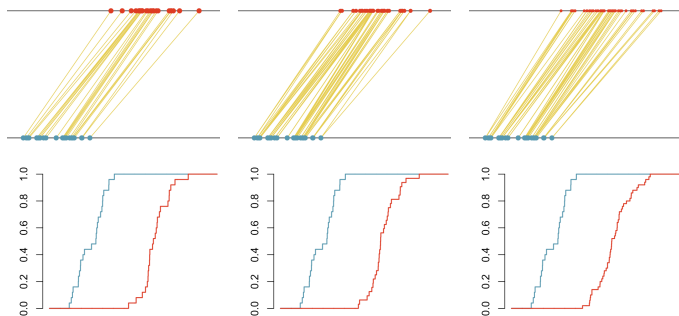
› If $c(x, y) = |x - y|$, the optimal transport solution might be non-unique.

Theorem 4.7: Optimal map for continuous univariate distributions

The optimal Monge map T^* such that $T_{\#}^* \mathbb{P}_A = \mathbb{P}_B$ is $T^* = F_B^{-1} \circ F_A$.

Optimal transport (discrete)

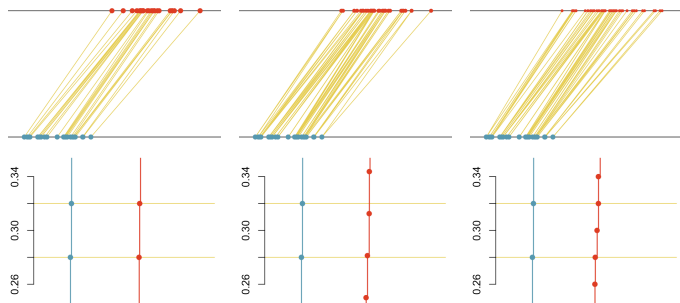
- Consider $n_A = 25$ and $n_B = 25$ points in \mathbb{R} , $n_B = 32$ and $n_B = 50$



$$\hat{F}_{n_A}(x) = \frac{1}{n_A} \sum_{i=1}^{n_A} \mathbf{1}(x_i \leq x) \quad \text{and} \quad \hat{F}_{n_B}(x) = \frac{1}{n_B} \sum_{i=1}^{n_B} \mathbf{1}(x_i \leq x)$$

Optimal transport (discrete)

- › Consider $n_A = 25$ and $n_B = 25$ points in \mathbb{R} , $n_B = 32$ and $n_B = 50$

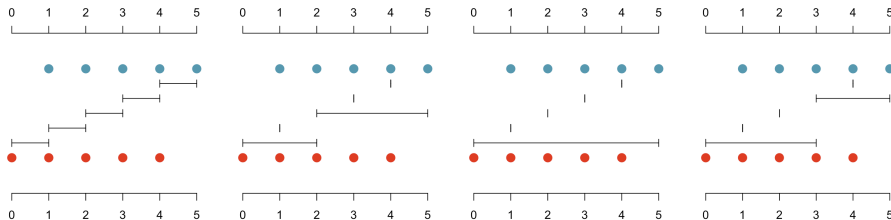


$$\hat{F}_{n_A}(x) = \frac{1}{n_A} \sum_{i=1}^{n_A} \mathbf{1}(x_i \leq x) \quad \text{and} \quad \hat{F}_{n_B}(x) = \frac{1}{n_B} \sum_{i=1}^{n_B} \mathbf{1}(x_i \leq x)$$

Optimal transport (discrete)

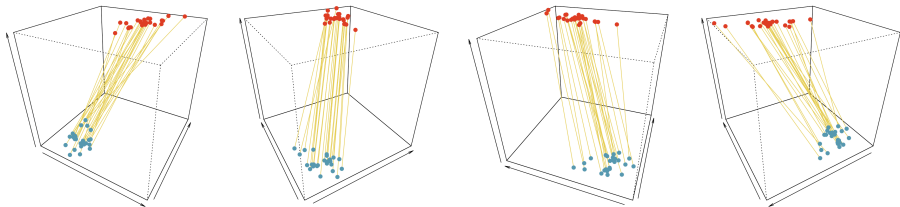
› In the univariate case, if $k = 1$,

$$W_1 = \frac{1}{n} \sum_{i=1}^n |x_i - y_{\sigma(i)}|$$

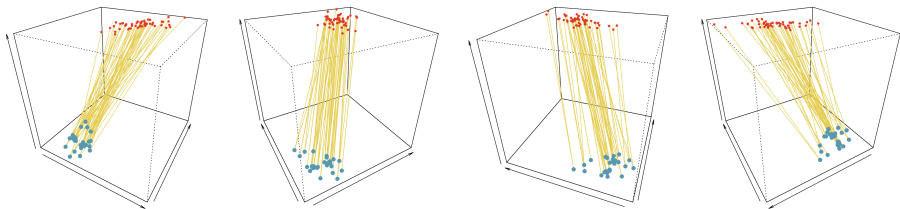


Multivariate Optimal Transport

- Consider n and n points in \mathbb{R}^2

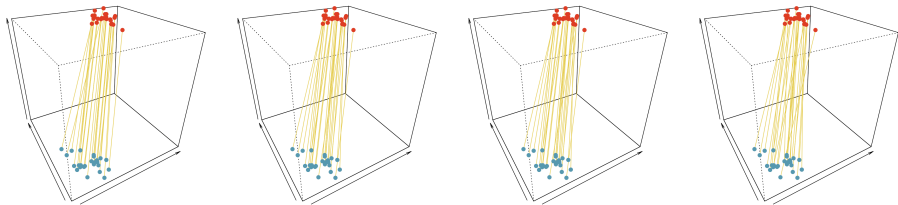


- Consider n and $2n$ points in \mathbb{R}^2

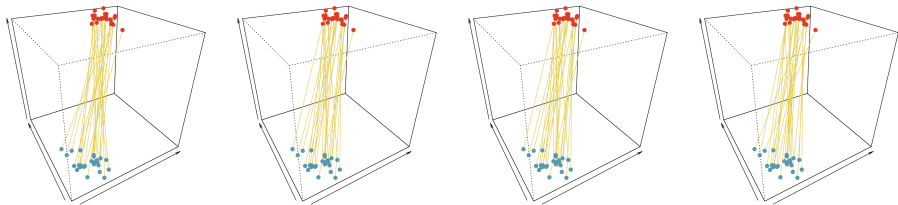


Multivariate Optimal Transport

- Consider n and n points in \mathbb{R}^2 , and $k = 1, 2, 3, 4$, $T_{\#}\mathbb{P}_A = \mathbb{P}_B$



- Consider n and n points in \mathbb{R}^2 , and $p = 1, 2, 3, 4$, $T_{\#}\mathbb{P}_B = \mathbb{P}_A$



Multivariate Optimal Transport

Theorem 4.8: Optimal map for continuous multivariate distributions, Brenier (1991)

With a quadratic cost, the optimal Monge map T^* is unique, and it is the gradient of a convex function, $T^* = \nabla\varphi$.

Example Multidimensional Gaussian

$$\mathbf{x}_B = \mathcal{T}^*(\mathbf{x}_A) = \boldsymbol{\mu}_B + \mathbf{A}(\mathbf{x}_A - \boldsymbol{\mu}_A),$$

where \mathbf{A} is a symmetric positive matrix that satisfies $\mathbf{A}\boldsymbol{\Sigma}_A\mathbf{A} = \boldsymbol{\Sigma}_B$, which has a unique solution given by $\mathbf{A} = \boldsymbol{\Sigma}_A^{-1/2}(\boldsymbol{\Sigma}_A^{1/2}\boldsymbol{\Sigma}_B\boldsymbol{\Sigma}_A^{1/2})^{1/2}\boldsymbol{\Sigma}_A^{-1/2}$, where $\mathbf{M}^{1/2}$ is the square root of the square (symmetric) positive matrix \mathbf{M} based on the Schur decomposition ($\mathbf{M}^{1/2}$ is a positive symmetric matrix), as described in Higham (2008).

Multivariate Optimal Transport

Proposition 4.33: W_2 for Gaussian vectors

Consider two Gaussian distributions, then

$$W_2(\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2))^2 = \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2 + \text{tr}(\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2 - 2(\boldsymbol{\Sigma}_1^{1/2} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_1^{1/2})^{1/2})$$

Proof: Let $\mathbf{X}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$, $\mathbf{X}_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$, and Γ define the covariance matrix of $(\mathbf{X}_1, \mathbf{X}_2)$,

$$\Gamma = \begin{pmatrix} \boldsymbol{\Sigma}_1 & C \\ C^\top & \boldsymbol{\Sigma}_2 \end{pmatrix}$$

where (generally), C is some $n_1 \times n_2$ matrix. Recall that $n_1 \times n_2$ matrices can have a pseudo-inverse, in the sense that (Penrose conditions)

$$\begin{cases} AA^-A = A & \{(AA^-)^\top = AA^- \\ A^-AA^- = A^-, & \{(A^-A)^\top = A^-A, \end{cases}$$

Multivariate Optimal Transport

► Observe that $\mathbb{E}(\|\mathbf{X}_1 - \mathbf{X}_2\|_{\ell_2}^2) = \text{tr}(\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2 - 2C)$. Recall that C must satisfy the Schur complement constraint, $\boldsymbol{\Sigma}_1 - C\boldsymbol{\Sigma}_2^{-1}C^\top \succeq 0$, so that we want to solve

$$C^* = \operatorname{argmin}\{-2\text{tr}(C)\} \text{ s.t. } \boldsymbol{\Sigma}_1 - C\boldsymbol{\Sigma}_2^{-1}C^\top \succeq 0,$$

as studied in [Olkin and Pukelsheim \(1982\)](#), where $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$ are positive ($\succeq 0$) matrices.

Let $\mathcal{G} = \{C, n_1 \times n_2 : \boldsymbol{\Sigma}_1 - C\boldsymbol{\Sigma}_2^{-1}C^\top \succeq 0\}$, $\mathcal{S} = \{S : SS^{-1}\boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}_2\}$, one can prove (standard duality and convexity arguments) that

$$\max_{C \in \mathcal{G}} \{2\text{tr}(C)\} = \max_{S \in \mathcal{S}} \{\text{tr}(\boldsymbol{\Sigma}_1 S + \boldsymbol{\Sigma}_2 S^{-1})\} = 2\text{tr}(\boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_2^{1/2})$$

with respective (unique) solutions

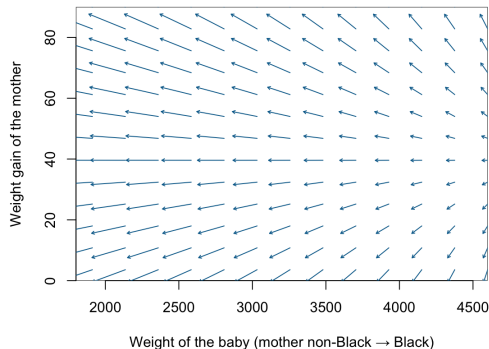
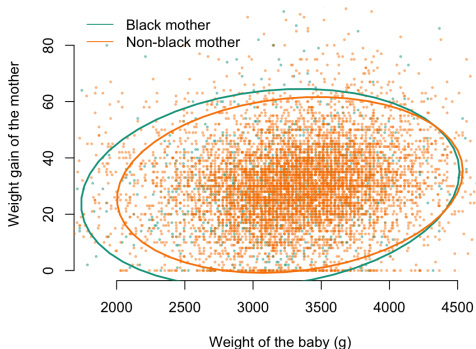
$$\begin{cases} C^* = \boldsymbol{\Sigma}_1 S^* \\ S^* = \boldsymbol{\Sigma}_2^{1/2} [(\boldsymbol{\Sigma}_2^{1/2} \boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_2^{1/2})^{1/2}] \boldsymbol{\Sigma}_2^{1/2} \end{cases}$$

Multivariate Optimal Transport

- See [Olkin and Pukelsheim \(1982\)](#), [Givens and Shortt \(1984\)](#) and [Knott and Smith \(1984\)](#), or more recently [Takatsu \(2008\)](#) and [Takatsu and Yokota \(2012\)](#), with more geometric interpretations.
- To illustrate, consider the previous example, with newborn weight and weight gain of mothers, in the U.S., with [Black](#) and [non-Black](#) mothers, with here a joint mapping $\mathbb{R}_2^+ \rightarrow \mathbb{R}_2^+$.

Multivariate Optimal Transport

$(x_1, x_2) \leftrightarrow (x_1, x_2)$ (newborn weight, weight gain of the mother)



– Part 3 –
Models

Generalized Linear Model

Definition 4.39: Exponential family, McCullagh and Nelder (1989)

The distribution of Y is in the **exponential family** if its density (with respect to some appropriate measure) is

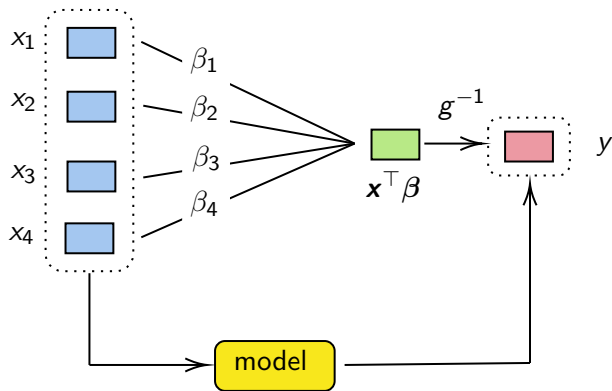
$$f_{\theta, \varphi}(y) = \exp\left(\frac{y\theta - b(\theta)}{\varphi} + c(y, \varphi)\right),$$

where θ is the canonical parameter, φ is a nuisance parameter, and $b : \mathbb{R} \rightarrow \mathbb{R}$ is some $\mathbb{R} \rightarrow \mathbb{R}$ function.

- Such as the binomial, Poisson, Gaussian, gamma distributions, etc.
- Also compound Poisson / Tweedie (from **Tweedie (1984)**).

Generalized Linear Model

- › Given some dataset (y_i, \mathbf{x}_i) , suppose that $\mu(\mathbf{x}) = g^{-1}(\mathbf{x}^\top \boldsymbol{\beta})$



- › OLS, $\mu(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta}$ and $\hat{\boldsymbol{\beta}}^{\text{ols}} = \operatorname{argmin} \left\{ \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.

Generalized Linear Model

- › Consider problems

$$\min_{\mathbf{x} \in \mathbb{R}^k} \{f(\mathbf{x})\} \quad \text{under constraint } g(\mathbf{x}) = \mathbf{0} \quad \text{or} \quad \min_{\mathbf{x} \in \mathbb{R}^k} \{f(\mathbf{x})\} \quad \text{under constraint } g(\mathbf{x}) \leq \mathbf{0}$$

- › Karush-Kuhn-Tucker condition is

$$\begin{cases} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \mathbf{z}^*) = \mathbf{0} \\ \nabla_{\mathbf{z}} \mathcal{L}(\mathbf{x}^*, \mathbf{z}^*) = \mathbf{0} \end{cases}$$

where

$$\mathcal{L}(\mathbf{x}, \mathbf{z}) = f(\mathbf{x}) + \mathbf{z}^\top g(\mathbf{x})$$

is the Lagrangian problem (parameter \mathbf{z} are multipliers)

Generalized Linear Model

Definition 4.40: Ridge Estimator (OLS), Hoerl and Kennard (1970)

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \underset{\beta \in \mathbb{R}^k}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^{\top} \beta)^2 + \lambda \sum_{j=1}^k \beta_j^2 \right\}.$$

$$\hat{\beta}_{\lambda}^{\text{ridge}} = (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^{\top} \mathbf{y}$$

Definition 4.41: Ridge Estimator (GLM)

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \underset{\beta \in \mathbb{R}^k}{\operatorname{argmin}} \left\{ - \sum_{i=1}^n \log f(y_i | \mu_i = g^{-1}(\mathbf{x}_i^{\top} \beta)) + \lambda \sum_{j=1}^k \beta_j^2 \right\}.$$

Definition 4.42: LASSO Estimator (OLS), Tibshirani (1996)

$$\hat{\beta}_{\lambda}^{\text{lasso}} = \operatorname{argmin} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^{\top} \beta)^2 + \lambda \sum_{j=1}^k |\beta_j| \right\}.$$

Definition 4.43: LASSO Estimator (GLM)

$$\hat{\beta}_{\lambda}^{\text{lasso}} = \operatorname{argmin} \left\{ - \sum_{i=1}^n \log f(y_i | \mu_i = g^{-1}(\mathbf{x}_i^{\top} \beta)) + \lambda \sum_{j=1}^k |\beta_j| \right\}.$$

Generalized Linear Model

```
1 > library(glmnet)
2 > fit_ridge = glmnet(x, y, alpha = 0)
3 > fit_lasso = glmnet(x, y, alpha = 1)
```

› Elastic net

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda_1 \sum_{j=1}^k |\beta_j| + \frac{\lambda_2}{2} \sum_{j=1}^k \beta_j^2 \right\},$$

e.g. $\lambda_1 = \alpha\lambda$ and $\lambda_2 = (1 - \alpha)\lambda$ (two parameters — one for the global regularization, one for the trade-off between Ridge (Tikhonov) vs. Lasso)

to go further → (for more details on penalization issues)

Definition 4.44: ROC curve

The ROC curve is the parametric curve

$$\{\mathbb{P}[m(\mathbf{X}) > t | Y = 0], \mathbb{P}[m(\mathbf{X}) > t | Y = 1]\} \text{ for } t \in [0, 1],$$

when the score $m(\mathbf{X})$ and Y evolve in the same direction (a high score indicates a high risk).

$$C(t) = \text{TPR} \circ \text{FPR}^{-1}(t),$$

where

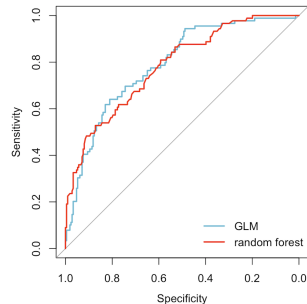
$$\begin{cases} \text{FRP}(t) = \mathbb{P}[m(\mathbf{X}) > t | Y = 0] = \mathbb{P}[m_0(\mathbf{X}) > t] \\ \text{TPR}(t) = \mathbb{P}[m(\mathbf{X}) > t | Y = 1] = \mathbb{P}[m_1(\mathbf{X}) > t]. \end{cases}$$

Accuracy

```
1 > library(ROCR)
2 > pred = prediction(df$yhat, df$y)
3 > roc = performance(pred,"tpr","fpr")
4 > plot(roc)
5 > auc = performance(pred,"auc")
```

see also

```
1 > library(pROC)
```

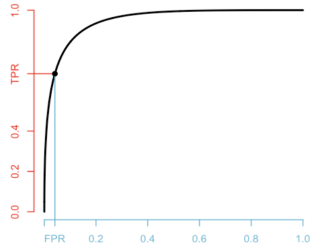
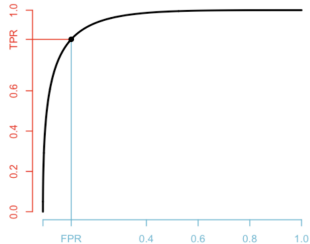
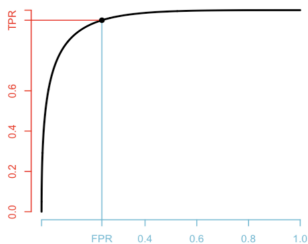
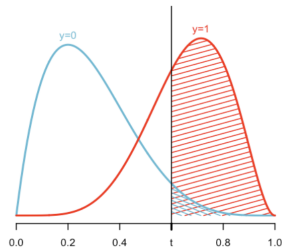
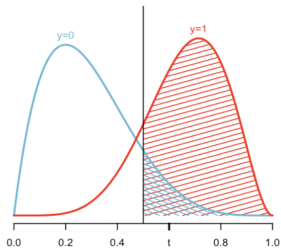
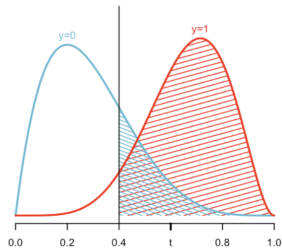


Definition 4.45: AUC, area under the ROC curve

The area under the curve is defined as the area below the ROC curve,

$$\text{AUC} = \int_0^1 C(t)dt = \int_0^1 \text{TPR} \circ \text{FPR}^{-1}(t)dt.$$

Accuracy



Calibration

- › Well-calibration was initially discussed in forecasting

Definition 4.46: Well-calibrated (1), Van Calster et al. (2019), Krüger and Ziegel (2021)

The forecast X of Y is a well-calibrated forecast of Y if $\mathbb{E}(Y|X) = X$ almost surely, or $\mathbb{E}[Y|X = x] = x$, for all x .

- › one can define “well-calibration” in prediction

Definition 4.47: Well-calibrated (2), Zadrozny and Elkan (2002); Cohen and Goldszmidt (2004)

The prediction $m(\mathbf{X})$ of Y is a well-calibrated prediction if $\mathbb{E}[Y|m(\mathbf{X}) = \hat{y}] = \hat{y}$, for all \hat{y} .

Calibration

“**Well calibrated** classifiers are probabilistic classifiers for which the output can be directly interpreted as a confidence level. For instance, a well calibrated (binary) classifier should classify the samples such that among the samples to which it gave a [predicted probability] value close to 0.8, approximately 80% actually belong to the positive class,” [scikit learn: Probability calibration](#)

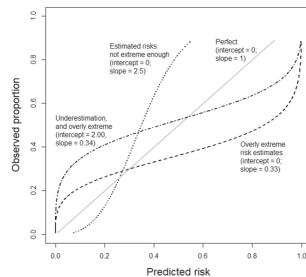
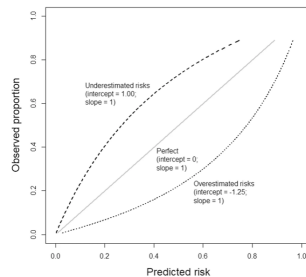
- “Suppose that a forecaster sequentially assigns probabilities to events. He is **well calibrated** if, for example, of those events to which he assigns a probability 30 percent, the long-run proportion that actually occurs turns out to be 30 percent,” [Dawid \(1982\)](#).
- “Out of all the times you said there was a 40 percent chance of rain, how often did rain actually occur? If, over the long run, it really did rain about 40 percent of the time, that means your forecasts were **well calibrated**,” [Silver \(2012\)](#),
- “we desire that the estimated class probabilities are reflective of the true underlying probability of the sample,” [Kuhn and Johnson \(2013\)](#)

Calibration

- See [Murphy and Epstein \(1967\)](#), [Roberts \(1968\)](#), [Gneiting and Raftery \(2005\)](#) on ensemble methods for weather forecasting, or more generally [Lichtenstein et al. \(1977\)](#), [Oakes \(1985\)](#), [Gneiting et al. \(2007\)](#).

Calibration

- As explained in [Van Calster et al. \(2019\)](#), "*among patients with an estimated risk of 20%, we expect 20 in 100 to have or to develop the event*",
 - ▶ If 40 out of 100 in this group are found to have the disease, the risk is **underestimated**
 - ▶ If we observe that in this group, 10 out of 100 have the disease, we have **overestimated** the risk.
- Hosmer-Lemeshow test, from [Hosmer Jr et al. \(2013\)](#) (logistic regression), and Bier score, from [Brier \(1950\)](#) and [Murphy \(1973\)](#)
- Function plotted in psychological papers [Keren \(1991\)](#)

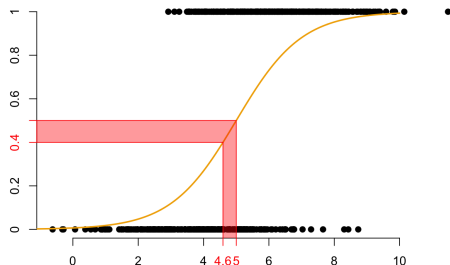
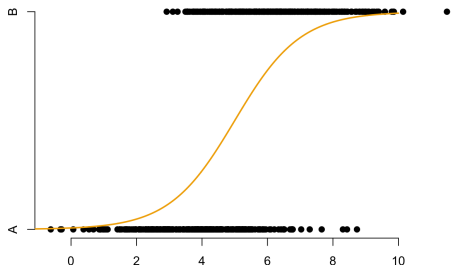


Calibration

› Consider a dataset (x_i, y_i) , $y_i \in \{A, B\}$, and consider model

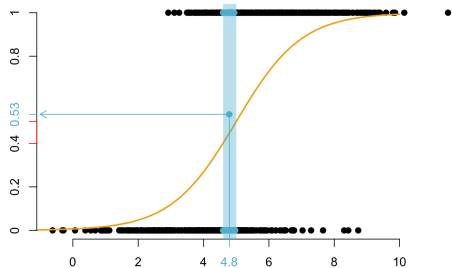
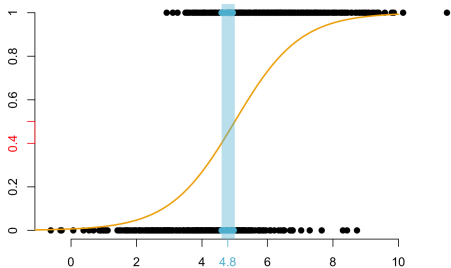
$\hat{m}(x) = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 x)}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 x)}$ to estimate $\mathbb{P}[Y = B|X = x]$ (logistic regression). Given

$[p_-, p_+] \subset [0, 1]$ (here $[0.4, 0.5]$), set $\mathcal{I} = \{i : m(x_i) \in [p_-, p_+]\}$.



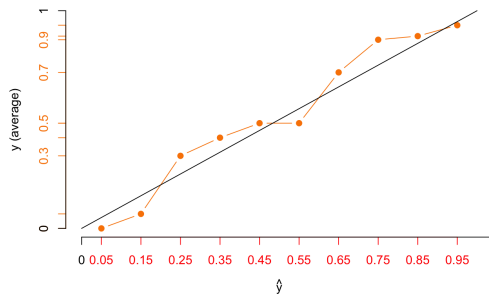
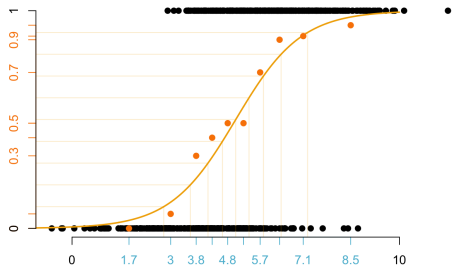
Calibration

- Given $\mathcal{I} = \{i : m(\mathbf{x}_i) \in [p_-, p_+]\}$, set $\bar{y}_{\mathcal{I}} = \frac{1}{n_{\mathcal{I}}} \sum_{i \in \mathcal{I}} y_i$



Calibration

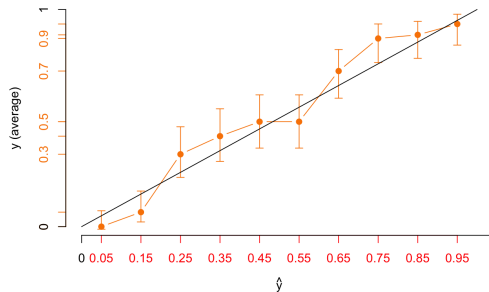
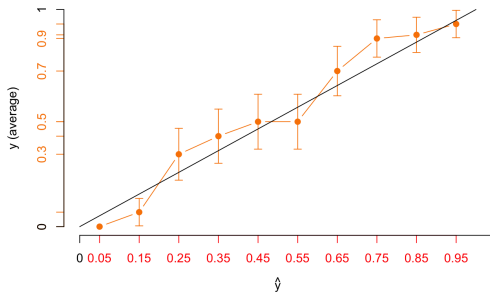
- Compute deciles $\bar{y}_1, \dots, \bar{y}_{10}$ associated with $[p_-, p_+]$ equal $[0, 0.1], [0.1, 0.2], \dots, [0.9, 1]$, with midpoints p_k . Then visualize $\{(p_k, \bar{y}_k)\}$, as in scikit-learn.



Calibration

- Asymptotic and Agresti and Coull (1998), $\tilde{n} = n + u_{1-\alpha/2}^2$ et $\tilde{p} = \frac{1}{\tilde{n}} \left(n\bar{x} + \frac{u_{1-\alpha/2}^2}{2} \right)$

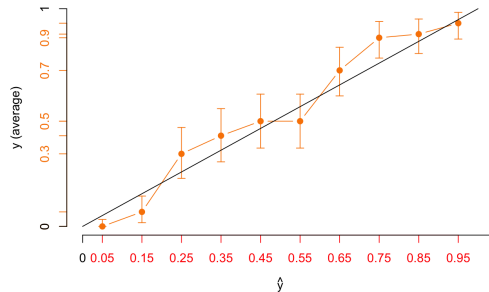
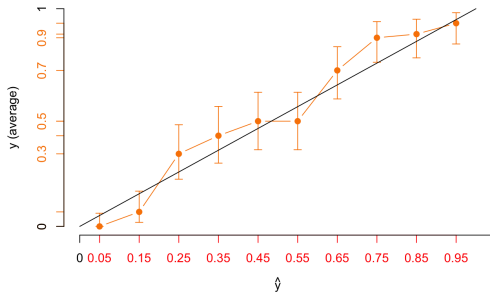
$$\left[\hat{p} \pm u_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right], \text{ and } \left[\tilde{p} \pm u_{1-\alpha/2} \sqrt{\frac{\tilde{p}(1-\tilde{p})}{\tilde{n}}} \right]$$



Calibration

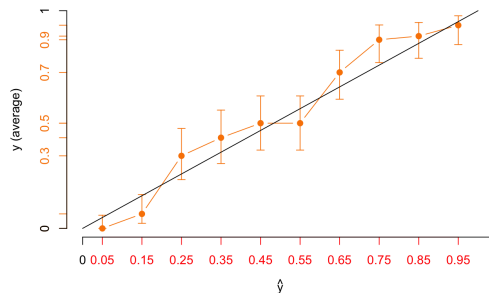
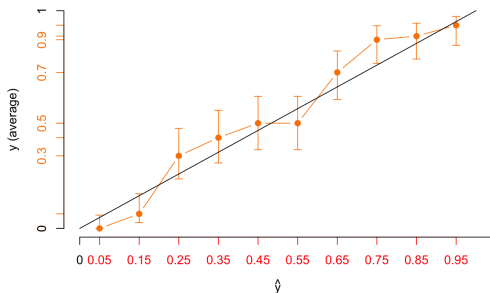
- Exact, and Wilson (1927)

$$\left[\frac{1}{1 + \frac{u_{1-\alpha/2}^2}{n}} \left(\hat{p} + \frac{u_{1-\alpha/2}^2}{2n} \right) \pm \frac{u_{1-\alpha/2}}{1 + \frac{u_{1-\alpha/2}^2}{n}} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n} + \frac{u_{1-\alpha/2}^2}{4n^2}} \right]$$



Calibration

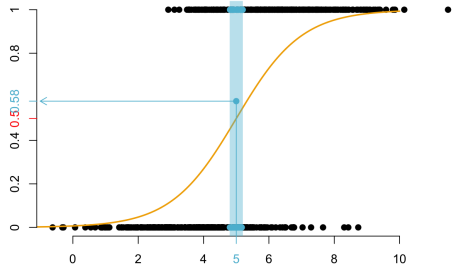
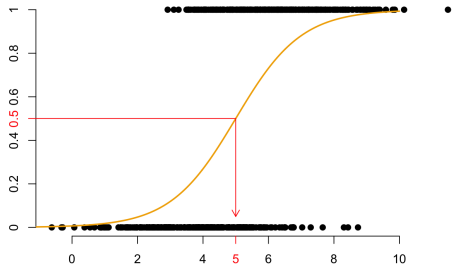
- › Bayes and Probit (via `binom.confint` in `binom` package)



Calibration

➤ Given $p \in (0, 1)$, consider $\mathcal{I}_p = \{i : \hat{m}(x_i) \in [p - h, p + h]\}$ for some $h > 0$, set

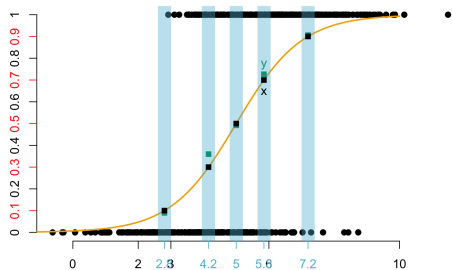
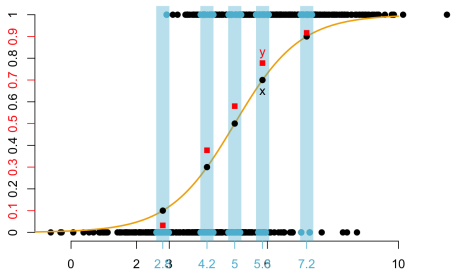
$$\bar{y}_p = \frac{1}{n_{\mathcal{I}_p}} \sum_{i \in \mathcal{I}_p} y_i$$



Calibration

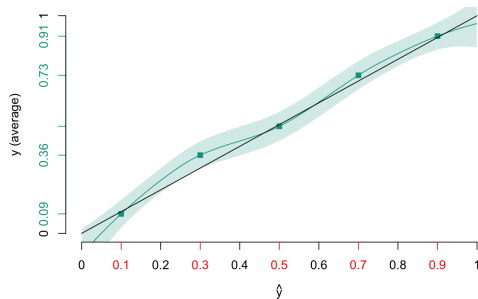
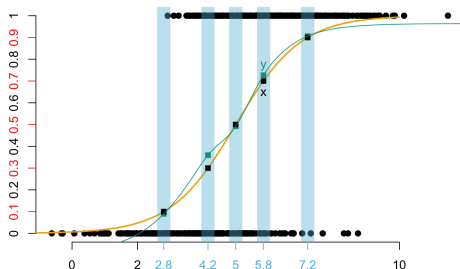
- Given $p \in (0, 1)$, compute $\bar{y}_p = \frac{1}{n_{\mathcal{I}_p}} \sum_{i \in \mathcal{I}_p} y_i$ (for appropriate bandwidth $h > 0$)

One could also consider some kernel based average... \bar{y}_p



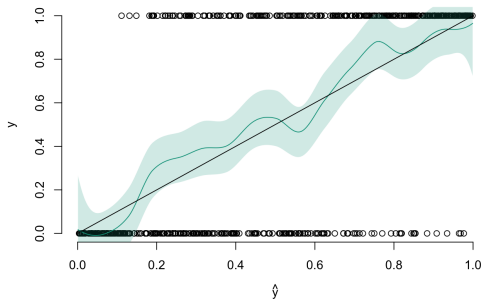
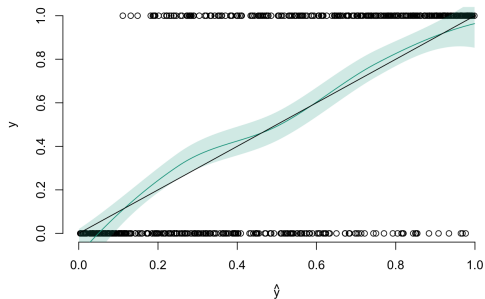
Calibration

- Compute \bar{y}_p for various p , and plot $p \mapsto \bar{y}_p$, that is an estimate of $\mathbb{E}[Y | \hat{m}(\mathbf{X}) = p]$.
- Add a confidence band around.
- ... but here, it works only because \hat{m} is smooth enough...



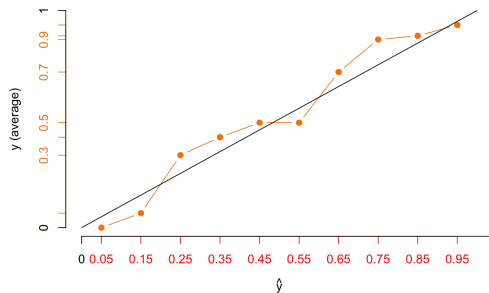
Calibration

- More generally, consider the local regression of y_i 's against $\hat{y}_i = \hat{m}(x_i)$'s.



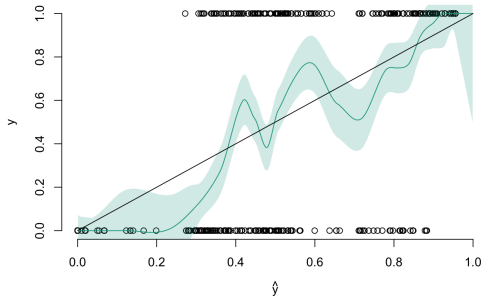
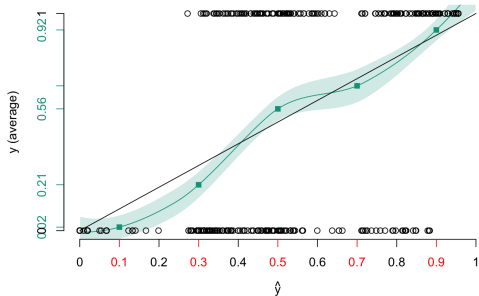
Calibration

- Consider e.g., some random forest model for \hat{m} , and then the calibration curve per decile,



Calibration

- Or, for our random forest model \hat{m} , some local regression approach



Definition 4.48: Calibration plot

The calibration plot associated with model m is the function $\hat{y} \mapsto \mathbb{E}(Y|m(\mathbf{X}) = \hat{y})$. The empirical version is some local regression on $\{y_i, m(\mathbf{x}_i)\}$.

Definition 4.49: Globally unbiased model m , [Denuit et al. \(2021\)](#)

Model m is globally unbiased if $\mathbb{E}[Y] = \mathbb{E}[m(\mathbf{X})]$.

Definition 4.50: Locally unbiased model m , [Denuit et al. \(2021\)](#)

Model m is locally unbiased at \hat{y} if $\mathbb{E}[Y|m(\mathbf{X}) = \hat{y}] = \hat{y}$.

Calibration

- For GLM, remember that

$$f(y_i) = \exp\left(\frac{y_i\theta_i - b(\theta_i)}{\varphi} + c(y_i, \varphi)\right),$$

$$\frac{\partial \log \mathcal{L}_i}{\partial \beta_j} = \frac{\partial \log \mathcal{L}_i}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_j} = \frac{\partial \log \mathcal{L}_i}{\partial \beta_j} = \frac{y_i - \mu_i}{\varphi} \cdot \frac{1}{V(\mu_i)} \cdot x_{i,j} \cdot \left(\frac{\partial \eta_i}{\partial \mu_i}\right)^{-1}$$

- When g is the canonical link ($g_* = b'^{-1}$ or $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta} = \theta_i$)

$$\nabla \log \mathcal{L} = \mathbf{X}^\top (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{0}$$

Proposition 4.34: Calibration of GLM

In the GLM framework with the canonical link function, $\hat{m}(\mathbf{x}) = g_*^{-1}(\mathbf{x}_i^\top \hat{\boldsymbol{\beta}})$ is globally unbiased (on the training dataset), but possibly locally biased.

Calibration

› Otherwise

$$\nabla \log \mathcal{L} = \mathbf{X}^\top \mathbf{\Omega} (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{0},$$

where $\mathbf{\Omega}$ is a diagonal matrix ($\mathbf{\Omega} = \mathbf{W}\mathbf{\Delta}$, where $\mathbf{W} = \text{diag}((V(\mu_i)g'(\mu_i)^2)^{-1})$ and $\mathbf{\Delta} = \text{diag}(g'(\mu_i))$), so that we recognize Fisher information - corresponding to the Hessian matrix (up to a negative sign) - $\mathbf{X}^\top \mathbf{W}\mathbf{X}$).

	training data					validation data				
	\bar{y}	GLM	CART	GAM	RF	\bar{y}	GLM	CART	GAM	RF
$\hat{m}(\mathbf{x}, s)$	8.73	8.73	8.73	8.73	8.27	8.55	9.05	9.03	8.84	8.70
$\hat{m}(\mathbf{x})$	8.73	8.73	8.73	8.73	8.29	8.55	9.05	9.03	8.84	8.73

Definition 4.51: Brier score (binary classifier) **Brier (1950)**

Brier score is the mean squared error of probability estimate,

$$\text{BS} = \frac{1}{n} \sum_{i=1}^n (\hat{m}(\mathbf{x}_i), y_i)^2$$

- › Let g be the calibration function, $g(\hat{m}(\mathbf{x})) \approx p(\mathbf{x})$.
- › Platt scaling (from **Platt et al. (1999)**), $g(s) = [1 + e^{-(ws+b)}]^{-1}$.
- › “confidence” value given by **Picpurify**, using pictures generate by a GAN (a generative adversarial network, used in **Hill and White (2020)**).

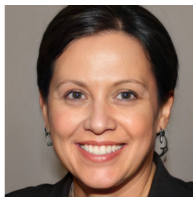
Calibration



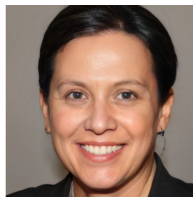
female (0.984)
male (0.016)



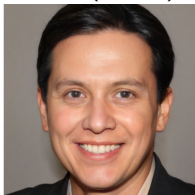
female (0.983)
male (0.017)



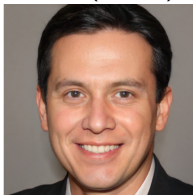
female (0.982)
male (0.018)



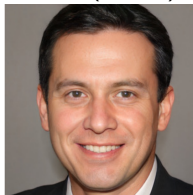
female (0.960)
male (0.040)



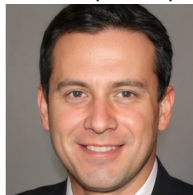
female (0.009)
male (0.991)



female (0.013)
male (0.987)

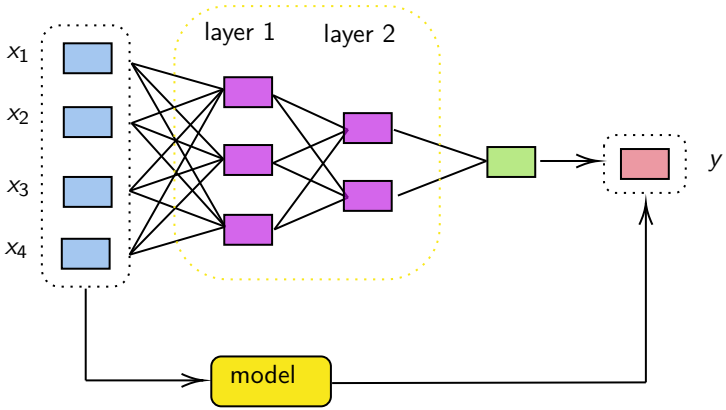


female (0.014)
male (0.986)

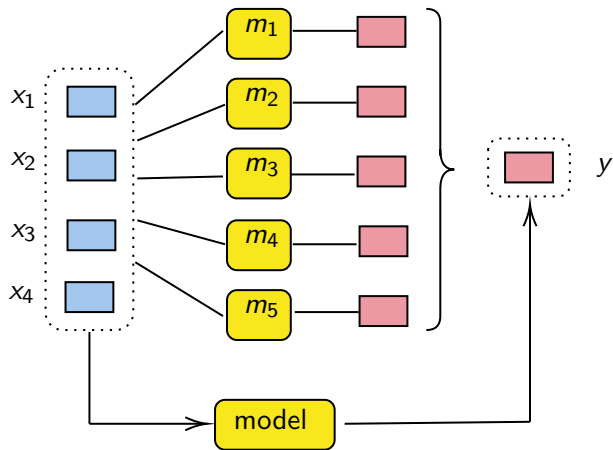


female (0.015)
male (0.985)

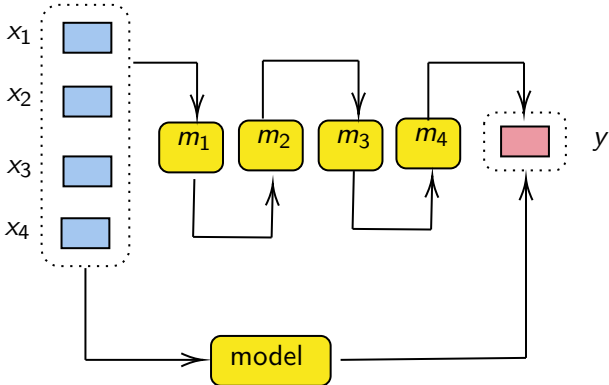
Standard modeling architecture



Standard modeling architecture



Standard modeling architecture

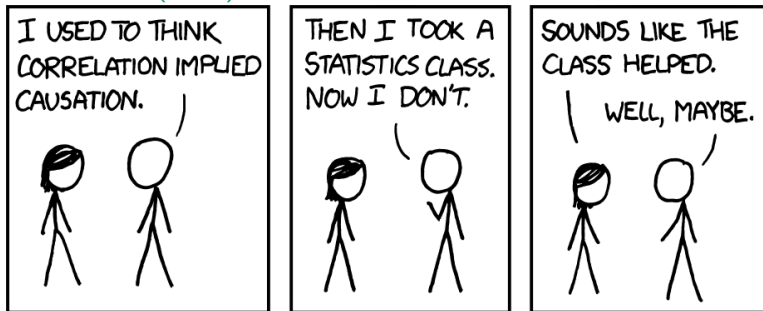


– Part 4 –
Data

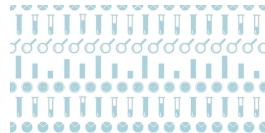
Data (the two types)

"It is often said, 'You cannot prove causality with statistics.' One of my professors, Frederick Mosteller, liked to counter, 'You can only prove causality with statistics.' (...) The title, 'Observation and Experiment,' marks the modern distinction between randomized experiments and observational studies."

Rosenbaum (2018)



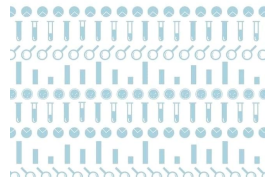
Correlation, Randall Munroe, 2009 <https://xkcd.com/552/>



Observation & Experiment

An Introduction to Causal Inference

PAUL R. ROSENBAUM



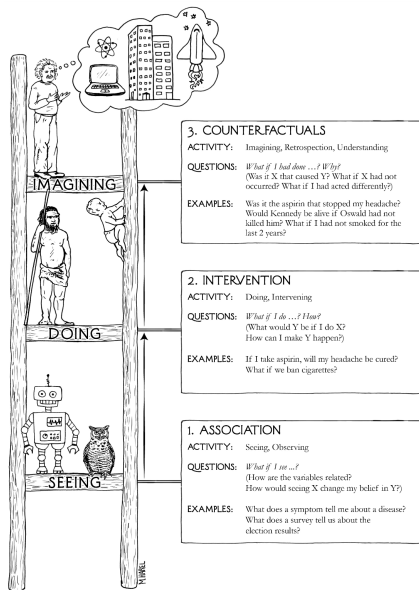
Data (the three rung ladder)

“Ladder of causation” from Pearl et al. (2009)

- 3. **Counterfactuals**
(Imagining, “*what if I had done...*”)
- 2. **Intervention**
(Doing, “*what if I do...*”)
- 1. **Association**
(Seeing, “*what if I see...*”)

Picture source: Pearl and Mackenzie (2018)

What would be the impact of a treatment T on a variable of interest Y ?



Proxy

- › “OK, let’s not use race, but should we use zip code, which of course is a proxy for race in our segregated society?,” O’Neil (2016).

Definition 4.52: Proxy, Merriam-Webster (2022)

A **proxy** is a person authorized to act for another (from a contracted form of the Middle English word *procuracie* (from French “procuration”)).

Definition 4.53: Perfect proxy, Datta et al. (2017)

A variable X is a perfect proxy for Z if there exist functions $\varphi : \mathcal{X} \rightarrow \mathcal{Z}$ and $\psi : \mathcal{Z} \rightarrow \mathcal{Y}$ such that

$$\mathbb{P}[X = \psi(Z)] = \mathbb{P}[\varphi(X) = Z] = 1.$$

Definition 4.54: Comonotonicity, [Hoeffding \(1940\)](#); [Fréchet \(1951\)](#)

Variables X and Y are comonotonic if $(X, Y) = (F_x^{-1}(U), F_y^{-1}(U))$ for some $U \sim \mathcal{U}([0, 1])$.

- › See also [Dhaene et al. \(2002a,b\)](#) on comonotonic vectors.
- › See also [Prince and Schwarcz \(2019\)](#), or [Tschantz \(2022\)](#) for discrimination by proxy.
- › Range of possible situation between [independence](#) and [perfect proxy](#).

Definition 4.55: Independence (dimension 2)

X and Y are independent, denoted $X \perp\!\!\!\perp Y$, if for any sets $\mathcal{A}, \mathcal{B} \subset \mathbb{R}$,

$$\mathbb{P}[X \in \mathcal{A}, Y \in \mathcal{B}] = \mathbb{P}[X \in \mathcal{A}] \cdot \mathbb{P}[Y \in \mathcal{B}].$$

Definition 4.56: Linear Independence (dimension 2)

Consider two random variables X and Y . $X \perp\!\!\!\perp Y$ if and only if $\text{Cov}[X, Y] = 0$.

Definition 4.57: Correlation (dimension 2), Pearson (1895)

X and Y are two random variables

$$\text{Corr}[X, Y] = \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X] \cdot \text{Var}[Y]}}.$$

where $\text{Cov}[X, Y] = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$.

Proposition 4.35: Correlation bounds (dimension 2)

For any random variables X and Y (with finite variances),
 $r_{\min} \leq \text{Corr}[X, Y] \leq r_{\max}$, where

$$r_{\min} = \frac{\text{Cov}[F_x^{-1}(U), F_y^{-1}(1 - U)]}{\sqrt{\text{Var}[X] \cdot \text{Var}[Y]}} \quad \text{and} \quad r_{\max} = \frac{\text{Cov}[F_x^{-1}(U), F_y^{-1}(U)]}{\sqrt{\text{Var}[X] \cdot \text{Var}[Y]}}$$

- Maximal correlation is obtained when X and Y are comonotonic (minimal correlation when X and $-Y$ are comonotonic).
- Related to optimal transport, see also [Knott and Smith \(1984\)](#).

Proposition 4.36

Consider two random variables X and Y . $X \perp\!\!\!\perp Y$ if and only if for any functions $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ and $\psi : \mathbb{R} \rightarrow \mathbb{R}$ (such that the expected values below exist and are well-defined) $\text{Cov}[\varphi(X), \psi(Y)] = 0$, i.e.,

$$\mathbb{E}[\varphi(X) \cdot \psi(Y)] = \mathbb{E}[\varphi(X)] \cdot \mathbb{E}[\psi(Y)].$$

Definition 4.58: Maximal Correlation, HGR

Consider two random variables X and Y ,

$$r^*(X, Y) = \max_{\varphi, \psi} \{ \text{Corr}[\varphi(X), \psi(Y)] \}.$$

Independence

- HGR because of Hirschfeld (1935), Gebelein (1941) and Rényi (1959) (also Sarmanov (1958a,b)).

$$r^*(X, Y) = \max_{\varphi \in \mathcal{F}_x, \psi \in \mathcal{G}_y} \mathbb{E}[\varphi(X)\psi(Y)],$$

where

$$\begin{cases} \mathcal{F}_x = \{\varphi : \mathcal{X} \rightarrow \mathbb{R} : \mathbb{E}[\varphi(X)] = 0 \text{ and } \mathbb{E}[\varphi^2(X)] = 1\} \\ \mathcal{G}_y = \{\psi : \mathcal{Y} \rightarrow \mathbb{R} : \mathbb{E}[\psi(Y)] = 0 \text{ and } \mathbb{E}[\psi^2(Y)] = 1\} \end{cases}$$

- See either `ccaPP` or `acepack` package,

```
1 > ccaPP::maxCorProj(x = x, y = y, method = "pearson")
2 > corstar = acepack::ace(x = x, y = y)
3 > cor(corstar$tx, corstar$ty)
```

Independence

Proposition 4.37

Consider two random variables X and Y . $X \perp\!\!\!\perp Y$ if and only if $r^*(X, Y) = 0$.

Proof: Given a random variable X , its characteristic function is $\phi_X(t) = \mathbb{E}[e^{itX}]$.

Recall that

$$\begin{cases} \phi_X(t) = \phi_Y(t), \quad \forall t \in \mathbb{R} \text{ if and only if } X \stackrel{\mathcal{L}}{=} Y \\ \phi_{X,Y}(s, t) = \mathbb{E}[e^{i(sX+tY)}] = \phi_X(s) \cdot \phi_Y(t), \quad \forall s, t \in \mathbb{R} \text{ if and only if } X \perp\!\!\!\perp Y \end{cases}$$

If $r^*(X, Y) = 0$, let $s, t \in \mathbb{R}$ and consider $\varphi(x) = \phi_X(x) = \mathbb{E}[e^{ix}]$ and $\psi(y) = \phi_Y(y) = \mathbb{E}[e^{iy}]$, then $\text{Cov}[e^{isX}, e^{itY}] = \text{Cov}[X'_s, Y'_t] = 0$, i.e. $\mathbb{E}[X'_s Y'_t] = \mathbb{E}[X'_s] \mathbb{E}[Y'_t]$,

$$\underbrace{\mathbb{E}[e^{i(sX+tY)}]}_{\phi_{XY}(s,t)} = \underbrace{\mathbb{E}[e^{isX}] \cdot \mathbb{E}[e^{itY}]}_{\phi_X(s) \cdot \phi_Y(t)}, \quad \forall s, t \in \mathbb{R} \text{ i.e. } X \perp\!\!\!\perp Y.$$

Proposition 4.38

Consider two random variables X and Y such that (X, Y) is a Gaussian vector. Then $r^*(X, Y) = |\text{Corr}[X, Y]|$.

- › See [Lancaster \(1957, 1958\)](#), and Gauss-Hermite decomposition

$$f(x, y) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(-\frac{x^2 - 2\rho xy + y^2}{2[1-\rho^2]}\right) = \phi(x)\phi(y) \cdot \sum_{i=0}^{\infty} r^i H_i(x)H_i(y)$$

where H_i 's are Hermite polynomial.

Independence

› Instead of

$$r^*(X, Y) = \max_{\varphi \in \mathcal{F}_x, \psi \in \mathcal{G}_y} \mathbb{E}[\varphi(X)\psi(Y)],$$

where

$$\begin{cases} \mathcal{F}_x = \{\varphi : \mathcal{X} \rightarrow \mathbb{R} : \mathbb{E}[\varphi(X)] = 0 \text{ and } \mathbb{E}[\varphi^2(X)] = 1\} \\ \mathcal{G}_y = \{\psi : \mathcal{Y} \rightarrow \mathbb{R} : \mathbb{E}[\psi(Y)] = 0 \text{ and } \mathbb{E}[\psi^2(Y)] = 1\} \end{cases}$$

Definition 4.59: Constrained Maximal Correlation, Bach and Jordan (2002), Gretton et al. (2005)

Consider two random variables X and Y , as well as some Hilbert spaces $\bar{\mathcal{F}}_x \subset \mathcal{F}_x$ and $\bar{\mathcal{G}}_y \subset \mathcal{G}_y$,

$$\bar{r}^*(X, Y) = \max_{\varphi \in \bar{\mathcal{F}}_x, \psi \in \bar{\mathcal{G}}_y} \{\text{Corr}[\varphi(X), \psi(Y)]\}.$$

Independence

- › Kimeldorf and Sampson (1978) and Kimeldorf et al. (1982) suggested to consider for $\bar{\mathcal{F}}_x$ and $\bar{\mathcal{G}}_y$ as subsets of monotone functions.

$$\begin{cases} \bar{\mathcal{F}}_x = \{\varphi \in \mathcal{F}_x : \varphi \text{ monotone}\} \\ \bar{\mathcal{G}}_y = \{\psi \in \mathcal{G}_y : \psi \text{ monotone}\} \end{cases}$$

- › See Mourier (1953), Hannan (1961), Jensen and Mayer (1977) and Lin (1987).

to go further → (for more details on RKHS issues)

Definition 4.60: Linear Independence

In a general context, consider two random vectors \mathbf{X} and \mathbf{Y} , in \mathbb{R}^{d_x} and \mathbb{R}^{d_y} , respectively. $\mathbf{X} \perp \mathbf{Y}$ if and only if for any $\mathbf{a} \in \mathbb{R}^{d_x}$ and $\mathbf{b} \in \mathbb{R}^{d_y}$

$$\text{Cov}[\mathbf{a}^\top \mathbf{X}, \mathbf{b}^\top \mathbf{Y}] = 0.$$

Definition 4.61: Independence

In a general context, consider two random vectors \mathbf{X} and \mathbf{Y} . $\mathbf{X} \perp\!\!\!\perp \mathbf{Y}$ if and only if for any $\mathcal{A} \subset \mathbb{R}^{d_x}$ and $\mathcal{B} \subset \mathbb{R}^{d_y}$,

$$\mathbb{P}[\{\mathbf{X} \in \mathcal{A}\} \cap \{\mathbf{Y} \in \mathcal{B}\}] = \mathbb{P}[\{\mathbf{X} \in \mathcal{A}\}] \cdot \mathbb{P}[\{\mathbf{Y} \in \mathcal{B}\}].$$

Proposition 4.39: Independence

Consider two random vectors \mathbf{X} and \mathbf{Y} . $\mathbf{X} \perp\!\!\!\perp \mathbf{Y}$ if and only if for any functions $\varphi : \mathbb{R}^{d_x} \rightarrow \mathbb{R}$ and $\psi : \mathbb{R}^{d_y} \rightarrow \mathbb{R}$ (such that the expected values below exist and are well-defined)

$$\mathbb{E}[\varphi(\mathbf{X})\psi(\mathbf{Y})] = \mathbb{E}[\varphi(\mathbf{X})] \cdot \mathbb{E}[\psi(\mathbf{Y})],$$

or equivalently

$$\text{Cov}[\varphi(\mathbf{X}), \psi(\mathbf{Y})] = 0.$$

Definition 4.62: Mutual Independence

Let $\mathbf{Y} = (Y_1, \dots, Y_k)$ denote some random vector. All components of \mathbf{Y} are (mutually) independent if for any $\mathcal{A}_1, \dots, \mathcal{A}_k \subset \mathbb{R}$

$$\mathbb{P} \left[\{(Y_1, \dots, Y_k) \in \bigcap_{i=1}^k \mathcal{A}_i\} \right] = \prod_{i=1}^k \mathbb{P}[\{Y_i \in \mathcal{A}_i\}].$$

Definition 4.63: Conditional Independence (dimension 2)

X and Y are independent conditionally on Z , denoted $X \perp\!\!\!\perp Y \mid Z$, if for any sets $\mathcal{A}, \mathcal{B}, \mathcal{C} \subset \mathbb{R}$,

$$\mathbb{P}[X \in \mathcal{A}, Y \in \mathcal{B} \mid Z \in \mathcal{C}] = \mathbb{P}[X \in \mathcal{A} \mid Z \in \mathcal{C}] \cdot \mathbb{P}[Y \in \mathcal{B} \mid Z \in \mathcal{C}].$$

Definition 4.64: Conditional Independence

In a general context, consider three random vectors \mathbf{X} , \mathbf{Y} and \mathbf{Z} . $(\mathbf{X} \perp\!\!\!\perp \mathbf{Y})|\mathbf{Z}$ if and only if for any $\mathcal{A} \subset \mathbb{R}^{d_x}$, $\mathcal{B} \subset \mathbb{R}^{d_y}$ and $\mathcal{C} \subset \mathbb{R}^{d_z}$,

$$\mathbb{P}[\{\mathbf{X} \in \mathcal{A}\} \cap \{\mathbf{Y} \in \mathcal{B}\} | \mathbf{Z} \in \mathcal{C}] = \mathbb{P}[\{\mathbf{X} \in \mathcal{A}\} | \mathbf{Z} \in \mathcal{C}] \cdot \mathbb{P}[\{\mathbf{Y} \in \mathcal{B}\} | \mathbf{Z} \in \mathcal{C}].$$

Proposition 4.40

Consider three random variables X , Y , and Z . If $X \perp Z$ and $Y \perp Z$, then $aX + bY \perp Z$, for any $a, b \in \mathbb{R}$.

Independence

Proposition 4.41: $X \perp Z, Y \perp Z \not\Rightarrow \psi(X, Y) \perp Z$

Consider three random variables X , Y , and Z . If $X \perp Z$ and $Y \perp Z$, it does not imply that $\psi(X, Y) \perp Z$, for any $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$.

$$(X, Y, Z) = \begin{cases} (0, 0, 0) & \text{with probability } 1/4, \\ (0, 1, 1) & \text{with probability } 1/4, \\ (1, 0, 1) & \text{with probability } 1/4, \\ (1, 1, 0) & \text{with probability } 1/4. \end{cases}$$

Proposition 4.42

Consider a random vector \mathbf{X} in \mathbb{R}^k , and a random variable Z . $\mathbf{X} \perp Z$ does not imply that $\psi(\mathbf{X}) \perp Z$, for any $\psi : \mathbb{R}^k \rightarrow \mathbb{R}$.

Proposition 4.43

Consider three random variables X , Y , and Z . Even if $X \perp\!\!\!\perp Z$ and $Y \perp\!\!\!\perp Z$, it does not imply either that $\psi(X, Y) \perp Z$ or that $\psi(X, Y) \perp\!\!\!\perp Z$, for any $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$.

Proposition 4.44

Consider a random vector \mathbf{X} in \mathbb{R}^k , and a random variable Z .

$\mathbf{X} \perp\!\!\!\perp Z$ does not imply either that $\psi(\mathbf{X}) \perp Z$ **or** $\psi(\mathbf{X}) \perp\!\!\!\perp Z$, for any $\psi : \mathbb{R}^k \rightarrow \mathbb{R}$.

Causation

Definition 4.65: Common cause, [Reichenbach \(1956\)](#)

If X and Y are non-independent, $X \not\perp\!\!\!\perp Y$, then, either

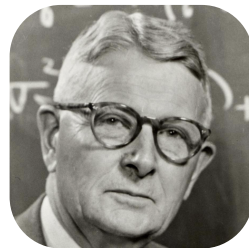
$$\left\{ \begin{array}{l} X \text{ causes } Y \\ Y \text{ causes } X \\ \text{there exists } Z \text{ such that } Z \text{ causes both } X \text{ and } Y. \end{array} \right.$$

- › See also [Bollen and Pearl \(2013\)](#)
- › SCM, [Goldberger \(1972\)](#), [Duncan \(1975\)](#) or [Bollen \(1989\)](#)
- › Bayesian network, [Pearl \(1985\)](#), [Henrion \(1988\)](#), [Charniak \(1991\)](#)
- › Causal path diagrams and probabilistic DAGs, [Spirtes et al. \(1993\)](#)

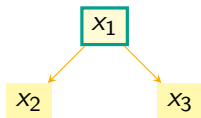


Causation

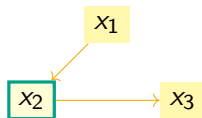
➤ Sewall Wright (see [Wright \(1921, 1934\)](#)) use **directed graphs** to represent probabilistic cause and effect relationships among a set of variables, and developed path diagrams and path analysis



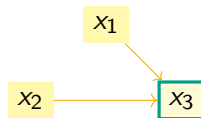
(a)
confounder



(b)
mediator



(c)
collider



Causation

Definition 4.66: Path

A path π from a node x_i to another node x_j is a sequence of nodes and edges starting at x_i and ending at x_j .

Definition 4.67: d -separation

A set of nodes \mathbf{x}_i is said to be d -separated with another set of nodes \mathbf{x}_j by \mathbf{x}_c whenever every path from any $x_i \in \mathbf{x}_i$ to any $x_j \in \mathbf{x}_j$ is blocked by \mathbf{x}_c . We will simply denote $\mathbf{x}_i \perp_{\mathcal{G}} \mathbf{x}_j \mid \mathbf{x}_c$.

Proposition 4.45

Two nodes x_i and x_j are d -separated by \mathbf{x}_c if and only members of \mathbf{x}_c block all paths from x_i to x_j .

Causation

- Chain rule :
$$\begin{cases} \mathbb{P}[x_1, x_2, x_3, x_4] = \mathbb{P}[x_1] \times \mathbb{P}[x_2|x_1] \times \mathbb{P}[x_3|x_1, x_2] \times \mathbb{P}[x_4|x_1, x_2, x_3] \\ \mathbb{P}[x_1, x_2, x_3, x_4] = \mathbb{P}[x_4] \times \mathbb{P}[x_3|x_4] \times \mathbb{P}[x_2|x_3, x_4] \times \mathbb{P}[x_1|x_2, x_3, x_4] \end{cases}$$

Definition 4.68: Directed acyclic graph, DAG (or causal graph)

A directed acyclic graph (DAG) \mathcal{G} is a directed graph with no directed cycles.

Definition 4.69: Markov Property

Given a causal graph \mathcal{G} with nodes \mathbf{x} , the joint distribution of \mathbf{X} satisfies the (global) Markov property with respect to \mathcal{G} if, for any disjoint \mathbf{x}_1 , \mathbf{x}_2 and \mathbf{x}_c

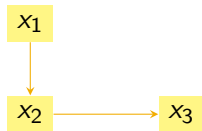
$$\mathbf{x}_1 \perp_{\mathcal{G}} \mathbf{x}_2 \mid \mathbf{x}_c \Rightarrow \mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2 \mid \mathbf{X}_c.$$

Proposition 4.46: Probabilistic graphical model

If \mathbf{X} satisfies the (global) Markov property with respect to \mathcal{G}

$$\mathbb{P}[x_1, \dots, x_n] = \prod_{i=1}^n \mathbb{P}[x_i | \text{parents}(x_i)]$$

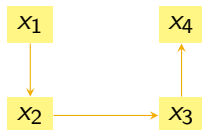
where $\text{parents}(x_i)$ are nodes with edges directed towards x_i



➤ Path from x_1 to x_3 is blocked by x_2 , i.e., $x_1 \perp_{\mathcal{G}} x_3 \mid x_2$, or $X_1 \perp\!\!\!\perp X_3 \mid X_2$. From the chain rule,

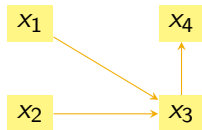
$$\mathbb{P}[x_1, x_2, x_3] = \mathbb{P}[x_1] \times \mathbb{P}[x_2 | x_1] \times \underbrace{\mathbb{P}[x_3 | x_2, x_1]}_{\mathbb{P}[x_3 | x_2]}$$

Causation



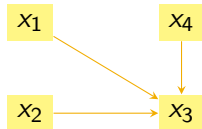
- From the chain rule, for the causal graph on the left (top),

$$\mathbb{P}[x_1, x_2, x_3, x_4] = \mathbb{P}[x_1] \times \mathbb{P}[x_2|x_1] \times \mathbb{P}[x_3|x_2] \times \mathbb{P}[x_4|x_3]$$



- From the chain rule, for the causal graph on the left (middle),

$$\mathbb{P}[x_1, x_2, x_3, x_4] = \mathbb{P}[x_1] \times \mathbb{P}[x_2] \times \mathbb{P}[x_3|x_1, x_2] \times \mathbb{P}[x_4|x_3]$$



- From the chain rule, for the causal graph on the left (bottom),

$$\mathbb{P}[x_1, x_2, x_3, x_4] = \mathbb{P}[x_1] \times \mathbb{P}[x_2] \times \mathbb{P}[x_3|x_1, x_2, x_4] \times \mathbb{P}[x_4]$$

Intervention

- › $\mathbb{P}[Y \in \mathcal{A}|X = x]$: how $Y \in \mathcal{A}$ is likely to occur if X happened to be equal to x
- › Therefore, it is an observational statement.
- › $P[Y \in \mathcal{A}|\text{do}(X = x)]$: how $Y \in \mathcal{A}$ is likely to occur if X is set to x
- › It is here an intervention statement.
- › Using causal graphs, intervention $\text{do}(X = x)$ means that all incoming edges to x are cut.
- › If $P[Y \in \mathcal{A}|\text{do}(X = x)] \neq \mathbb{P}[Y \in \mathcal{A}|X = x]$, it means that X and Y are confounded, see [Pearl \(2009\)](#).

Definition 4.70: Structural Causal Models (SCM)

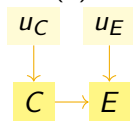
In a simple causal graph, with two nodes C (the cause) and E (the effect), the causal graph is $C \rightarrow E$, and the mathematical interpretation can be summarized in two assignments

$$\begin{cases} C = h_c(U_C) \\ E = h_e(C, U_E), \end{cases}$$

where U_C and U_E are two independent random variables, $U_C \perp\!\!\!\perp U_E$.

(a)

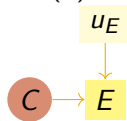
observation



$$\begin{cases} C = h_c(U_C) \\ E = h_e(C, U_E) \end{cases}$$

(b)

intervention

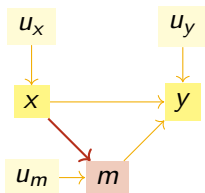


$$\begin{cases} C = c \text{ (or do}(C = c)) \\ E_c^* = h_e(c, U_E) \end{cases}$$

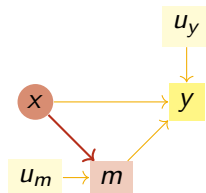
Intervention

(a)

m mediator variable

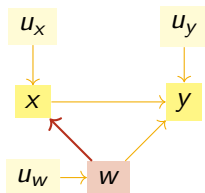


(b)

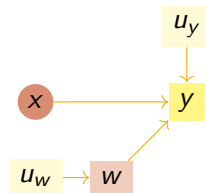


(c)

w confounding variable



(d)



$$\begin{cases} \text{mediator : } & \mathbb{P}[Y_x^* = 1] = \mathbb{P}[Y = 1 | \text{do}(X = x)] = \mathbb{P}[Y = 1 | X = x] \\ \text{confusion : } & \mathbb{P}[Y_x^* = 1] = \mathbb{P}[Y = 1 | \text{do}(X = x)] \neq \mathbb{P}[Y = 1 | X = x]. \end{cases}$$

Intervention

› In fact, in the presence of a confounding factor, $\mathbb{P}[Y_x^* = 1]$ which corresponds to $\mathbb{P}[Y = 1|\text{do}(X = x)]$ should be written

$$\sum_w \mathbb{P}[Y = 1|W = w, X = x] \cdot \mathbb{P}[W = w] = \mathbb{E}(\mathbb{P}[Y = 1|W, X = x]).$$

Causal Inference and counterfactuals

› Define **potential outcomes** to quantify the treatment effect, $TE = y_{i,T \leftarrow 1}^* - y_{i,T \leftarrow 0}^*$

$\left\{ \begin{array}{l} \text{observation} : y_{i,T \leftarrow 1}^* \text{ when } t_i = 1 \text{ is observed, and } \mathbf{x}_i \\ \text{counterfactual} : y_{i,T \leftarrow 0}^* \text{ when } t_i = 1 \text{ is observed, and } \mathbf{x}_i \end{array} \right.$

› Here we want to observe counterfactuals $y_{i,T \leftarrow t'}^*$ at the individual level.

	Gender	Name	Treatment	Outcome (Weight)				Height	...
		t_i		y_i	$y_{i,T \leftarrow 0}^*$	$y_{i,T \leftarrow 1}^*$	TE	x_i	...
1	H	Alex	0	75	75	64	11	172	...
2	F	Betty	1	52	67	52	15	161	...
3	F	Beatrix	1	57	71	57	14	163	...
4	H	Ahmad	0	78	78	61	17	183	...

› Different notations are used $y(1)$ and $y(0)$ in Imbens and Rubin (2015), y^1 and y^0 in Cunningham (2021), or $y_{t=1}$ and $y_{t=0}$ in Pearl and Mackenzie (2018).

Causal Inference and counterfactuals

› Define **potential outcomes** to quantify the treatment effect, $TE = y_{i,T \leftarrow 1}^* - y_{i,T \leftarrow 0}^*$

$\left\{ \begin{array}{l} \text{observation} : y_{i,T \leftarrow 1}^* \text{ when } t_i = 1 \text{ is observed, and } \mathbf{x}_i \\ \text{counterfactual} : y_{i,T \leftarrow 0}^* \text{ when } t_i = 1 \text{ is observed, and } \mathbf{x}_i \end{array} \right.$

› Here we want to observe counterfactuals $y_{i,T \leftarrow t'}^*$ at the individual level.

	Gender	Name	Treatment	Outcome (Weight)				Height	...
		t_i		y_i	$y_{i,T \leftarrow 0}^*$	$y_{i,T \leftarrow 1}^*$	TE	x_i	...
1	H	Alex	0	75	75	?	?	172	...
2	F	Betty	1	52	?	52	?	161	...
3	F	Beatrix	1	57	?	57	?	163	...
4	H	Ahmad	0	78	78	?	?	183	...

› Different notations are used $y(1)$ and $y(0)$ in Imbens and Rubin (2015), y^1 and y^0 in Cunningham (2021), or $y_{t=1}$ and $y_{t=0}$ in Pearl and Mackenzie (2018).

Causal Inference and counterfactuals

Definition 4.71: Average Treatment Effect, **Holland (1986)**

Given a treatment T , the average treatment effect on outcome Y is

$$\tau = \text{ATE} = \mathbb{E}[Y_{t \leftarrow 1}^* - Y_{t \leftarrow 0}^*].$$

Definition 4.72: Conditional Average Treatment Effect, **Wager and Athey (2018)**

Given a treatment T , the conditional average treatment effect on outcome Y , given some covariates \mathbf{X} , is

$$\tau(\mathbf{x}) = \text{CATE}(\mathbf{x}) = \mathbb{E}[Y_{t \leftarrow 1}^* - Y_{t \leftarrow 0}^* | \mathbf{X} = \mathbf{x}].$$

Definition 4.73: Individual Average Treatment Effect

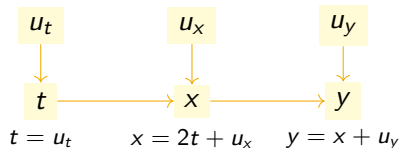
Given a treatment T , the conditional average treatment effect on outcome Y , for individual i , given covariates \mathbf{X}_i , is

$$\text{IATE}(i) = \mathbb{E}[Y_{i,t \leftarrow (1-t_i)}^* - Y_{i,t \leftarrow t_i}^*].$$

Causal Inference and counterfactuals

› Pearl (2009) suggest to use a **twin network representation of the counterfactual**

› Start with a simple structural causal model, e.g.,

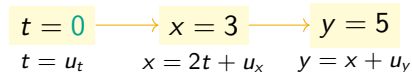


$$\begin{cases} t = u_t \\ x = 2t + u_x \\ y = x + u_y \end{cases}$$

› Suppose we were able to estimate that model

Causal Inference and counterfactuals

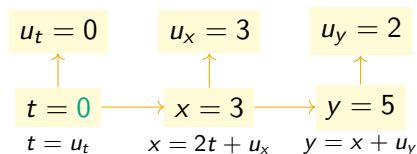
- Pearl (2009) suggest to use a twin network representation of the counterfactual



- Consider a single observation
- i.e. a triplet (t_i, x_i, y_i)

Causal Inference and counterfactuals

› Pearl (2009) suggest to use a **twin network representation of the counterfactual**



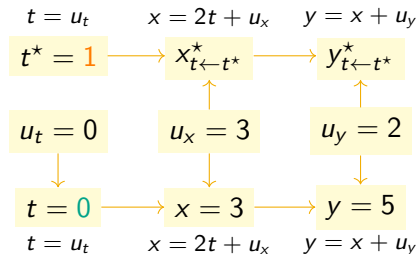
› Inverse the SCM

$$\begin{cases} u_t = t \\ u_x = x - 2t \\ u_y = y - x \end{cases}$$

- › From that triplet (t_i, x_i, y_i)
- › Derive unobserved u 's.

Causal Inference and counterfactuals

➤ Pearl (2009) suggest to use a **twin network representation of the counterfactual**

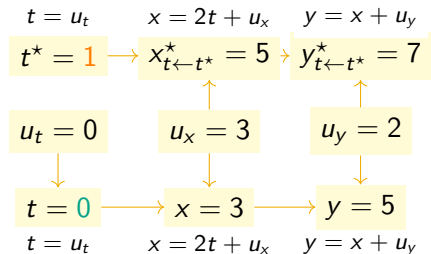


➤ Suppose that the same SCM holds in the counterfactual world

$$\begin{cases} t^* \\ x_{t^* \leftarrow t^*}^* = 2t^* + u_t \\ y_{t^* \leftarrow t^*}^* = x_{t^* \leftarrow t^*}^* + u_y \end{cases}$$

Causal Inference and counterfactuals

- Pearl (2009) suggest to use a twin network representation of the counterfactual

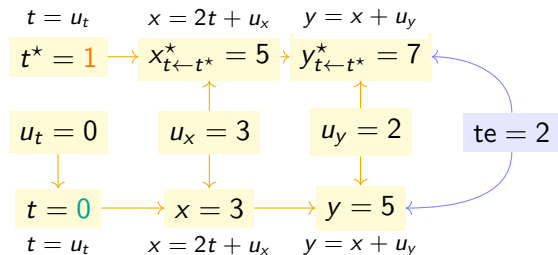


- Plugin u 's obtained from (t_i, \mathbf{x}_i, y_i)

$$\begin{cases} t^* \\ x^*_{t^* \leftarrow t^*} = 2t^* + u_t \\ y^*_{t^* \leftarrow t^*} = x^*_{t^* \leftarrow t^*} + u_y \end{cases}$$

Causal Inference and counterfactuals

► Pearl (2009) suggest to use a **twin network representation of the counterfactual**

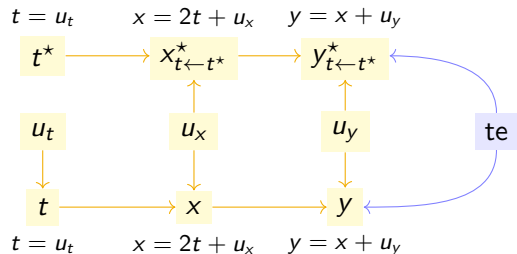


► We can compute the treatment effect

$$te = y - y_{t^*}^*$$

Causal Inference and counterfactuals

- Pearl (2009) suggest to use a **twin network representation of the counterfactual**



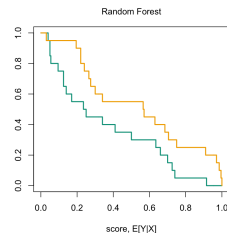
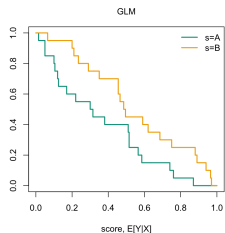
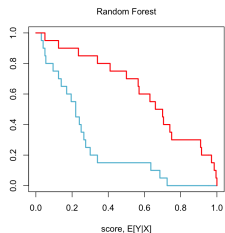
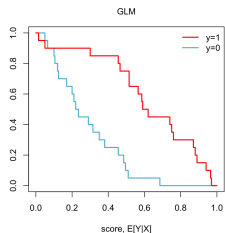
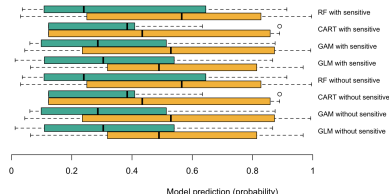
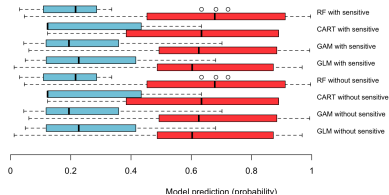
- Works well only if the SCM can be
- inverted (to derive the u 's)

– Part 5 –

Group Fairness

Group Fairness

➤ Back on **toydata2**, distributions of scores, $\hat{m}(x_i)$'s conditional on y_i and s_i



Group Fairness

Definition 5.1: Fairness through unawareness, [Dwork et al. \(2012\)](#)

A model m satisfies the fairness through unawareness criteria, with respect to sensitive attribute $s \in \mathcal{S}$ if $m : \mathcal{X} \rightarrow \mathcal{Y}$.

by Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold and Richard Zemel,



Group Fairness

› See introduction about the gender directive,

“institutional messages of color blindness may therefore artificially depress formal reporting of racial injustice. Color-blind messages may thus appear to function effectively on the surface even as they allow explicit forms of bias to persist,”

Apfelbaum et al. (2010)

Definition 5.2: Aware and unaware regression functions μ

The aware regression function is $\mu(\mathbf{x}, s) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}, S = s]$
and the unaware regression function is $\mu(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$.

Historical Perspective: "Cultural Fairness" and "Statistical Discrimination"

Definition 5.3: Four definitions of cultural fairness, Darlington (1971)

A test (\hat{y}) is considered "culturally fair" if it fits the appropriate equation

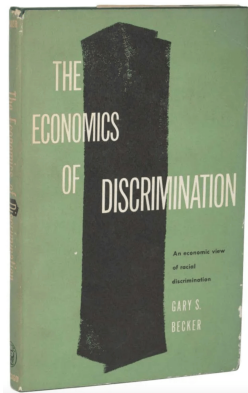
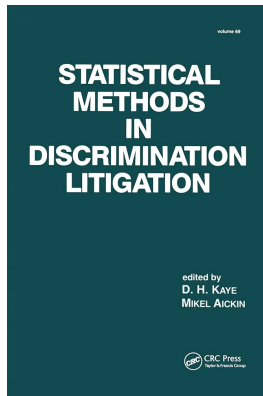
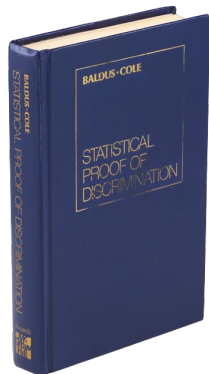
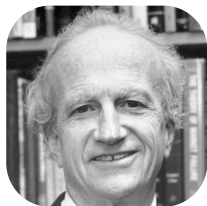
$$\left\{ \begin{array}{l} \text{Cor}[S, \hat{Y}] = \text{Cor}[S, Y] / \text{Cor}[Y, \hat{Y}] \\ \text{Cor}[S, \hat{Y}] = \text{Cor}[S, Y] \\ \text{Cor}[S, \hat{Y}] = \text{Cor}[S, Y] \cdot \text{Cor}[Y, \hat{Y}] \\ \text{Cor}[S, \hat{Y}] = 0 \end{array} \right.$$



See also [Thorndike \(1971\)](#), [Linn and Werts \(1971\)](#), following [Cleary \(1968\)](#).

"Economics of Discrimination" and "Statistical Discrimination"

- See [Becker \(1957\)](#) or [Baldus and Cole \(1980\)](#), among (many) others.



Historical Perspective: Decomposition

$$\begin{cases} y_{A:i} = \mathbf{x}_{A:i}^\top \boldsymbol{\beta}_A + \varepsilon_{A:i} \text{ (group A)}, & \bar{y}_A = \bar{\mathbf{x}}_A^\top \hat{\boldsymbol{\beta}}_A \\ y_{B:i} = \mathbf{x}_{B:i}^\top \boldsymbol{\beta}_B + \varepsilon_{B:i} \text{ (group B)}, & \bar{y}_B = \bar{\mathbf{x}}_B^\top \hat{\boldsymbol{\beta}}_B. \end{cases}$$

› Using ordinary least squares estimates

Definition 5.4: Kitagawa (1955), Oaxaca (1973),
Blinder (1973)

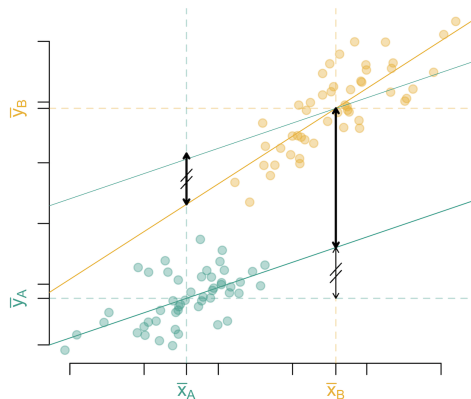
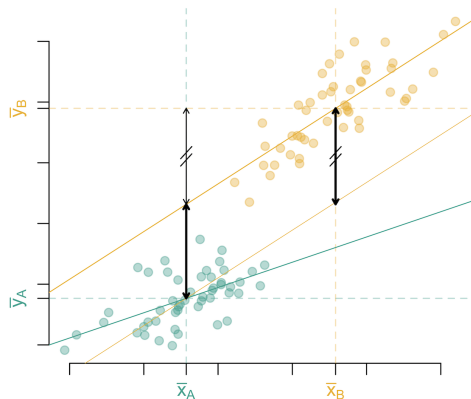
$$\bar{y}_A - \bar{y}_B = \underbrace{(\bar{\mathbf{x}}_A - \bar{\mathbf{x}}_B)^\top \hat{\boldsymbol{\beta}}_B}_{\text{characteristics}} + \underbrace{\bar{\mathbf{x}}_A^\top (\hat{\boldsymbol{\beta}}_A - \hat{\boldsymbol{\beta}}_B)}_{\text{coefficients}}, \quad (7)$$

$$\bar{y}_A - \bar{y}_B = \underbrace{(\bar{\mathbf{x}}_A - \bar{\mathbf{x}}_B)^\top \hat{\boldsymbol{\beta}}_A}_{\text{characteristics}} + \underbrace{\bar{\mathbf{x}}_B^\top (\hat{\boldsymbol{\beta}}_A - \hat{\boldsymbol{\beta}}_B)}_{\text{coefficients}}. \quad (8)$$

› Also Brown et al. (1980) and Conway and Roberts (1983).



Historical Perspective: Decomposition



$x_A(\hat{\beta}_A - \hat{\beta}_B)$ and $(\bar{x}_A - \bar{x}_B)\hat{\beta}_B$ (as in Equation 7) on the left

$x_B(\hat{\beta}_A - \hat{\beta}_B)$ and $(\bar{x}_A - \bar{x}_B)\hat{\beta}_A$ (as in Equation 8) on the right.

Independence and Demographic Parity

Definition 5.5: Independence, [Barocas et al. \(2017\)](#)

A model m satisfies the independence property if $m(\mathbf{Z}) \perp\!\!\!\perp S$, with respect to the distribution \mathbb{P} of the triplet (\mathbf{X}, S, Y) .

by Solon Barocas, Moritz Hardt and Arvind Narayanan



➤ For classifiers, one might ask for independence $\hat{Y} \perp\!\!\!\perp S$ (where \hat{y} is a class), as [Darlington \(1971\)](#).

Independence and Demographic Parity

Definition 5.6: Demographic Parity, Calders and Verwer (2010), Corbett-Davies et al. (2017)

A decision function \hat{y} – or a classifier m_t , taking values in $\{0, 1\}$ – satisfies demographic parity, with respect to some sensitive attribute S if (equivalently)

$$\begin{cases} \mathbb{P}[\hat{Y} = 1 | S = \mathbf{A}] = \mathbb{P}[\hat{Y} = 1 | S = \mathbf{B}] = \mathbb{P}[\hat{Y} = 1] \\ \mathbb{E}[\hat{Y} | S = \mathbf{A}] = \mathbb{E}[\hat{Y} | S = \mathbf{B}] = \mathbb{E}[\hat{Y}] \\ \mathbb{P}[m_t(\mathbf{Z}) = 1 | S = \mathbf{A}] = \mathbb{P}[m_t(\mathbf{Z}) = 1 | S = \mathbf{B}] = \mathbb{P}[m_t(\mathbf{Z}) = 1]. \end{cases}$$

by Toon Calders, Sizzo Verwer, Sam Corbett-Davies, Emma Pierson, Sharad Goel, etc



Independence and Demographic Parity

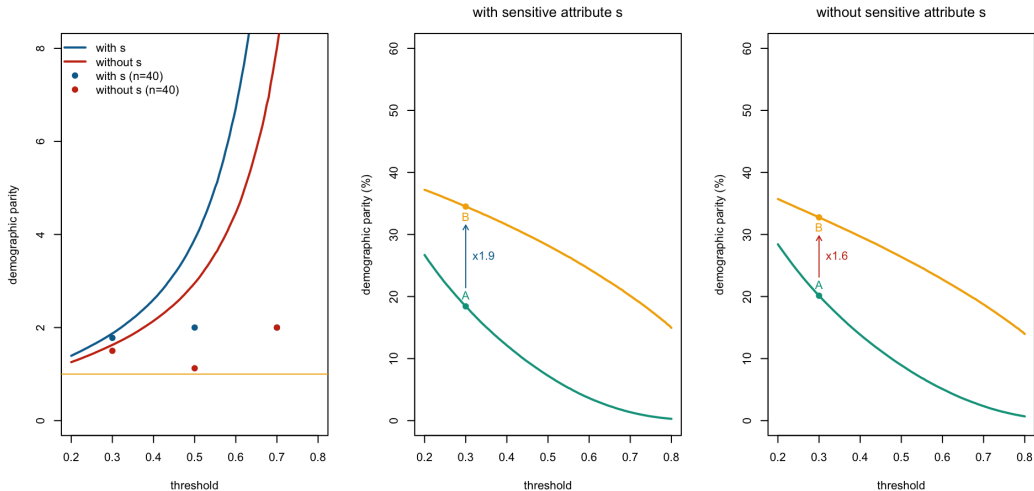
	unaware (without s)				aware (with s)			
	GLM	GAM	CART	RF	GLM	GAM	CART	RF
$n = 1000$, various t , ratio $\mathbb{P}[\hat{Y} = 1 S = \text{B}]/\mathbb{P}[\hat{Y} = 1 S = \text{A}]$								
$t = 30\%$	1.652	1.519	1.235	1.559	1.918	1.714	1.235	1.798
$t = 50\%$	1.877	2.451	2.918	2.404	2.944	3.457	2.918	2.180
$t = 70\%$	6.033	8.711	26.000	4.621	7.917	19.333	26.000	4.578

(`dem_parity` from R package `fairness`)

- On the left-hand side, evolution of the ratio $\mathbb{P}[\hat{Y} = 1|S = \text{B}]/\mathbb{P}[\hat{Y} = 1|S = \text{A}]$. The horizontal line (at $y = 1$) corresponds to perfect demographic parity.

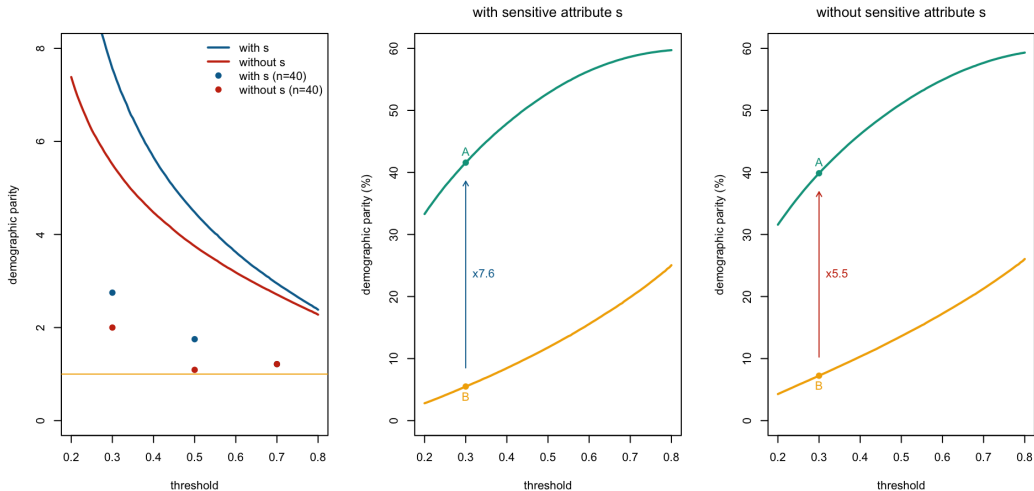
In the middle $t \mapsto \mathbb{P}[m_t(\mathbf{X}) > t|S = \text{B}]$ and $t \mapsto \mathbb{P}[m_t(\mathbf{X}) > t|S = \text{A}]$ on the model with s , and on the right-hand side without s .

Independence and Demographic Parity



On the left-hand side, evolution of the ratio ratio $\mathbb{P}[\hat{Y} = 1|S = B]/\mathbb{P}[\hat{Y} = 1|S = A]$.

Independence and Demographic Parity



➤ On the left-hand side, evolution of the ratio ratio $\mathbb{P}[\hat{Y} = 0|S = A]/\mathbb{P}[\hat{Y} = 0|S = B]$

Independence and Demographic Parity

Definition 5.7: Weak Demographic Parity

A decision function \hat{y} satisfies weak demographic parity if

$$\mathbb{E}[\hat{Y}|S = \mathbf{A}] = \mathbb{E}[\hat{Y}|S = \mathbf{B}].$$

Definition 5.8: Strong Demographic Parity

A decision function \hat{y} satisfies demographic parity if $\hat{Y} \perp\!\!\!\perp S$, i.e., for all A ,

$$\mathbb{P}[\hat{Y} \in \mathcal{A}|S = \mathbf{A}] = \mathbb{P}[\hat{Y} \in \mathcal{A}|S = \mathbf{B}], \quad \forall \mathcal{A} \subset \mathcal{Y}.$$

Independence and Demographic Parity

Proposition 5.1

A model m satisfies the strong demographic parity property if and only if

$$d_{\text{TV}}(\mathbb{P}_{m|\mathbf{A}}, \mathbb{P}_{m|\mathbf{B}}) = d_{\text{TV}}(\mathbb{P}_{\mathbf{A}}, \mathbb{P}_{\mathbf{B}}) = 0.$$

➤ $d_{\text{TV}}(\mathbb{P}_{m|\mathbf{A}}, \mathbb{P}_{m|\mathbf{B}})$ could be seen as a measure of “unfairness”, but for a non-binary sensitive attribute, a more general definition is necessary (see [Denis et al. \(2021\)](#)).

Independence and Demographic Parity

Definition 5.9: Conditional demographic parity, Corbett-Davies et al. (2017)

We will have a conditional demographic parity if (at choice) for all \mathbf{x} ,

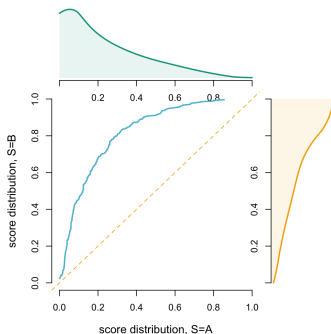
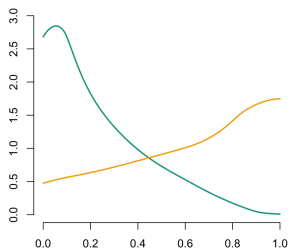
$$\begin{cases} \mathbb{P}[\hat{Y} = 1 | \mathbf{X}_L = \mathbf{x}, S = \mathbf{A}] = \mathbb{P}[\hat{Y} = 1 | \mathbf{X}_L = \mathbf{x}, S = \mathbf{B}], \quad \forall y \in \{0, 1\} \\ \mathbb{E}[\hat{Y} | \mathbf{X}_L = \mathbf{x}, S = \mathbf{A}] = \mathbb{E}[\hat{Y} | \mathbf{X}_L = \mathbf{x}, S = \mathbf{B}], \\ \mathbb{P}[\hat{Y} \in \mathcal{A} | \mathbf{X}_L = \mathbf{x}, S = \mathbf{A}] = \mathbb{P}[\hat{Y} \in \mathcal{A} | \mathbf{X}_L = \mathbf{x}, S = \mathbf{B}], \quad \forall \mathcal{A} \subset \mathcal{Y}, \end{cases}$$

where L denotes a “legitimate” subset of unprotected covariates.

Independence and Demographic Parity

Proposition 5.2

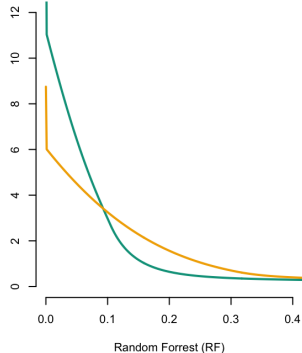
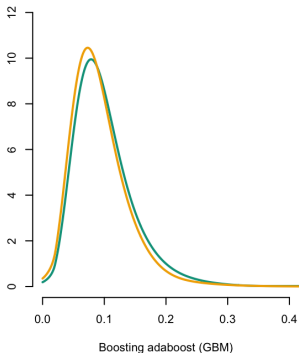
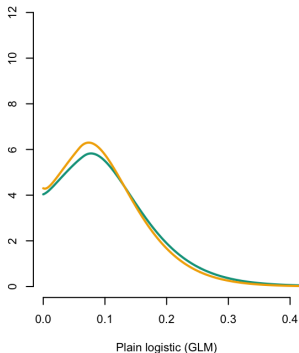
A model m satisfies is strongly fair if and only if $W_2(\mathbb{P}_A, \mathbb{P}_B) = 0$.



```
1 > model_glm = glm(y~x1
+ x2+x3, data=
toydata2, family=
binomial)
2 > pred_y_glm = predict
(model_glm, type="
response")
3 > sA = pred_y_glm[
toydata2$sensitive
=="A"]
4 > library(transport)
5 > wasserstein1d(sA,sB)
6 [1] 0.3860795
```

Independence and Demographic Parity

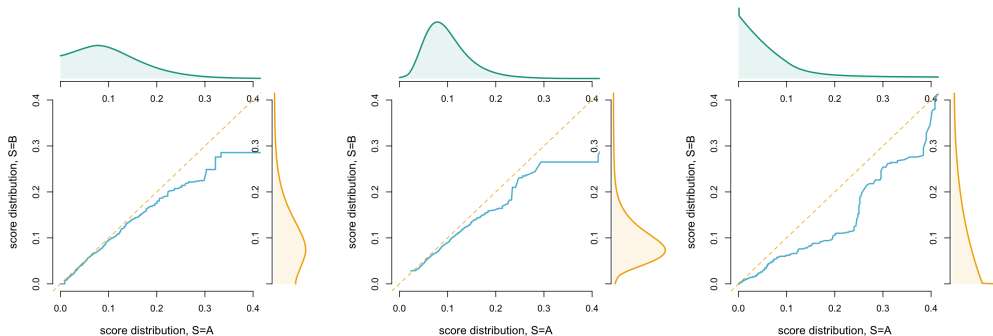
- On the **FrenchMotor** dataset, consider GLM, GBM and RF for claim occurrence



```
1 > wasserstein1d(lA,lB) 1 > wasserstein1d(bA,bB) 1 > wasserstein1d(fA,fB)
2 [1] 0.007220468      2 [1] 0.008895917      2 [1] 0.01001088
```


Independence and Demographic Parity

```
1 > wasserstein1d(lA,lB) 1 > wasserstein1d(bA,bB) 1 > wasserstein1d(fA,fB)
2 [1] 0.007220468      2 [1] 0.008895917      2 [1] 0.01001088
```



Independence and Demographic Parity

Definition 5.10: Unfairness, Denis et al. (2021); Chzhen and Schreuder (2022)

Given a model m , let \mathbb{P}_m denote the distribution of $m(\mathbf{X}, S)$ and $\mathbb{P}_{m|s}$ denote the conditional distribution of $m(\mathbf{X}, S)$ given $S = s$, define

$$\left\{ \begin{array}{l} \mathcal{U}_{\text{TV}}(m) = \max_{s \in \{A, B\}} \{d_{\text{TV}}(\mathbb{P}_m, \mathbb{P}_{m|s})\} \text{ or } \sum_{s \in \{A, B\}} d_{\text{TV}}(\mathbb{P}_m, \mathbb{P}_{m|s}) \\ \mathcal{U}_{\text{KS}}(m) = \max_{s \in \{A, B\}} \{d_{\text{KS}}(\mathbb{P}_m, \mathbb{P}_{m|s})\} \text{ or } \sum_{s \in \{A, B\}} d_{\text{KS}}(\mathbb{P}_m, \mathbb{P}_{m|s}) \\ \mathcal{U}_{W_k}(m) = \max_{s \in \{A, B\}} \{W_k(\mathbb{P}_m, \mathbb{P}_{m|s})\} \text{ or } \sum_{s \in \{A, B\}} W_k(\mathbb{P}_m, \mathbb{P}_{m|s}) \end{array} \right.$$

➤ In the original version, Chzhen and Schreuder (2022) suggested to use the one on the right.

Independence and Demographic Parity

- › Those measures characterize strong demographic parity,

Proposition 5.3: Strong Demographic Parity

A model m is strongly fair if and only if $\mathcal{U}(m) = 0$.

Separation and Equalized Odds

Definition 5.11: Separation, [Barocas et al. \(2017\)](#)

A model $m : \mathcal{Z} \rightarrow \mathcal{Y}$ satisfies the separation property if $m(\mathbf{Z}) \perp\!\!\!\perp S \mid Y$, with respect to the distribution \mathbb{P} of the triplet (\mathbf{X}, S, Y) .

by Solon Barocas, Moritz Hardt and Arvind Narayanan



Separation and Equalized Odds

Definition 5.12: True positive equality, (Weak) Equal Opportunity, [Hardt et al. \(2016\)](#)

A decision function \hat{y} – or a classifier $m_t(\cdot)$, taking values in $\{0, 1\}$ – satisfies equal opportunity, with respect to some sensitive attribute S if

$$\begin{cases} \mathbb{P}[\hat{Y} = 1 | S = \mathbf{A}, Y = 1] = \mathbb{P}[\hat{Y} = 1 | S = \mathbf{B}, Y = 1] = \mathbb{P}[\hat{Y} = 1 | Y = 1] \\ \mathbb{P}[m_t(\mathbf{Z}) = 1 | S = \mathbf{A}, Y = 1] = \mathbb{P}[m_t(\mathbf{Z}) = 1 | S = \mathbf{B}, Y = 1] = \mathbb{P}[m_t(\mathbf{Z}) = 1 | Y = 1], \end{cases}$$

which corresponds to parity of true positives, in the two groups, $\{\mathbf{A}, \mathbf{B}\}$.

Definition 5.13: Strong Equal Opportunity

A classifier $m(\cdot)$, taking values in $\{0, 1\}$, satisfies equal opportunity, with respect to some sensitive attribute S if

$$\mathbb{P}[m(\mathbf{X}, S) \in \mathcal{A} | S = \mathbf{A}, Y = 1] = \mathbb{P}[m(\mathbf{X}, S) \in \mathcal{A} | S$$

for all $\mathcal{A} \subset [0, 1]$.

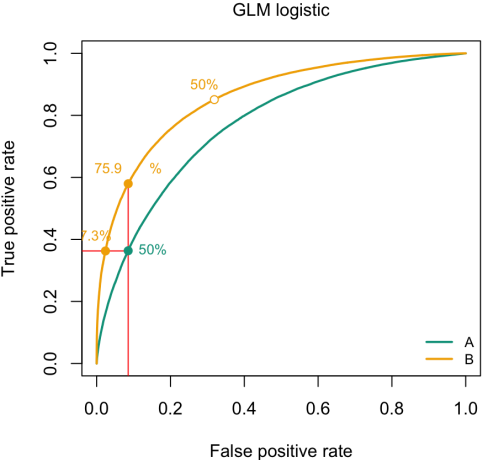
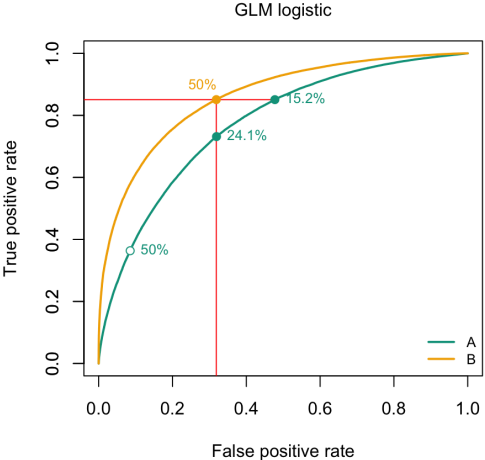
Separation and Equalized Odds

Definition 5.14: False positive equality, [Hardt et al. \(2016\)](#)

A decision function \hat{y} – or a classifier $m_t(\cdot)$, taking values in $\{0, 1\}$ – satisfies parity of false positives, with respect to some sensitive attribute s , if

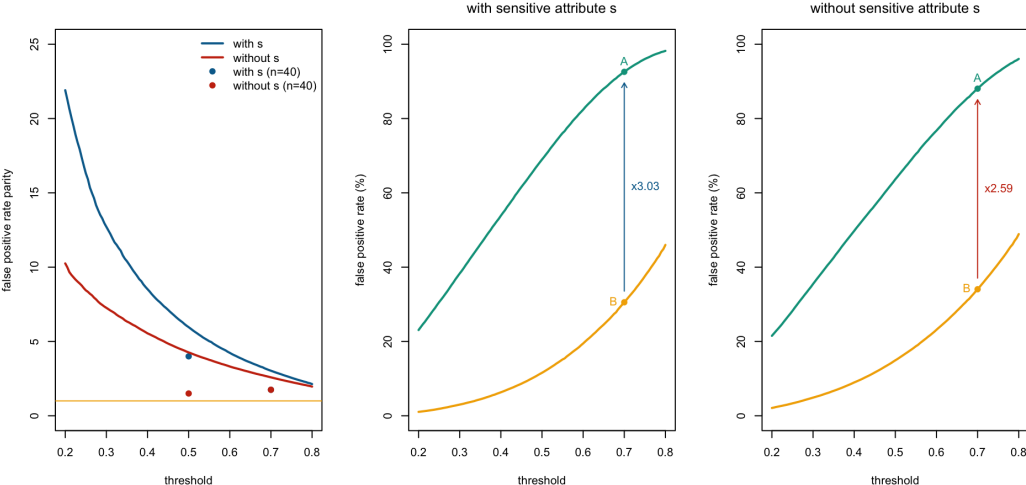
$$\begin{cases} \mathbb{P}[\hat{Y} = 1 | S = \mathbf{A}, Y = 0] = \mathbb{P}[\hat{Y} = 1 | S = \mathbf{B}, Y = 0] = \mathbb{P}[\hat{Y} = 1 | Y = 0] \\ \mathbb{P}[m_t(\mathbf{Z}) = 1 | S = \mathbf{A}, Y = 0] = \mathbb{P}[m_t(\mathbf{Z}) = 1 | S = \mathbf{B}, Y = 0] = \mathbb{P}[m_t(\mathbf{Z}) = 1 | Y = 0]. \end{cases}$$

Separation and Equalized Odds



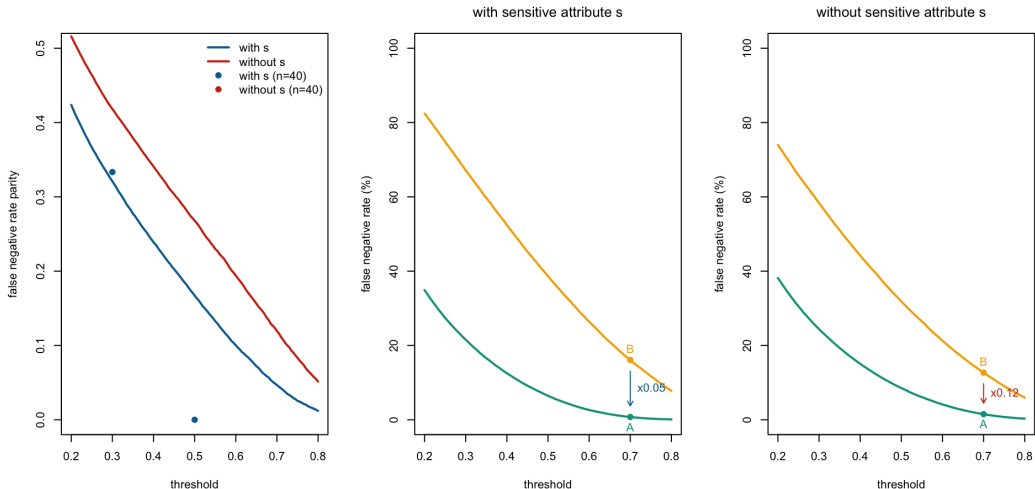
➤ ROC curves (TPR against FPR) for the logistic regression on **toydata2**.

Separation and Equalized Odds



➤ Evolution of the false positive rates, **fpr_parity** from **fairness**.

Separation and Equalized Odds



➤ Evolution of the false negative rates, `fnr_parity` from `fairness`.

Separation and Equalized Odds

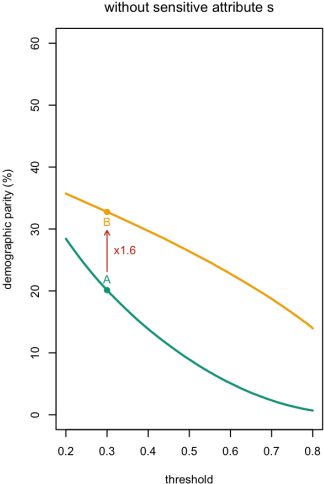
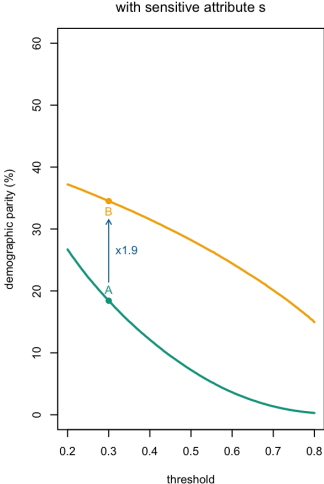
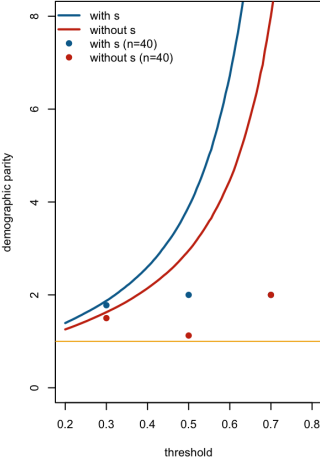
Definition 5.15: Equalized Odds, [Hardt et al. \(2016\)](#)

A decision function \hat{y} – or a classifier $m_t(\cdot)$ taking values in $\{0, 1\}$ – satisfies equal odds constraint, with respect to some sensitive attribute S , if

$$\begin{cases} \mathbb{P}[\hat{Y} = 1 | S = \mathbf{A}, Y = y] = \mathbb{P}[\hat{Y} = 1 | S = \mathbf{B}, Y = y] = \mathbb{P}[\hat{Y} = 1 | Y = y], \forall y \in \{0, 1\} \\ \mathbb{P}[m_t(\mathbf{Z}) = 1 | S = \mathbf{A}, Y = y] = \mathbb{P}[m_t(\mathbf{Z}) = 1 | S = \mathbf{B}, Y = y], \forall y \in \{0, 1\}, \end{cases}$$

which corresponds to parity of true positive and false positive, in the two groups.

Separation and Equalized Odds



➤ Evolution of the equalized odds metrics

Separation and Equalized Odds

- One can also consider any kind of standard metrics on confusion matrices, such as ϕ (introduced in [Yule \(1912\)](#)), usually named "Matthews correlation coefficient"

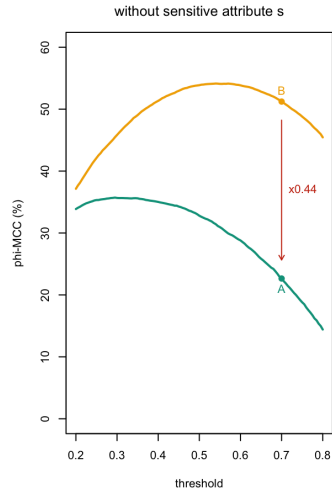
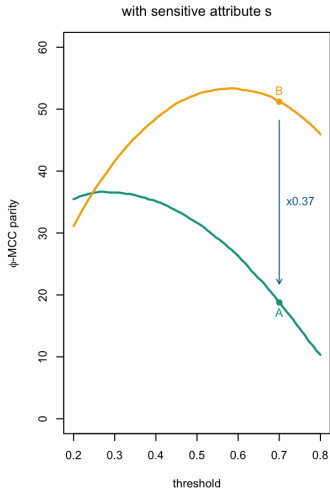
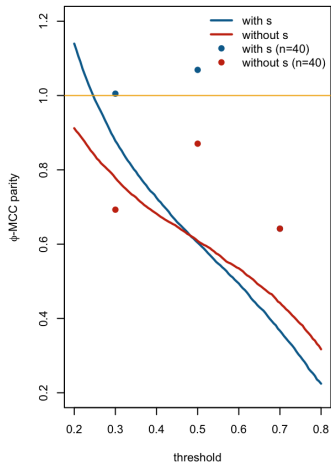
Definition 5.16: ϕ -fairness, [Chicco and Jurman \(2020\)](#)

We will have ϕ -fairness if $\phi_A = \phi_B$, where ϕ_s denotes Matthews correlation coefficient for the s group,

$$\phi_s = \frac{TP_s \cdot TN_s - FP_s \cdot FN_s}{\sqrt{(TP_s + FP_s)(TP_s + FN_s) \cdot (TN_s + FP_s)(TN_s + FN_s)}}, \quad s \in \{A, B\}.$$

- but one could consider the F_1 -score (as defined in [Van Rijsbergen \(1979\)](#)), Fowlkes–Mallows or Jaccard indices (in [Fowlkes and Mallows \(1983\)](#) or [Jaccard \(1901\)](#)).
- .. or AUC as we will considered later on.

Separation and Equalized Odds



➤ Evolution of the ϕ -fairness metric

Definition 5.17: Class Balance, Kleinberg et al. (2016)

We will have class balance in the weak sense if

$$\mathbb{E}[m(\mathbf{X})|Y = y, S = \mathbf{A}] = \mathbb{E}[m(\mathbf{X})|Y = y, S = \mathbf{B}], \quad \forall y \in \{0, 1\},$$

or in the strong sense if

$$\mathbb{P}[m(\mathbf{X}) \in \mathcal{A}|Y = y, S = \mathbf{A}] = \mathbb{P}[m(\mathbf{X}) \in \mathcal{A}|Y = y, S = \mathbf{B}], \quad \forall \mathcal{A} \subset [0, 1], \quad \forall y \in \{0, 1\}.$$

Separation and Equalized Odds

Definition 5.18: Similar Mistreatment, [Zafar et al. \(2019\)](#)

We will have similar mistreatment, or “*lack of disparate mistreatment*,” if

$$\begin{cases} \mathbb{P}[\hat{Y} = Y | S = \mathbf{A}] = \mathbb{P}[\hat{Y} = Y | S = \mathbf{B}] = \mathbb{P}[\hat{Y} = Y] \\ \mathbb{P}[m_t(\mathbf{X}) = Y | S = \mathbf{A}] = \mathbb{P}[m_t(\mathbf{X}) = Y | S = \mathbf{B}] = \mathbb{P}[m_t(\mathbf{X}) = Y]. \end{cases}$$

Definition 5.19: Equality of ROC curves, [Vogel et al. \(2021\)](#)

Let $\text{FRP}_s(t) = \mathbb{P}[m(\mathbf{X}) > t | Y = 0, S = s]$ and $\text{TPR}_s(t) = \mathbb{P}[m(\mathbf{X}) > t | Y = 1, S = s]$, where $s \in \{\mathbf{A}, \mathbf{B}\}$. Set $\Delta_{\text{TPR}}(t) = \text{TPR}_{\mathbf{B}} \circ \text{TPR}_{\mathbf{A}}^{-1}(t) - t$ et $\Delta_{\text{FRP}}(t) = \text{FPR}_{\mathbf{B}} \circ \text{FPR}_{\mathbf{A}}^{-1}(t) - t$. We will have fairness with respect to ROC curves if $\|\Delta_{\text{TPR}}\|_{\infty} = \|\Delta_{\text{FRP}}\|_{\infty} = 0$.

Separation and Equalized Odds

Definition 5.20: AUC Fairness, [Borkan et al. \(2019\)](#)

We will have AUC fairness if $AUC_A = AUC_B$, where AUC_s is the AUC associated with model m within the s group.

	unaware (without s)				aware (with s)			
	GLM	GAM	CART	RF	GLM	GAM	CART	RF
ratio of AUC	0.837	0.839	0.913	0.768	0.857	0.860	0.913	0.763

Sufficiency and Calibration

- › Inspired by [Cleary \(1968\)](#), define

Definition 5.21: Sufficiency, [Barocas et al. \(2017\)](#)

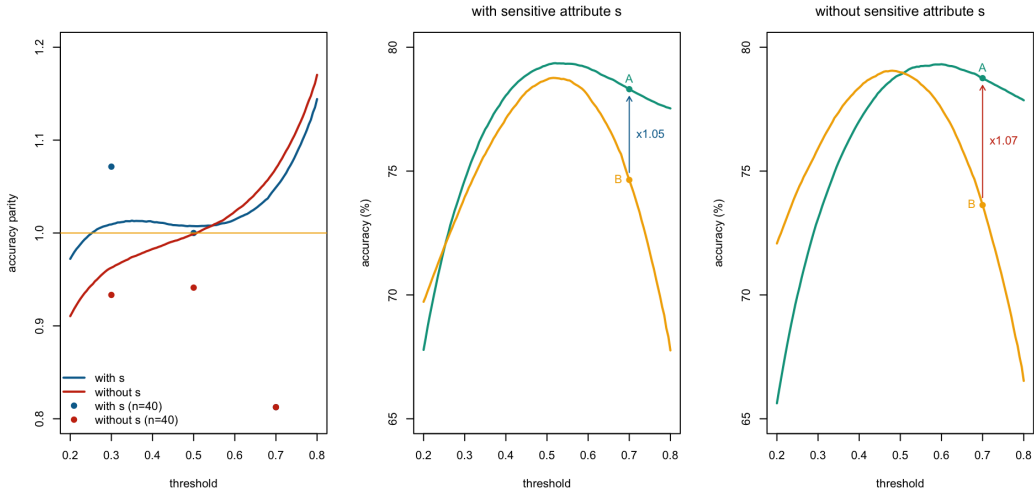
A model $m : \mathcal{Z} \rightarrow \mathcal{Y}$ satisfies the sufficiency property if $Y \perp\!\!\!\perp S \mid m(\mathbf{Z})$, with respect to the distribution \mathbb{P} of the triplet (\mathbf{X}, S, Y) .

Definition 5.22: Calibration Parity, Accuracy Parity, [Kleinberg et al. \(2016\)](#), [Zafar et al. \(2019\)](#)

Calibration parity is met if

$$\mathbb{P}[Y = 1 | m(\mathbf{X}) = t, S = \mathbf{A}] = \mathbb{P}[Y = 1 | m(\mathbf{X}) = t, S = \mathbf{B}], \quad \forall t \in [0, 1].$$

Sufficiency and Calibration



➤ Evolution of accuracy, in groups A and B.

Sufficiency and Calibration

Definition 5.23: Good Calibration, Kleinberg et al. (2017), Verma and Rubin (2018)

Fairness of good calibration is met if

$$\mathbb{P}[Y = 1 | m(\mathbf{X}) = t, S = \mathbf{A}] = \mathbb{P}[Y = 1 | m(\mathbf{X}) = t, S = \mathbf{B}] = t, \forall t \in [0, 1].$$

Definition 5.24: Non-Reconstruction of Protected Attribute, Kim (2017)

If we cannot tell from the result $(\mathbf{x}, m(\mathbf{x}), y$ and $\hat{y})$ whether the subject was a member of a protected group or not, we will talk about fairness by non-reconstruction of the protected attribute

$$\mathbb{P}[S = \mathbf{A} | \mathbf{X}, m(\mathbf{X}), \hat{Y}, Y] = \mathbb{P}[S = \mathbf{B} | \mathbf{X}, m(\mathbf{X}), \hat{Y}, Y].$$

Relaxation and Approximate Fairness

Definition 5.25: Disparate Impact, [Feldman et al. \(2015\)](#)

A decision function \hat{Y} has a disparate impact, for a given threshold τ , if,

$$\min \left\{ \frac{\mathbb{P}[\hat{Y} = 1 | S = \text{A}]}{\mathbb{P}[\hat{Y} = 1 | S = \text{B}]}, \frac{\mathbb{P}[\hat{Y} = 1 | S = \text{B}]}{\mathbb{P}[\hat{Y} = 1 | S = \text{A}]} \right\} < \tau \text{ (usually 80\%).}$$

➤ The **80% rule** was suggested by the "Technical Advisory Committee on Testing", from the State of California Fair Employment Practice Commission (FEPC) in 1971, or the 1978 "Uniform Guidelines on Employee Selection Procedures", a document used by the U.S. Equal Employment Opportunity Commission (EEOC), see [Biddle \(2017\)](#).

Relaxation and Approximate Fairness

- › We have defined (Definition 5.10) unfairness as

$$\mathcal{U}_k(m) = \max_{s \in \{A, B\}} \{W_k(\mathbb{P}_m, \mathbb{P}_{m|s})\},$$

so that m is (strongly) fair if and only if $\mathcal{U}_k(m) = 0$.

- › [Chzhen and Schreuder \(2022\)](#) introduced the notion of Relative Improvement

Definition 5.26: ε -Approximate Fairness

Model m is ε -approximately fair if $\mathcal{U}_k(m) \leq \varepsilon \cdot \mathcal{U}_k(m^*)$, where m^* is Bayes regressor, for some $\varepsilon \geq 0$.

Three different concepts ?

$$\left\{ \begin{array}{l} \text{Independence (Definition 5.5)} : m(\mathbf{Z}) \perp\!\!\!\perp S \\ \text{Separation (Definition 5.11)} : m(\mathbf{Z}) \perp\!\!\!\perp S \mid Y \\ \text{Sufficiency (Definition 5.21)} : Y \perp\!\!\!\perp S \mid m(\mathbf{Z}) \end{array} \right.$$

- ▶ Independence assumes no differences among groups, regardless of accuracy
- ▶ Separation minimizes differences among groups by not trying to maximize accuracy
- ▶ Sufficiency maximizes accuracy by not trying to minimize differences among groups

See [Kleinberg et al. \(2016\)](#) or [Chouldechova \(2017\)](#).

Impossibility theorems

- Unless very specific properties are assumed on \mathbb{P} , there is no prediction function $m(\cdot)$ that can satisfy at the same time two fairness criteria.

$$\left\{ \begin{array}{l} \text{Independence (Definition 5.5)} : m(\mathbf{Z}) \perp\!\!\!\perp S \\ \text{Separation (Definition 5.11)} : m(\mathbf{Z}) \perp\!\!\!\perp S \mid Y \\ \text{Sufficiency (Definition 5.21)} : Y \perp\!\!\!\perp S \mid m(\mathbf{Z}) \end{array} \right.$$

Proposition 5.4

Suppose that a model m satisfies the independence condition (5.5) and the sufficiency property (5.21), with respect to a sensitive attribute s , then necessarily, $Y \perp\!\!\!\perp S$.

- Therefore, unless the sensitive attribute s has no impact on the outcome y , there is no model m which satisfies independence and sufficiency simultaneously.

Impossibility theorems

- From the sufficiency property , $S \perp\!\!\!\perp Y \mid m(\mathbf{Z})$, then, for $s \in \mathcal{S}$ and $\mathcal{A} \subset \mathcal{Y}$,

$$\mathbb{P}[S = s, Y \in \mathcal{A}] = \mathbb{E}[\mathbb{P}[S = s, Y \in \mathcal{A} \mid m(\mathbf{Z})]],$$

can be written

$$\mathbb{P}[S = s, Y \in \mathcal{A}] = \mathbb{E}[\mathbb{P}[S = s \mid m(\mathbf{Z})] \cdot \mathbb{P}[Y \in \mathcal{A} \mid m(\mathbf{Z})]].$$

And from the independence property (5.21), $m(\mathbf{Z}) \perp\!\!\!\perp S$, we can write the first component $\mathbb{P}[S = s \mid m(\mathbf{Z})] = \mathbb{P}[S = s]$, almost surely, and therefore

$$\mathbb{P}[S = s, Y \in \mathcal{A}] = \mathbb{E}[\mathbb{P}[S = s] \cdot \mathbb{P}[Y \in \mathcal{A} \mid m(\mathbf{Z})]] = \mathbb{P}[S = s] \cdot \mathbb{P}[Y \in \mathcal{A}],$$

for all $s \in \mathcal{S}$ and $\mathcal{A} \subset \mathcal{Y}$, corresponding to the independence between S and Y .

Impossibility theorems

Proposition 5.5

Consider a classifier m_t taking values in $\mathcal{Y} = \{0, 1\}$. Suppose that m_t satisfies the independence condition (5.5) and the separation property (5.11), with respect to a sensitive attribute s , then necessarily either $m_t(\mathbf{Z}) \perp\!\!\!\perp Y$ or $Y \perp\!\!\!\perp S$ (possibly both).

➤ Because m_t satisfies the independence condition (5.5), $m_t(\mathbf{Z}) \perp\!\!\!\perp S$, and the separation property (5.11), $m_t(\mathbf{Z}) \perp\!\!\!\perp S \mid Y$, then, for $\hat{y} \in \mathcal{Y}$ and for $s \in \mathcal{S}$,

$$\mathbb{P}[m_t(\mathbf{Z}) = \hat{y}] = \mathbb{P}[m_t(\mathbf{Z}) = \hat{y} \mid S = s] = \mathbb{E}[\mathbb{P}[m_t(\mathbf{Z}) = \hat{y} \mid Y, S = s]],$$

that we can write

$$\mathbb{P}[m_t(\mathbf{Z}) = \hat{y}] = \sum_y \mathbb{P}[m_t(\mathbf{Z}) = \hat{y} \mid Y = y, S = s] \cdot \mathbb{P}[Y = y \mid S = s],$$

Impossibility theorems

or

$$\mathbb{P}[m_t(\mathbf{Z}) = \hat{y}] = \sum_y \mathbb{P}[m_t(\mathbf{Z}) = \hat{y} | Y = y] \cdot \mathbb{P}[Y = y | S = s],$$

almost surely. Furthermore, we can also write

$$\mathbb{P}[m_t(\mathbf{Z}) = \hat{y}] = \sum_y \mathbb{P}[m_t(\mathbf{Z}) = \hat{y} | Y = y] \cdot \mathbb{P}[Y = y],$$

so that, if we combine the two expressions, we get

$$\sum_y \mathbb{P}[m_t(\mathbf{Z}) = \hat{y} | Y = y] \cdot \left(\mathbb{P}[Y = y | S = s] - \mathbb{P}[Y = y] \right) = 0,$$

almost surely. And since we assumed that y was a binary variable, $\mathbb{P}[Y = 0] = 1 - \mathbb{P}[Y = 1]$, as well as $\mathbb{P}[Y = 0 | S = s] = 1 - \mathbb{P}[Y = 1 | S = s]$, and therefore

$$\mathbb{P}[m_t(\mathbf{Z}) = \hat{y} | Y = 1] \cdot \left(\mathbb{P}[Y = 1 | S = s] - \mathbb{P}[Y = 1] \right)$$

Impossibility theorems

or

$$-\mathbb{P}[m_t(\mathbf{Z}) = \hat{y} | Y = 0] \cdot (\mathbb{P}[Y = 0 | S = s] - \mathbb{P}[Y = 0])$$

can be written

$$\mathbb{P}[m_t(\mathbf{Z}) = \hat{y} | Y = 0] \cdot (\mathbb{P}[Y = 1 | S = s] - \mathbb{P}[Y = 1]).$$

Thus, either $\mathbb{P}[Y = 1 | S = s] - \mathbb{P}[Y = 1]$ almost surely, or

$\mathbb{P}[m_t(\mathbf{Z}) = \hat{y} | Y = 0] = \mathbb{P}[m_t(\mathbf{Z}) = \hat{y} | Y = 1]$ (or both).

› Of course, the previous proposition holds only when y is a binary variable.

Proposition 5.6

Consider a classifier m_t taking values in $\mathcal{Y} = \{0, 1\}$. Suppose that m_t satisfies the sufficiency condition (5.21) and the separation property (5.11), with respect to a sensitive attribute s , then necessarily either $\mathbb{P}[m_t(\mathbf{Z}) = 1 | Y = 1] = 0$ or $Y \perp\!\!\!\perp S$ or $m_t(\mathbf{Z}) \perp\!\!\!\perp Y$.

➤ Suppose that m_t satisfies the sufficiency condition (5.21) and the separation property (5.11), respectively $Y \perp\!\!\!\perp S | m_t(\mathbf{Z})$ and $m_t(\mathbf{Z}) \perp\!\!\!\perp S | Y$. For all $s \in \mathcal{S}$, we can write, using Bayes formula

$$\mathbb{P}[Y = 1 | S = s, m_t(\mathbf{Z}) = 1] = \frac{\mathbb{P}[m_t(\mathbf{Z}) = 1 | Y = 1, S = s] \cdot \mathbb{P}[Y = 1 | S = s]}{\mathbb{P}[m_t(\mathbf{Z}) = 1 | S = s]},$$

Impossibility theorems

i.e.,

$$\mathbb{P}[Y = 1|S = s, m_t(\mathbf{Z}) = 1] = \frac{\mathbb{P}[m_t(\mathbf{Z}) = 1|Y = 1] \cdot \mathbb{P}[Y = 1|S = s]}{\sum_{y \in \{0,1\}} \mathbb{P}[m_t(\mathbf{Z}) = 1|Y = y] \cdot \mathbb{P}[Y = 1|S = s]},$$

that should not depend on s (from the sufficiency property). So a similar property holds if $S = s'$. Observe further that $\mathbb{P}[m_t(\mathbf{Z}) = 1|Y = 1]$ is the *true positive rate* (TPR) while $\mathbb{P}[m_t(\mathbf{Z}) = 1|Y = 0]$ is the *false positive rate* (FPR). Let $p_s = \mathbb{P}[Y = 1|S = s]$, so that

$$\mathbb{P}[Y = 1|S = s, m_t(\mathbf{Z}) = 1] = \frac{\text{TPR}}{p_s \cdot \text{TPR} + (1 - p_s) \cdot \text{FPR}}.$$

Impossibility theorems

➤ Suppose that Y and S are not independent (otherwise $Y \perp\!\!\!\perp S$ as stated in the proposition), i.e., there are s and s' such that $p_s = \mathbb{P}[Y = 1|S = s] \neq \mathbb{P}[Y = 1|S = s'] = p_{s'}$. Hence, $p_s \neq p_{s'}$, but at the same time

$$\frac{\text{TPR}}{p_s \cdot \text{TPR} + (1 - p_s) \cdot \text{FPR}} = \frac{\text{TPR}}{p_{s'} \cdot \text{TPR} + (1 - p_{s'}) \cdot \text{FPR}}.$$

Supposes that $\text{TPR} \neq 0$ (otherwise $\text{TPR} = \mathbb{P}[m_t(\mathbf{Z}) = 1|Y = 1] = 0$ as stated in the proposition), then

$$(p_s - p_{s'}) \cdot \text{TPR} = (p_s - p_{s'}) \cdot \text{FPR} \neq 0,$$

and therefore $m_t(\mathbf{Z}) \perp\!\!\!\perp Y$.

Group fairness, wrap-up

independence, $\hat{Y} \perp\!\!\!\perp S$, (Definition 5.5)

<i>statistical parity</i>	Dwork et al. (2012)	$\mathbb{P}[\hat{Y} = 1 S = s] = \text{cst}, \forall s$
<i>conditional statistical parity</i>	Corbett-Davies et al. (2017)	$\mathbb{P}[\hat{Y} = 1 S = s, X = x] = \text{cst}_x, \forall s, y$

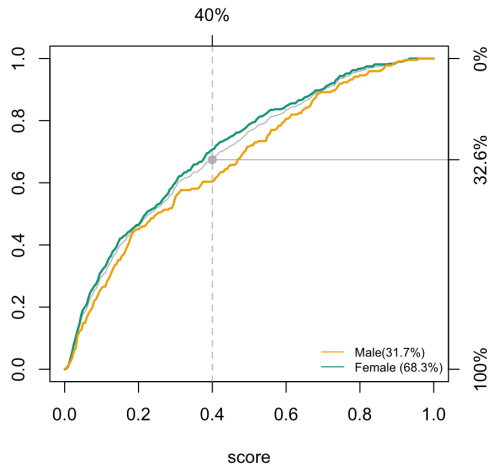
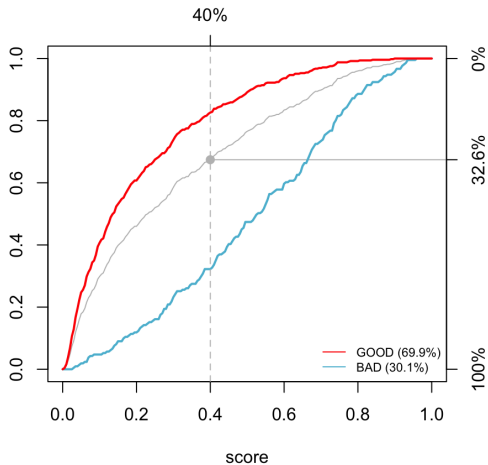
separation, $\hat{Y} \perp\!\!\!\perp S | Y$, (Definition 5.11)

<i>equalized odds</i>	Hardt et al. (2016)	$\mathbb{P}[\hat{Y} = 1 S = s, Y = y] = \text{cst}_y, \forall s, y$
<i>equalized opportunity</i>	Hardt et al. (2016)	$\mathbb{P}[\hat{Y} = 1 S = s, Y = 1] = \text{cst}, \forall s$
<i>predictive equality</i>	Corbett-Davies et al. (2017)	$\mathbb{P}[\hat{Y} = 1 S = s, Y = 0] = \text{cst}, \forall s$
<i>balance</i>	Kleinberg et al. (2016)	$\mathbb{E}[M S = s, Y = 1] = \text{cst}_y, \forall s, y$

sufficiency, $Y \perp\!\!\!\perp S | \hat{Y}$, (Definition 5.21)

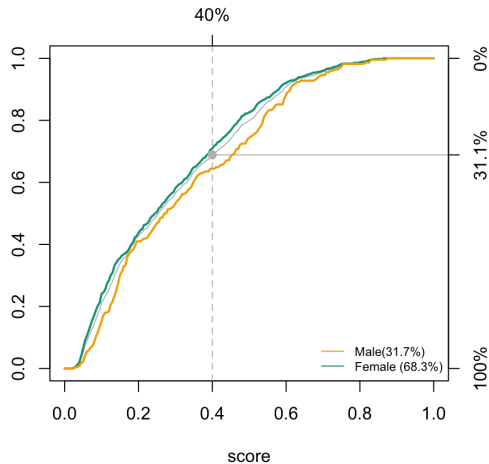
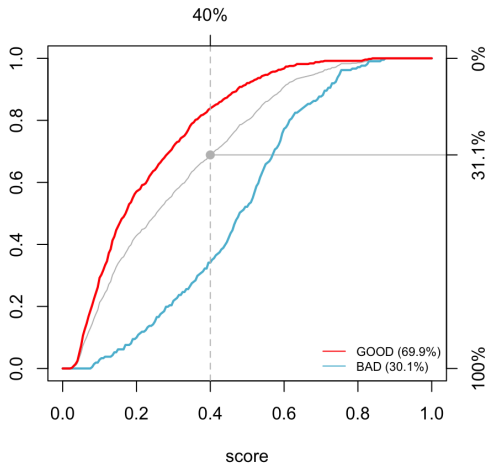
<i>disparate mistreatment</i>	Zafar et al. (2019)	$\mathbb{P}[Y = y S = s, \hat{Y} = y] = \text{cst}_y, \forall s, y$
<i>predictive parity</i>	Chouldechova (2017)	$\mathbb{P}[Y = 1 S = s, \hat{Y} = 1] = \text{cst}, \forall s$
<i>calibration</i>	Chouldechova (2017)	$\mathbb{P}[Y = 1 M = m, S = s] = \text{cst}_m, \forall m, s$
<i>well-calibration</i>	Chouldechova (2017)	$\mathbb{P}[Y = 1 M = m, S = s] = s, \forall m, s$

Numerical examples



Conditional distributions of scores on [GermanCredit](#), logistic regression.

Numerical examples



Conditional distributions of scores on **GermanCredit**, boosting model.

Numerical examples

	with sensitive				without sensitive			
	GLM	tree	boosting	bagging	GLM	tree	boosting	bagging
$\mathbb{P}[m(\mathbf{X}) > t]$	51.7%	28.0%	54.7%	61.7%	50.7%	28.0%	56.0%	60.7%
Predictive Rate Parity	0.992	1.190	0.992	1.050	0.957	1.190	1.041	1.037
Demographic Parity	0.998	1.091	1.159	1.027	1.213	1.091	1.112	1.208
FNR Parity	1.398	0.740	1.078	1.124	1.075	0.740	1.064	0.970
Proportional Parity	0.922	1.008	1.071	0.949	1.121	1.008	1.027	1.116
Equalized odds	0.816	1.069	0.947	0.888	0.956	1.069	0.953	1.031
Accuracy Parity	0.843	1.181	0.912	0.904	0.896	1.181	0.943	0.966
FPR Parity	1.247	0.683	1.470	0.855	2.004	0.683	0.962	1.069
NPV Parity	0.676	1.141	0.763	0.772	0.735	1.141	0.799	0.823
Specificity Parity	0.941	1.439	0.930	1.028	0.851	1.439	1.007	0.990
ROC AUC Parity	0.928	1.162	0.997	1.108	0.926	1.162	1.004	1.090
MCC Parity	0.604	2.013	0.744	0.851	0.639	2.013	0.884	0.930

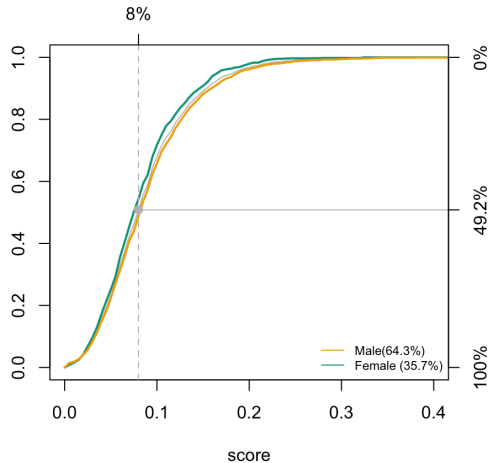
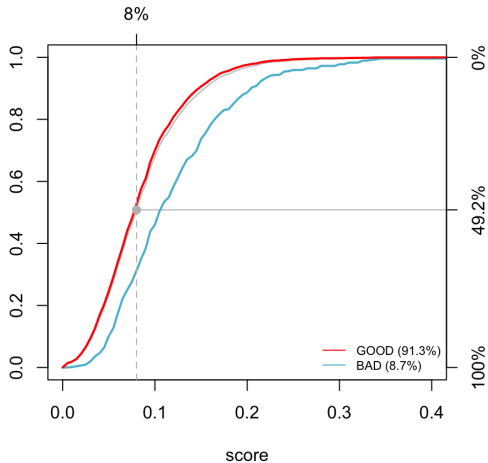
Fairness metrics on **GermanCredit**, with threshold at 20%.

Numerical examples

	with sensitive				without sensitive			
	GLM	tree	boosting	bagging	GLM	tree	boosting	bagging
$\mathbb{P}[m(\mathbf{X}) > t]$	30.3%	26.0%	27.7%	25.7%	30.7%	26.0%	28.0%	27.0%
Predictive Rate Parity	1.030	1.179	1.110	1.182	1.034	1.179	1.111	1.200
Demographic Parity	1.090	1.062	1.074	1.069	1.108	1.062	1.044	1.019
FNR Parity	1.533	0.851	1.110	0.781	1.342	0.851	1.322	0.962
Proportional Parity	1.007	0.981	0.992	0.987	1.024	0.981	0.964	0.942
Equalized odds	0.925	1.032	0.982	1.041	0.944	1.032	0.955	1.008
Accuracy Parity	0.949	1.154	1.054	1.164	0.963	1.154	1.038	1.159
FPR Parity	1.118	0.703	0.820	0.653	1.118	0.703	0.784	0.641
NPV Parity	0.738	1.080	0.890	1.108	0.766	1.080	0.848	1.082
Specificity Parity	0.935	1.470	1.169	1.480	0.935	1.470	1.203	1.652
ROC AUC Parity	0.928	1.162	0.997	1.108	0.926	1.162	1.004	1.090
MCC Parity	0.745	1.817	1.105	1.754	0.779	1.817	1.056	2.055

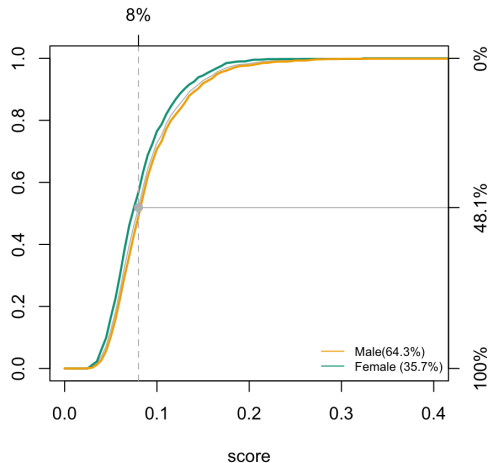
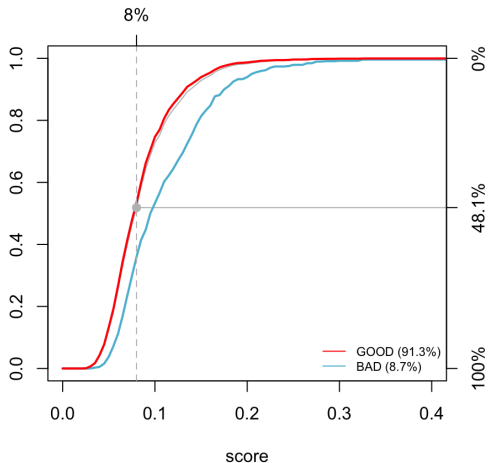
Fairness metrics on **GermanCredit**, with threshold at 40%.

Numerical examples



Conditional distributions of scores on **FrenchMotor**, from the logistic regression.

Numerical examples



Conditional distributions of scores on **FrenchMotor**, from a boosting classification.

– Part 6 –

Individual Fairness

Definition 6.1: Similarity Fairness, Luong et al. (2011), Dwork et al. (2012)

Consider two metrics, one on $\mathcal{Y} \times \mathcal{Y}$ (or for a classifier $[0, 1]$ and not $\{0, 1\}$) noted D_y , and one on \mathcal{X} noted D_x , such that we will have similarity fairness on a database of size n if we have the following property (called Lipschitz property)

$$D_y(m(\mathbf{x}_i, s_i), m(\mathbf{x}_j, s_j)) \leq L \cdot D_x(\mathbf{x}_i, \mathbf{x}_j), \quad \forall i, j = 1, \dots, n,$$

for some $L < \infty$.

Definition 6.2: Local individual fairness, [Petersen et al. \(2021\)](#)

Consider two metrics, one on \mathcal{Y} ($[0, 1]$ for a classifier and not $\{0, 1\}$) noted D_y , and one on \mathcal{X} noted D_x , model m is locally individually fair if

$$\mathbb{E}_{(\mathbf{X}, S)} \left[\limsup_{\mathbf{x}': D_x(\mathbf{X}, \mathbf{x}') \rightarrow 0} \frac{D_y(m(\mathbf{X}, S), m(\mathbf{x}', S))}{D_x(\mathbf{X}, \mathbf{x}')} \right] \leq L < \infty.$$

Individual Fairness

Definition 6.3: Proxy Based Fairness, [Kilbertus et al. \(2017\)](#)

A decision making process \hat{y} exhibits no proxy discrimination with respect to sensitive attribute s if

$$\mathbb{E}[\hat{Y}|\text{do}(S = \text{A})] = \mathbb{E}[\hat{Y}|\text{do}(S = \text{B})].$$

Definition 6.4: Fairness on Average Treatment Effect, [Kusner et al. \(2017\)](#)

We achieve fairness on average treatment effect (counterfactual fairness on average)

$$\text{ATE} = \mathbb{E}[Y_{S \leftarrow \text{A}}^* - Y_{S \leftarrow \text{B}}^*] = 0.$$

Individual Fairness

- A decision satisfies counterfactual fairness if "*had the protected attributes (e.g., race) of the individual been different, other things being equal, the decision would have remained the same.*"

Definition 6.5: Counterfactual Fairness, [Kusner et al. \(2017\)](#)

We achieve counterfactual fairness for an individual with characteristics \mathbf{x} if

$$\text{CATE}(\mathbf{x}) = \mathbb{E}[Y_{S \leftarrow \mathbf{A}}^* - Y_{S \leftarrow \mathbf{B}}^* | \mathbf{X} = \mathbf{x}] = 0.$$

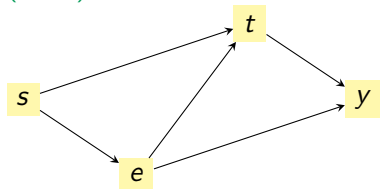
Individual Fairness

Definition 6.6: Path-Specific Counterfactual Effect, Wu et al. (2019)

Given a causal diagram, a factual condition (denoted \mathcal{F}), and a path π from s to y , the π -effect of a change of s from B to A on y is

$$\text{PCE}_{\pi}(B \rightarrow A|\mathcal{F}) = \mathbb{E}[Y|\text{do}_{\pi}(S = A), \mathcal{F}] - \mathbb{E}[Y|S = B, \mathcal{F}].$$

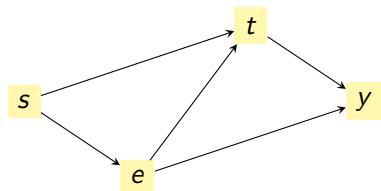
- A counterfactual value replaces the cause of interest, then propagates "downstream" the causal graph via the structural equations, *ceteris paribus*.
- **fairadapt**, based on **twin network**, and "recursive substitution", Ma and Koenker (2006).



$\left\{ \begin{array}{l} s : \text{gender} \\ e : \text{previous educational achievement} \\ t : \text{admission score} \\ y : \text{admission indicator} \end{array} \right.$

Individual Fairness

- Following [Ma and Koenker \(2006\)](#), consider some structural model



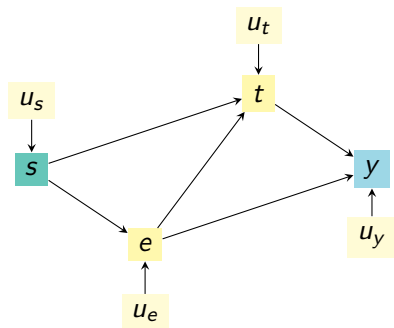
$$\begin{cases} s = U_s \\ e = h_e(s, U_e) \\ t = h_t(s, e, U_t) \\ y = h_y(e, t, U_y) \end{cases}$$

with U_s, U_e, U_t, U_y independent.

- Consider observation (s_i, e_i, t_i, y_i) .
- Can we get a counterfactual of that observation, where $s_i = A$?
- Quantile regression ([Koenker and Bassett Jr \(1978\)](#), [Koenker et al. \(2017\)](#)), possibly quantile forest ([Meinshausen and Ridgeway \(2006\)](#)).

Individual Fairness

Individual Fairness

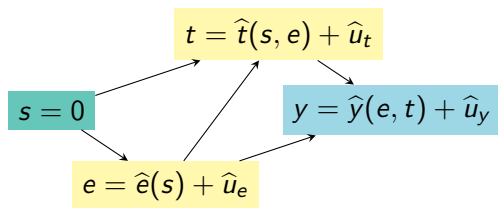


the DAG can be expressed as an additive SCM

$$\begin{cases} s = U_s \\ e = h_e(s) + U_e \\ t = h_t(s, e) + U_t \\ y = h_y(e, t) + U_y \end{cases}$$

with U_s, U_e, U_t, U_y independent.

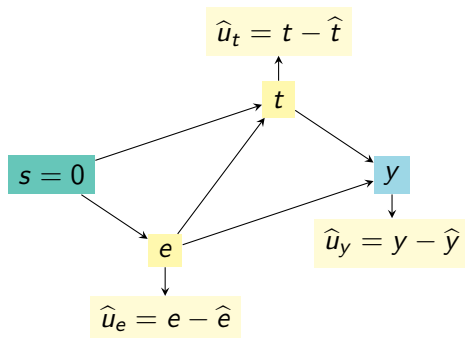
Individual Fairness



Consider some individual in group $s = 0$

$$\begin{cases} s = 0 \\ e = \hat{h}_e(s) + \hat{u}_e \\ t = \hat{h}_t(s, e) + \hat{u}_t \\ y = \hat{h}_y(e, t) + \hat{u}_y \end{cases}$$

Individual Fairness



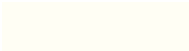
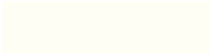
For that individuals, residuals are

$$\begin{cases} s = 0 \\ \hat{u}_e = e - \hat{h}_e(s) \\ \hat{u}_t = t - \hat{h}_t(s, e) \\ \hat{u}_y = y - \hat{h}_y(e, t) \end{cases}$$

To generate a counterfactual, suppose that residuals remain unchanged.

Individual Fairness

$$s = 1$$



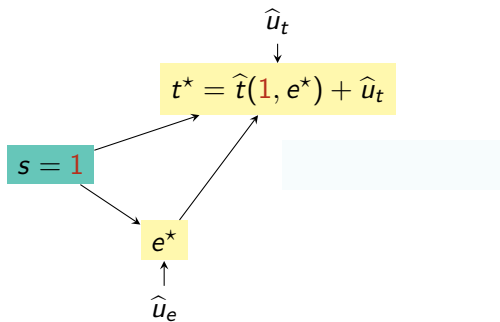
Individual Fairness

$s = 1$

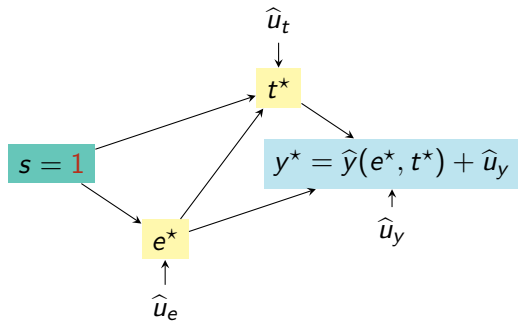
$$e^* = \hat{e}(s = 1) + \hat{u}_e$$

\hat{u}_e

Individual Fairness



Individual Fairness



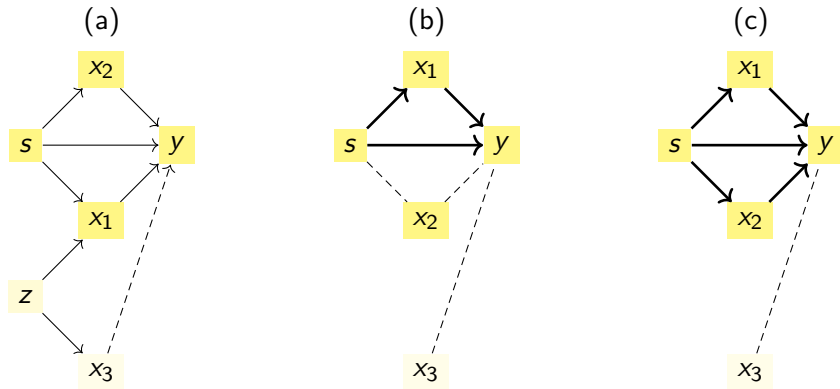
Individual Fairness

Individual Fairness

Individual Fairness

Similarity Fairness (<i>Lipschitz</i>)	Dwork et al. (2012)	$D_Y(\hat{y}_i, \hat{y}_j) \leq D_X(\mathbf{x}_i, \mathbf{x}_j), \forall i, j$
Proxy Based Fairness,	Kilbertus et al. (2017)	$\mathbb{E}[Y \text{do}(S = \mathbf{A})] = \mathbb{E}[Y \text{do}(S = \mathbf{B})]$
Fairness on Average Treatment Effect	Kusner et al. (2017)	$\mathbb{E}[Y_{S \leftarrow \mathbf{A}}^*] = \mathbb{E}[Y_{S \leftarrow \mathbf{B}}^*]$
Counterfactual Fairness,	Kusner et al. (2017)	$\mathbb{E}[Y_{S \leftarrow \mathbf{A}}^* \mathbf{X} = \mathbf{x}] = \mathbb{E}[Y_{S \leftarrow \mathbf{B}}^* \mathbf{X} = \mathbf{x}]$
Path-Specific Effect	Avin et al. (2005)	$\mathbb{E}[Y \text{do}_\pi(S = \mathbf{A})] = \mathbb{E}[Y \text{do}_\pi(S = \mathbf{B})]$
Path-Specific Counterfactual Effect	Wu et al. (2019)	$\mathbb{E}[Y \text{do}_\pi(S = \mathbf{A}), \mathcal{F}] = \mathbb{E}[Y \text{do}_\pi(S = \mathbf{B}), \mathcal{F}]$
Mutatis Mutandis Counterfactual	Kusner et al. (2017)	$\mathbb{E}[Y_{S \leftarrow \mathbf{A}}^* \mathbf{X} = \mathcal{T}(\mathbf{x})] = \mathbb{E}[Y_{S \leftarrow \mathbf{B}}^* \mathbf{X} = \mathcal{T}(\mathbf{x})]$

Individual Fairness

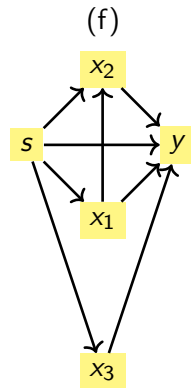
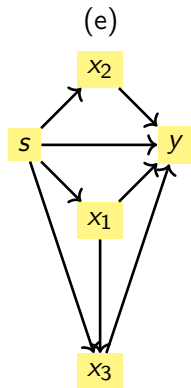
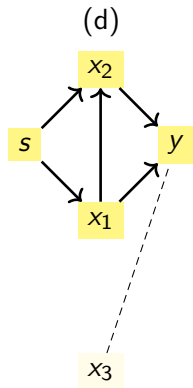


(a) Causal graph used to generate variables in `toydata2`.

(b) Causal graph, where s might cause y , either directly, or indirectly, through x_1 .

(c) Causal graph, where s might cause y , either directly or indirectly, via with two possible paths and two mediator variables, x_1 and x_2 .

Individual Fairness



(d) Causal graph with no direct impact of s on y , but two mediators, and possibly, x_1 might cause x_2 .

(e) similar to (c) with an additional indirect connection from x_1 to y , via mediator x_3 .

(f) similar to (d) with an additional indirect connection from x_1 to y , via mediator x_3 .

Individual Fairness

Original data

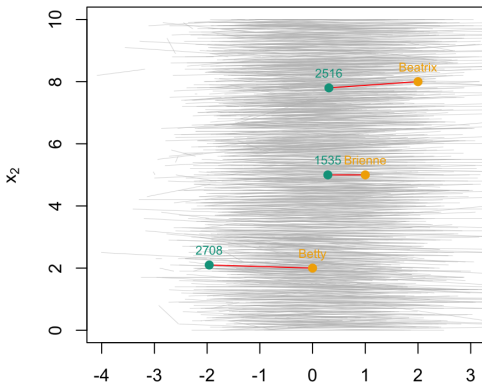
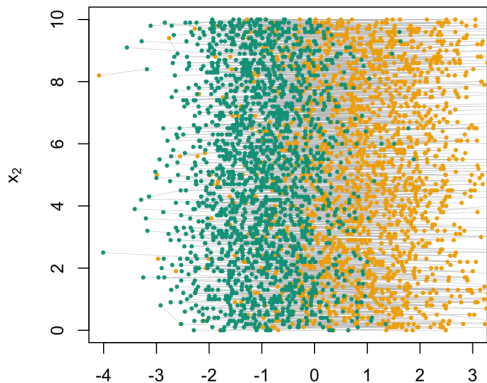
	s	x_1	x_2	x_3	$\hat{m}_{\text{glm}}(\mathbf{x})$	$\hat{m}_{\text{glm}}(\mathbf{x}, s)$	$\hat{m}_{\text{gam}}(\mathbf{x})$	$\hat{m}_{\text{gam}}(\mathbf{x}, s)$	$\hat{m}_{\text{rf}}(\mathbf{x})$	$\hat{m}_{\text{rf}}(\mathbf{x}, s)$
Betty	B	0	2	0	18.22%	24.06%	13.23%	17.63%	17.4%	29.6%
Brienne	B	1	5	1	67.19%	70.47%	66.18%	67.09%	63.60%	61.80%
Beatrix	B	2	8	2	94.95%	94.73%	97.53%	97.58%	96.60%	98.40%
Alex	A	0	2	0	18.22%	13.71%	13.23%	10.05%	17.40%	9.20%
Ahmad	A	1	5	1	67.19%	54.48%	66.18%	50.49%	63.60%	64.40%
Anthony	A	2	8	2	94.95%	90.02%	97.53%	90.51%	96.60%	68.20%

Individual Fairness

Counterfactual

	s	x_1	x_2	x_3	$\hat{m}_{\text{glm}}(\mathbf{x})$	$\hat{m}_{\text{glm}}(\mathbf{x}, s)$	$\hat{m}_{\text{gam}}(\mathbf{x})$	$\hat{m}_{\text{gam}}(\mathbf{x}, s)$	$\hat{m}_{\text{rf}}(\mathbf{x})$	$\hat{m}_{\text{rf}}(\mathbf{x}, s)$
adjusted data, using marginal quantiles										
Betty	A	-1.68	2.1	-1.68	3.51%	3.58%	4.78%	4.85%	10.40%	10.80%
Brienne	A	-0.98	5.1	-0.96	19.39%	17.65%	16.64%	16.13%	29.00%	41.00%
Beatrix	A	-0.27	7.9	-0.26	59.83%	53.65%	51.89%	46.37%	53.60%	49.00%
adjusted data, using optimal transport, (c)										
Betty	A	-1.96	2.1	-1.9	2.62%	2.82%	4.65%	4.81%	0.00%	0.00%
Brienne	A	0.29	5	0.25	48.24%	38.92%	40.04%	32.14%	21.40%	12.20%
Beatrix	A	0.31	7.8	0.21	72.83%	65.1%	67.5%	58.83%	20.80%	15%
adjusted data, using Gaussian transport, (c)										
Betty	A	-1.58	2.15	-1.59	3.95%	3.96%	4.96%	4.99%	0.40%	0.40%
Brienne	A	-0.98	4.96	-0.99	18.47%	16.84%	15.84%	15.40%	19.80%	27.20%
Beatrix	A	-0.38	7.79	-0.38	55.71%	50.05%	47.86%	43.16%	51.80%	63.60%

Individual Fairness



Optimal matching, of individuals in group B to individuals in group A, on right, where points ● are Betty, Brienne and Beatrix, and ● their counterfactual version in group A.

Individual Fairness

Counterfactual

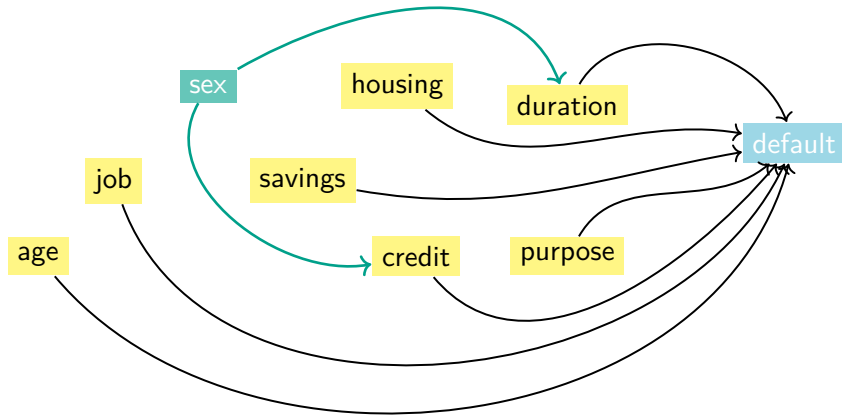
	s	x_1	x_2	x_3	$\hat{m}_{\text{glm}}(\mathbf{x})$	$\hat{m}_{\text{glm}}(\mathbf{x}, s)$	$\hat{m}_{\text{gam}}(\mathbf{x})$	$\hat{m}_{\text{gam}}(\mathbf{x}, s)$	$\hat{m}_{\text{rf}}(\mathbf{x})$	$\hat{m}_{\text{rf}}(\mathbf{x}, s)$
adjusted data, with fairAdapt, Figure (e)										
Betty	A	-1.65	2	-1.32	3.63%	3.54%	4.72%	4.60%	14.60%	8.00%
Brienne	A	-0.97	4.55	-0.94	16.57%	14.96%	13.96%	13.51%	2.20%	5.20%
Beatrix	A	-0.33	7.72	-0.44	56.3%	50.71%	48.49%	43.74%	70.60%	74.80%
adjusted data, with fairAdapt, Figure (f)										
Betty	A	-1.75	2.28	-1.68	3.5%	3.6%	5.03%	5.13%	7.20%	7.00%
Brienne	A	-0.96	5.3	-0.96	20.9%	19.05%	17.91%	17.34%	5.80%	8.40%
Beatrix	A	-0.24	8.12	-0.34	62.31%	56.43%	54.8%	49.3%	45.60%	39.20%

Numerical illustrations

Original data

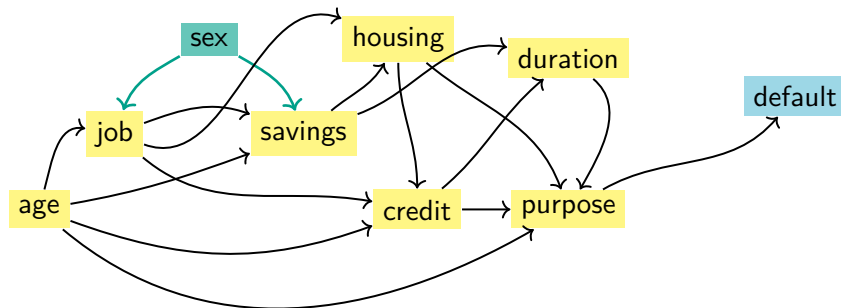
	s	x_1	x_2	x_3	$\hat{m}_{\text{glm}}(\mathbf{x})$	$\hat{m}_{\text{glm}}(\mathbf{x}, s)$	$\hat{m}_{\text{gam}}(\mathbf{x})$	$\hat{m}_{\text{gam}}(\mathbf{x}, s)$	$\hat{m}_{\text{rf}}(\mathbf{x})$	$\hat{m}_{\text{rf}}(\mathbf{x}, s)$
Betty	B	0	2	0	18.22%	24.06%	13.23%	17.63%	17.4%	29.6%
Brienne	B	1	5	1	67.19%	70.47%	66.18%	67.09%	63.60%	61.80%
Beatrix	B	2	8	2	94.95%	94.73%	97.53%	97.58%	96.60%	98.40%
Alex	A	0	2	0	18.22%	13.71%	13.23%	10.05%	17.40%	9.20%
Ahmad	A	1	5	1	67.19%	54.48%	66.18%	50.49%	63.60%	64.40%
Anthony	A	2	8	2	94.95%	90.02%	97.53%	90.51%	96.60%	68.20%

Numerical illustrations



Simple causal graph on the `GermanCredit` dataset,

Numerical illustrations



Causal graph on the `germancredit` dataset, from [Watson et al. \(2021\)](#)

Numerical illustrations

	Alex	Ahmad	Anthony	Betty	Brienne	Beatrix
s (gender)	M	M	M	F	F	F
x_1 Duration	12	18	30	12	18	30
$u = F_{1 s}(x_1)$	36%	57%	86%	34%	50%	79%
$T(x_1) = F_{1 s=M}^{-1}(u)$	12	18	30	12	18	24
x_2 Credit	1262	2319	4720	1262	2319	4720
$u = F_{2 s}(x_2)$	25%	55%	82%	17%	45%	76%
$T(x_2) = F_{2 s=M}^{-1}(u)$	1262	2319	4720	1074	1855	3854

Numerical illustrations

On the **GermanCredit** dataset

Firstname	s	Firstname	s	Job	Savings	Housing
Alex	M	Betty	F	highly qualified employee	100 DM	rent
Ahmad	M	Brienne	F	skilled employee	100<=...<500 DM	own
Anthony	M	Beatrix	F	unskilled - resident	no savings	for free

Original data

	s	Age	Duration	Credit	$\hat{m}_{\text{glm}}(\mathbf{x})$	$\hat{m}_{\text{glm}}(\mathbf{x}, s)$	$\hat{m}_{\text{gbm}}(\mathbf{x})$	$\hat{m}_{\text{gbm}}(\mathbf{x}, s)$	$\hat{m}_{\text{cart}}(\mathbf{x})$	$\hat{m}_{\text{cart}}(\mathbf{x}, s)$
Betty	F	26	12	1262	39.69%	36.66%	42.30%	43.26%	31.75%	31.75%
Brienne	F	33	18	2320	24.30%	22.61%	23.88%	21.08%	21.31%	21.31%
Beatrix	F	45	30	4720	30.88%	30.08%	28.49%	30.42%	15.38%	15.38%
Alex	M	26	12	1262	39.69%	42.10%	42.30%	44.86%	31.75%	31.75%
Ahmad	M	33	18	2320	24.30%	26.84%	23.88%	22.18%	21.31%	21.31%
Anthony	M	45	30	4720	30.88%	35.08%	28.49%	31.82%	15.38%	15.38%

Numerical illustrations

Original data

	s	Age	Duration	Credit	$\hat{m}_{\text{glm}}(\mathbf{x})$	$\hat{m}_{\text{glm}}(\mathbf{x}, s)$	$\hat{m}_{\text{gbm}}(\mathbf{x})$	$\hat{m}_{\text{gbm}}(\mathbf{x}, s)$	$\hat{m}_{\text{cart}}(\mathbf{x})$	$\hat{m}_{\text{cart}}(\mathbf{x}, s)$
Betty	M	26	12	1074	39.51%	41.90%	40.69%	44.86%	31.75%	31.75%
Brienne	M	33	18	1855	23.95%	26.46%	23.88%	22.18%	21.31%	21.31%
Beatrix	M	45	24	3854	24.91%	28.58%	20.55%	20.31%	21.31%	21.31%

adjusted data, with fairAdapt, causal graph from Figure ??

Betty	M	26	12	1110	42.73%	45.18%	44.24%	46.64%	31.75%	31.75%
Brienne	M	33	18	1787	23.90%	26.40%	23.88%	22.18%	21.31%	21.31%
Beatrix	M	45	24	3990	25.01%	28.70%	22.17%	23.60%	21.31%	21.31%

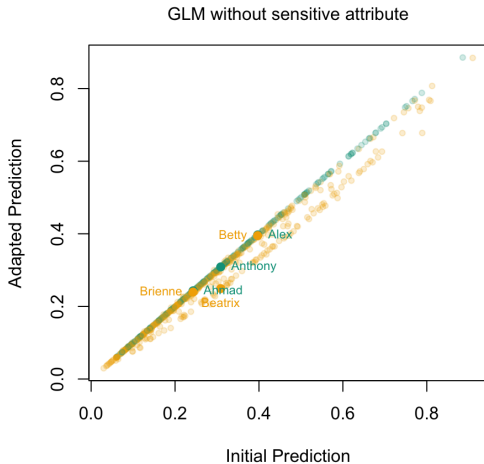
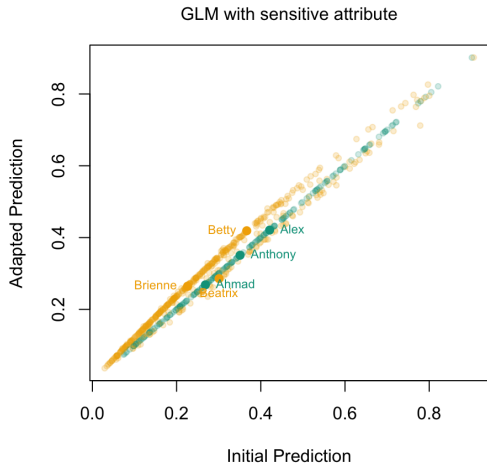
adjusted data, with fairAdapt, causal graph from Figure ??

Betty	M	26	18	1778	52.23%	54.03%	40.05%	46.81%	21.31%	21.31%
Brienne	M	33	15	1864	32.25%	35.85%	31.60%	25.97%	21.31%	21.31%
Beatrix	M	45	21	3599	39.70%	43.16%	28.36%	28.90%	21.31%	21.31%

adjusted data, with fairAdapt, causal graph from Figure ??

Betty	M	26	15	1882	49.05%	50.86%	35.32%	40.12%	21.31%	21.31%
Brienne	M	33	18	1881	50.76%	53.49%	43.00%	38.77%	21.31%	21.31%
Beatrix	M	45	24	3234	24.20%	26.23%	14.63%	16.84%	21.31%	21.31%

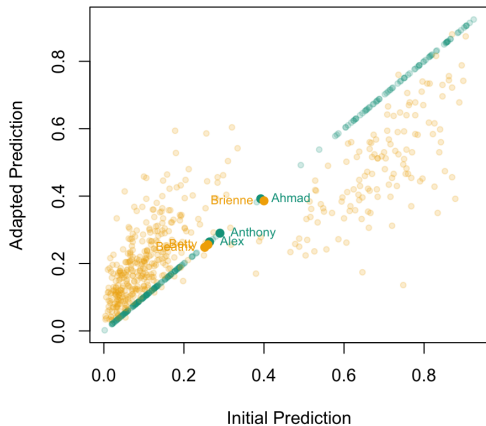
Numerical illustrations



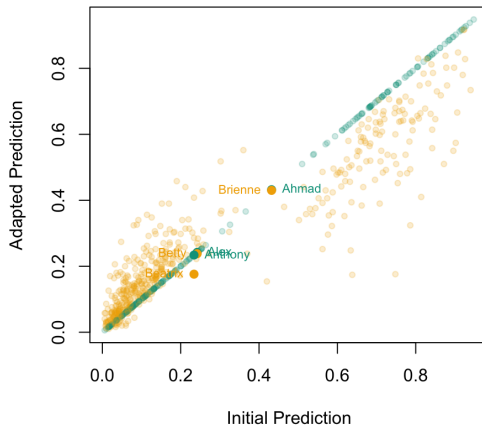
Scatterplot $(m(\mathbf{x}_i), m(\mathcal{T}^*(\mathbf{x}_i)))$ for individuals in groups **M** and **F**.

Numerical illustrations

RF with sensitive attribute

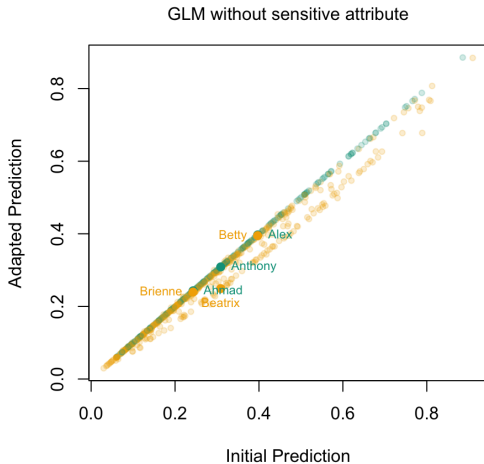
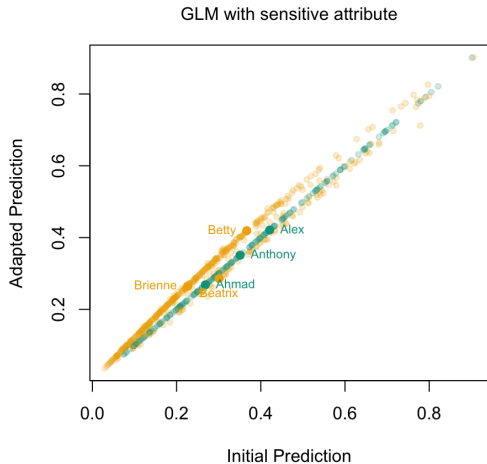


RF without sensitive attribute



Scatterplot $(m(\mathbf{x}_i), m(\mathcal{T}^*(\mathbf{x}_i)))$ for individuals in groups **M** and **F**.

Numerical illustrations



Scatterplot $(m(\mathbf{x}_i), m(\mathcal{T}^*(\mathbf{x}_i)))$ for individuals in groups **M** and **F**.

– Part 7 –

Mitigating Discrimination

Achieving a Fair Prediction

› Mitigating discrimination is usually seen as paradoxical, because in order to avoid discrimination, we must create another discrimination. More precisely, Supreme Court Justice Harry Blackmun stated, in 1978, “*in order to get beyond racism, we must first take account of race. There is no other way. And in order to treat some persons equally, we must treat them differently,*” cited in Knowlton (1978), as mentioned in Lippert-Rasmussen (2020)).



› More formally, an argument in favor of affirmative action – called “*the present-oriented anti-discrimination argument*” – is simply that justice requires that we eliminate or at least mitigate (present) discrimination by the best morally permissible means of doing so, which corresponds to affirmative action. Freeman (2007) suggested a “*time-neutral anti-discrimination argument,*” in order to mitigate past, present, or future discrimination.

Achieving a Fair Prediction

➤ But there are also arguments against affirmative action, corresponding to “*the reverse discrimination objection*,” as defined in [Goldman \(1979\)](#): some might consider that there is an absolute ethical constraint against unfair discrimination (including affirmative action). To quote another Supreme Court Justice, in 2007, John G. Roberts of the US Supreme Court submits: “*The way to stop discrimination on the basis of race is to stop discriminating on the basis of race*” ([Turner \(2015\)](#) and [Sabbagh \(2007\)](#)).



➤ The arguments against affirmative action are usually based on two theoretical moral claims, according to [Pojman \(1998\)](#). The first denies that groups have moral status (or at least meaningful status). According to this view, individuals are only responsible for the acts they perform as specific individuals and, as a corollary, we should only compensate individuals for the harms they have specifically suffered. The second asserts that a society should distribute its goods according to merit.

Achieving a Fair Prediction

- › We have defined the risk of a model $m \in \mathcal{M}$ as $\mathcal{R}(m) = \mathbb{E}[\ell(Y, m(\mathbf{X}))]$.
- › Define the classes of fair models,

$$\begin{cases} \mathcal{M}_{\text{DP}} = \{m \in \mathcal{M} \text{ s.t. } m(\mathbf{X}) \perp\!\!\!\perp S\} \\ \mathcal{M}_{\text{EO}} = \{m \in \mathcal{M} \text{ s.t. } m(\mathbf{X}) \perp\!\!\!\perp S \mid Y\} \end{cases}$$

- › Fairness is achieved by projection onto a fair subspace

$$\hat{m}_{\text{fair}} \in \underset{m \in \mathcal{M}_{\text{fair}}}{\operatorname{argmin}} \{ \hat{\mathcal{R}}_n(m) \}$$

Definition 7.1: Price of fairness

Given a risk \mathcal{R} , a class \mathcal{M} and the fair subclass $\mathcal{M}_{\text{fair}}$, the price of fairness

$$\mathcal{E}_{\text{fair}}(\mathcal{M}) = \min_{m \in \mathcal{M}_{\text{fair}}} \{ \mathcal{R}(m) \} - \min_{m \in \mathcal{M}} \{ \mathcal{R}(m) \}.$$

Achieving a Fair Prediction

- › Recall that Bayes estimator is the best model,

$$\mu(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}] \text{ and set } \begin{cases} \mu_{\mathbf{A}}(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}, S = \mathbf{A}] \\ \mu_{\mathbf{B}}(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}, S = \mathbf{B}] \end{cases}$$

- › From the definition of Wasserstein distance,

$$W_2(p, q) = \left(\inf_{\pi \in \Pi(p, q)} \int |x - y|^2 d\pi(x, y) \right)^{1/2}$$

Thus,

$$\mathbb{E}[|m(\mathbf{X}, S) - \mu_S(\mathbf{X})|^2 | S = s] \geq W_2(\mathbb{P}_m, \mathbb{P}_s)^2$$

Achieving a Fair Prediction

Proposition 7.1: Price of fairness and Wasserstein Barycenter

$$\mathcal{E}_{\text{fair}}(\mathcal{M}) = \min_{m \in \mathcal{M}_{\text{fair}}} \{\mathcal{R}(m)\} - \min_{m \in \mathcal{M}} \{\mathcal{R}(m)\} \geq \min_{g \in \mathcal{M}} \{\mathbb{E} \left(W_2(\mathbb{P}_S, \mathbb{P}_{S,g})^2 \right) \}$$

where \mathbb{P}_S is the condition distribution of $\mu(\mathbf{X}, S)$, given S , and $\mathbb{P}_{S,g}$ is the condition distribution of $g(\mathbf{X}, S)$, given S . Moreover, if $\mathcal{M}_{\text{fair}} = \mathcal{M}_{\text{DP}}$, and if \mathbb{P}_s is absolutely continuous (w.r.t. Lebesgue measure),

$$\mathcal{E}_{\text{DP}}(\mathcal{M}) = \min_{g \in \mathcal{M}} \{\mathbb{E} \left(W_2(\mathbb{P}_S, \mathbb{P}_{S,g})^2 \right) \} = \min_{g \in \mathcal{M}} \left\{ \sum_s \mathbb{P}[S = s] \cdot W_2(\mathbb{P}_s, \mathbb{P}_{s,g})^2 \right\}$$

See [Gouic et al. \(2020\)](#).

- The minimum is reached at the **Wasserstein barycenter** of \mathbb{P}_S 's.

Pre-Processing

- Write the $n \times k$ matrix \mathbf{S} as a collection of k vectors in \mathbb{R}^n , $\mathbf{S} = (\mathbf{s}_1 \cdots \mathbf{s}_k)$, that will correspond to k sensitive attributes. The orthogonal projection on variables $\{\mathbf{s}_1, \dots, \mathbf{s}_k\}$ is associate to matrix $\Pi_{\mathbf{S}} = \mathbf{S}(\mathbf{S}^T \mathbf{S})^{-1} \mathbf{S}^T$, while the projection on the orthogonal of \mathbf{S} is $\Pi_{\mathbf{S}^\perp} = \mathbb{I} - \Pi_{\mathbf{S}}$ (Gram-Schmidt orthogonalization,).
- Let $\tilde{\mathbf{S}}$ denote the collection of centered vectors (using matrix notations, $\tilde{\mathbf{S}} = \mathbf{H}\mathbf{S}$ where $\mathbf{H} = \mathbb{I} - (\mathbf{1}\mathbf{1}^T)/n$).
- Write the $n \times p$ matrix \mathbf{X} as a collection of p vectors in \mathbb{R}^n , $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_p)$. For any \mathbf{x}_j , define

$$\mathbf{x}_j^\perp = \Pi_{\tilde{\mathbf{S}}^\perp} \mathbf{x}_j = \mathbf{x}_j - \tilde{\mathbf{S}}(\tilde{\mathbf{S}}^T \tilde{\mathbf{S}})^{-1} \tilde{\mathbf{S}}^T \mathbf{x}_j.$$

One can easily prove that \mathbf{x}_j^\perp is then orthogonal to any \mathbf{s} , since

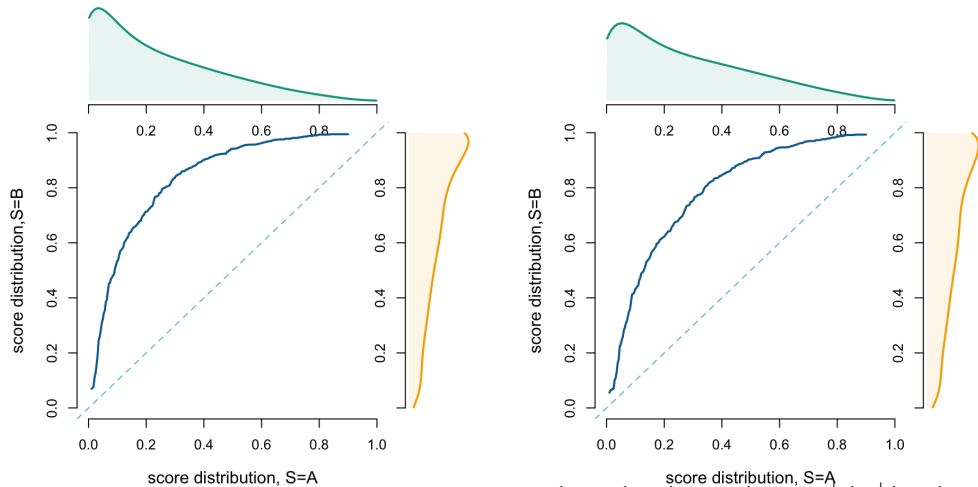
$$\text{Cov}(\mathbf{s}, \mathbf{x}_j^\perp) = \frac{1}{n} \mathbf{s}^T \mathbf{H} \mathbf{x}_j^\perp = \frac{1}{n} \tilde{\mathbf{s}}^T \Pi_{\tilde{\mathbf{S}}^\perp} \mathbf{x}_j = 0.$$

Pre-Processing

- And similarly the centered version of \mathbf{x}_j^\perp is then also orthogonal to any \mathbf{s} . From an econometric perspective, \mathbf{x}_j^\perp can be seen as the residual of the regression of \mathbf{x}_j against \mathbf{s} 's, obtained from least square estimation

$$\mathbf{x}_j = \tilde{\mathbf{s}}^\top \hat{\beta}_j + \mathbf{x}_j^\perp.$$

Pre-Processing



Optimal transport between distributions of $\hat{m}(\mathbf{x}_i, s_i)$'s (x-axis) to $\hat{m}^\perp(\mathbf{x}_i^\perp)$'s (y-axis), for individuals in group **A** on the left-hand side, and in group **B** on the right-hand side.

Pre-Processing

- › Consider the linear model $\mathbf{y} = \mathbf{S}\boldsymbol{\alpha} + \mathbf{X}^\perp\boldsymbol{\beta} + \varepsilon$
- › Consider the fairness constraint

$$R_{\text{fair}}^2(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\text{Var}[\mathbf{S}\boldsymbol{\alpha}]}{\text{Var}[\mathbf{S}\boldsymbol{\alpha} + \mathbf{X}^\perp\boldsymbol{\beta}]} = \frac{\boldsymbol{\alpha}^\top \text{Var}[\mathbf{S}]\boldsymbol{\alpha}}{\boldsymbol{\alpha}^\top \text{Var}[\mathbf{S}]\boldsymbol{\alpha} + \boldsymbol{\beta}^\top \text{Var}[\mathbf{X}^\perp]\boldsymbol{\beta}}$$

- › Then solve

$$\min_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \left\{ \mathbb{E}[\|\mathbf{y} - \mathbf{S}\boldsymbol{\alpha} - \mathbf{X}^\perp\boldsymbol{\beta}\|^2] \right\} \text{ s.t. } R_{\text{fair}}^2(\boldsymbol{\alpha}, \boldsymbol{\beta}) \leq r^2 \quad (r \in \mathbb{R}_+).$$

Pre-Processing

- An alternative was considered in [Komiyama and Shimao \(2017\)](#), with a Ridge penalty

$$\min_{\alpha, \beta} \left\{ \mathbb{E}[\|\mathbf{y} - \mathbf{S}\alpha - \mathbf{X}^\perp\beta\|_{\ell_2}^2] + \lambda\|\alpha\|_{\ell_2}^2 \right\}$$

- The penalty is on α only because (by construction) there is no discriminating information in \mathbf{X}^\perp . There is a closed form solution

$$\begin{pmatrix} (\mathbf{S}^\top \mathbf{S} + \lambda \mathbb{I})^{-1} \mathbf{S}^\top \mathbf{y} \\ (\mathbf{X}^{\perp\top} \mathbf{X}^\perp)^{-1} \mathbf{X}^\perp \mathbf{y} \end{pmatrix}$$

- In a linear regression problem, $\mathbf{y} = \mathbf{X}\beta + \varepsilon$. [Zafar et al. \(2017\)](#) suggested

$$\beta^* = \min_{\beta} \left\{ \mathbb{E}[\|\mathbf{y} - \mathbf{X}\beta\|^2] \right\} \text{ s.t. } |\text{Cov}[\mathbf{X}\beta, S]| \leq c \text{ (} \in \mathbb{R}_+ \text{)}.$$

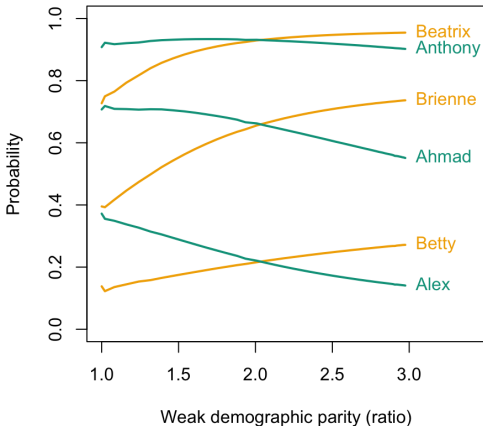
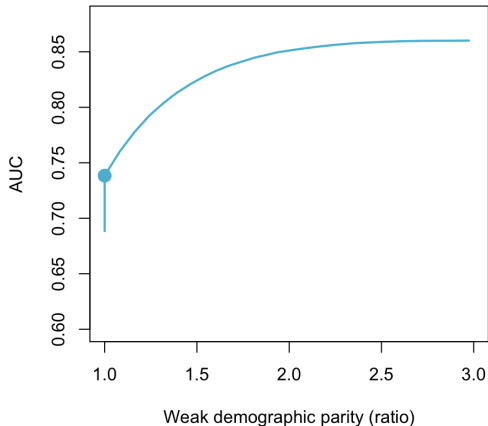
Pre-Processing

	$\widehat{m}(\mathbf{x}, s)$, aware					$\widehat{m}(\mathbf{x})$, unaware			
	← less fair		more fair →			← less fair		more fair →	
$\widehat{\beta}_0$ (Intercept)	-2.55	-2.29	-1.97	-1.51	-1.03	-2.14	-1.98	-1.78	-1.63
$\widehat{\beta}_1$ (x_1)	0.88	0.88	0.85	0.77	0.62	1.01	0.84	0.57	0.26
$\widehat{\beta}_2$ (x_2)	0.37	0.37	0.35	0.32	0.25	0.37	0.35	0.31	0.24
$\widehat{\beta}_3$ (x_3)	0.02	0.02	0.02	0.02	0.03	0.15	0.02	-0.15	-0.29
$\widehat{\beta}_B$ ($\mathbf{1}_B$)	0.82	0.44	-0.03	-0.70	-1.31	-	-	-	-

Pre-Processing

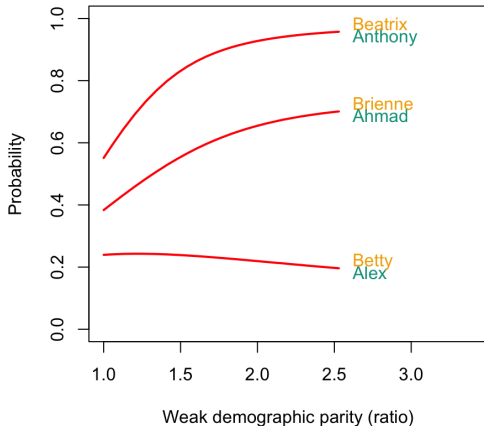
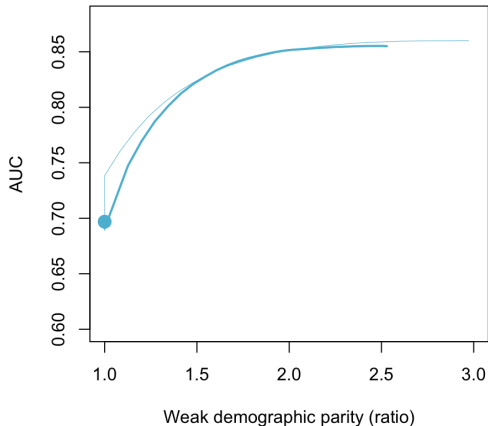
	$\hat{m}(\mathbf{x}, s)$, aware					$\hat{m}(\mathbf{x})$, unaware			
	← less fair		more fair →			← less fair		more fair →	
Betty	0.27	0.25	0.22	0.17	0.14	0.20	0.22	0.24	0.24
Brienne	0.74	0.71	0.66	0.54	0.40	0.70	0.66	0.55	0.38
Beatrix	0.95	0.95	0.93	0.87	0.73	0.96	0.93	0.82	0.55
Alex	0.14	0.17	0.22	0.29	0.37	0.20	0.22	0.24	0.24
Ahmad	0.55	0.61	0.66	0.70	0.71	0.70	0.66	0.55	0.38
Anthony	0.90	0.92	0.93	0.93	0.91	0.96	0.93	0.82	0.55
$\mathbb{E}[\hat{m}(\mathbf{x}_i, s_i) S = \text{A}]$	0.23	0.26	0.31	0.36	0.42	0.25	0.30	0.37	0.41
$\mathbb{E}[\hat{m}(\mathbf{x}_i, s_i) S = \text{B}]$	0.67	0.65	0.61	0.53	0.42	0.64	0.61	0.54	0.41
(ratio)	×2.97	×2.49	×2.01	×1.46	×1.00	×2.53	×2.02	×1.48	×1.00
AUC	0.86	0.86	0.85	0.82	0.74	0.86	0.85	0.82	0.70

Pre-Processing



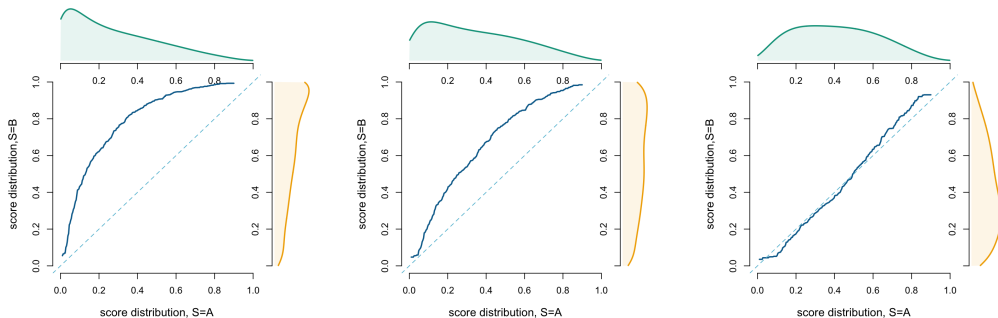
AUC of $\hat{m}_{\hat{\beta}_\lambda}$ and evolution of $\hat{m}_{\hat{\beta}_\lambda}(\mathbf{x}_i, s_i)$ (with a logistic regression)

Pre-Processing



AUC of $\hat{m}_{\hat{\beta}_\lambda}$ and evolution of $\hat{m}_{\hat{\beta}_\lambda}(\mathbf{x}_i)$ (with a logistic regression)

Pre-Processing



Optimal transport between distributions of $\widehat{m}_{\widehat{\beta}_\lambda}(\mathbf{x}_i, s_i)$'s from individuals in group **A** and in **B**, for different values of λ (low value on the left-hand side and high value on the right-hand side), associated with a demographic parity penalty criteria

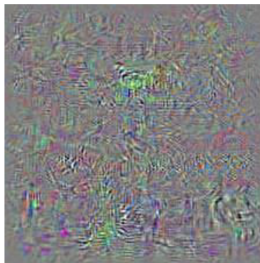
Pre-Processing

- **Adversarial learning** has to do with robustness of learning algorithm, Szegedy et al. (2013) (“*are neural network stables?*”).
- “*Adversarial examples are inputs to machine learning models that an attacker has intentionally designed to cause the model to make a mistake*”, Bengio et al. (2017)



Schoolbus

+



Perturbation

=



Ostrich

Pre-Processing

- Adversarial learning deals with the problem that the distributions we obtain IRL are not the ones we train the model on... and we try to quantify what can go wrong
- Popular in pictures (what happens if we rotate an object, add glasses to people, etc). Brittleness of ML algorithms...
- Problem of data pollution (add outliers) and problems of adversarial examples.
- Machine learning perspective

$$\min_{\theta} \{ \mathbb{E}_{(\mathbf{X}, Y) \sim \mathbb{P}} [\ell(m_{\theta}(\mathbf{X}), Y)] \}$$

- Adversarial perspective

$$\max_{\varepsilon \in \mathcal{E}} \{ \mathbb{E}_{(\mathbf{X}, Y) \sim \mathbb{P}} [\ell(m_{\theta}(\mathbf{X} + \varepsilon), Y)] \}$$

leads to robust learning...

Pre-Processing

$$\min_{\theta} \left\{ \underbrace{\max_{\varepsilon \in \mathcal{E}} \left\{ \mathbb{E}_{(\mathbf{X}, Y) \sim \mathbb{P}} [\ell(m_{\theta}(\mathbf{X} + \varepsilon), Y)] \right\}}_{\text{creating an adversarial example}} \right\}$$

training a robust classifier

- Approaches based on **robust optimization**, [Ben-Tal et al. \(2009\)](#), e.g., Danskin's Theorem, [Danskin \(1967\)](#),

$$\nabla_{\theta} \max_{\varepsilon \in \mathcal{E}} \{ \ell(m_{\theta}(\mathbf{X} + \varepsilon), Y) \} = \nabla_{\theta} \ell(m_{\theta}(\mathbf{X} + \varepsilon^*), Y)$$

where $\varepsilon^* = \operatorname{argmax}_{\varepsilon \in \mathcal{E}} \{ \ell(m_{\theta}(\mathbf{X} + \varepsilon), Y) \}$.

- Recall the **minimax** theorem from [von Neumann \(1928\)](#)

Proposition 7.2: Nash equilibrium and Minimax

Let A be some $m \times n$ real-valued matrix, there is a Nash equilibrium $(\mathbf{x}_*, \mathbf{y}_*)$ associated with A if

$$\mathbf{y}_*^\top A \mathbf{x}_* = \max_{\mathbf{x} \in \mathcal{S}_m} \min_{\mathbf{y} \in \mathcal{S}_n} \{\mathbf{y}^\top A \mathbf{x}\} = \min_{\mathbf{y} \in \mathcal{S}_n} \max_{\mathbf{x} \in \mathcal{S}_m} \{\mathbf{y}^\top A \mathbf{x}\}.$$

Pre-Processing

- Consider a **Minimax games**: given that the discriminator will try to do the best job it can, the generator is set to make the discriminator as wrong as possible

$$\min_{\theta_g} \max_{\theta_d} \{ \mathbb{E}_{\mathbf{X} \sim \mathbb{P}} [\log(m_{\theta_d}(\mathbf{x}))] + \mathbb{E}_{\mathbf{Z} \sim \mathbb{Q}} [\log(1 - m_{\theta_d}(G_{\theta_d}(\mathbf{z})))] \}$$

where $\mathbf{X} \sim \mathbb{P}$ denotes data sampled from the training data, while $\mathbf{Z} \sim \mathbb{Q}$ are sampled by the opponent

- See [Wadsworth et al. \(2018\)](#), [Xu et al. \(2021\)](#), [Lima et al. \(2022\)](#) for achieving fairness through adversarial learning

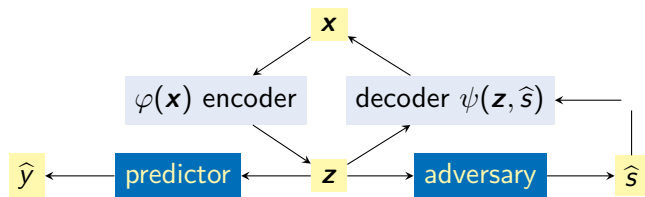
Pre-Processing

- FairGAN, [Xu et al. \(2018\)](#)

Pre-processing approach actually, with demographic parity (DP)

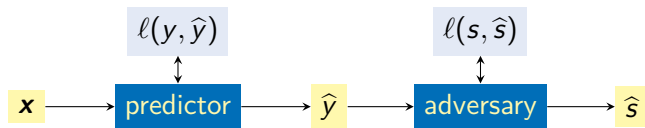
Other algorithms are in-processing approaches, with demographic parity (DP) and equalized odds (EO)

- Learning adversarially fair and transferable representations, [Madras et al. \(2018\)](#)



- Adversarially learning fair representations, [Beutel et al. \(2017\)](#)
- Fair Adversarial Debiasing Approach, [Zhang et al. \(2018\)](#)

Pre-Processing



Following [Zhang et al. \(2018\)](#)

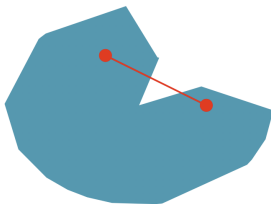
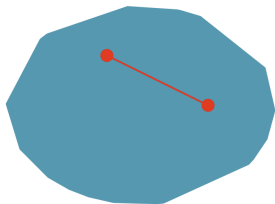
- › the predictor predicts y given x ,
- › the adversary tries to predict s based on the output of the predictor
- › the predictor targets to increase its prediction accuracy
- › and tries to increase the adversary's loss

Barycenter

Several approaches can be considered to define **means**, **averages**, **centroids**, **barycenters** (etc.), as discussed in [Fréchet \(1948\)](#) and [Grove and Karcher \(1973\)](#),

- ▶ convex properties (from [Möbius \(1827\)](#) and [Rockafellar \(1970\)](#))
- ▶ axiomatization (from [Nagumo \(1930\)](#), [Kolmogorov \(1930\)](#) and [Aczél \(1948\)](#))
- ▶ optimization (from [Hey \(1814\)](#), [Nathan \(1952\)](#) and [Agueh and Carlier \(2011\)](#))

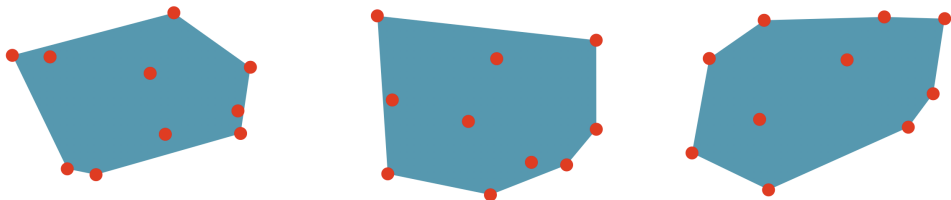
$C \subset \mathbb{R}^n$ is convex if $\mathbf{x}, \mathbf{y} \in C \implies t\mathbf{x} + (1 - t)\mathbf{y} \in C$ for all $t \in [0, 1]$



Barycenter

Let $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^n$, then a convex combination is any linear combination $\omega_1 \mathbf{x}_1 + \dots + \omega_k \mathbf{x}_k$ with $(\omega_1, \dots, \omega_k) \in \mathcal{S}_k \subset \mathbb{R}_+$.

The convex hull of a set C is the set of all convex combinations of elements of C .



The geometric centroid of a convex object always lies in the object.

Barycenter

Define the barycenter for two points, with equal weights as a function $M : E \times E \rightarrow E$

- ▶ Reflexivity: $M(\mathbf{x}, \mathbf{x}) = \mathbf{x}$,
- ▶ Symmetry: $M(\mathbf{x}_1, \mathbf{x}_2) = M(\mathbf{x}_2, \mathbf{x}_1)$,
- ▶ Continuity: $M(\cdot, \cdot)$ is continuous,
- ▶ Bisymmetry: $M(M(\mathbf{x}_{11}, \mathbf{x}_{12}), M(\mathbf{x}_{21}, \mathbf{x}_{22})) = M(M(\mathbf{x}_{11}, \mathbf{x}_{21}), M(\mathbf{x}_{12}, \mathbf{x}_{22}))$

Then (see [Aczél \(1948\)](#)), there is f such that

$$M(\mathbf{x}_1, \mathbf{x}_2) = f^{-1} \left(\frac{1}{2}f(\mathbf{x}_1) + \frac{1}{2}f(\mathbf{x}_2) \right).$$

If $E \subset \mathbb{R}^k$, consider means on each coordinate axis independently.

› A natural extension is

$$B_f(\mathbf{x}, \boldsymbol{\omega}) = f^{-1} \left(\sum_{i=1}^n \omega_i f(x_i) \right).$$

Barycenter

- › For the optimisation approach, given a distance d on E , set

$$B_d(\mathbf{x}, \omega) = \operatorname{argmin}_{z \in E} \left\{ \sum_{i=1}^n \omega_i d(x_i, z) \right\}$$

- › Consider some points $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ in a metric space \mathbb{R}^2
- › The **mean** is

$$\bar{\mathbf{x}} = \frac{\mathbf{x}_1 + \mathbf{x}_2 + \dots + \mathbf{x}_k}{k} = \frac{1}{k} \sum_{i=1}^k \mathbf{x}_i,$$

or equivalently

$$\bar{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \left\{ \frac{1}{k} \sum_{i=1}^k \|\mathbf{x} - \mathbf{x}_i\|_{\ell_2}^2 \right\}.$$

- › But they can be defined using any distance/divergence/discrepancy

Barycenter

- › Instead of points $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ in the metric space \mathbb{R}^2 , we can consider some measures $\{\mathbb{P}_1, \mathbb{P}_2, \dots, \mathbb{P}_k\}$.
- › The **Euclidean mean** is

$$\bar{\mathbb{Q}} = \operatorname{argmin}_{\mathbb{Q}} \left\{ \frac{1}{k} \sum_{i=1}^k \Delta^2(\mathbb{Q}, \mathbb{P}_i) \right\},$$

where $\Delta^2(\mathbb{Q}, \mathbb{P}_i) = \int_{\mathbb{R}^2} (d\mathbb{Q} - d\mathbb{P}_i)^2$.

But any discrepancy function can be considered

- › One can consider Wasserstein discrepancy

Definition 8.1: Wasserstein W_2 Barycenter, [Agueh and Carlier \(2011\)](#)

$$\bar{\mathbb{Q}} = \underset{\mathbb{Q}}{\operatorname{argmin}} \left\{ \sum_{i=1}^k \omega_i W_2(\mathbb{Q}, \mathbb{P}_i)^2 \right\},$$

- This can be seen as a multi-marginal optimal transport problem.
- Recall that the “*push-forward*” measure is

$$\mathbb{P}_1(\mathcal{A}) = \mathcal{T}_{\#} \mathbb{P}_0(\mathcal{A}) = \mathbb{P}_0(\mathcal{T}^{-1}(\mathcal{A})), \quad \forall \mathcal{A} \subset \mathbb{R}.$$

An optimal transport \mathcal{T}^* (in Brenier’s sense, from [Brenier \(1991\)](#), see [Villani \(2009\)](#) or [Galichon \(2016\)](#)) from \mathbb{P}_0 towards \mathbb{P}_1 will be solution of

$$\mathcal{T}^* \in \underset{\mathcal{T}: \mathcal{T}_{\#} \mathbb{P}_0 = \mathbb{P}_1}{\operatorname{arginf}} \left\{ \int_{\mathbb{R}^k} \ell(\mathbf{x}, \mathcal{T}(\mathbf{x})) d\mathbb{P}_0(\mathbf{x}) \right\},$$

Barycenter

and for univariate distributions, the optimal transport \mathcal{T}^* is the monotone transformation.

$$\mathcal{T}^* : x_0 \mapsto x_1 = F_1^{-1} \circ F_0(x_0).$$

➤ Given a reference measure, say \mathbb{P}_1 , it is possible to write the barycenter as the "*average push-forward*" transformation of \mathbb{P}_1 : if $\mathbb{P}_i = \mathcal{T}_{\#}^{1 \rightarrow i} \mathbb{P}_1$ (with the convention that $\mathcal{T}_{\#}^{1 \rightarrow 1}$ is the identity),

Proposition 8.1: Wasserstein W_2 Barycenter,

$$\bar{\mathbb{Q}} = \left(\sum_{i=1}^k \omega_i \mathcal{T}^{1 \rightarrow i} \right)_{\#} \mathbb{P}_1.$$

Proposition 8.2: Wasserstein W_2 Barycenter,

$$\bar{\mathbb{Q}} = \left(\sum_{i=1}^k \omega_i \mathcal{T}^{1 \rightarrow i} \right)_{\#} \mathbb{P}_1.$$

- › Computation of barycenters can be computationally difficult, [Altschuler and Boix-Adsera \(2021\)](#)
- › For univariate distributions, there is a simple expression, $\mathcal{T}^{1 \rightarrow i}$ is simply a rearrangement, defined as $\mathcal{T}^{1 \rightarrow i} = F_i^{-1} \circ F_1$, where $F_i(t) = \mathbb{P}_i((-\infty, t])$ and F_i^{-1} is its generalized inverse

Proposition 8.3: Wasserstein W_2 Barycenter, univariate distributions

.... $\mathcal{T}^{1 \rightarrow i}$ is simply a rearrangement, defined as $\mathcal{T}^{1 \rightarrow i} = F_i^{-1} \circ F_1$, where $F_i(t) = \mathbb{P}_i((-\infty, t])$...

$$\bar{\mathbb{Q}} = \left(\sum_{i=1}^n k\omega_i \mathcal{T}^{1 \rightarrow i} \right)_{\#} \mathbb{P}_1.$$

Proposition 8.4: Wasserstein W_2 Barycenter, univariate distributions

Given two scores $m(\mathbf{x}, s = \text{A})$ and $m(\mathbf{x}, s = \text{B})$, the “fair barycenter score” is

$$\begin{cases} m^*(\mathbf{x}, s = \text{A}) = \mathbb{P}[S = \text{A}] \cdot m(\mathbf{x}, s = \text{A}) + \mathbb{P}[S = \text{B}] \cdot F_{\text{B}}^{-1} \circ F_{\text{A}}(m(\mathbf{x}, s = \text{A})) \\ m^*(\mathbf{x}, s = \text{B}) = \mathbb{P}[S = \text{A}] \cdot F_{\text{A}}^{-1} \circ F_{\text{B}}(m(\mathbf{x}, s = \text{B})) + \mathbb{P}[S = \text{B}] \cdot m(\mathbf{x}, s = \text{B}). \end{cases}$$

Barycenter

- that is generally numerically intractable (computing one subgradient requires solving k optimal transports)
- In the discrete case, if we consider a fixed grid (so that C can be computed once only)

$$\min_{\mathbf{a}} \left\{ \sum_{i=1}^k \min_{P_i \in U_{\mathbf{a}, \mathbf{b}_i}} \{ \langle P_i, C \rangle \} \right\},$$

- We can write this as a large linear program
$$\min_{\mathbb{Q}} \left\{ \min_{P_1, \dots, P_k, \mathbf{a}} \sum_{i=1}^k \{ \langle P_i, C \rangle \} \right\}, \text{ where } \begin{cases} P_1^\top \mathbf{1}_n = \mathbf{b}_1 \\ \vdots \\ P_k^\top \mathbf{1}_n = \mathbf{b}_k \\ P_1 \mathbf{1}_n = \dots = P_k \mathbf{1}_n = \mathbf{a} \end{cases}$$

Barycenter

- › In the Gaussian case, some linear algebra might help
- › A variance matrix Σ is a positive-semidefinite ($\mathbf{a}^\top \Sigma \mathbf{a} \geq 0$ for all $\mathbf{a} \in \mathbb{R}^n$) symmetric ($\Sigma^\top = \Sigma$) matrix.

Definition 8.2: Matrix Exponential

Consider some $n \times n$ matrix \mathbf{X}

$$\exp(\mathbf{X}) = \sum_{k=0}^{\infty} \frac{1}{k!} \mathbf{X}^k = \lim_{k \rightarrow \infty} \left(\mathbb{I} + \frac{1}{k} \mathbf{X} \right)^k$$

- › The matrix exponential of a real symmetric matrix is positive definite.

Definition 8.3: Matrix Logarithm

Consider some $n \times n$ matrix \mathbf{X} such that $\|\mathbf{X} - \mathbb{I}\| < 1$

$$\log(\mathbf{X}) = \sum_{k=0}^{\infty} \frac{(-1)^{k+1}}{k} (\mathbf{X} - \mathbb{I})^k$$

- › As expected, $\exp(\log(\mathbf{X})) = \mathbf{X}$. In R, use `expm::expm` and `expm::logm`.
- › See Theorem 7.2.6 in [Horn and Johnson \(2012\)](#)

Definition 8.4: Square root of a symmetric matrix

Let \mathbf{X} be a real positive semidefinite matrix, then there is exactly one positive semidefinite matrix \mathbf{Z} such that $\mathbf{X} = \mathbf{Z}\mathbf{Z} = \mathbf{Z}^2$,

$$\mathbf{Z} = \mathbf{X}^{\frac{1}{2}} = \sum_{n=0}^{\infty} (-1)^n \binom{\frac{1}{2}}{n} (\mathbb{I} - \mathbf{X})^n.$$

- From power series $(1 - z)^{\frac{1}{2}} = \sum_{n=0}^{\infty} (-1)^n \binom{1/2}{n} z^n$, if $z = \mathbb{I} - \mathbf{X}$.
- In R, use `expm::sqrtm`.

$$\begin{pmatrix} 1 & 1.2 \\ 1.2 & 2 \end{pmatrix}^{\frac{1}{2}} = \begin{pmatrix} 0.8244771 & 0.5658953 \\ 0.5658953 & 1.2960565 \end{pmatrix}$$

Barycenter

► If \mathbb{P}_i is a Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, $\bar{\mathbb{Q}} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma}$ is some variance matrix (positive definite)

$$\boldsymbol{\mu} = \sum_{i=1}^k \omega_i \boldsymbol{\mu}_i \text{ and } \boldsymbol{\Sigma} \text{ satisfies } \boldsymbol{\Sigma} = \sum_{i=1}^k \omega_i (\boldsymbol{\Sigma}^{1/2} \boldsymbol{\Sigma}_i \boldsymbol{\Sigma}^{1/2})^{1/2}$$

```
1 > library(T4transport)
2 > par_mean = rbind(mu1, mu2)
3 > par_vars = array(0,c(2,2,2))
4 > par_vars[, ,1] = S1; par_vars[, ,2] = S2
5 > gmean = gaussbarypd(par_mean, par_vars)
```

Definition 8.5: Generalized Inverse, Rao and Mitra (1972)

Consider some $m \times n$ matrix \mathbf{X} . Matrix $n \times m$ \mathbf{Y} is said to be a generalized inverse of \mathbf{X} if $\mathbf{X}\mathbf{Y}\mathbf{X} = \mathbf{X}$. And \mathbf{Y} is denoted \mathbf{Y}^- .

- › From [Shurbet et al. \(1974\)](#), we can solve “quadratic matrix equations”

Theorem 8.1: Matrix Equation $XAX = B$

1. Equation $XAX = B$ has a solution if and only if

$$\begin{cases} (\mathbf{AB})^{1/2} \text{ exists,} \\ \mathbf{AA}^{-1}(\mathbf{AB})^{\frac{1}{2}} = (\mathbf{AB})^{\frac{1}{2}}, \\ \mathbf{B}(\mathbf{AB})^{\frac{1}{2}} - (\mathbf{AB})^{\frac{1}{2}} = \mathbf{B}. \end{cases}$$

2. If \mathbf{A} and \mathbf{B} are positive definite, the unique positive definite solution of $XAX = B$ is

$$\mathbf{X} = \mathbf{A}^{-\frac{1}{2}} \left(\mathbf{A}^{\frac{1}{2}} \mathbf{B} \mathbf{A}^{\frac{1}{2}} \right)^{\frac{1}{2}} \mathbf{A}^{-\frac{1}{2}}.$$

Theorem 8.2: Variance Σ

If $k = 2$, Σ satisfies

$$\Sigma = \omega_1 (\Sigma^{1/2} \Sigma_1 \Sigma^{1/2})^{1/2} + \omega_2 (\Sigma^{1/2} \Sigma_2 \Sigma^{1/2})^{1/2}$$

and the explicit expression is

$$\Sigma = \omega_1^2 \Sigma_1 + \omega_2^2 \Sigma_2 + \omega_1 \omega_2 \left(\Sigma_1^{1/2} (\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2})^{1/2} \Sigma_1^{-1/2} + \Sigma_1^{-1/2} \Sigma_2 \Sigma_1^{1/2} \right)^{1/2} \Sigma_1^{1/2}$$

Proposition 8.5: Variance Σ

$\sum_{i=1}^k \omega_i \Sigma_i - \Sigma$ is a positive matrix.

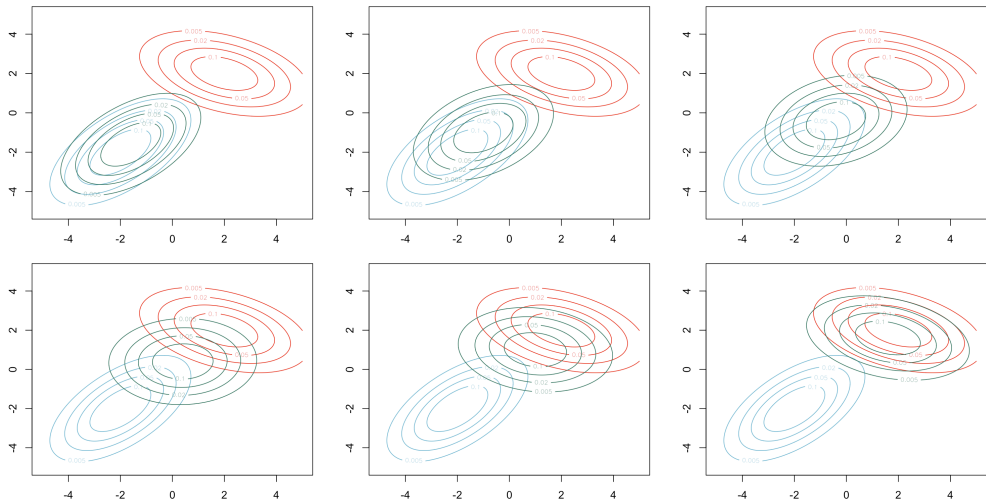
Proposition 8.6: Variance Σ

$$\text{If } \Sigma_i = \mathbf{P}\Delta_i\mathbf{P}^\top, \text{ then } \Sigma = \mathbf{P} \left(\sum_{i=1}^k \omega_i \Delta_i^{\frac{1}{2}} \right)^2 \mathbf{P}^\top$$

Consider two Gaussian distributions, $\mathcal{N}(\mu_A, \Sigma_A)$ and $\mathcal{N}(\mu_B, \Sigma_B)$, and weights $\omega_A = t$ and $\omega_B = 1 - t$, with $t \in [0, 1]$.

Barycenter

➤ Barycenter of two bivariate Gaussian distribution ($t = 0.1, 0.25, 0.4, 0.6, 0.75, 0.9$)



Computational Issues

- › Barycenters and k -coupling
- › To simplify, suppose that $\omega_1 = \dots = \omega_k$, the uniform barycenter $\bar{\mathbb{Q}}$ of (\mathbb{P}_i) is any solution of

$$\inf_{\mathbb{Q}} \left\{ \sum_{i=1}^k W_2(\mathbb{P}_i, \mathbb{Q})^2 \right\}$$

- › It is related to the k -coupling problem, defined in [Rüschendorf and Uckelmann \(2002\)](#), inspired by [Knott and Smith \(1994\)](#), who reduced 3-coupling to 2-coupling

$$\sup \left\{ \mathbb{E} \left[\left\| \sum_{i=1}^k X_i \right\|^2 \right] ; \text{ where } X_i \sim \mathbb{P}_i, \forall i \right\} \text{ and then } Y^* \sim \mathbb{Q}^*$$

as discussed in [Agueh and Carlier \(2011\)](#), under continuity assumptions.

Computational Issues

If $\bar{X} = \frac{1}{k} \sum_{i=1}^k X_i$, recall that for all Y ,

$$\sum_{i=1}^k \|X_i - Y\|^2 \geq \sum_{i=1}^k \|X_i - \bar{X}\|^2$$

and we can write, if $Z \sim \mathbb{Q}$,

$$\sum_{i=1}^k W_2(\mathbb{P}_i, \mathbb{Q})^2 = \sum_{i=1}^k \|X_i - Y\|^2 \geq \sum_{i=1}^k \|X_i - \bar{X}\|^2 \geq \sum_{i=1}^k W_2(\mathbb{P}_i, \bar{\mathbb{Q}})^2$$

and

$$\inf_{\mathbb{Q}} \left\{ \sum_{i=1}^k W_2(\mathbb{P}_i, \mathbb{Q})^2 \right\} = \inf \left\{ \sum_{i=1}^k \|X_i - \bar{X}\|^2; \text{ where } X_i \sim \mathbb{P}_i, \forall i \right\}$$

and simple arguments yield to the fact that this corresponds to an optimal k -coupling.

Computational Issues

- $\mathbf{X} = (X_1, \dots, X_k)$ is an optimal k -coupling if and only if the distribution of \bar{X} is a barycenter of \mathbb{P}_i 's.
- Following [Puccetti \(2017\)](#), observe that an equivalent representation of k coupling is

$$\sup \{ \mathbb{E}[\psi(X_1, \dots, X_k)] \} \text{ where } \psi(x_1, \dots, x_k) = \sum_{i=1}^k \langle x_i, \sum_{j \neq i} x_j \rangle.$$

that can be solve using an **iterative swapping algorithm** (ISA), as in [Puccetti et al. \(2020\)](#)

- [Oberman and Ruan \(2015\)](#) suggested multi-assignment problem

$$\max \left\{ \frac{1}{n} \sum_{j=1}^n \psi(\mathbf{x}_{1:\sigma_1(i)}, \dots, \mathbf{x}_{k:\sigma_k(i)}) \right\} \text{ where } \sigma_1, \dots, \sigma_k \text{ are permutations of } \{1, 2, \dots, n\}.$$

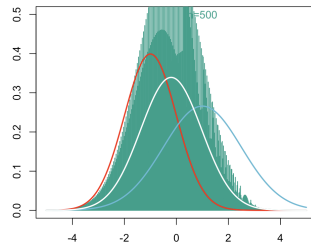
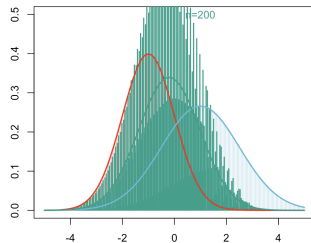
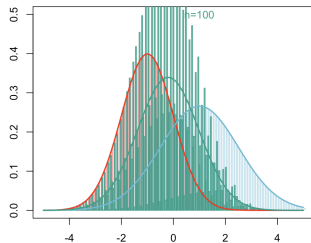
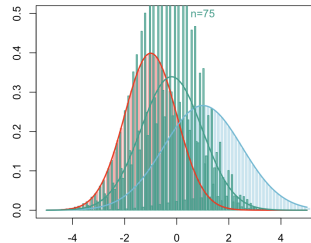
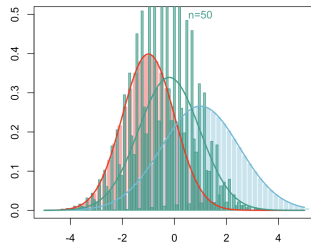
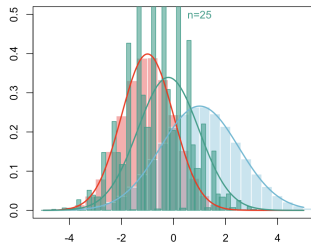
- Usually this problem is a lower bound the k -coupling problem, but they are equivalent in dimension $k = 2$.

Computational Issues

- › Barycenters based on histograms, [Cuturi and Doucet \(2014\)](#), computed in `bary14C` and `histbary14C`, in package `T4transport`
- › Barycenters based on histograms, [Benamou et al. \(2015\)](#), computed in `bary15C` and `histbary15C`, in package `T4transport`
- › We can also compute barycenters of pictures

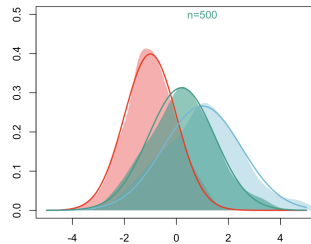
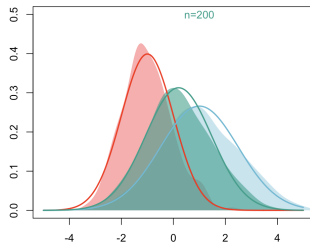
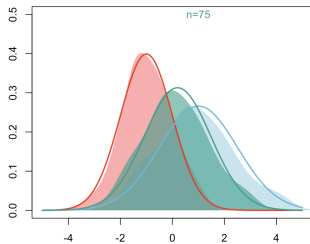
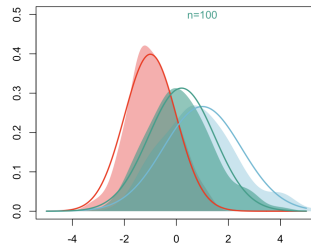
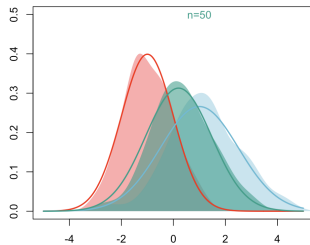
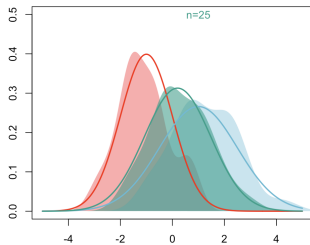
Computational Issues

➤ Numerical computation, Gaussian case



Computational Issues

➤ Numerical computation, Gaussian case (simulation and kernel density estimation)



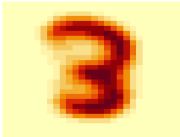
Barycenter of Several Pictures



Barycenter of Several Pictures



Barycenter of Several Pictures



- A geodesic is a curve representing in some sense the shortest[a] path (arc) between two points in a surface, or more generally in a Riemannian manifold. It is the general extension of a segment.

Definition 8.6: Constant Speed Geodesic

Consider some metric space, (E, d) . A constant speed geodesic between two points $x_0, x_1 \in E$ is a continuous curve $x : [0, 1] \rightarrow E$ such that for every $s, t \in (0, 1)$, $d(x_s, x_t) = |s - t|d(x_0, x_1)$.

Proposition 8.7

Consider μ_0 and μ_1 two measures on $X \subset \mathbb{R}^k$, compact and convex. Let $\pi^* \in \Pi(\mu_0, \mu_1)$ be an optimal transport plan. Define

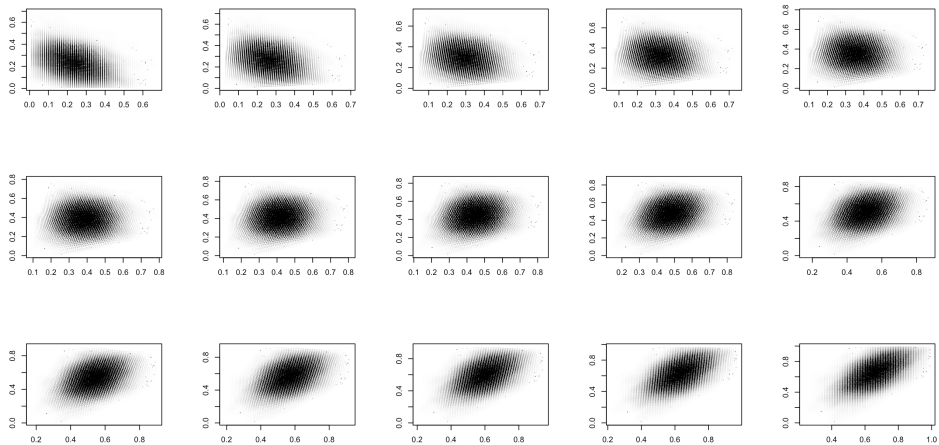
$$\mu_t = (\gamma_t)_\# \pi^*, \text{ where } \gamma_t(x, y) = (1 - t)x + ty, \ t \in (0, 1).$$

(μ_t) is a constant speed geodesic between μ_0 and μ_1 .

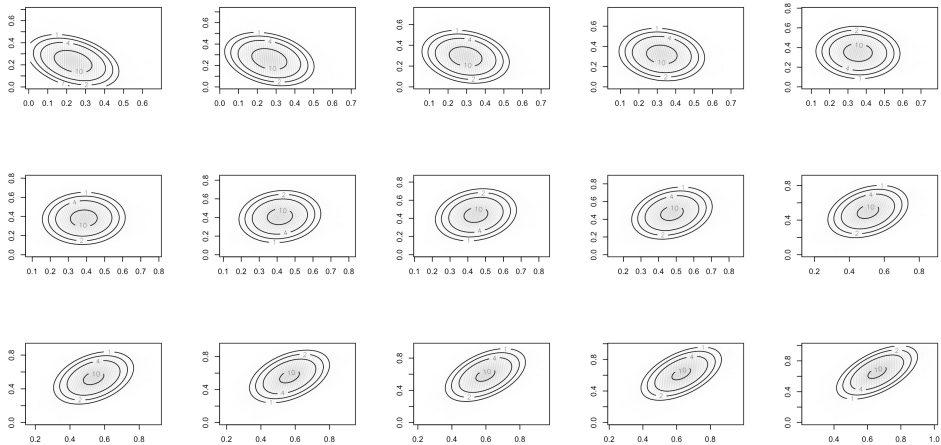
And if there exist an optimal transport map T^* ,

$$\mu_t = ((1 - t)I + tT^*)_\# \mu_0, \text{ where } t \in (0, 1).$$

Wasserstein Geodesic



Wasserstein Geodesic

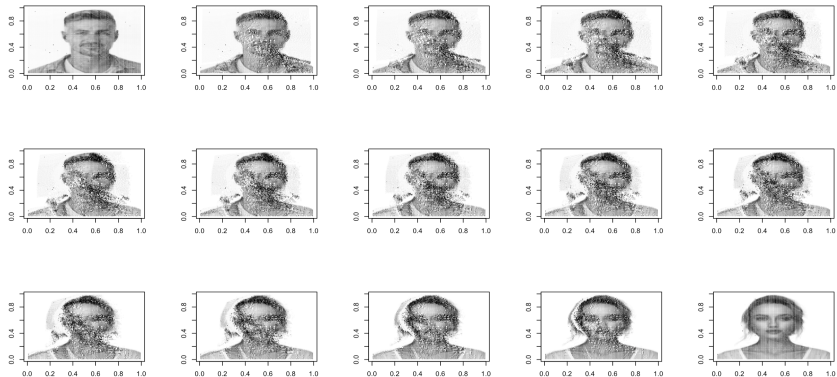


Wasserstein Geodesic



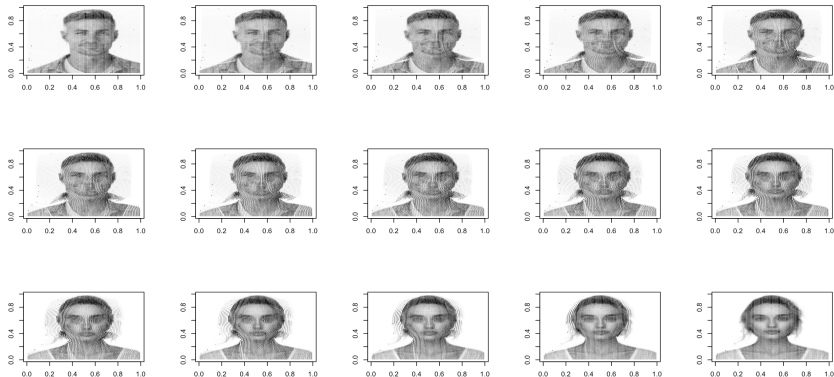
Wasserstein Geodesic

with the ℓ_1 cost function



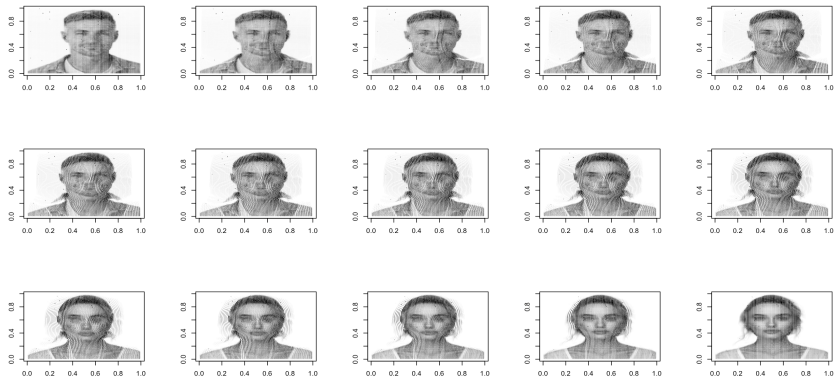
Wasserstein Geodesic

with the ℓ_2 cost function

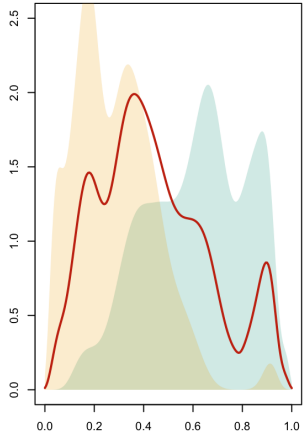
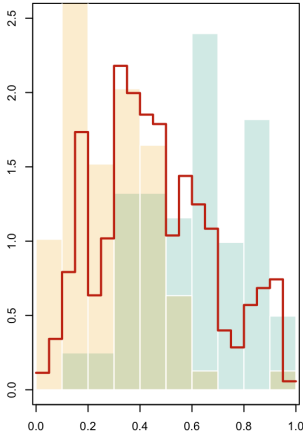
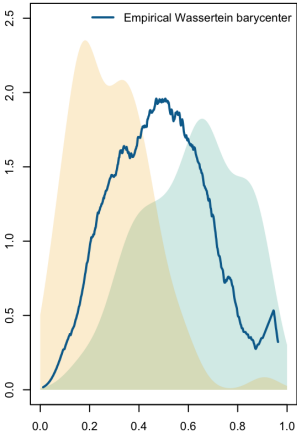


Wasserstein Geodesic

with the ℓ_3 cost function



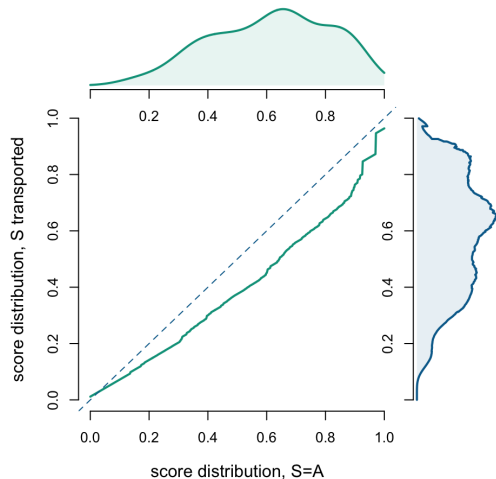
Application to toydata1



Application to toydata1

➤ Given scores $m(\mathbf{x}, s = \text{A})$ and $m(\mathbf{x}, s = \text{B})$, the “fair barycenter score” is

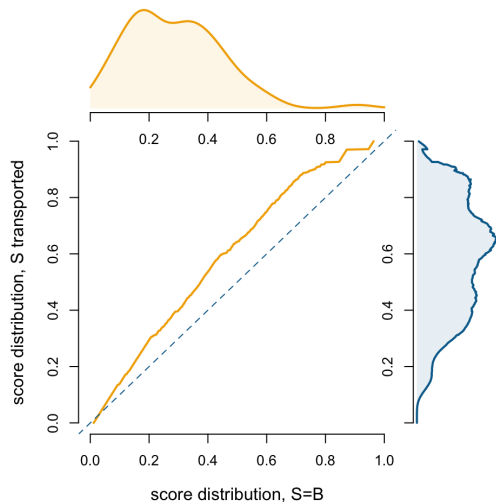
$$\begin{aligned} & m^*(\mathbf{x}, s = \text{A}) \\ &= \mathbb{P}[S = \text{A}] \cdot m(\mathbf{x}, s = \text{A}) \\ &+ \mathbb{P}[S = \text{B}] \cdot F_{\text{B}}^{-1} \circ F_{\text{A}}(m(\mathbf{x}, s = \text{A})) \end{aligned}$$



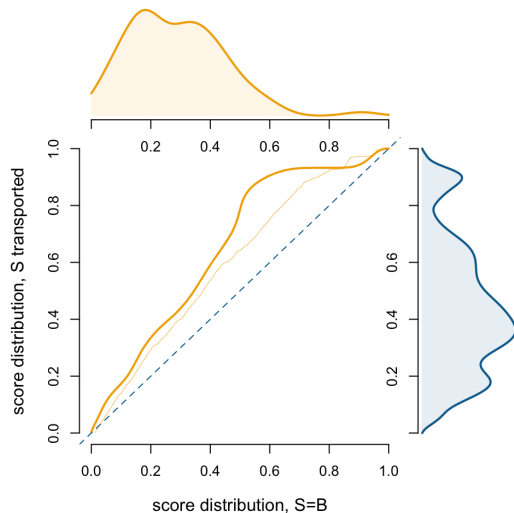
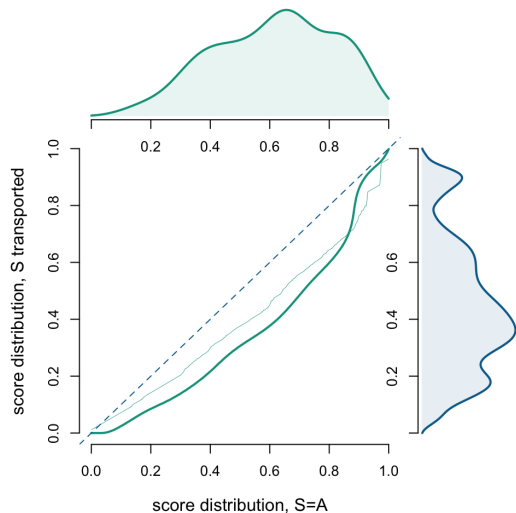
Application to toydata1

➤ Given scores $m(\mathbf{x}, s = \text{A})$ and $m(\mathbf{x}, s = \text{B})$, the “fair barycenter score” is

$$\begin{aligned} & m^*(\mathbf{x}, s = \text{B}) \\ &= \mathbb{P}[S = \text{A}] \cdot F_{\text{A}}^{-1} \circ F_{\text{B}}(m(\mathbf{x}, s = \text{B})) \\ &+ \mathbb{P}[S = \text{B}] \cdot m(\mathbf{x}, s = \text{B}) \end{aligned}$$



Application to toydata1

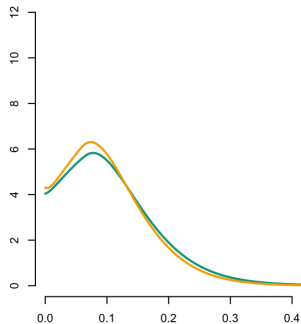


Application to toydata1

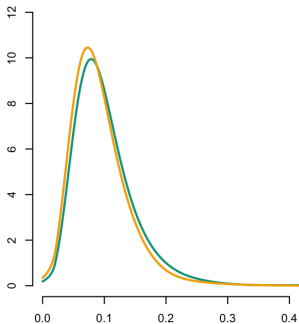
	x	s	\bar{y}	$\hat{m}(x, s)$	$\hat{m}(x)$	$\hat{m}_w^*(x)$	$\hat{m}_{jkl}^*(x)$
Alex	-1	A	0.475	0.250	0.219	0.154	0.094
Betty	-1	B	0.475	0.205	0.219	0.459	0.357
Ahmad	0	A	0.475	0.490	0.465	0.341	0.279
Brienne	0	B	0.475	0.426	0.465	0.719	0.692
Anthony	+1	A	0.475	0.734	0.730	0.571	0.521
Beatrix	+1	B	0.475	0.681	0.730	0.842	0.932

Application to FrenchMotor

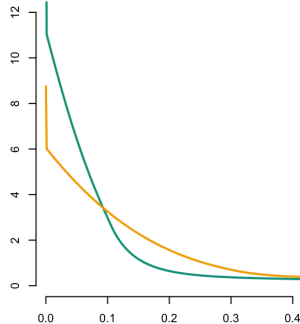
- If the two models are balanced, m^* is also balanced.
- Annual claim occurrence (motor insurance, Charpentier et al. (2023))
- Three models (plain GLM, GBM, Random Forest)



Plain logistic (GLM)



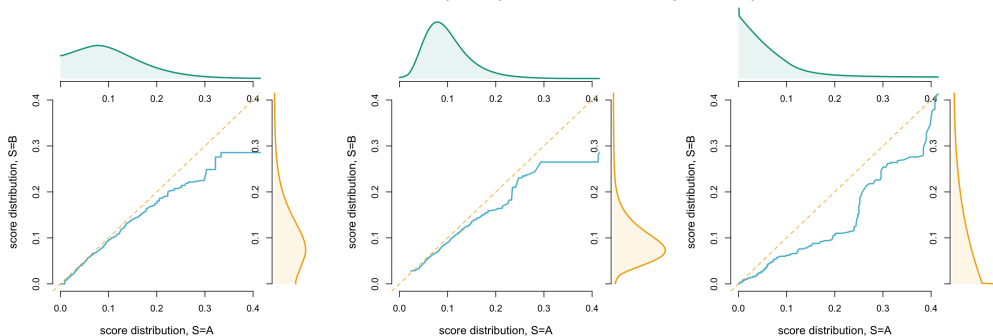
Boosting adaboost (GBM)



Random Forrest (RF)

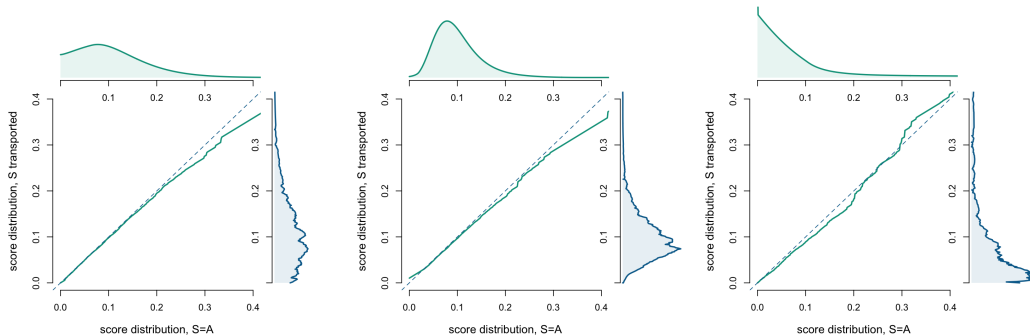
Application to FrenchMotor

- Predictions are different for men (= A) and women (S = B)



- since $W_2 \neq 0$ consider post processing mitigation

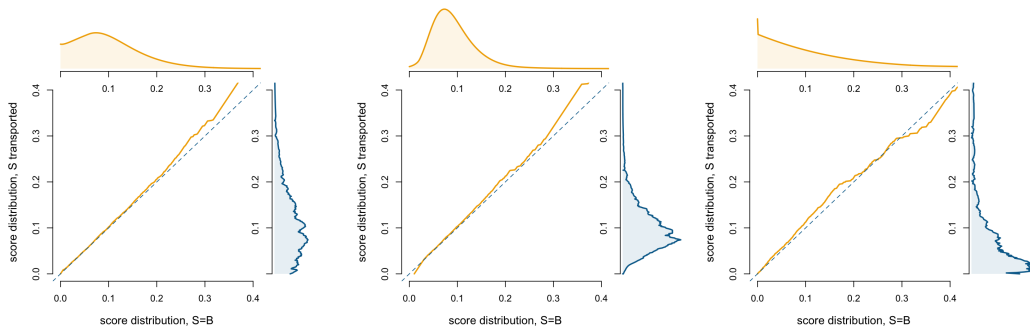
Application to FrenchMotor



➤ Given scores $m(\mathbf{x}, s = \mathbf{A})$ and $m(\mathbf{x}, s = \mathbf{B})$, the “fair barycenter score” is

$$m^*(\mathbf{x}, s = \mathbf{A}) = \mathbb{P}[S = \mathbf{A}] \cdot m(\mathbf{x}, s = \mathbf{A}) + \mathbb{P}[S = \mathbf{B}] \cdot F_{\mathbf{B}}^{-1} \circ F_{\mathbf{A}}(m(\mathbf{x}, s = \mathbf{A}))$$

Application to FrenchMotor

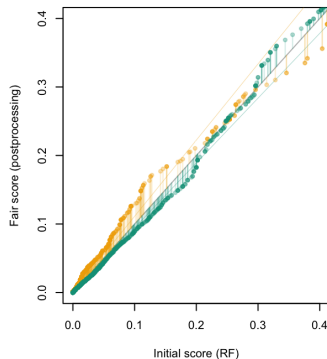
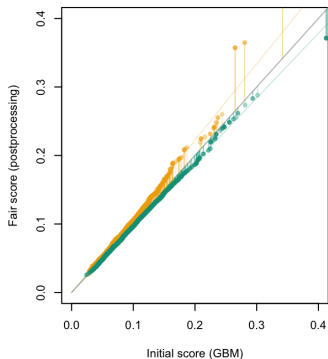
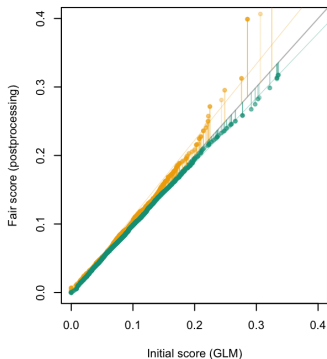


➤ Given scores $m(\mathbf{x}, s = \text{A})$ and $m(\mathbf{x}, s = \text{B})$, the “fair barycenter score” is

$$m^*(\mathbf{x}, s = \text{B}) = \mathbb{P}[S = \text{A}] \cdot F_{\text{A}}^{-1} \circ F_{\text{B}}(m(\mathbf{x}, s = \text{B})) + \mathbb{P}[S = \text{B}] \cdot m(\mathbf{x}, s = \text{B})$$

Application to FrenchMotor

- We can plot $\{(m(\mathbf{x}_i, \mathbf{A}), m^*(\mathbf{x}_i, \mathbf{A}))\}$ and $\{(m(\mathbf{x}_i, \mathbf{B}), m^*(\mathbf{x}_i, \mathbf{B}))\}$



Application to FrenchMotor

➤ Numerical values, for initial occurrence probability of 5%, 10% and 20%, we have

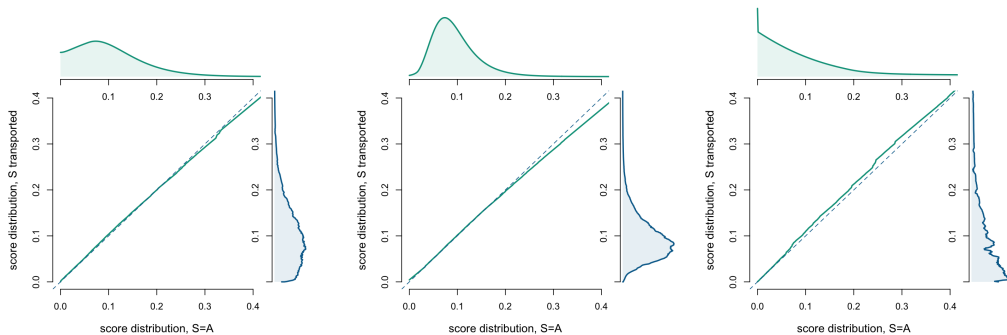
	A (men)				B (women)			
	×0.94	GLM	GBM	RF	×1.11	GLM	GBM	RF
$m(\mathbf{x}) = 5\%$	4.73%	4.94%	4.80%	4.42%	5.56%	5.16%	5.25%	6.15%
$m(\mathbf{x}) = 10\%$	9.46%	9.83%	9.66%	8.92%	11.12%	10.38%	10.49%	12.80%
$m(\mathbf{x}) = 20\%$	18.91%	19.50%	18.68%	18.26%	22.25%	20.77%	21.63%	21.12%

Application to FrenchMotor

We can do the same for discrimination against "old" drivers.

	A (younger < 65)				B (old > 65)			
	×1.01	GLM	GBM	RF	×0.94	GLM	GBM	RF
$m(\mathbf{x}) = 5\%$	5.05%	5.17%	5.10%	5.27%	4.71%	3.84%	3.84%	3.96%
$m(\mathbf{x}) = 10\%$	10.09%	10.37%	10.16%	11.00%	9.42%	7.81%	9.10%	6.88%
$m(\mathbf{x}) = 20\%$	20.19%	19.98%	19.65%	21.26%	18.85%	19.78%	23.79%	12.54%

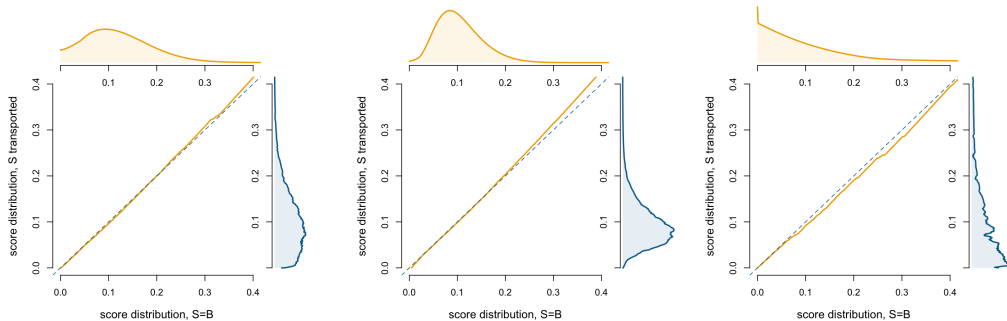
Application to FrenchMotor



➤ Given scores $m(\mathbf{x}, s = \mathbf{A})$ and $m(\mathbf{x}, s = \mathbf{B})$, the “fair barycenter score” is

$$m^*(\mathbf{x}, s = \mathbf{A}) = \mathbb{P}[S = \mathbf{A}] \cdot m(\mathbf{x}, s = \mathbf{A}) + \mathbb{P}[S = \mathbf{B}] \cdot F_{\mathbf{B}}^{-1} \circ F_{\mathbf{A}}(m(\mathbf{x}, s = \mathbf{A}))$$

Application to FrenchMotor

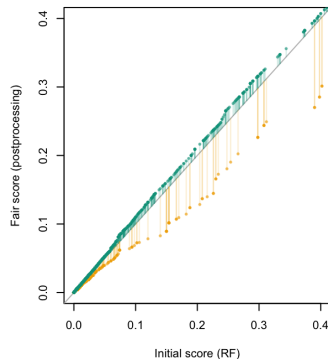
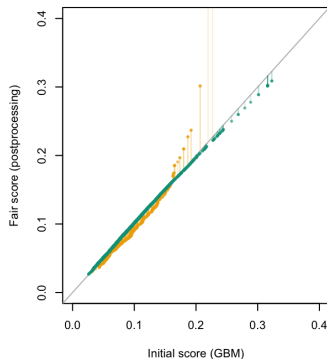
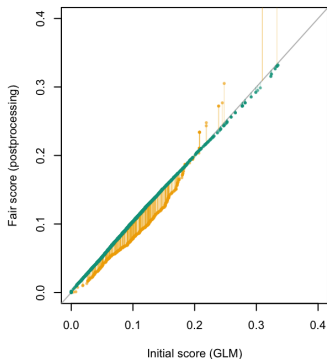


➤ Given scores $m(\mathbf{x}, s = \mathbf{A})$ and $m(\mathbf{x}, s = \mathbf{B})$, the “fair barycenter score” is

$$m^*(\mathbf{x}, s = \mathbf{B}) = \mathbb{P}[S = \mathbf{A}] \cdot F_{\mathbf{A}}^{-1} \circ F_{\mathbf{B}}(m(\mathbf{x}, s = \mathbf{B})) + \mathbb{P}[S = \mathbf{B}] \cdot m(\mathbf{x}, s = \mathbf{B})$$

Application to FrenchMotor

- We can plot $\{(m(\mathbf{x}_i, \mathbf{A}), m^*(\mathbf{x}_i, \mathbf{A}))\}$ and $\{(m(\mathbf{x}_i, \mathbf{B}), m^*(\mathbf{x}_i, \mathbf{B}))\}$



What's next ?

- While algorithmic fairness may result in greater short-term benefits, studies indicate that common fairness criteria may not promote improvement over time, as claimed by [Liu et al. \(2018\)](#).
- For example, requiring a bank to give out loans to individuals who are less likely to repay them ultimately impoverishes the individuals who end up defaulting as a result. This is an issue that algorithmic fairness cannot address on its own.



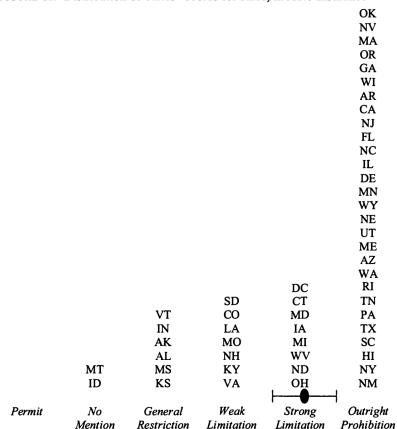
– Appendix –
Additional Results

Appendix: Sensitive attributes in insurance (in the U.S.)

← to come back From Avraham et al. (2013),

- › **Expressly Permit (-1)** - The state has a statute expressly or impliedly permitting insurers to take the characteristic into account.
- › **No Law on Point (0)** - The state laws are silent with respect to the particular characteristic.
- › **General Restriction (1)** - The state has a statute that generally prohibits "unfair discrimination," either across all lines of insurance or in some lines of insurance, but that statute does not provide any explanation as to what constitutes unfair discrimination and does not single out any particular trait for limitation.

FIGURE 1a. Distribution of States' Scores for Race, in Auto Insurance



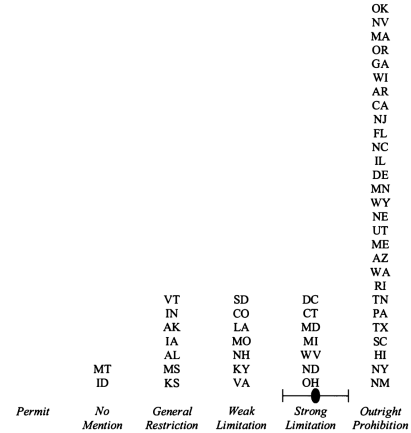
Appendix: Sensitive attributes in insurance (in the U.S.)

➤ **Characteristic-Specific Weak Limitation (2)** - The state has a statute that limits the use of a particular characteristic in either issuance, renewal, or cancellation.

➤ **Characteristic-Specific Strong Limitation (3)** - The state has a statute that prohibits the use of a particular characteristic when the policy is either issued, renewed, or cancelled, or the state has a statute that limits but does not completely prohibit the use of a particular characteristic in rate setting.

➤ **Characteristic-Specific Prohibition (4)** - The state has a statute that expressly prohibits insurers from taking into account a specific characteristic in setting rates.

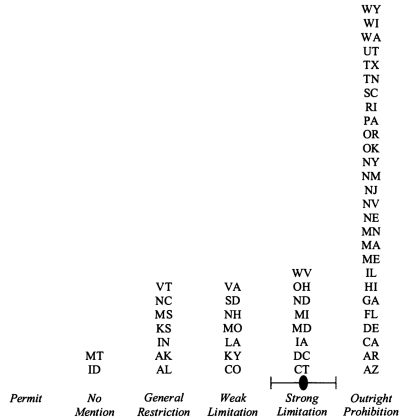
FIGURE 1b. Distribution of States' Scores for National Origin, in Auto Insurance



Appendix: Sensitive attributes in insurance (in the U.S.)

"Race, national origin, and religion have a special place in this country's history; and, as discussed above, discrimination on the basis of these three characteristics has been subject to stricter scrutiny in American law than have other characteristics," Avraham et al. (2013)

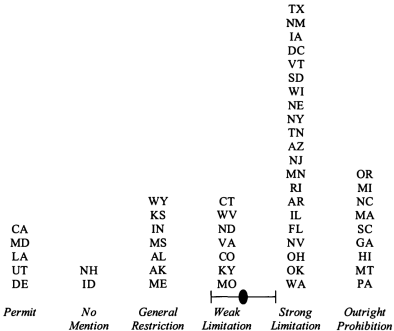
FIGURE 1c. Distribution of States' Scores for Religion, in Auto Insurance



Appendix: Sensitive attributes in insurance (in the U.S.)

”Gender-based discrimination in insurance has long been controversial. And differential treatment on the basis of gender is, of course, in many contexts widely considered unacceptable or illegal. Nevertheless, there does not seem to be the same level of agreement-as there is for race, religion, and national origin-that drawing gender-based distinctions is always wrong. Federal constitutional law treats gender as only a quasi-suspect classification; as a result, laws that discriminate on the basis of gender are subject to an intermediate level of scrutiny.” Avraham et al. (2013)

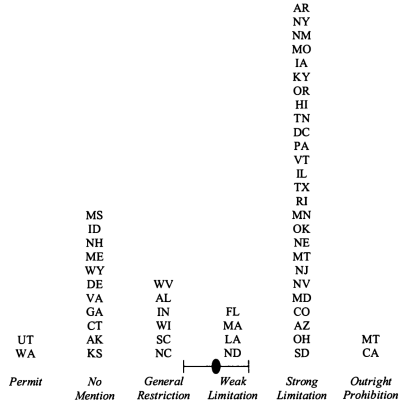
FIGURE 3a. Distribution of States' Scores for Gender, in Auto Insurance



Appendix: Sensitive attributes in insurance (in the U.S.)

FIGURE 3c. Distribution of States' Scores for Gender, in Disability Insurance

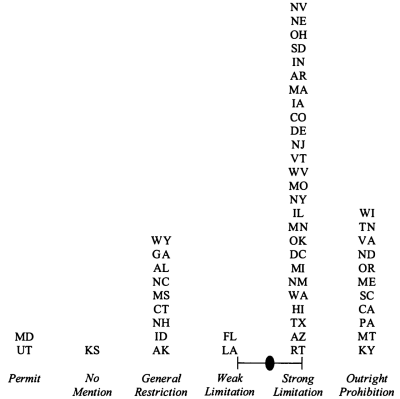
"With respect to life insurance, we predict that the laws regulating gender discrimination will be on average relatively weak, since adverse selection in the life insurance market is especially problematic." Avraham et al. (2013)



Appendix: Sensitive attributes in insurance (in the U.S.)

FIGURE 3d. Distribution of States' Scores for Gender, in Property/Casualty Insurance

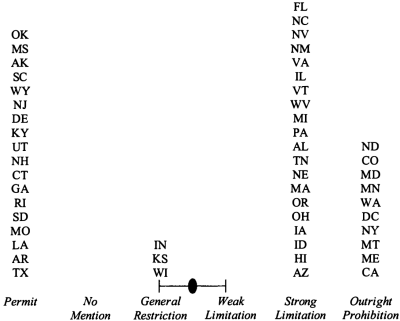
"Regarding property/casualty insurance, as there seems to be no conceivable correlation between those risks and gender, we predict either states will cluster around no regulation, or, alternatively, states will cluster around forbidding the use of gender in property/casualty insurance on symbolic or expressive grounds." Avraham et al. (2013)



Appendix: Sensitive attributes in insurance (in the U.S.)

"The gender discrimination will be more strictly regulated on average for health insurance (where gender-rated policies often result in higher premiums for women) than for auto insurance (where gender-rated policies result in higher premiums for men)." Avraham et al. (2013)

FIGURE 3e. Distribution of States' Scores for Gender, in Health Insurance

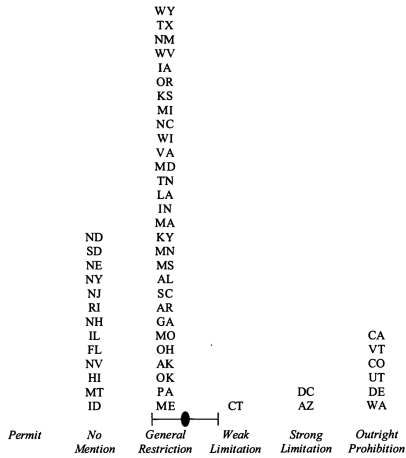


Appendix: Sensitive attributes in insurance (in the U.S.)

FIGURE 4a. Distribution of States' Scores for Sexual Orientation, in Auto Insurance

"Unlike with race, national origin, religion, and gender, legal classifications on the basis of an individual's sexual orientation have not clearly been identified by the Supreme Court as deserving special scrutiny. In addition, unlike race, national origin, and gender, there are no federal laws forbidding discrimination on the basis of sexual orientation in employment."

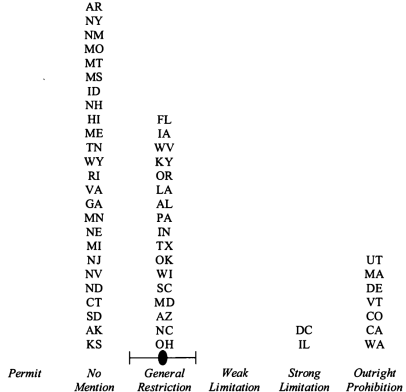
Avraham et al. (2013)



Appendix: Sensitive attributes in insurance (in the U.S.)

"However, there are state laws that forbid discrimination on the basis of sexual orientation, and some lower courts have held that sexual orientation should be a suspect or quasi-suspect characterisation." Avraham et al. (2013)

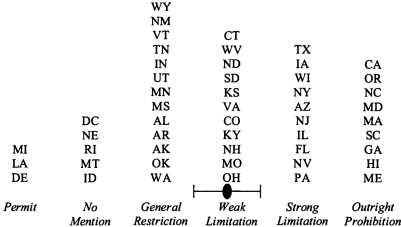
FIGURE 4c. Distribution of States' Scores for Sexual Orientation, in Disability Insurance



Appendix: Sensitive attributes in insurance (in the U.S.)

"We expect that *age* will have the lowest average regulatory score of all the risk characteristics we are studying. First, age is not a suspect classification, at least not by constitutional standards. Second, age tends to correlate causally with several important areas of risk (mortality, health, and perhaps disability risks), thereby increasing the perceived fairness of rating on that basis." Avraham et al. (2013)

FIGURE 5a. Distribution of States' Scores for Age, in Auto Insurance

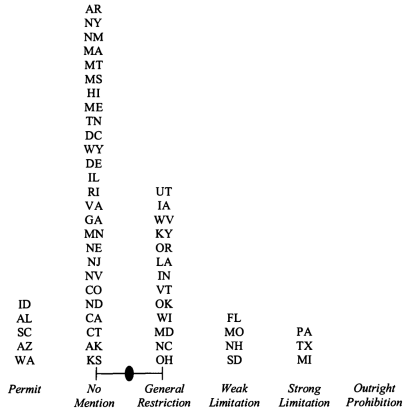


Appendix: Sensitive attributes in insurance (in the U.S.)

”Third, age can present serious adverse selection problems for insurers if they are forbidden from taking it into account, since individual insureds know their own age and the associated risks. Fourth, social solidarity arguments with respect to age are relatively weak, since individuals can spread risk over their lifetime through various income smoothing products.”

Avraham et al. (2013)

FIGURE 5c. Distribution of States’ Scores for Age, in Disability Insurance



Appendix: Sensitive attributes in insurance (in the U.S.)

Avraham et al. (2013) suggested to visualize the distribution of scores (Expressly Permit (-1) / No Law on Point (0) / General Restriction (1) / ... / Characteristic-Specific Prohibition (4))

FIGURE 6. Distribution of States' Scores for Age, by Line of Insurance

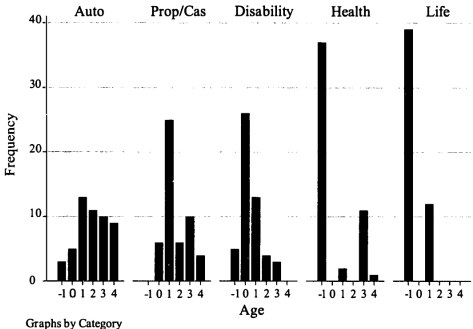
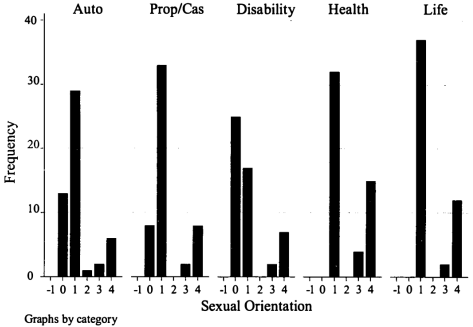


FIGURE 7. Distribution of States' Scores for Sexual Orientation, by Line of Insurance



Appendix: Sensitive attributes in insurance (in the U.S.)

FIGURE 8. Distribution of States' Scores for Zip Code, by Line of Insurance

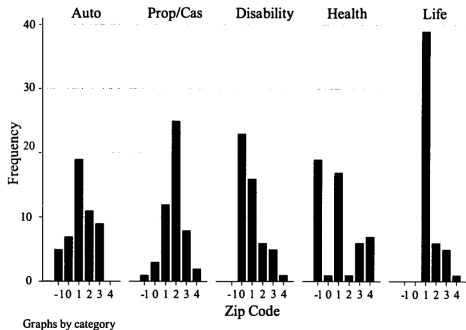
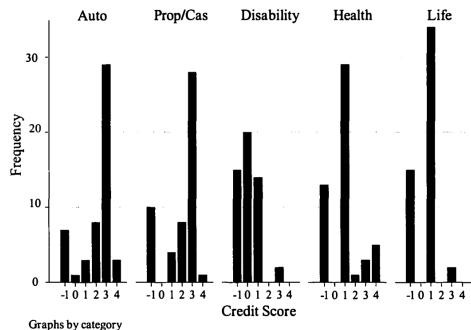


FIGURE 9. Distribution of States' Scores for Credit Score, by Line of Insurance



“Credit score and zip code are not, by themselves, socially suspect characteristics. However, some commentators have argued that credit score and zip code are used by auto and home insurers as proxies for potentially socially suspect characteristics.”

Appendix: L^k and L^2

← to come back

Note: A Hilbert space is an abstract vector space possessing the structure of an inner product.

If \mathcal{H} is finite, $\mathcal{H} = \{h_1 \cdots, h_d\}$, $\langle x, y \rangle_{\mathcal{H}}$ takes value $K_{i,j}$ if $x = h_i$ and $y = h_j$. Let $\mathbf{K} = [K_{i,j}]$.

\mathbf{K} is a symmetric $d \times d$ matrix, $\mathbf{K} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{\top}$ for some orthogonal matrix \mathbf{V} where columns are eigenvectors, and $\mathbf{\Lambda} = \text{diag}[\lambda_i]$ (positive values). Let

$$\Phi(x) = (\sqrt{\lambda_1}V_{1,i}, \sqrt{\lambda_2}V_{2,i}, \cdots, \sqrt{\lambda_d}V_{d,i}) \text{ if } x = h_i$$

Note that

$$K_{i,j} = [\mathbf{K}]_{i,j} = [\mathbf{V}\mathbf{\Lambda}\mathbf{V}^{\top}]_{i,j} = \sum_{l=1}^d \lambda_l V_{l,i} V_{l,j} = \langle \Phi(h_i), \Phi(h_j) \rangle$$

Appendix: L^k and L^2

Matrix \mathbf{K} defines an inner product, it is called a **kernel**. It is symmetric, associated with a positive semi-definite matrix.

Then $K(u, u) \geq 0$ and $K(u, v) \leq \sqrt{K(u, u) \cdot K(v, v)}$.

Let $\Phi : u \mapsto K(\cdot, u)$, then $K(x, y) = \langle \Phi(x), \Phi(y) \rangle$

Proposition 8.8: ℓ_p norms for random variables

- ▶ $\|X\|_k \geq 0$,
- ▶ $\|X\|_k = 0$ if and only if $X = 0$ (a.s., i.e. $\mathbb{P}[X \neq 0] = 0$),
- ▶ $\|X + Y\|_k \leq \|X\|_k + \|Y\|_k$ (Minkowski inequality)

Proof for Minkowski inequality, $(x, y) \mapsto g(x, y) = (x^{1/k} + y^{1/k})^k$ is concave, and from Jensen inequality, for positive variables U and V (\mathbb{R}_+^2 is a convex set),

$$\mathbb{E} \left[(U^{1/k} + V^{1/k})^k \right] \leq \left([\mathbb{E}(U)]^{1/k} + [\mathbb{E}(V)]^{1/k} \right)^k$$

(then use $U = |X|^k$ and $V = |Y|^k$)

Appendix: L^k and L^2

Definition 8.7: L^k (Lebesgue space)

For $k \in [1, \infty)$, $L^k = \{X \in \mathcal{V} : |X|^k < \infty\}$.

Proposition 8.9: Lyapunov inequality

If $1 \leq j < k < \infty$, $\|X\|_j \leq \|X\|_k$, or $L^k \subset L^j$.

$d_k(X, Y) = \|X - Y\|_k = \left[\mathbb{E}(|X - Y|^k) \right]^{1/k}$ is a distance on \mathcal{V} .

Proposition 8.10: Inner product on L^2

$\langle X, Y \rangle = \mathbb{E}(XY)$ is an inner product on L^2 ,

- ▶ $\langle X, Y \rangle = \langle Y, X \rangle$ (symmetric property),
- ▶ $\langle X, X \rangle \geq 0$ and $\langle X, X \rangle = 0$ if and only if $X = 0$ a.s. (positive property)
- ▶ $\langle aX, Y \rangle = a\langle X, Y \rangle$ (scaling property),
- ▶ $\langle X + Y, Z \rangle = \langle X, Z \rangle + \langle Y, Z \rangle$, the additive property.

Proposition 8.11: Hölder inequality

If $1 \leq j < k < \infty$ with $j^{-1} + k^{-1} = 1$, $X \in L^j$ and $Y \in L^k$, then

$$\langle |X|, |Y| \rangle \leq \|X\|_j \cdot \|Y\|_k.$$

Appendix: L^k and L^2

Observe that $\text{Cov}[X, Y] = \langle X - \mathbb{E}(X), Y - \mathbb{E}(Y) \rangle$.

When $j = k = 2$, Hölder inequality is **Cauchy-Schwarz** inequality

$$\underbrace{\mathbb{E}(|X - \mathbb{E}(X)| \cdot |Y - \mathbb{E}(Y)|)}_{\text{Cov}[|X|, |Y|]} \leq \underbrace{\sqrt{\mathbb{E}([X - \mathbb{E}(X)]^2)}}_{\text{sd}[X]} \cdot \underbrace{\sqrt{\mathbb{E}([Y - \mathbb{E}(Y)]^2)}}_{\text{sd}[Y]}$$

Suppose that \mathcal{H} is a subspace of L^2 . Let $Y \in L^2$, then the projection of Y onto \mathcal{H} (if it exists) is the vector $X \in \mathcal{H}$ such that $\langle X - Y, Z \rangle = 0, \forall Z \in \mathcal{H}$. It will be denoted $\mathbb{E}[Y|\mathcal{H}]$.

Appendix: L^k and L^2

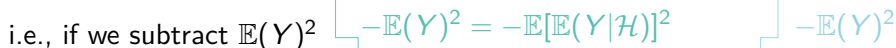
Let $A = \mathbb{E}[Y|\mathcal{H}]$, $B = Y - \mathbb{E}[Y|\mathcal{H}]$ and $C = A + B = Y$. Then, since $\langle A, B \rangle = 0$,

$$\|A\|_2^2 + \|B\|_2^2 = \|C\|_2^2,$$

from [Pythagorean theorem](#), or,

$$\mathbb{E}[\mathbb{E}(Y|\mathcal{H})^2] + \mathbb{E}[(Y - \mathbb{E}[Y|\mathcal{H}])^2] = \mathbb{E}[Y^2]$$

i.e., if we subtract $\mathbb{E}(Y)^2$



$$\text{var}(\mathbb{E}[Y | \mathcal{H}]) + \mathbb{E}[\text{var}(Y | \mathcal{H})] = \text{var}(Y)$$

$$= \text{var}(Y - \mathbb{E}[Y | \mathcal{H}])$$


Appendix: Optimisation Issues

← to come back

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if $\text{dom}(f)$ is convex and

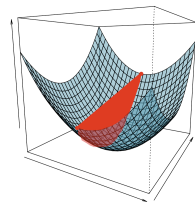
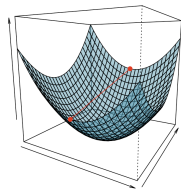
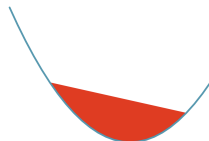
$$f(t\mathbf{x} + (1-t)\mathbf{y}) \leq tf(\mathbf{x}) + (1-t)f(\mathbf{y})$$

for all $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$.

Proposition 8.12

If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable, then f is convex if $\text{dom}(f)$ is convex and for all $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$



Appendix: Optimisation Issues

If f is convex and non-differentiable, for all \mathbf{x} , there is \mathbf{g} such that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x})$$

for all $\mathbf{y} \in \text{dom}(f)$.

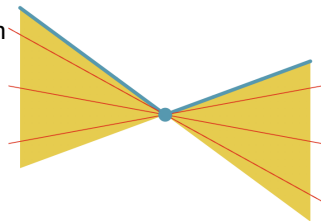
\mathbf{g} is called subgradient at point \mathbf{x} .

If f differentiable at \mathbf{x} , \mathbf{g} is unique and $\mathbf{g} = \nabla f(\mathbf{x})$

The set of subgradients of a convex function f is the subdifferential,

$$\partial f(\mathbf{x}) = \{\mathbf{g} \in \mathbb{R}^n : \mathbf{g} \text{ is a subgradient at } \mathbf{x}\}$$

Note that $\partial f(\mathbf{x})$ is a convex set, and if f is differentiable at point \mathbf{x} , $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$



Appendix: Optimisation Issues

Consider the case of the median ($\alpha = 1/2$). Given $\{y_1, \dots, y_n\}$, we want to solve

$$\min_{\mu} \left\{ \sum_{i=1}^n |y_i - \mu| \right\}$$

This problem is equivalent to

$$\min_{\mu, \mathbf{a}, \mathbf{b}} \left\{ \sum_{i=1}^n a_i + b_i \right\} \text{ s.t. } \begin{cases} a_i, b_i \geq 0 \\ y_i - \mu = a_i - b_i \end{cases}, \quad \forall i = 1, \dots, n$$

More generally, for a $\alpha \in (0, 1)$, solve

$$\min_{\mu, \mathbf{a}, \mathbf{b}} \left\{ \sum_{i=1}^n \tau a_i + (1 - \tau) b_i \right\} \text{ s.t. } \begin{cases} a_i, b_i \geq 0 \\ y_i - \mu = a_i - b_i \end{cases}, \quad \forall i = 1, \dots, n$$

Appendix: Optimisation Issues

And for the quantile regression is

$$\min_{\gamma, \mathbf{a}, \mathbf{b}} \left\{ \sum_{i=1}^n \tau a_i + (1 - \tau) b_i \right\} \text{ s.t. } \begin{cases} a_i, b_i \geq 0 \\ y_i - \mathbf{x}_i^\top \gamma = a_i - b_i \end{cases}$$

Introduce **slack variables** to turn inequality constraints into equality constraints with positive unknowns : any inequality $a_1 x_1 + \dots + a_n x_n \leq c$ can be replaced by $a_1 x_1 + \dots + a_n x_n + u = c$ with $u \geq 0$.

A linear programming problem written in a **standard form** is

$$\text{Primal problem: } \min_{\mathbf{x} \in \mathbb{R}^n} \{ \mathbf{c}^\top \mathbf{x} \} \text{ s.t. } \begin{cases} A\mathbf{x} = \mathbf{b}, \\ \mathbf{x} \geq \mathbf{0}. \end{cases}$$

for some $m \times n$ A matrix, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{c} \in \mathbb{R}^n$, and

$$\text{Dual problem: } \max_{\mathbf{y} \in \mathbb{R}^m} \{ \mathbf{y}^\top \mathbf{b} \} \text{ s.t. } \mathbf{y}^\top A \leq \mathbf{c}^\top.$$

Appendix: Optimisation Issues

The minimal cost and the maximal income coincide, i.e., the two problems are equivalent.

Theorem 8.3: Strong duality theorem, Dantzig and Thapa (1997)

The primal problem has a nondegenerate solution \mathbf{x}^* if and only if the dual problem has a nondegenerate solution \mathbf{y}^* . And in this case $\mathbf{y}^{*\top} \mathbf{b} = \mathbf{c}^\top \mathbf{x}^*$.

In quantile regression,

$$\text{Primal problem: } \min_{\beta, \mathbf{u}, \mathbf{v}} \left\{ \tau \mathbf{1}^\top \mathbf{u} + (1 - \tau) \mathbf{1}^\top \mathbf{v} \right\} \text{ s.t. } \mathbf{y} = \mathbf{X}\beta + \mathbf{u} - \mathbf{v}, \text{ with } \mathbf{u}, \mathbf{v} \in \mathbb{R}_+^n$$

and

$$\text{Dual problem: } \max_{\mathbf{d}} \left\{ \mathbf{y}^\top \mathbf{d} \right\} \text{ s.t. } \mathbf{X}^\top \mathbf{d} = (1 - \tau) \mathbf{X}^\top \mathbf{1} \text{ with } \mathbf{d} \in [0, 1]^n$$

Appendix: Optimisation Issues

For the median, Heuristically, the idea is to write $y_i = \mu + \varepsilon_i$, and then define a_i 's and b_i 's so that $\varepsilon_i = a_i - b_i$ and $|\varepsilon_i| = a_i + b_i$, i.e.

$$\begin{cases} a_i = (\varepsilon_i)_+ = \max\{0, \varepsilon_i\} = |\varepsilon_i| \cdot \mathbf{1}_{\varepsilon_i > 0} \\ b_i = (-\varepsilon_i)_+ = \max\{0, -\varepsilon_i\} = |\varepsilon_i| \cdot \mathbf{1}_{\varepsilon_i < 0} \end{cases}$$

Thus, set $\mathbf{z} = (\mu^+; \mu^-; \mathbf{a}, \mathbf{b})^\top \in \mathbb{R}_+^{2n+2}$, and then write the constraint as $\mathbf{Az} = \mathbf{b}$ with $\mathbf{b} = \mathbf{y}$ and $\mathbf{A} = [\mathbf{1}_n; -\mathbf{1}_n; \mathbb{I}_n; -\mathbb{I}_n]$

For the objective function $\mathbf{c} = (\mathbf{0}, \mathbf{1}_n, \mathbf{1}_n)^\top \in \mathbb{R}_+^{2n+2}$ and our program is

$$\min_{\mathbf{z}} \{ \mathbf{c}^\top \mathbf{z} \} \text{ s.t. } \mathbf{Az} = \mathbf{b}, \mathbf{z} \geq \mathbf{0}.$$

```
1 > n = 101
2 > set.seed(1)
3 > y = rlnorm(n)
4 > median(y)
5 [1] 1.077415
6 > library(lpSolve)
```

```
1 > X = rep(1, n)
2 > A = cbind(X, -X, diag(n), -diag(n))
3 > b = y
4 > c = c(rep(0, 2), rep(1, n), rep(1, n))
5 > r = lp("min", c, A, rep("=", n), b)$solution
6 [1] 1.077415
```

Appendix: Optimisation Issues

More generally, if the quantile of order τ is a solution of $\tau \in (0, 1)$,

$$\min_q \left\{ \sum_{i=1}^n \max\{\tau(y_i - \mu), (1 - \tau)(y_i - \mu)\} \right\}$$

The linear program is now

$$\min_{q^+, q^-, \mathbf{a}, \mathbf{b}} \left\{ \sum_{i=1}^n \tau a_i + (1 - \tau) b_i \right\}$$

with $a_i, b_i, q^+, q^- \geq 0$ and $y_i = q^+ - q^- + a_i - b_i, \forall i = 1, \dots, n$.

```
1 > tau = .3
2 > quantile(y, tau)
3     30%
4 0.6741586
```

```
1 > c = c(rep(0, 2), tau*rep(1, n), (1-tau)*rep(1, n))
2 > r = lp("min", c, A, rep("=", n), b)$solution
3 [1] 0.6741586
```

Appendix: Optimisation Issues

In a regression, we use $\mathbf{x}_i^\top \boldsymbol{\beta}$ instead of $\boldsymbol{\mu}$. The linear program is

$$\min_{\boldsymbol{\beta}^+, \boldsymbol{\beta}^-, \mathbf{a}, \mathbf{b}} \left\{ \sum_{i=1}^n \tau a_i + (1 - \tau) b_i \right\}$$

with $a_i, b_i \geq 0$ and $y_i = \mathbf{x}_i^\top [\boldsymbol{\beta}^+ - \boldsymbol{\beta}^-] + a_i - b_i, \forall i = 1, \dots, n$ and $\beta_j^+, \beta_j^- \geq 0$
 $\forall j = 0, \dots, k$.

Appendix: Optimisation Issues

Let \mathbf{Q} be a positive semidefinite matrix, the following problem corresponds to **quadratic programming**,

$$\begin{aligned} \min_{\mathbf{x}} \left\{ \mathbf{c}^\top \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} \right\} \\ \text{subject to } \begin{cases} \mathbf{D} \mathbf{x} \geq \mathbf{d} \\ \mathbf{A} \mathbf{x} = \mathbf{b} \end{cases} \end{aligned}$$

More generally, consider a convex problem

$$\begin{aligned} \min_{\mathbf{x} \in D} \{ f(\mathbf{x}) \} \\ \text{subject to } \begin{cases} g_i(\mathbf{x}) \leq 0, \forall i = 1, \dots, m \\ \mathbf{A} \mathbf{x} = \mathbf{b} \end{cases} \end{aligned}$$

for some convex function f, g_1, \dots, g_m

Definition 8.8: Dual vector space

Given a vector space E , its dual E^* is the set of all linear forms on E . Define $\langle \cdot, \cdot \rangle : E^* \times E \rightarrow \mathbb{R}$ as $\langle p, x \rangle = p(x)$.

Definition 8.9: Convex set

$C \subset E$ is a convex set if for all $x, y \in C$, $[x, y] = \{(1-t)x + ty, t \in [0, 1]\} \subset C$.

Definition 8.10: Convex function

$f : E \rightarrow \mathbb{R}$ is a convex function if

$$f((1-t)x + ty) \leq (1-t)f(x) + tf(y), \quad \forall x, y \in E, t \in [0, 1].$$

More generally (if f is not proper, i.e. taking values in $[-\infty, +\infty]$), f is convex if its epigraph is convex, where $\text{epi}(f) = \{(x, y) : y \geq f(x)\} \subset E \times \mathbb{R}$.

Proposition 8.13

Let $f_i : E \rightarrow \mathbb{R}$ be convex functions, then $f^+(x) = \sup\{f_i(x)\}$ is a convex function.

Proposition 8.14

If $f : E \rightarrow \mathbb{R}$ is convex, for every $x_0 \in E$, there exists an affine function $\ell_{x_0}(x) = a_0x + y$ such that $\ell_{x_0}(x_0) = f(x_0)$ and f lies above ℓ_{x_0} .

Proof:

Proposition 8.15: Hahn-Banach

A convex lower semi continuous (for every $x_n \rightarrow x$, $\lim f(x_n) \geq f(x)$) function f is equal to the supremum of affine functions.

Convexity

Proof: Let \mathcal{L} denote the set of the affine functions from Proposition 8.14, $\mathcal{L} = \{\ell_{x_0}, x_0 \in E\}$. Then, for all $x \in E$,

$$\begin{cases} \sup_{\ell \in \mathcal{L}} \{\ell(x)\} \geq \ell_x(x) = f(x) \text{ because } \ell_x \text{ passes through } (x, f(x)), \\ \sup_{\ell \in \mathcal{L}} \{\ell(x)\} \leq f(x) \text{ because all } \ell \text{ lies below } f \end{cases}$$

therefore $f(x) = \sup_{\ell \in \mathcal{L}} \{\ell(x)\}$.

Definition 8.11: (Fenchel) convex conjugate

If $f : E \rightarrow \mathbb{R}$ is a function, its convex conjugate $f^* : E^* \rightarrow \mathbb{R}$ is defined as

$$f^*(p) = \sup_{x \in E} \langle p, x \rangle - f(x).$$

Convexity

Let E denote some nonempty subset of \mathbb{R}^k , and define the indicator function of E ,

$$\mathbf{1}_E(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{x} \notin E \\ +\infty & \text{if } \mathbf{x} \in E \end{cases}$$

Then $\mathbf{1}_E^*(\mathbf{s}) = \sup_{\mathbf{x} \in E} \{\mathbf{s}^\top \mathbf{x}\} = \sup_{\mathbf{x} \in E} \{\langle \mathbf{s}, \mathbf{x} \rangle\}$ which is the support function of E .

Definition 8.12: Subgradient

If $f : E \rightarrow \mathbb{R}$ is a convex function, the subgradient of f at point $x \in E$ is the set of elements in E^* defined as

$$\partial f(x) = \{p \in E^* : f(y) \geq f(x) + \langle p, y - x \rangle \text{ for all } y \in E\}.$$

Convexity

Therefore, with a little abuse of notations,

$$\langle \partial f(x) - \partial f(y), x - y \rangle \geq 0$$

Proposition 8.16: Legendre-Fenchel identity

If $f : E \rightarrow \mathbb{R}$ is a convex function,

$$p \in \partial f(x) \iff x \in \partial f^*(p) \iff f^*(p) + f(x) = \langle p, x \rangle.$$

If $f : \mathbb{R}^k \rightarrow \mathbb{R}$ is convex, twice differentiable, then $\nabla f^*(p) = [\nabla f]^{-1}(p)$, or $\nabla f^* \circ \nabla f$ and $\nabla f \circ \nabla f^*$ are the identity (on E and E^* respectively).

Let X be a random variable with c.d.f. F_X and quantile function Q_X . The Fenchel-Legendre transform of

$$\Psi(x) = \mathbb{E}[(x - X)_+] = \int_{-\infty}^x F_X(z) dz$$

Convexity

is

$$\Psi^*(y) = \sup_{x \in \mathbb{R}} \{xy - \Psi(x)\} = \int_0^y Q_X(t) dt$$

on $[0, 1]$.

Indeed, from Fubini,

$$\Psi(x) = \int_{-\infty}^{\infty} x \mathbb{P}(X \leq z) dz = \int_{-\infty}^{\infty} x \mathbb{E}(\mathbf{1}_{X \leq z}) dz = \mathbb{E} \left(\int_{-\infty}^{\infty} x \mathbf{1}_{X \leq z} dz \right)$$

i.e.

$$\Psi(x) = \mathbb{E}([x - X]_+) = \int_0^1 [x - Q_X(t)]_+ dt$$

Observe that

$$\Psi^*(1) = \sup_{x \in \mathbb{R}} \{x - \Psi(x)\} = \lim_{x \uparrow \infty} \int_0^1 [x - (x - Q_X(t))_+] dt = \int_0^1 Q_X(t) dt$$

Convexity

and $\Psi^*(0) = 0$. Now, the proof of the result when $y \in (0, 1)$ can be obtained since

$$\frac{\partial_x y - \Psi(x)}{\partial x} = y - F_X(x)$$

The optimum is then obtained when $y = F_X(x)$, or $x = Q_Y(y)$.

Definition 8.13: Proximal Operator (Moreau-Yosida)

If $f : E \rightarrow \mathbb{R}$ is a convex function, the proximal operator is defined as

$$\text{proximal}_{\epsilon, f}(x) = \underset{y}{\operatorname{argmin}} \left\{ f(y) + \frac{1}{2\epsilon} \|x - y\|^2 \right\} = (\operatorname{Id} + \epsilon \partial f)^{-1} x$$

Note that $\text{proximal}_{\epsilon, f}(x)$ is uniquely defined.

Convex Optimization Problem

$$\min_{\mathbf{x}} \{f(\mathbf{x})\}$$

with f convex, and differentiable.

Algorithm 1: Gradient Descent

- 1 initialization : $\mathbf{x}^{(0)}$;
- 2 **for** $t=1,2,\dots$ **do**
- 3 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} - h_t \nabla f(\mathbf{x}^{(t-1)})$

Heuristics: Taylor expansion

$$f(\mathbf{y}) \sim f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2h} \|\mathbf{y} - \mathbf{x}\|^2$$

Then one can prove that $|f(\mathbf{x}^{(t)}) - f^*| \leq O(t^{-1})$.

Convex Optimization Problem

$$\min_{\mathbf{x}} \{f(\mathbf{x})\}$$

with f convex, and differentiable.

Algorithm 2: Accelerated Gradient Descent (Nesterov)

- 1 initialization : $\mathbf{x}^{(0)}, \mathbf{x}^{(-1)}$;
- 2 **for** $t=1,2,\dots$ **do**
- 3 $\mathbf{y}^{(t-1)} \leftarrow \mathbf{x}^{(t-1)} + \frac{t-1}{t+2}(\mathbf{x}^{(t-1)} - \mathbf{x}^{(t-2)});$
- 4 $\mathbf{x}^{(t)} \leftarrow \mathbf{y}^{(t-1)} - h_t \nabla f(\mathbf{y}^{(t-1)});$

Then one can prove that $|f(\mathbf{x}^{(t)}) - f^*| \leq O(t^{-2})$.

From Gradient Descent to Newton's Method

$$\min_{\mathbf{x}} \{f(\mathbf{x})\}$$

with f convex, twice differentiable.

Algorithm 3: Newton's Method

- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 $\mathbf{H}_t \leftarrow \nabla^2 f(\mathbf{x}^{(t-1)})$;
 - 4 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} - \mathbf{H}_t^{-1} \nabla f(\mathbf{x}^{(t-1)})$
-

Use a better quadratic approximation in Taylor expansion $-\frac{1}{h}\mathbb{I} \rightarrow H$,

$$f(\mathbf{y}) \sim f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^\top H(\mathbf{y} - \mathbf{x})$$

Convex Optimization Problem

$$\min_{\mathbf{x}} \{f(\mathbf{x})\}$$

with f convex, but non-differentiable.

Algorithm 4: Subgradient 'Descent'

- 1 initialization : $\mathbf{x}^{(0)}$;
- 2 **for** $t=1,2,\dots$ **do**
- 3 $\mathbf{g}^{(t-1)} \in \partial f(\mathbf{x}^{(t-1)})$;
- 4 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} - h_t \mathbf{g}^{(t-1)}$

Note that it is not necessarily a descent, so pick

$$\mathbf{x}^* = \operatorname{argmin}\{f(\mathbf{x}^{(0)}), f(\mathbf{x}^{(1)}), f(\mathbf{x}^{(2)}), \dots\}$$

Convex Optimization Problem

$$\min_{\mathbf{x}} \{f(\mathbf{x})\} \text{ or } \min_{\mathbf{x}} \{f_1(\mathbf{x}) + f_2(\mathbf{x})\}$$

with f_1 and f_2 convex, but f_2 non-differentiable.

Algorithm 5: Proximal Gradient 'Descent'

- 1 initialization : $\mathbf{x}^{(0)}$;
- 2 **for** $t=1,2,\dots$ **do**
- 3 $\gamma_h(\mathbf{x}) = \frac{1}{h}(\mathbf{x} - \text{proximal}_{h,f_2}(\mathbf{x} - h\nabla f_1(\mathbf{x})));$
- 4 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} - h_t \gamma_{h_{t-1}}(\mathbf{x}^{(t-1)}) = \text{proximal}_{h,f_2}(\mathbf{x}^{(t-1)} - h\nabla f_1(\mathbf{x}^{(t-1)}))$

γ_h is the "generalized gradient of f ".

The trick is that $\text{proximal}_{h,f_2}(\cdot)$ usually has a closed form in most applications.

Convex Optimization Problem

E.g. consider the LASSO objective, $f(\beta) = \underbrace{\frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2}_{=f_1(\beta)} + \underbrace{\lambda \|\beta\|_{\ell_1}}_{=f_2(\beta)}$. The proximal

mapping is the soft-thresholding operator

$$\text{proximal}_h(\beta) = \underset{\mathbf{z}}{\text{argmin}} \left\{ \frac{1}{2h} \|\beta - \mathbf{z}\|_2^2 + \lambda \|\mathbf{z}\|_{\ell_1} \right\} = S_{\lambda h}(\beta)$$

where

$$S_{\lambda h}(\beta_j) = \begin{cases} \beta_j - \lambda & \text{if } \beta_j > \lambda \\ \beta_j + \lambda & \text{if } \beta_j < -\lambda \\ 0 & \text{otherwise.} \end{cases}$$

Hence

$$\begin{cases} \nabla f_1(\beta) = -\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta) \\ \gamma_h(\beta) = S_{\lambda h}(\beta + h \cdot \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta)) \end{cases}$$

see [Beck and Teboulle \(2009\)](#), "iterative soft-thresholding algorithm".

Convex Optimization Problem

$$\min_{\mathbf{x}} \{f(\mathbf{x})\} \text{ or } \min_{\mathbf{x}} \{f_1(\mathbf{x}) + f_2(\mathbf{x})\}$$

with f_1 and f_2 convex, but f_2 non-differentiable.

Algorithm 6: Accelerated Proximal Gradient 'Descent' (Nesterov)

- 1 initialization : $\mathbf{x}^{(-1)}, \mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 $\mathbf{y}^{(t-1)} \leftarrow \mathbf{x}^{(t-1)} + \frac{t-2}{t+1} [\mathbf{x}^{(t-1)} - \mathbf{x}^{(t-2)}];$
 - 4 $\mathbf{x}^{(t)} \leftarrow \text{proximal}_{h,f_2}(\mathbf{y}^{(t-1)} - h_t \nabla f_2(\mathbf{y}^{(t-1)}))$
-

Coordinate Descent

Let $\{\vec{e}_1, \dots, \vec{e}_n\}$ denote the standard basis in \mathbb{R}^n ,

$$\vec{e}_i = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^n$$

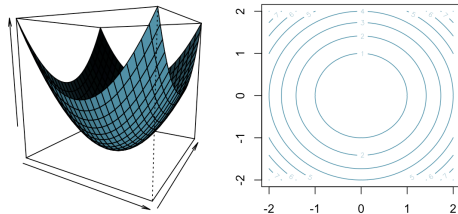
Proposition 8.17

If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, differentiable,

$$f(\mathbf{x}) \leq f(\mathbf{x} + \delta \vec{e}_i), \forall i \implies f(\mathbf{x}) = \min\{f\}$$

i.e. if we are at a point \mathbf{x} such that $f(\mathbf{x})$ is minimized along each coordinate axis, then we have found a global minimizer.

Coordinate Descent



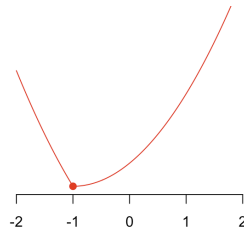
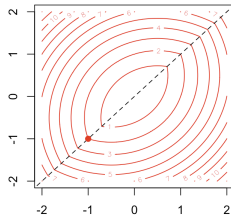
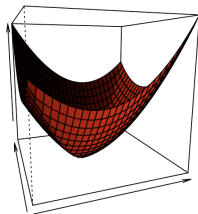
Proposition 8.18

If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, but **not differentiable**,

$$f(\mathbf{x}) \leq f(\mathbf{x} + \delta \vec{e}_i), \forall i \not\Rightarrow f(\mathbf{x}) = \min\{f\}$$

i.e. if we are at a point \mathbf{x} such that $f(\mathbf{x})$ is minimized along each coordinate axis, then we have **not** found a global minimizer.

Coordinate Descent



Proposition 8.19

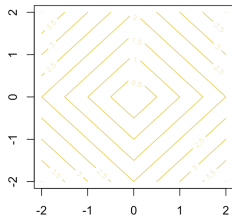
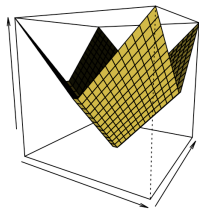
If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ can be written

$$f(\mathbf{x}) = f_0(\mathbf{x}) + \underbrace{\sum_{i=1}^n f_i(\mathbf{x}_i)}_{\text{separable}}, \quad \text{where } \begin{cases} f_0 \text{ convex and differentiable} \\ f_i \text{ convex and non-differentiable} \end{cases}$$

$$f(\mathbf{x}) \leq f(\mathbf{x} + \delta \vec{\mathbf{e}}_i), \quad \forall i \implies f(\mathbf{x}) = \min\{f\}$$

i.e. if we are at a point \mathbf{x} such that $f(\mathbf{x})$ is minimized along each coordinate axis, then we have found a global minimizer.

Coordinate Descent



Coordinate Descent

If we want to solve $\min\{f(\mathbf{x})\}$ for some $f : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$f(\mathbf{x}) = f_0(\mathbf{x}) + \underbrace{\sum_{i=1}^n f_i(\mathbf{x}_i)}_{\text{separable}}, \quad \text{where } \begin{cases} f_0 \text{ convex and differentiable} \\ f_i \text{ convex and non-differentiable} \end{cases}$$

we can use a **coordinate descent algorithm**

Algorithm 7: Coordinate Descent

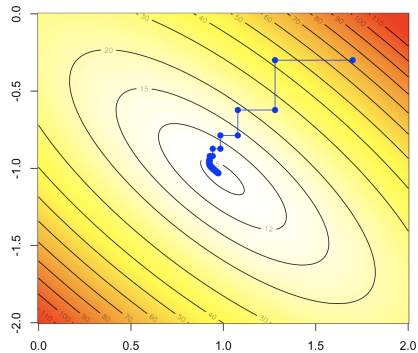
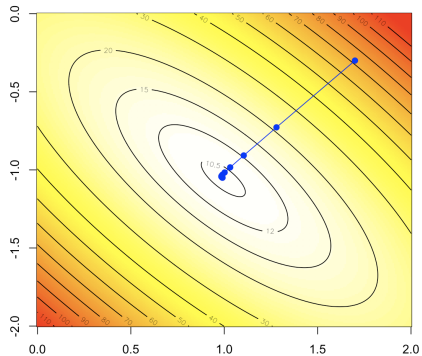
- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 **for** $j=1,2,\dots,n$ **do**
 - 4 $\mathbf{x}_j^{(t)} \leftarrow \operatorname{argmin}\{f(\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_{j-1}^{(t)}, \mathbf{x}_j, \mathbf{x}_{j+1}^{(t-1)}, \dots, \mathbf{x}_n^{(t-1)})\}$
-

Gradient vs. Coordinate Descent

Consider the problem $\min\{f(\boldsymbol{\beta})\}$ where $f(\boldsymbol{\beta}) = \frac{1}{2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$

- ▶ Gradient descent, $\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} + h\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$
- ▶ Coordinate descent, $\beta_j \leftarrow \beta_j + \frac{1}{\mathbf{x}_j^\top \mathbf{x}_j} \mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$

Gradient vs. Coordinate Descent



← to come back

Appendix: Constrained Optimization and Penalization



Let $C \subset \mathbb{R}^n$ denote a convex set, with f convex, and differentiable.

$$\min_{\mathbf{x} \in C} \{f(\mathbf{x})\} \iff \min_{\mathbf{x}} \{f(\mathbf{x}) + I_C(\mathbf{x})\} \text{ where } I_C(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{x} \in C \\ \infty & \text{if } \mathbf{x} \notin C \end{cases}$$

$$\text{proximal}_h(\mathbf{x}) = \underset{\mathbf{z}}{\text{argmin}} \left\{ \frac{1}{2h} \|\mathbf{x} - \mathbf{z}\|_2^2 + \lambda I_C(\mathbf{z}) \right\} = \underset{\mathbf{z} \in C}{\text{argmin}} \left\{ \frac{1}{2h} \|\mathbf{x} - \mathbf{z}\|_2^2 \right\}$$

i.e. $\text{proximal}_h(\mathbf{x})$ is the projection operator onto C , $\Pi_C(\mathbf{x})$

Algorithm 8: Projected Gradient Descent

- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 $\mathbf{x}^{(t)} \leftarrow \Pi_C(\mathbf{x}^{(t-1)} - h_t \nabla f(\mathbf{x}^{(t-1)}))$
-

Appendix: Constrained Optimization and Penalization

Consider a very general problem

$$\min_{\mathbf{x}} \{f(\mathbf{x})\} \text{ s.t. } \begin{cases} g_i(\mathbf{x}) \leq 0, \forall i \\ h_j(\mathbf{x}) = 0, \forall j \end{cases}$$

We can define the **Lagrangian** as

$$\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f(\mathbf{x}) + \sum_i u_i \cdot g_i \mathbf{x} + \sum_j v_j \cdot h_j \mathbf{x},$$

where $\mathbf{u} \in \mathbb{R}_+^m$ and $\mathbf{v} \in \mathbb{R}^n$.

Observe that, for feasible \mathbf{x} ,

$$\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f(\mathbf{x}) + \sum_i u_i \cdot \underbrace{g_i \mathbf{x}}_{\leq 0} + \sum_j v_j \cdot \underbrace{h_j \mathbf{x}}_{=0} \leq f(\mathbf{x}).$$

Appendix: Constrained Optimization and Penalization

Let F denote the set of feasible \mathbf{x} , define

$$\psi(\mathbf{u}, \mathbf{v}) = \min_{\mathbf{x}} \{\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v})\} \leq \min_{\mathbf{x} \in F} \{\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v})\}$$

called **Lagrange dual function**. One can easily get that

$$\psi(\mathbf{u}, \mathbf{v}) \leq f^*.$$

For example xxxx

$$\min_{\mathbf{x}} \left\{ \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \gamma^\top \mathbf{x} \right\} \text{ s.t. } \begin{cases} \mathbf{x} \geq \mathbf{0} \\ A\mathbf{x} = \mathbf{b} \end{cases}$$

then the Lagrangian is

$$\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v}) = \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \gamma^\top \mathbf{x} - \mathbf{u}^\top \mathbf{x} + \mathbf{v}^\top (A\mathbf{x} - \mathbf{b})$$

Appendix: Constrained Optimization and Penalization

and Lagrange dual function is

$$\psi(\mathbf{u}, \mathbf{v}) = \min_{\mathbf{x}} \{\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v})\} = -\frac{1}{2}(\boldsymbol{\gamma} - \mathbf{u} + A^T \mathbf{v})^T Q^{-1}(\boldsymbol{\gamma} - \mathbf{u} + A^T \mathbf{v}) - \mathbf{b}^T \mathbf{v}$$

Consider Lagrange dual problem,

$$\max_{\mathbf{u}, \mathbf{v}} \{\psi(\mathbf{u}, \mathbf{v})\} \text{ s.t. } \mathbf{u} \geq \mathbf{0}.$$

Proposition 8.20: Weak duality

If g^* is the optimal value of this problem,

$$f^* \geq \psi^* \quad (\text{weak duality}).$$

One can prove that this dual problem is a convex optimization problem (as a pointwise maximum of convex functions, in \mathbf{u}, \mathbf{v}).

Appendix: Constrained Optimization and Penalization

In some cases,

$$f^* = \psi^* \quad (\text{strong duality}).$$

Proposition 8.21: Slater's condition

If the primal problem is convex, i.e. f and g_i 's are convex, and h_j 's are affine, and if there exists at least one strictly feasible \mathbf{x} (in the sense that $g_i(\mathbf{x}) < 0$, $\forall i$), then

$$f^* = \psi^* \quad (\text{strong duality}).$$

Definition 8.14: Karush-Kuhn-Tucker (KKT) Conditions

xxx

$$\begin{cases} 0 \in \partial_{\mathbf{x}}(\mathcal{L}(\mathbf{x}, \mathbf{u}, \mathbf{v})) \\ u_i \cdot g_i(\mathbf{x}) = 0, \forall i \text{ (slackness)} \\ g_i(\mathbf{x}) \leq 0 \text{ and } h_j(\mathbf{x}) = 0, \forall i, j \text{ (primal feasibility)} \\ u_i \geq 0, \forall i \text{ (dual feasibility)} \end{cases}$$

Proposition 8.22: Karush-Kuhn-Tucker Conditions

For a problem with strong duality (e.g., assume Slater's condition, convex problem and there exists x strictly satisfying non-affine inequality constraints), \mathbf{x}^* and $(\mathbf{u}^*, \mathbf{v}^*)$ are primal and dual solutions if and only if \mathbf{x}^* and $(\mathbf{u}^*, \mathbf{v}^*)$ satisfy KKT conditions (Definition 8.14).

See [Boyd and Vandenberghe \(2004\)](#).

Appendix: Constrained Optimization and Penalization

$$\min_{Ax=b} \{f(\mathbf{x})\}$$

The dual problem is

$$\max_{\mathbf{u}} \underbrace{\{-f^*(-A^\top \mathbf{u}) - b^\top \mathbf{u}\}}_{=\psi(\mathbf{u})}$$

Observe that $\partial\psi(\mathbf{u}) = A \underbrace{\partial f^*(-A^\top \mathbf{u})}_{=\mathbf{x}} - b$. Hence,

$$\partial\psi(\mathbf{u}) = A\mathbf{x} - b \text{ where } \mathbf{x} \in \operatorname{argmin}_{\mathbf{z}} \{f(\mathbf{z}) + \mathbf{u}^\top A\mathbf{z}\}$$

Use a coordinate descent/ascent approach,

Appendix: Constrained Optimization and Penalization

Algorithm 9: Coordinate Descent (dual subgradient method)

- 1 initialization : $\mathbf{x}^{(0)}, \mathbf{u}^{(0)}$;
- 2 **for** $t=1,2,\dots$ **do**
- 3 $\mathbf{x}^{(t)} \in \operatorname{argmin}_{\mathbf{x}} \{f(\mathbf{x}) + \mathbf{u}^{(t-1)\top} A\mathbf{x}\}$;
- 4 $\mathbf{u}^{(t)} \leftarrow \mathbf{u}^{(t-1)} + \gamma(A\mathbf{x}^{(t)} - b)$;

$\mathbf{x}^{(t)}$ is unique if f is strictly convex.

Appendix: Constrained Optimization and Penalization

$$\min_{A\mathbf{x}+B\mathbf{y}=c} \{f(\mathbf{x}) + g(\mathbf{y})\}$$

$$\min_{\mathbf{x}, \mathbf{y}} \left\{ \sup_{\mathbf{z}} \{f(\mathbf{x}) + g(\mathbf{y}) + \langle \mathbf{z}, c - (A\mathbf{x} + B\mathbf{y}) \rangle + \frac{\gamma}{2} \|c - (A\mathbf{x} + B\mathbf{y})\|^2\} \right\}$$

Use a coordinate descent/ascent approach,

Algorithm 10: Coordinate Descent

- 1 initialization : $\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{z}^{(0)}$;
 - 2 **for** $t=1, 2, \dots$ **do**
 - 3 $\mathbf{x}^{(t)} \leftarrow \underset{\mathbf{x}}{\operatorname{argmin}} \{f(\mathbf{x}) - \langle \mathbf{z}^{(t-1)}, A\mathbf{x} \rangle + \frac{\gamma}{2} \|c - (A\mathbf{x} + B\mathbf{y}^{(t-1)})\|^2\}$;
 - 4 $\mathbf{y}^{(t)} \leftarrow \underset{\mathbf{y}}{\operatorname{argmin}} \{g(\mathbf{y}) - \langle \mathbf{z}^{(t-1)}, B\mathbf{y} \rangle + \frac{\gamma}{2} \|c - (A\mathbf{x}^{(t)} + B\mathbf{y})\|^2\}$;
 - 5 $\mathbf{z}^{(t)} \leftarrow \mathbf{z}^{(t-1)} + \gamma(c - (A\mathbf{x}^{(t)} + B\mathbf{y}^{(t)}))$;
-

Appendix: Constrained Optimization and Penalization

REMETTRE LE LIEN AVEC NORME 0, RAJOUTER UNE CONTRAINTE, ETC

Define $\|\mathbf{a}\|_{\ell_0} = \sum \mathbf{1}(|a_i| > 0)$ (pseudo-norm)

Here $\dim(\boldsymbol{\beta}) = k$ but $\|\boldsymbol{\beta}\|_{\ell_0} = s$.

We wish we could solve

$$\hat{\boldsymbol{\beta}}^{\text{selec}} = \underset{\boldsymbol{\beta} \in \{\|\boldsymbol{\beta}\|_{\ell_0} = s\}}{\text{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^T \boldsymbol{\beta}\|_{\ell_2}^2 \}$$

We might **convexify the ℓ_0 "norm"**, $\|\cdot\|_{\ell_0}$.

On $[-1, +1]^k$, the convex hull of $\|\boldsymbol{\beta}\|_{\ell_0}$ is $\|\boldsymbol{\beta}\|_{\ell_1}$

On $[-a, +a]^k$, the convex hull of $\|\boldsymbol{\beta}\|_{\ell_0}$ is $a^{-1} \|\boldsymbol{\beta}\|_{\ell_1}$

Hence, why not solve

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}; \|\boldsymbol{\beta}\|_{\ell_1} \leq \tilde{s}}{\text{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^T \boldsymbol{\beta}\|_{\ell_2} \},$$

Appendix: Constrained Optimization and Penalization

which is equivalent (Kuhn-Tucker theorem) to the Lagrangian optimization problem

$$\hat{\beta} = \operatorname{argmin}\{\|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2}^2 + \lambda \|\beta\|_{\ell_1}\}$$

Recall that the sub-differential of $x \mapsto |x|$ is

$$\partial|x| = \begin{cases} \{-1\} & \text{if } x < 0 \\ [-1, +1] & \text{if } x = 0 \\ \{+1\} & \text{if } x > 0 \end{cases}$$

Here, we want to find $\min\{\|\mathbf{y} - \mathbf{X}\beta\|_{\ell_2}^2 + \lambda\|\beta\|_{\ell_1}\}$, the **first order condition** is

$$\mathbf{0} \in -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \beta^* + \lambda \partial\|\beta^*\|_{\ell_1}$$

Appendix: Constrained Optimization and Penalization

i.e., for the (univariate) j -th condition, if all variables are orthogonal

$$0 \in -\hat{\beta}_j^{\text{ols}} + \beta_j^* + \frac{\lambda}{2} \partial |\beta_j^*|.$$

i.e.

$$\beta_j^* = \begin{cases} \hat{\beta}_j^{\text{ols}} + \lambda/2 & \text{if } \beta_j^* < 0 \\ \hat{\beta}_j^{\text{ols}} - \lambda/2 & \text{if } \beta_j^* > 0 \end{cases}$$

Penalized OLS

Let us define the **soft-thresholding** function,

$$S_\gamma(z) = \text{sign}(z) \cdot (|z| - \gamma)_+$$

then $\beta_j^* = S_{\lambda/2}(\hat{\beta}_j^{\text{ols}})$. xx

Penalized OLS

In the (OLS) Ridge regression

$$\min_{\beta \in \mathbb{R}^k} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2 + \lambda \sum_{j=1}^k \beta_j^2 \right\}$$

..

$$f(x) \approx u(x) = \sum_{i=1}^m \alpha_i K(x, x_i), \quad (9)$$

..

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \|Y - K\alpha\|_{\mathbb{R}^n}^2 + \lambda \alpha^\top K \alpha, \quad (10)$$

Tikhonov (1943)

← to come back

Definition 8.15: Wasserstein, [Wasserstein \(1969\)](#)

Consider two measures p and q on \mathbb{R}^k , with a norm $\|\cdot\|$ (on \mathbb{R}^k). Then define [Wasserstein distance](#)

$$W_k(p, q) = \left(\inf_{\pi \in \Pi(p, q)} \int_{\mathbb{R}^k \times \mathbb{R}^k} \|\mathbf{x} - \mathbf{y}\|^k d\pi(\mathbf{x}, \mathbf{y}) \right)^{1/k},$$

where $\Pi(p, q)$ is the set of all couplings of p and q .

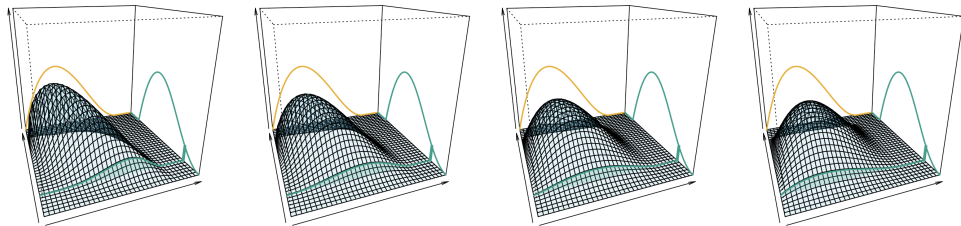
Appendix: Optimal Transport

► One can prove that $W_1(p, q) \leq W_k(p, q)$ for all $k \geq 1$. And because of Jensen's inequality, if $k \leq k'$,

$$\left(\inf_{\pi \in \Pi(p, q)} \int_{\mathbb{R}^k \times \mathbb{R}^k} \|\mathbf{x} - \mathbf{y}\|^k d\pi(\mathbf{x}, \mathbf{y}) \right)^{1/k} \leq \left(\inf_{\pi \in \Pi(p, q)} \int_{\mathbb{R}^k \times \mathbb{R}^k} \|\mathbf{x} - \mathbf{y}\|^{k'} d\pi(\mathbf{x}, \mathbf{y}) \right)^{1/k'}$$

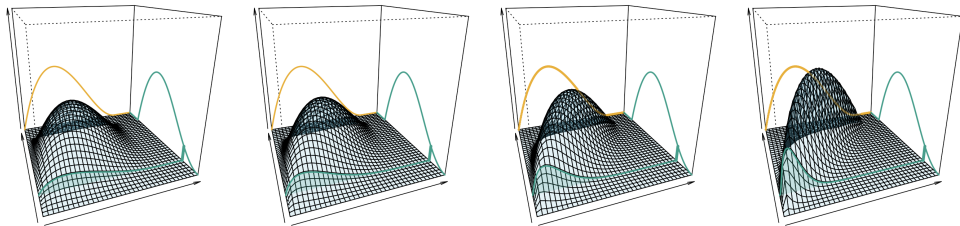
Appendix: Optimal Transport

- Wasserstein/Monge-Kantorovich distance between probability measures \mathbb{P}_A and \mathbb{P}_B : how much (kinetic) energy does it require to move a mass from \mathbb{P}_A to \mathbb{P}_B ?
- $\Pi(\mathbb{P}_A, \mathbb{P}_B)$ denotes the set of **joint probabilities** with “marginals” \mathbb{P}_A and \mathbb{P}_B



Appendix: Optimal Transport

- Wasserstein/Monge-Kantorovich distance between probability measures \mathbb{P}_A and \mathbb{P}_B : how much (kinetic) energy does it require to move a mass from \mathbb{P}_A to \mathbb{P}_B ?
- $\Pi(\mathbb{P}_A, \mathbb{P}_B)$ denotes the set of **joint probabilities** with “marginals” \mathbb{P}_A and \mathbb{P}_B



Appendix: Optimal Transport

› If $P \in \Pi(\mathbb{P}_A, \mathbb{P}_B)$, for all \mathcal{A} and \mathcal{B} ,

$$P(\mathcal{A} \times \mathcal{Y}) = \mathbb{P}_A(\mathcal{A}) \text{ and } P(\mathcal{X} \times \mathcal{B}) = \mathbb{P}_B(\mathcal{B})$$

or, put differently, for all ψ and φ ,

$$\int_{\mathcal{X} \times \mathcal{Y}} (\psi(\mathbf{x}) + \varphi(\mathbf{y})) dP(\mathbf{x}, \mathbf{y}) = \int_{\mathcal{X}} \psi(\mathbf{x}) d\mathbb{P}_A(\mathbf{x}) + \int_{\mathcal{Y}} \varphi(\mathbf{y}) d\mathbb{P}_B(\mathbf{y}),$$

or equivalently

$$\sup_{\psi, \varphi} \left\{ \int_{\mathcal{X}} \psi(\mathbf{x}) d\mathbb{P}_A(\mathbf{x}) + \int_{\mathcal{Y}} \varphi(\mathbf{y}) d\mathbb{P}_B(\mathbf{y}) - \int_{\mathcal{X} \times \mathcal{Y}} (\psi(\mathbf{x}) + \varphi(\mathbf{y})) dP(\mathbf{x}, \mathbf{y}) \right\} = 0.$$

Wasserstein- k distance (with $k \geq 1$) is

$$W_k(\mathbb{P}_A, \mathbb{P}_B) = \left(\inf_{P \in \Pi(\mathbb{P}_A, \mathbb{P}_B)} \mathbb{E}_P [c(X, Y)^k] \right)^{1/k} \text{ where } (X, Y) \sim P.$$

Appendix: Optimal Transport

Definition 8.16: Kantorovich Problem

Kantorovich Problem is defined as

$$W_c(\mathbb{P}_A, \mathbb{P}_B) = \inf_{P \in \Pi(\mathbb{P}_A, \mathbb{P}_B)} \int_{\mathcal{X} \times \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) dP(\mathbf{x}, \mathbf{y}),$$

for cost function c (or loss function).



› We can rewrite using a Lagrangian form

$$W_c(\mathbb{P}_A, \mathbb{P}_B) = \inf_P \left\{ \sup_{\psi, \varphi} \{L(P, \psi, \varphi)\} \right\},$$

where

$$L(P, \psi, \varphi) = \int_{\mathcal{X} \times \mathcal{Y}} (c(\mathbf{x}, \mathbf{y}) - (\psi(\mathbf{x}) + \varphi(\mathbf{y}))) dP(\mathbf{x}, \mathbf{y}) + \int_{\mathcal{X}} \psi(\mathbf{x}) d\mathbb{P}_A(\mathbf{x}) + \int_{\mathcal{Y}} \varphi(\mathbf{y}) d\mathbb{P}_B(\mathbf{y}).$$

Appendix: Optimal Transport

$$L(P, \psi, \varphi) = \int_{\mathcal{X} \times \mathcal{Y}} (c(\mathbf{x}, \mathbf{y}) - (\psi(\mathbf{x}) + \varphi(\mathbf{y}))) dP(\mathbf{x}, \mathbf{y}) + \int_{\mathcal{X}} \psi(\mathbf{x}) d\mathbb{P}_{\mathbf{A}}(\mathbf{x}) + \int_{\mathcal{Y}} \varphi(\mathbf{y}) d\mathbb{P}_{\mathbf{B}}(\mathbf{y}).$$

➤ The dual problem becomes

$$\sup_{\psi, \varphi} \left\{ \inf_P \{L(P, \psi, \varphi)\} \right\},$$

clearly

$$\inf_P \{L(P, \psi, \varphi)\} = \int_{\mathcal{X}} \psi(\mathbf{x}) d\mathbb{P}_{\mathbf{A}}(\mathbf{x}) + \int_{\mathcal{Y}} \varphi(\mathbf{y}) d\mathbb{P}_{\mathbf{B}}(\mathbf{y}) \text{ if } c(\mathbf{x}, \mathbf{y}) \geq \psi(\mathbf{x}) + \varphi(\mathbf{y}).$$

Appendix: Optimal Transport

Theorem 8.4: Minimax theorem, von Neumann (1928)

Let $\mathcal{X} \subset \mathbb{R}^n$ and $\mathcal{Y} \subset \mathbb{R}^m$ be compact convex sets. If $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a continuous function that is concave-convex, i.e.

$$\begin{cases} f(\cdot, y) : \mathcal{X} \rightarrow \mathbb{R} \text{ is concave for fixed } y \\ f(x, \cdot) : \mathcal{Y} \rightarrow \mathbb{R} \text{ is convex for fixed } x, \end{cases}$$

Then we have that

$$\max_{x \in \mathcal{X}} \min_{y \in \mathcal{Y}} f(x, y) = \min_{y \in \mathcal{Y}} \max_{x \in \mathcal{X}} f(x, y).$$

➤ For example, $f(x, y) = x^\top Ay$ for a finite matrix $A \in \mathbb{R}^{n \times m}$,

$$\max_{x \in \mathcal{X}} \min_{y \in \mathcal{Y}} x^\top Ay = \min_{y \in \mathcal{Y}} \max_{x \in \mathcal{X}} x^\top Ay.$$

Appendix: Optimal Transport

- Using a generalized minimax theorem, [Parthasarathy \(1970\)](#),

$$W_c(\mathbb{P}_A, \mathbb{P}_B) = \sup_{\psi, \varphi} \left\{ \int_{\mathcal{X}} \psi(\mathbf{x}) d\mathbb{P}_A(\mathbf{x}) + \int_{\mathcal{Y}} \varphi(\mathbf{y}) d\mathbb{P}_B(\mathbf{y}) \right\}, \text{ s.t. } c(\mathbf{x}, \mathbf{y}) \geq \psi(\mathbf{x}) + \varphi(\mathbf{y}),$$

the constraint can be rewritten

$$\psi(\mathbf{x}) \leq \overbrace{\min_{\mathbf{z}} \{c(\mathbf{x}, \mathbf{z}) - \psi(\mathbf{z})\}}^{=\varphi^c(\mathbf{x})}$$

where φ^c is the c -transform of φ , and

$$W_c(\mathbb{P}_A, \mathbb{P}_B) = \sup_{\varphi} \left\{ \int_{\mathcal{X}} \varphi^c(\mathbf{x}) d\mathbb{P}_A(\mathbf{x}) + \int_{\mathcal{Y}} \varphi(\mathbf{y}) d\mathbb{P}_B(\mathbf{y}) \right\}$$

called **Kantorovich duality formula**.

Appendix: Hilbert Spaces (and RKHS)

← to come back

Definition 8.17: Inner Product

Let \mathcal{H} denote some real-valued vector space. An **inner product** on \mathcal{H} is the application $(f, g) \mapsto \langle f, g \rangle_{\mathcal{H}}$ (taking value in \mathbb{R}) bilinear, symmetric, definite positive:

- $\langle f, g \rangle_{\mathcal{H}} = \langle g, f \rangle_{\mathcal{H}}$
- $\langle \alpha f + \beta g, h \rangle_{\mathcal{H}} = \alpha \langle f, h \rangle_{\mathcal{H}} + \beta \langle g, h \rangle_{\mathcal{H}}$
- $\langle f, f \rangle_{\mathcal{H}} \geq 0$ and $\langle f, f \rangle_{\mathcal{H}} = 0$ if and only if $f = 0$.

Appendix: Hilbert Spaces (and RKHS)

Definition 8.18: $L^2(\mathbb{P})$

Define $L^2(\mathbb{P})$ as the Hilbert space

$$L^2(\mathbb{P}) = \{h : \mathcal{X} \rightarrow \mathbb{R} : \mathbb{E}[h^2(X)] < \infty \text{ where } X \sim \mathbb{P}\},$$

with inner-product

$$(f, g) \mapsto \langle f, g \rangle_{L^2} = \mathbb{E}[f(X)g(X)], \quad X \sim \mathbb{P}.$$

Appendix: Hilbert Spaces (and RKHS)

Most properties are related to the geometry the Euclidean space of \mathbb{R}^n .

Example : $\mathcal{H} = \mathbb{R}^n$, $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$

Example : $\mathcal{H} = \ell_2 = \left\{ u : \sum_{i=1}^{\infty} u_i^2 < \infty \right\}$, $\langle u, v \rangle = \sum_{i=1}^{\infty} u_i v_i$

Example : $\mathcal{H} = L_2(\mu) = \left\{ f : \int f(x)^2 d\mu(x) < \infty \right\}$, $\langle f, g \rangle = \int f(x)g(x)d\mu(x)$

If \mathcal{H} is finite, $\mathcal{H} = \{h_1 \cdots, h_d\}$, $\langle x, y \rangle_{\mathcal{H}}$ takes value $K_{i,j}$ if $x = h_i$ and $y = h_j$. Let $\mathbf{K} = [K_{i,j}]$

\mathbf{K} is a symmetric $d \times d$ matrix, $\mathbf{K} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$ for some orthogonal matrix \mathbf{V} where columns are eigenvectors, and $\mathbf{\Lambda} = \text{diag}[\lambda_i]$ (positive values). Let

$$\Phi(x) = (\sqrt{\lambda_1}V_{i,1}, \sqrt{\lambda_2}V_{2,i}, \cdots, \sqrt{\lambda_d}V_{d,i}) \text{ if } x = h_i$$

Note that

$$K_{i,j} = [\mathbf{K}]_{i,j} = [\mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top]_{i,j} = \sum_{l=1}^k \lambda_l V_{i,l} V_{l,j} = \langle \Phi(h_i), \Phi(h_j) \rangle$$

Appendix: Hilbert Spaces (and RKHS)

Matrix \mathbf{K} defines an inner product, it is called a **kernel**. It is symmetric, associated with a positive semi-definite matrix.

Then $K(u, u) \geq 0$ and $K(u, v) \leq \sqrt{K(u, u) \cdot K(v, v)}$.

Let $\Phi : u \mapsto K(\cdot, u)$, then $K(x, y) = \langle \Phi(x), \Phi(y) \rangle$

Definition 8.19: Kernel

A function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called **kernel** if there exists \mathcal{H} and a map $\varphi : \mathcal{X} \rightarrow \mathcal{H}$ such that

$$k(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}, \quad \forall x, y \in \mathcal{X}$$

If k is a kernel, and $\alpha > 0$, so is αk .

If k_1 and k_2 are kernels, so is $k_1 + k_2$.

In a general setting, let $\|f\|_{\mathcal{H}} = \langle f, f \rangle_{\mathcal{H}}^{1/2}$, and define the distance from f to $\mathcal{G} \subset \mathcal{H}$

$$d(f, \mathcal{G}) = \inf_{g \in \mathcal{G}} \{\|f - g\|_{\mathcal{H}}\} = d(f, g^*) \text{ where } g^* \in \mathcal{G}$$

Appendix: Hilbert Spaces (and RKHS)

Note that $\langle g, f - g^* \rangle_{\mathcal{H}} = 0, \forall g \in \mathcal{G}$. And $\mathcal{H} = \mathcal{G} \oplus \mathcal{G}^{\perp}$.

Proposition 8.23: Riesz representation theorem

For any continuous linear functionals L from \mathcal{H} into the field \mathbb{R} , there exists a unique $g_L \in \mathcal{H}$ such that $\forall f \in \mathcal{H}, \langle g_L, f \rangle_{\mathcal{H}} = Lf$.

Consider the case where $\mathcal{H} = \mathbb{R}^n$. Let Σ denote some symmetric $n \times n$ positive definite matrix. Then

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\Sigma} = \mathbf{x}^{\top} \Sigma^{-1} \mathbf{y} \text{ is an inner product on } \mathbb{R}^n.$$

Note that if σ_i denote columns of Σ $\langle \sigma_i, \sigma_j \rangle_{\Sigma} = \sigma_i^{\top} \Sigma^{-1} \sigma_j = \Sigma_{i,j}$, and more generally, $\langle \sigma_i, \mathbf{x} \rangle_{\Sigma} = x_i$

The space \mathcal{H} of functions $\mathbb{R}^p \rightarrow \mathbb{R}$ is a Reproducing Kernel Hilbert Space (**RKHS**) if there is an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ such that \mathcal{H} with $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is an Hilbert space, and for all $\mathbf{x} \in \mathbb{R}^p$, linear functional $\delta_{\mathbf{x}} : \mathcal{H} \rightarrow \mathbb{R}$ defined as $\delta_{\mathbf{x}}(f) = f(\mathbf{x})$ is bounded.

Appendix: Hilbert Spaces (and RKHS)

Thus, \mathcal{H} is a RKHS if and only if $\forall f \in \mathcal{H}$ and $\mathbf{x} \in \mathbb{R}^p$, there exists $M_{\mathbf{x}}$ such that $|f(\mathbf{x})| \leq M_{\mathbf{x}} \cdot \|f\|_{\mathcal{H}}$.

From Riesz theorem, there exists a unique $\zeta_{\mathbf{x}} \in \mathcal{H}$ associated with $\delta_{\mathbf{x}}$, i.e.

$$\langle \zeta_{\mathbf{x}}, f \rangle_{\mathcal{H}} = f(\mathbf{x})$$

Definition 8.20: Reproducing Kernel of \mathcal{H}

Function $\mathbf{x} \mapsto \zeta_{\mathbf{x}}$ is called reproducing function in \mathbf{x} and $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ defined as $K(\mathbf{x}, \mathbf{y}) = \zeta_{\mathbf{x}}(\mathbf{y})$ is the reproducing kernel of \mathcal{H} .

Observe that $\langle K(\mathbf{x}, \cdot), K(\mathbf{y}, \cdot) \rangle_{\mathcal{H}} = K(\mathbf{x}, \mathbf{y})$.

The kernel is unique, and is (semi-)definite positive.

If \mathcal{H} is a closed subspace of Hilbert space \mathcal{X} . For any function $f \in \mathcal{X}$, $\mathbf{x} \mapsto \langle f, K(\cdot, \mathbf{x}) \rangle_{\mathcal{X}}$ is the projection of f on \mathcal{H} .

Note that conversely, **Moore-Aronszajn's theorem** allows to create a RKHS from a definite positive kernel K .

Appendix: Hilbert Spaces (and RKHS)

Proposition 8.24: Mercer's kernel

Let μ denote some measure on \mathbb{R}^p and $\mathcal{H} = L^2(\mathbb{R}^p, \mu)$, define

$$(L_K f)(\mathbf{x}) = \int K(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mu(\mathbf{y})$$

which is a compact bounded linear operator, self-adjoint and positive. Let $\lambda_1 \geq \lambda_2 \geq \dots$ denote eigenvalues of L_K , with (orthonormal) eigenvectors ψ_1, ψ_2, \dots , then

$$K(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^p \lambda_k \psi_k(\mathbf{x}) \psi_k(\mathbf{y}) = \Psi(\mathbf{x})^\top \Psi(\mathbf{y}) = \langle \Psi(\mathbf{x}), \Psi(\mathbf{y}) \rangle_{L^2}$$

where $\Psi(\mathbf{x}) = (\sqrt{\lambda_k} \psi_k(\mathbf{x}))$.

Appendix: Hilbert Spaces (and RKHS)

Example : Consider the space \mathcal{H} defined as

$$\mathcal{H}_1 = \{f : [0, 1] \rightarrow \mathbb{R} \text{ continuously differentiable, with } f' \in L^2([0, 1]) \text{ and } f(0) = 0\}$$

\mathcal{H}_1 is an Hilbert space on $[0, 1]$ with inner product

$$\langle f, g \rangle_{\mathcal{H}_1} = \int_0^1 f'(t)g'(t)dt$$

with (definite positive) kernel $K_1(x, y) = \min\{x, y\}$:

$$\langle f, K(x, \cdot) \rangle_{\mathcal{H}_1} = \int_0^1 f'(t) \underbrace{\frac{\partial K_1(t, x)}{\partial x}}_{= \mathbf{1}_{[0, x]}(t)} dt = \int_0^x f'(t) dt = f(x)$$

Example : Consider the Sobolev space $W^1([0, 1])$ defined as

$$W^1([0, 1]) = \{f : [0, 1] \rightarrow \mathbb{R} \text{ continuously differentiable, with } f' \in L^2([0, 1])\}$$

Appendix: Hilbert Spaces (and RKHS)

Observe that $W^1([0, 1]) = \mathcal{H}_0 \oplus \mathcal{H}_1$ where

$$\mathcal{H}_0 = \{f : [0, 1] \rightarrow \mathbb{R} \text{ continuously differentiable, with } f' = 0\}$$

The later is an Hilbert space with kernel $K_0(x, y) = 1$.

One can consider kernel $K(x, y) = K_0(x, y) + K_1(x, y)$ (related to linear splines). More generally, consider

$$\mathcal{H}_2 = \{f : [0, 1] \rightarrow \mathbb{R} \text{ twice cont. diff., with } f'' \in L^2([0, 1] \text{ with } f'(0) = 0)\}$$

Then $\langle f, g \rangle_{\mathcal{H}_2} = \int_0^1 f''(t)g''(t)dt$ is an inner product, with kernel

$$K_2(x, y) = \int_0^1 (x - t)_+(y - t)_+ dt$$

Appendix: Hilbert Spaces (and RKHS)

Consider some Hilbert space \mathcal{H} with kernel K and some functional $\Psi : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$, increasing in its last argument.

Proposition 8.25: Representation theorem

Given $\mathbf{x}_1, \dots, \mathbf{x}_n$, $\min_{f \in \mathcal{H}} \{ \Psi(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \|f\|_{\mathcal{H}}) \}$ admits solution

$$\forall \mathbf{x}, f(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x})$$

A classical expression for Ψ is, for some convex function ψ ,

$$\Psi(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \|f\|_{\mathcal{H}}) = \psi(\mathbf{y}, f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) + \lambda \|f\|_{\mathcal{H}}$$

$$\Psi(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n), \|f\|_{\mathcal{H}}) = \sum_{i=1}^n \ell(\mathbf{y}, f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{H}}$$

Appendix: Hilbert Spaces (and RKHS)

Assume that $y_i = m(x_i) + \varepsilon_i$, where $m \in W_2([0, 1])$, then polynomial splines of degree 2 is the solution of

$$\min_{m \in W_2} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - m(x_i))^2 + \nu \int_0^1 [m''(t)]^2 dt \right\}$$

then $m^*(x) = \beta_0 + \beta_1 x + \sum_{i=1}^n \gamma_i K_2(x_i, x)$ Note that one can use a matrix representation

$$\min \{ (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Q}\boldsymbol{\gamma})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Q}\boldsymbol{\gamma}) + n\nu \boldsymbol{\gamma}^\top \mathbf{Q}\boldsymbol{\gamma} \}$$

where $\mathbf{Q} = [K_1(x_i, x_j)]$. If $\mathbf{M} = \mathbf{Q} + n\nu \mathbb{I}$,

$$\boldsymbol{\beta}^* = (\mathbf{X}^\top \mathbf{M}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}^{-1} \mathbf{y} \text{ and } \boldsymbol{\gamma}^* = \mathbf{M}^{-1} (\mathbb{I} - \mathbf{X}(\mathbf{X}^\top \mathbf{M}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}^{-1}) \mathbf{y}$$

Appendix: Hilbert Spaces (and RKHS)

Proposition 8.26: **Kimeldorf and Wahba (1971)** Representation Theorem

Consider a kernel K and \mathcal{H}_K the associated RKHS. For any (convex) loss function $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$, the solution

$$m^* \in \operatorname{argmin}_{m \in \mathcal{H}_K} \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i)) + \|m\|_{\mathcal{H}_K}^2$$

can be expressed

$$m^*(\cdot) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \cdot)$$

Appendix: Hilbert Spaces (and RKHS)

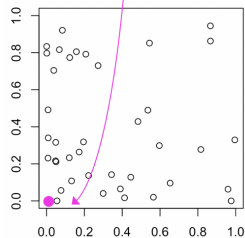
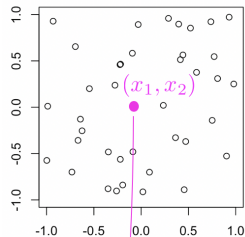
What's going on, here ?

For more technical (mathematical) results, see Wahba (1990, [Spline Models for Observational Data](#))

We've seen that in many cases, $K(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})^\top \varphi(\mathbf{y})$ for some $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^q$ (here $q = p$)

Consider for example $\varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1^2 \\ x_2^2 \end{pmatrix}$

so that $K(\mathbf{x}, \mathbf{y}) = x_1^2 y_1^2 + x_2^2 y_2^2$



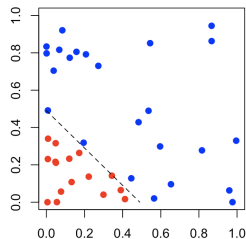
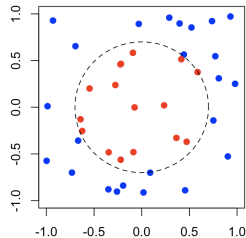
Appendix: Hilbert Spaces (and RKHS)

What's going on, here ?

Consider $\varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1^2 \\ x_2^2 \end{pmatrix}$

$$K(\mathbf{x}, \mathbf{y}) = x_1^2 y_1^2 + x_2^2 y_2^2$$

From data (y_i, \mathbf{x}_i) , transform the covariates into $(y_i, \phi(\mathbf{x}_i))$, and use a (classical) linear model

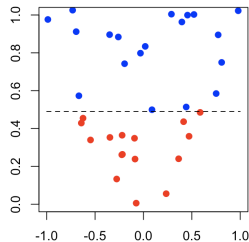
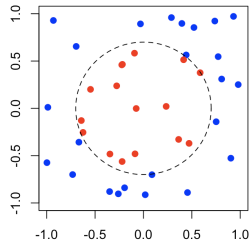


Appendix: Hilbert Spaces (and RKHS)

What's going on, here ?

$$\text{Consider } \varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1 \\ x_1^2 + x_2^2 \end{pmatrix}$$

$$K(\mathbf{x}, \mathbf{y}) = x_1 y_1 + (x_1^2 + x_2^2)(y_1^2 + y_2^2)$$



Appendix: Hilbert Spaces (and RKHS)

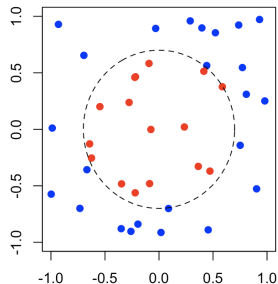
Appendix: Hilbert Spaces (and RKHS)

What's going on, here ?

But we can have $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^q$ (with $q \neq p$)

$$\text{Consider } \varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1 \\ x_2 \\ x_1 \cdot x_2 \end{pmatrix}$$

A classical idea with SVM will be to consider $q > p$ to be able to find a linear separator of points red and blue



Appendix: Hilbert Spaces (and RKHS)

What's going on, here ?

$$\text{Consider } \varphi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix}$$

$K(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^\top \mathbf{y})^2$ is called **polynomial kernel** (of order 2). More generally

$K(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^\top \mathbf{y})^d$, with $d \in \mathbb{N}$.

For any degree $d \geq 2$, $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^q$ with $q = \binom{p+d}{d}$ (default in R, $d = 3$)

Appendix: Hilbert Spaces (and RKHS)

How to get those kernels ?

If K_1 and K_2 are two kernels, so are

$$K(\mathbf{x}, \mathbf{y}) = a_1 K_1(\mathbf{x}, \mathbf{y}) + a_2 K_2(\mathbf{x}, \mathbf{y}) \text{ and } K(\mathbf{x}, \mathbf{y}) = K_1(\mathbf{x}, \mathbf{y}) \cdot K_2(\mathbf{x}, \mathbf{y})$$

If h is some $\mathbb{R}^n \rightarrow \mathbb{R}^n$ function, $K(\mathbf{x}, \mathbf{y}) = K_1(h(\mathbf{x}), h(\mathbf{y}))$

If h is some $\mathbb{R}^n \rightarrow \mathbb{R}$ function, $K(\mathbf{x}, \mathbf{y}) = h(\mathbf{x}) \cdot h(\mathbf{y})$

If P is a polynomial with positive coefficients, $K(\mathbf{x}, \mathbf{y}) = P(K(\mathbf{x}, \mathbf{y}))$ as well as

$$K(\mathbf{x}, \mathbf{y}) = \exp[\mathbf{x}^\top \mathbf{y}]$$

← to come back

Residuals and Independence

In a standard regression setting, $Y = \mu(\mathbf{X}) + \varepsilon$ where

$$\begin{cases} \mu(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}] \text{ is the regression function (conditional mean)} \\ \varepsilon = Y - \mu(\mathbf{X}) \text{ is the residual term} \end{cases}$$

where ideally we want residuals ε to be independent of \mathbf{X} .

Set $\epsilon = g(Y, \mathbf{X})$ so that

$$\begin{cases} \epsilon \perp\!\!\!\perp \mathbf{X}, \\ \sigma(\{Y, \mathbf{X}\}) = \sigma(\{\epsilon, \mathbf{X}\}) \end{cases}$$

i.e., $Y = h(\mathbf{X}, g(Y, \mathbf{X}))$.

There is no unique g , suppose

$$\begin{cases} y \mapsto g(y, \mathbf{x}) \text{ is strictly increasing, for every } \mathbf{x}, \\ g(Y, \mathbf{X}) \mid \mathbf{X} = \mathbf{x} \text{ is uniformly distributed, for every } \mathbf{x}. \end{cases}$$

Residuals and Independence

Then g is unique, and $g(y, \mathbf{x}) = F_{Y, \mathbf{X}}(y, \mathbf{x})$, and $h(\mathbf{X}, u) = F_{Y, \mathbf{X}}^{-1}(u | \mathbf{x})$.

Proof: From the uniform property, $\mathbb{P}[g(Y, \mathbf{X}) \leq u | \mathbf{X} = \mathbf{x}] = u$, for any $u \in [0, 1]$. Set $g_{\mathbf{x}}(y) = g(y, \mathbf{x})$, so that

$$\mathbb{P}[g(Y, \mathbf{X}) \leq u | \mathbf{X} = \mathbf{x}] = \mathbb{P}[g_{\mathbf{x}}(Y) \leq u | \mathbf{X} = \mathbf{x}] = \mathbb{P}[Y \leq g_{\mathbf{x}}^{-1}(u) | \mathbf{X} = \mathbf{x}] = F_{Y | \mathbf{X}}(g_{\mathbf{x}}^{-1}(u) | \mathbf{x}),$$

thus, $F_{Y | \mathbf{X}}(g_{\mathbf{x}}^{-1}(u) | \mathbf{x}) = u$, for every $u \in [0, 1]$, i.e., $g_{\mathbf{x}}(y) = F_{Y | \mathbf{X}}(y | \mathbf{x})$, i.e.,

$$h(\mathbf{x}, g(y, \mathbf{x})) = F_{Y | \mathbf{X}}^{-1}(g(y, \mathbf{x}) | \mathbf{x}) = F_{Y | \mathbf{X}}^{-1}(F_{Y | \mathbf{X}}(y | \mathbf{x}) | \mathbf{x}) = y.$$

References

- Abraham, K. (1986). *Distributing risk: Insurance, legal theory and public policy*. Yale University Press,.
- Aczél, J. (1948). On mean values. *Bulletin of the American Mathematical Society*, 54(4):392–400.
- Agresti, A. and Coull, B. A. (1998). Approximate is better than “exact” for interval estimation of binomial proportions. *The American Statistician*, 52(2):119–126.
- Agueh, M. and Carlier, G. (2011). Barycenters in the wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924.
- Ali, S. M. and Silvey, S. D. (1966). A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society: Series B (Methodological)*, 28(1):131–142.
- Altschuler, J. M. and Boix-Adsera, E. (2021). Wasserstein barycenters can be computed in polynomial time in fixed dimension. *The Journal of Machine Learning Research*, 22(1):2000–2018.
- Amari, S.-i. (2016). *Information geometry and its applications*, volume 194. Springer.
- Andrus, M., Spitzer, E., Brown, J., and Xiang, A. (2021). What we can't measure, we can't understand: Challenges to demographic data procurement in the pursuit of fairness. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 249–260.
- Angrist, J. D. and Pischke, J.-S. (2009). *Mostly harmless econometrics: An empiricist's companion*. Princeton university press.

References

- Angwin, J., Larson, J., Mattu, S., and Kirchner, L. (2016). Machine bias: There's software used across the country to predict future criminals and it's biased against blacks. *ProPublica*, May 23.
- Apfelbaum, E. P., Pauker, K., Sommers, S. R., and Ambady, N. (2010). In blind pursuit of racial equality? *Psychological science*, 21(11):1587–1592.
- Austin, R. (1983). The insurance classification controversy. *University of Pennsylvania Law Review*, 131(3):517–583.
- Avin, C., Shpitser, I., and Pearl, J. (2005). Identifiability of path-specific effects. *IJCAI International Joint Conference on Artificial Intelligence*, pages 357–363.
- Avraham, R. (2017). Discrimination and insurance. In Lippert-Rasmussen, K., editor, *Handbook of the Ethics of Discrimination*, pages 335–347. Routledge.
- Avraham, R., Logue, K. D., and Schwarcz, D. (2013). Understanding insurance antidiscrimination law. *Southern California Law Review*, 87:195.
- Bach, F. R. and Jordan, M. I. (2002). Kernel independent component analysis. *Journal of machine learning research*, 3(Jul):1–48.
- Bailey, R. A. and Simon, L. J. (1959). An actuarial note on the credibility of experience of a single private passenger car. *Proceedings of the Casualty Actuarial Society*, XLVI:159.
- Baldus, D. C. and Cole, J. W. (1980). *Statistical proof of discrimination*. McGraw-Hill.

References

- Banerjee, A., Merugu, S., Dhillon, I. S., Ghosh, J., and Lafferty, J. (2005). Clustering with bregman divergences. *Journal of machine learning research*, 6(10).
- Barbour, V. (1911). Privateers and pirates of the west indies. *The American Historical Review*, 16(3):529–566.
- Barocas, S., Hardt, M., and Narayanan, A. (2017). Fairness in machine learning. *Nips tutorial*, 1:2017.
- Barry, L. and Charpentier, A. (2020). Personalization as a promise: Can big data change the practice of insurance? *Big Data & Society*, 7(1):2053951720935143.
- Bartlett, P. L., Jordan, M. I., and McAuliffe, J. D. (2006). Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156.
- Bauschke, H. H., Borwein, J. M., et al. (1997). Legendre functions and the method of random bregman projections. *Journal of convex analysis*, 4(1):27–67.
- Bauschke, H. H. and Lewis, A. S. (2000). Dykstras algorithm with bregman projections: A convergence proof. *Optimization*, 48(4):409–427.
- Beck, A. and Teboulle, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202.
- Becker, G. S. (1957). *The economics of discrimination*. University of Chicago press.

References

- Bellemare, M. G., Dabney, W., and Munos, R. (2017a). A distributional perspective on reinforcement learning. *arXiv:1707.06887*.
- Bellemare, M. G., Danihelka, I., Dabney, W., Mohamed, S., Lakshminarayanan, B., Hoyer, S., and Munos, R. (2017b). The cramer distance as a solution to biased wasserstein gradients. *arXiv:1705.10743*.
- Ben-Tal, A., El Ghaoui, L., and Nemirovski, A. (2009). *Robust optimization*, volume 28. Princeton university press.
- Benamou, J.-D., Carlier, G., Cuturi, M., Nenna, L., and Peyré, G. (2015). Iterative Bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138.
- Bengio, Y., Goodfellow, I., and Courville, A. (2017). *Deep learning*, volume 1. MIT press Cambridge, MA, USA.
- Beutel, A., Chen, J., Zhao, Z., and Chi, E. H. (2017). Data decisions and theoretical implications when adversarially learning fair representations. *arXiv*, 1707.00075.
- Bickel, P. J., Hammel, E. A., and O'Connell, J. W. (1975). Sex bias in graduate admissions: Data from Berkeley. *Science*, 187(4175):398–404.

References

- Biddle, D. (2017). *Adverse impact and test validation: A practitioner's guide to valid and defensible employment testing*. Routledge.
- Billingsley, P. (2017). *Probability and measure*. John Wiley & Sons.
- Blanpain, N. (2018). L'espérance de vie par niveau de vie-méthode et principaux résultats. *INSEE Document de Travail*, F1801.
- Blinder, A. S. (1973). Wage discrimination: Reduced form and structural estimates. *The Journal of Human Resources*, 8(4):436–455.
- Bollen, K. A. (1989). *Structural equations with latent variables*, volume 210. John Wiley & Sons.
- Bollen, K. A. and Pearl, J. (2013). Eight myths about causality and structural equation models. In *Handbook of causal analysis for social research*, pages 301–328. Springer.
- Borkan, D., Dixon, L., Sorensen, J., Thain, N., and Vasserman, L. (2019). Nuanced metrics for measuring unintended bias with real data for text classification. In *Companion proceedings of the 2019 world wide web conference*, pages 491–500.
- Bourdieu, P. (2018). Distinction a social critique of the judgement of taste. In *Inequality Classic Readings in Race, Class, and Gender*, pages 287–318. Routledge.
- Box, G. E., Luceño, A., and del Carmen Paniagua-Quinones, M. (2011). *Statistical control by monitoring and adjustment*, volume 700. John Wiley & Sons.

References

- Boyd, S. P. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.
- Boyle, J. P. and Dykstra, R. L. (1986). A method for finding projections onto the intersection of convex sets in hilbert spaces. In *Advances in Order Restricted Statistical Inference: Proceedings of the Symposium on Order Restricted Statistical Inference held in Iowa City, Iowa, September 11–13, 1985*, pages 28–47. Springer.
- Bregman, L. M. (1967). The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics*, 7(3):200–217.
- Brenier, Y. (1991). Polar factorization and monotone rearrangement of vector-valued functions. *Communications on pure and applied mathematics*, 44(4):375–417.
- Brier, G. W. (1950). Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3.
- Brilmayer, L., Hekeler, R. W., Laycock, D., and Sullivan, T. A. (1979). Sex discrimination in employer-sponsored insurance plans: A legal and demographic analysis. *University of Chicago Law Review*, 47:505.
- Britz, G. (2008). *Einzelfallgerechtigkeit versus Generalisierung: verfassungsrechtliche Grenzen statistischer Diskriminierung*. Mohr Siebeck.

References

- Brown, R. S., Moon, M., and Zoloth, B. S. (1980). Incorporating occupational attainment in studies of male-female earnings differentials. *Journal of Human Resources*, pages 3–28.
- Brualdi, R. A. (2006). *Combinatorial matrix classes*, volume 13. Cambridge University Press.
- Calders, T. and Verwer, S. (2010). Three naive Bayes approaches for discrimination-free classification. *Data mining and knowledge discovery*, 21(2):277–292.
- Casella, G. and Berger, R. L. (1990). *Statistical Inference*. Duxbury Advanced Series.
- Casey, B., Pezier, J., and Spetzler, C. (1976). *The Role of Risk Classification in Property and Casualty Insurance: A Study of the Risk Assessment Process : Final Report*. Stanford Research Institute.
- Censor, Y. and Reich, S. (1998). The dykstra algorithm with bregman projections. *Communications in Applied Analysis*, 2(3):407–420.
- Cha, S.-H. (2007). Comprehensive survey on distance/similarity measures between probability density functions, intl. *Journal of Mathematical Models and Methods in Applied Sciences*, Issue, 4.
- Chambert-Loir, A. (2023). *Information Theory: Three Theorems by Claude Shannon*, volume 144. Springer Nature.
- Charniak, E. (1991). Bayesian networks without tears. *AI magazine*, 12(4):50–50.

References

- Charpentier, A. (2014). Mesures de risque. In Dreesbeke, J.-J. and Saporta, G., editors, *Approches statistiques du risque*. Éditions Technip.
- Charpentier, A. (2023). *Insurance: biases, discrimination and fairness*. Springer Verlag.
- Charpentier, A., Hu, F., and Ratz, P. (2023). Mitigating discrimination in insurance with wasserstein barycenters. *BIAS, 3rd Workshop on Bias and Fairness in AI, International Workshop of ECML PKDD*.
- Cheney-Lippold, J. (2017). We are data. In *We Are Data*. New York University Press.
- Chicco, D. and Jurman, G. (2020). The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation. *BMC genomics*, 21(1):1–13.
- Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163.
- Chzhen, E. and Schreuder, N. (2022). A minimax framework for quantifying risk-fairness trade-off in regression. *The Annals of Statistics*, 50(4):2416–2442.
- Cleary, T. A. (1968). Test bias: Prediction of grades of negro and white students in integrated colleges. *Journal of Educational Measurement*, 5(2):115–124.

References

- Cohen, I. and Goldszmidt, M. (2004). Properties and benefits of calibrated classifiers. In *8th European Conference on Principles and Practice of Knowledge Discovery in Databases*, volume 3202, pages 125–136. Springer.
- Conway, D. A. and Roberts, H. V. (1983). Reverse regression, fairness, and employment discrimination. *Journal of Business & Economic Statistics*, 1(1):75–85.
- Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., and Huq, A. (2017). Algorithmic decision making and the cost of fairness. *arXiv*, 1701.08230.
- Cramér, H. (1928a). On the composition of elementary errors: First paper: Mathematical deductions. *Scandinavian Actuarial Journal*, 1928(1):13–74.
- Cramér, H. (1928b). On the composition of elementary errors: second paper: statistical applications. *Scandinavian Actuarial Journal*, 1928(1):141–180.
- Crossney, K. B. (2016). Redlining. <https://philadelphiaencyclopedia.org/essays/redlining/>.
- Csiszár, I. (1964). Eine informationstheoretische ungleichung und ihre anwendung auf beweis der ergodizitaet von markoffschen ketten. *Magyar Tudományos Akadémia Matematikai Kutató Intézetének Közleményei*, 8:85–108.
- Csiszár, I. (1967). On information-type measure of difference of probability distributions and indirect observations. *Studia Scientiarum Mathematicarum Hungarica*, 2:299–318.

References

- Cunningham, S. (2021). *Causal inference*. Yale University Press.
- Cuturi, M. and Doucet, A. (2014). Fast computation of wasserstein barycenters. In *International conference on machine learning*, pages 685–693. PMLR.
- Da Silva, N. (2023). *La bataille de la Sécu: une histoire du système de santé*. La fabrique éditions.
- Dall’Aglio, G. (1956). Sugli estremi dei momenti delle funzioni di ripartizione doppia. *Annali della Scuola Normale Superiore di Pisa-Classe di Scienze*, 10(1-2):35–74.
- Danskin, J. M. (1967). *The theory of max-min and its application to weapons allocation problems*. Springer.
- Dantzig, G. B. and Thapa, M. N. (1997). *Linear programming: Introduction*, volume 1. Springer.
- Darlington, R. B. (1971). Another look at “cultural fairness” 1. *Journal of educational measurement*, 8(2):71–82.
- Datta, A., Fredrikson, M., Ko, G., Mardziel, P., and Sen, S. (2017). Proxy non-discrimination in data-driven systems. *arXiv*, 1707.08120.
- Dawid, A. P. (1982). The well-calibrated bayesian. *Journal of the American Statistical Association*, 77(379):605–610.

References

- De Baere, G. and Goessens, E. (2011). Gender differentiation in insurance contracts after the judgment in case c-236/09, *Association Belge des Consommateurs Test-Achats asbl v. conseil des ministres*. *Colum. J. Eur. L.*, 18:339.
- de La Fontaine, J. (1668). *Fables*. Barbin.
- De Pril, N. and Dhaene, J. (1996). Segmentering in verzekeringen. *DTEW Research Report 9648*, pages 1–56.
- De Wit, G. and Van Eeghen, J. (1984). Rate making and society's sense of fairness. *ASTIN Bulletin: The Journal of the IAA*, 14(2):151–163.
- Dedecker, J. and Merlevède, F. (2007). The empirical distribution function for dependent variables: asymptotic and nonasymptotic results in. *ESAIM: Probability and Statistics*, 11:102–114.
- Denis, C., Elie, R., Hebiri, M., and Hu, F. (2021). Fairness guarantee in multi-class classification. *arXiv*, 2109.13642.
- Denuit, M. and Charpentier, A. (2004). *Mathématiques de l'assurance non-vie: Tome I Principes fondamentaux de théorie du risque*. Economica.
- Denuit, M., Charpentier, A., and Trufin, J. (2021). Autocalibration and tweedie-dominance for insurance pricing with machine learning. *Insurance: Mathematics & Economics*.

References

- Devroye, L., Mehrabian, A., and Reddad, T. (2018). The total variation distance between high-dimensional gaussians with the same mean. *arXiv*, 1810.08693.
- Dhaene, J., Denuit, M., Goovaerts, M. J., Kaas, R., and Vyncke, D. (2002a). The concept of comonotonicity in actuarial science and finance: applications. *Insurance: Mathematics and Economics*, 31(2):133–161.
- Dhaene, J., Denuit, M., Goovaerts, M. J., Kaas, R., and Vyncke, D. (2002b). The concept of comonotonicity in actuarial science and finance: theory. *Insurance: Mathematics and Economics*, 31(1):3–33.
- Dieterich, W., Mendoza, C., and Brennan, T. (2016). Compas risk scales: Demonstrating accuracy equity and predictive parity. *Northpointe Inc*, 7(7.4):1.
- Dressel, J. and Farid, H. (2018). The accuracy, fairness, and limits of predicting recidivism. *Science advances*, 4(1):eaao5580.
- Duncan, O. D. (1975). *Introduction to structural equation models*. Academic Press.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. (2012). Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pages 214–226.

References

- Elliott, M. N., Morrison, P. A., Fremont, A., McCaffrey, D. F., Pantoja, P., and Lurie, N. (2009). Using the census bureau's surname list to improve estimates of race/ethnicity and associated disparities. *Health Services and Outcomes Research Methodology*, 9:69–83.
- Endres, D. M. and Schindelin, J. E. (2003). A new metric for probability distributions. *IEEE Transactions on Information theory*, 49(7):1858–1860.
- Feeley, M. and Simon, J. (1994). Actuarial justice: The emerging new criminal law. *The futures of criminology*, 173:174.
- Feeley, M. M. and Simon, J. (1992). The new penology: Notes on the emerging strategy of corrections and its implications. *Criminology*, 30(4):449–474.
- Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C., and Venkatasubramanian, S. (2015). Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 259–268.
- Feller, A., Pierson, E., Corbett-Davies, S., and Goel, S. (2016). A computer program used for bail and sentencing decisions was labeled biased against blacks. it's actually not that clear. *The Washington Post*, October 17.
- Fourcade, M. and Healy, K. (2013). Classification situations: Life-chances in the neoliberal era. *Accounting, Organizations and Society*, 38(8):559–572.

References

- Fowlkes, E. B. and Mallows, C. L. (1983). A method for comparing two hierarchical clusterings. *Journal of the American Statistical Association*, 78(383):553–569.
- Fox, E. T. (2013). *'Piratical Schemes and Contracts': Pirate Articles and Their Society 1660-1730*. PhD Thesis, University of Exeter.
- François, P. (2022). Catégorisation, individualisation. retour sur les scores de crédit. *hal*, 03508245.
- Fréchet, M. (1948). Les éléments aléatoires de nature quelconque dans un espace distancié. In *Annales de l'institut Henri Poincaré*, volume 10, pages 215–310.
- Fréchet, M. (1951). Sur les tableaux de corrélation dont les marges sont données. *Annales de l'Université de Lyon, 3^e série, Sciences, Sect. A*, 14 : 53 – –77.
- Freeman, S. (2007). *Rawls*. Routledge.
- Galichon, A. (2016). *Optimal transport methods in economics*. Princeton University Press.
- Gandy, O. H. (2016). *Coming to terms with chance: Engaging rational discrimination and cumulative disadvantage*. Routledge.
- Gangbo, W. (1999). The monge mass transfer problem and its applications. *Contemporary Mathematics*, 226:79–104.
- Garrioch, D. (2011). Mutual aid societies in eighteenth-century paris. *French History & Civilization*, 4.

References

- Gebelein, H. (1941). Das statistische problem der korrelation als variations- und eigenwertproblem und sein zusammenhang mit der ausgleichsrechnung. *ZAMM - Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik*, 21(6):364–379.
- Ginsburg, M. (1940). Roman military clubs and their social functions. In *Transactions and Proceedings of the American Philological Association*, volume 71, pages 149–156. JSTOR.
- Givens, C. R. and Shortt, R. M. (1984). A class of wasserstein metrics for probability distributions. *Michigan Mathematical Journal*, 31(2):231–240.
- Glenn, B. J. (2000). The shifting rhetoric of insurance denial. *Law and Society Review*, pages 779–808.
- Glenn, B. J. (2003). Postmodernism: the basis of insurance. *Risk Management and Insurance Review*, 6(2):131–143.
- Gneiting, T., Balabdaoui, F., and Raftery, A. E. (2007). Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(2):243–268.
- Gneiting, T. and Raftery, A. E. (2005). Weather forecasting with ensemble methods. *Science*, 310(5746):248–249.
- Gneiting, T. and Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378.

References

- Goldberger, A. S. (1972). Structural equation methods in the social sciences. *Econometrica*, pages 979–1001.
- Goldman, A. (1979). *Justice and Reverse Discrimination*. Princeton University Press.
- Good, I. J. (1952). Rational decisions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 14(1):107–114.
- Gouic, T. L., Loubes, J.-M., and Rigollet, P. (2020). Projection to fairness in statistical learning. *arXiv*, 2005.11720.
- Gowri, A. (2014). *The Irony of Insurance: Community and Commodity*. PhD thesis, University of Southern California.
- Gretton, A., Smola, A., Bousquet, O., Herbrich, R., Belitski, A., Augath, M., Murayama, Y., Pauls, J., Schölkopf, B., and Logothetis, N. (2005). Kernel constrained covariance for dependence measurement. In *International Workshop on Artificial Intelligence and Statistics*, pages 112–119. PMLR.
- Grove, K. and Karcher, H. (1973). How to conjugate c 1-close group actions. *Mathematische Zeitschrift*, 132(1):11–20.
- Hacking, I. (1990). *The taming of chance*. Number 17. Cambridge University Press.

References

- Hannan, E. J. (1961). The general theory of canonical correlation and its relation to functional analysis. *Journal of the Australian Mathematical Society*, 2(2):229–242.
- Harari, Y. N. (2018). *21 Lessons for the 21st Century*. Random House.
- Harcourt, B. E. (2015). Risk as a proxy for race: The dangers of risk assessment. *Federal Sentencing Reporter*, 27(4):237–243.
- Hardt, M., Price, E., and Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29:3315–3323.
- Hardy, G. H., Littlewood, J. E., Pólya, G., Pólya, G., et al. (1952). *Inequalities*. Cambridge university press.
- Havens, H. V. (1979). Issues and needed improvements in state regulation of the insurance business. *U.S. General Accounting Office*.
- He, X. D., Kou, S., and Peng, X. (2022). Risk measures: robustness, elicibility, and backtesting. *Annual Review of Statistics and Its Application*, 9:141–166.
- Heimer, C. A. (1985). *Reactive Risk and Rational Action*. University of California Press.
- Hellinger, E. (1909). Neue begründung der theorie quadratischer formen von unendlichvielen veränderlichen. *Journal für die reine und angewandte Mathematik*, 1909(136):210–271.

References

- Henrion, M. (1988). Propagating uncertainty in bayesian networks by probabilistic logic sampling. In *Machine intelligence and pattern recognition*, volume 5, pages 149–163. Elsevier.
- Hey, R. (1814). Xviii. propositions containing some properties of tangents to circles; and of trapeziums inscribed in circles, and non-inscribed. together with propositions on the elliptic representations of circles, upon a plane surface, by perspective. *Philosophical Transactions of the Royal Society of London*, (104):348–396.
- Higham, N. J. (2008). *Functions of matrices: theory and computation*. SIAM.
- Hill, K. and White, J. (2020). Designed to deceive: do these people look real to you? *The New York Times*, 11(21).
- Hirschfeld, H. O. (1935). A connection between correlation and contingency. *Mathematical Proceedings of the Cambridge Philosophical Society*, 31(4):520–524.
- Hoeffding, W. (1940). Masstabinvariante korrelationstheorie. *Schriften des Mathematischen Instituts und Instituts für Angewandte Mathematik der Universität Berlin*, 5:181–233.
- Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67.
- Hoffman, F. L. (1896). *Race traits and tendencies of the American Negro*, volume 11. American Economic Association.

References

- Hoffman, F. L. (1918). *Mortality from respiratory diseases in dusty trades (inorganic dusts)*. Number 231. US Government Printing Office.
- Hoffman, F. L. (1931). Cancer and smoking habits. *Annals of surgery*, 93(1):50.
- Holland, P. W. (1986). Statistics and causal inference. *Journal of the American statistical Association*, 81(396):945–960.
- Horn, R. A. and Johnson, C. R. (2012). *Matrix analysis*. Cambridge university press.
- Hosmer Jr, D. W., Lemeshow, S., and Sturdivant, R. X. (2013). *Applied logistic regression*, volume 398. John Wiley & Sons.
- Hubbard, G. N. (1852). *De l'organisation des sociétés de bienfaisance ou de secours mutuels et des bases scientifiques sur lesquelles elles doivent être établies*. Paris, Guillaumin.
- Huttegger, S. M. (2013). In defense of reflection. *Philosophy of Science*, 80(3):413–433.
- Imbens, G. W. and Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- Ismay, P. (2018). *Trust among strangers: friendly societies in modern Britain*. Cambridge University Press.

References

- Jaccard, P. (1901). Étude comparative de la distribution florale dans une portion des alpes et des jura. *Bulletin de la Société Vaudoise de Sciences Naturelles*, 37:547–579.
- Jeffreys, H. (1946). An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 186(1007):453–461.
- Jensen, D. and Mayer, L. (1977). Some variational results and their applications in multiple inference. *The Annals of Statistics*, pages 922–931.
- Jordan, C. (1881). Sur la serie de fourier. *Camptes Rendus Hebdomadaires de l'Academie des Sciences*, 92:228–230.
- Kantorovich, L. and Rubinstein, G. (1958). On the space of completely additive functions. *Vestnic Leningrad Univ., Ser. Mat. Mekh. i Astron.*, 13(7):52–59. In Russian.
- Kantorovich, L. V. (1942). On the translocation of masses. In *Doklady Akademii Nauk USSR*, volume 37, pages 199–201.
- Kearns, M. and Roth, A. (2019). *The ethical algorithm: The science of socially aware algorithm design*. Oxford University Press.
- Keren, G. (1991). Calibration and probability judgements: Conceptual and methodological issues. *Acta psychologica*, 77(3):217–273.

References

- Kilbertus, N., Rojas-Carulla, M., Parascandolo, G., Hardt, M., Janzing, D., and Schölkopf, B. (2017). Avoiding discrimination through causal reasoning. *arXiv*, 1706.02744.
- Kim, P. T. (2017). Auditing algorithms for discrimination. *University of Pennsylvania Law Review*, 166:189.
- Kimeldorf, G., May, J. H., and Sampson, A. R. (1982). Concordant and discordant monotone correlations and their evaluation by nonlinear optimization. *Studies in the Management Sciences*, 19:117–130.
- Kimeldorf, G. and Sampson, A. R. (1978). Monotone dependence. *The Annals of Statistics*, pages 895–903.
- Kimeldorf, G. and Wahba, G. (1971). Some results on tchebycheffian spline functions. *Journal of mathematical analysis and applications*, 33(1):82–95.
- Kitagawa, E. M. (1955). Components of a difference between two rates. *Journal of the american statistical association*, 50(272):1168–1194.
- Kleinberg, J., Lakkaraju, H., Leskovec, J., Ludwig, J., and Mullainathan, S. (2017). Human Decisions and Machine Predictions. *The Quarterly Journal of Economics*, 133(1):237–293.
- Kleinberg, J., Mullainathan, S., and Raghavan, M. (2016). Inherent trade-offs in the fair determination of risk scores. *arXiv*, 1609.05807.

References

- Knott, M. and Smith, C. S. (1984). On the optimal mapping of distributions. *Journal of Optimization Theory and Applications*, 43:39–49.
- Knott, M. and Smith, C. S. (1994). On a generalization of cyclic monotonicity and distances among random vectors. *Linear algebra and its applications*, 199:363–371.
- Knowlton, R. E. (1978). Regents of the university of california v. bakke. *Arkansas Law Review*, 32:499.
- Koenker, R. and Bassett Jr, G. (1978). Regression quantiles. *Econometrica: journal of the Econometric Society*, pages 33–50.
- Koenker, R., Chernozhukov, V., He, X., and Peng, L. (2017). *Handbook of quantile regression*. CRC press.
- Kolmogorov, A. (1933). Sulla determinazione empirica di una legge di distribuzione. *Giornale dell'Istituto Italiano degli Attuari*, 4:83–91.
- Kolmogorov, A. N. (1930). *Sur la notion de la moyenne*. G. Bardi, tip. della R. Accad. dei Lincei.
- Komiyama, J. and Shimao, H. (2017). Two-stage algorithm for fairness-aware machine learning. *arXiv*, 1710.04924.
- Kranzberg, M. (1986). Technology and history: "kranzberg's laws". *Technology and culture*, 27(3):544–560.

References

- Krüger, F. and Ziegel, J. F. (2021). Generic conditions for forecast dominance. *Journal of Business & Economic Statistics*, 39(4):972–983.
- Kuhn, H. W. (1955). The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97.
- Kuhn, H. W. (1956). Variants of the hungarian method for assignment problems. *Naval research logistics quarterly*, 3(4):253–258.
- Kuhn, M. and Johnson, K. (2013). *Applied Predictive Modeling*. Springer.
- Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *The Annals of Mathematical Statistics*, 22(1):79–86.
- Kusner, M. J., Loftus, J., Russell, C., and Silva, R. (2017). Counterfactual fairness. In *Advances in Neural Information Processing Systems*, pages 4066–4076.
- Lancaster, H. O. (1957). Some properties of the bivariate normal distribution considered in the form of a contingency table. *Biometrika*, 44(1/2):289–292.
- Lancaster, H. O. (1958). The Structure of Bivariate Distributions. *The Annals of Mathematical Statistics*, 29(3):719 – 736.
- Larson, J., Mattu, S., Kirchner, L., and Angwin, J. (2016). How we analyzed the compas recidivism algorithm. *ProPublica*, 23-05.

References

- Leeson, P. T. (2009). The calculus of piratical consent: the myth of the myth of social contract. *Public Choice*, 139:443–459.
- Levin, D. A. and Peres, Y. (2017). *Markov chains and mixing times*, volume 107. American Mathematical Soc.
- Li, K. C.-W. (1996). The private insurance industry's tactics against suspected homosexuals: redlining based on occupation, residence and marital status. *American Journal of Law & Medicine*, 22(4):477–502.
- Lichtenstein, S., Fischhoff, B., and Phillips, L. D. (1977). Calibration of probabilities: The state of the art. *Decision making and change in human affairs*, pages 275–324.
- Lima, L. F. F. P. d., Ricarte, D. R. D., and Siebra, C. d. A. (2022). An overview on the use of adversarial learning strategies to ensure fairness in machine learning models. In *XVIII Brazilian Symposium on Information Systems*, pages 1–8.
- Lin, J. (1991). Divergence measures based on the shannon entropy. *IEEE Transactions on Information theory*, 37(1):145–151.
- Lin, P.-E. (1987). Measures of association between vectors. *Communications in Statistics-Theory and Methods*, 16(2):321–338.

References

- Linn, R. L. and Werts, C. E. (1971). Considerations for studies of test bias. *Journal of Educational Measurement*, 8(1):1–4.
- Lippert-Rasmussen, K. (2020). *Making sense of affirmative action*. Oxford University Press.
- Liu, L. T., Dean, S., Rolf, E., Simchowitz, M., and Hardt, M. (2018). Delayed impact of fair machine learning. In *International Conference on Machine Learning*, pages 3150–3158. PMLR.
- Luong, B. T., Ruggieri, S., and Turini, F. (2011). k -nn as an implementation of situation testing for discrimination discovery and prevention. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 502–510.
- Ma, L. and Koenker, R. (2006). Quantile regression methods for recursive structural equation models. *Journal of Econometrics*, 134(2):471–506.
- Madras, D., Creager, E., Pitassi, T., and Zemel, R. (2018). Learning adversarially fair and transferable representations. In *International Conference on Machine Learning*, pages 3384–3393. PMLR.
- Massey, D. S. (2007). *Categorically unequal: The American stratification system*. Russell Sage Foundation.
- McCullagh, P. and Nelder, J. (1989). *Generalized linear models*. Chapman & Hall.
- McDonald, S. (2015). Indirect gender discrimination and the ‘test-achats ruling’: an examination of the uk motor insurance market. In *Royal Economic Society Conf., Manchester*.

References

- Meinshausen, N. and Ridgeway, G. (2006). Quantile regression forests. *Journal of machine learning research*, 7(6).
- Merriam-Webster (2022). *Dictionary*. .
- Möbius, A. F. (1827). *Der barycentrische Calcul, ein Hülfsmittel zur analytischen Behandlung der Geometrie (etc.)*. Leipzig: J.A. Barth.
- Mourier, E. (1953). Éléments aléatoires dans un espace de banach. In *Annales de l'institut Henri Poincaré*, volume 13, pages 161–244.
- Mowbray, A. (1921). Classification of risks as the basis of insurance rate making with special reference to workmen's compensation. *Proceedings of the Casualty Actuarial Society*.
- Müller, A. (1997). Integral probability metrics and their generating classes of functions. *Advances in applied probability*, 29(2):429–443.
- Murphy, A. H. (1973). A new vector partition of the probability score. *Journal of Applied Meteorology and Climatology*, 12(4):595–600.
- Murphy, A. H. and Epstein, E. S. (1967). Verification of probabilistic predictions: A brief review. *Journal of Applied Meteorology and Climatology*, 6(5):748–755.
- Nagumo, M. (1930). Über eine klasse der mittelwerte. In *Japanese journal of mathematics*, volume 7, pages 71–79. The Mathematical Society of Japan.

References

- Nathan, A. (1952). *College Geometry: An Introduction to the Modern Geometry of the Triangle and the Circle*. Barnes & Noble.
- Neumann, J. v. and Morgenstern, O. (1947). *Theory of games and economic behavior*. Princeton University Press.
- Newey, W. K. and Powell, J. L. (1987). Asymmetric least squares estimation and testing. *Econometrica: Journal of the Econometric Society*, pages 819–847.
- Nielsen, F. (2022). The many faces of information geometry. *Notices of the American Mathematical Society*, 69(1):36–45.
- Nielsen, F. and Nock, R. (2013). On the chi square and higher-order chi distances for approximating f-divergences. *IEEE Signal Processing Letters*, 21(1):10–13.
- Oakes, D. (1985). Self-calibrating priors do not exist. *Journal of the American Statistical Association*, 80(390):339–339.
- Oaxaca, R. (1973). Male-female wage differentials in urban labor markets. *International Economic Review*, 14(3):693–709.
- Oberman, A. M. and Ruan, Y. (2015). An efficient linear programming method for optimal transportation. *arXiv preprint arXiv:1509.03668*.

References

- Olkin, I. and Pukelsheim, F. (1982). The distance between two random vectors with given dispersion matrices. *Linear Algebra and its Applications*, 48:257–263.
- O’Neil, C. (2016). *Weapons of math destruction: How big data increases inequality and threatens democracy*. Crown.
- Pardo, L. (2018). *Statistical inference based on divergence measures*. CRC press.
- Parlett, B. and Landis, T. (1982). Methods for scaling to doubly stochastic form. *Linear Algebra and its Applications*, 48:53–79.
- Parthasarathy, T. (1970). On games over the unit square. *SIAM Journal on Applied Mathematics*, 19(2):473–476.
- Pearl, J. (1985). Bayesian networks: A model of self-activated memory for evidential reasoning. In *Proceedings of the 7th conference of the Cognitive Science Society, University of California, Irvine, CA, USA*, pages 15–17.
- Pearl, J. (2009). *Causality*. Cambridge university press.
- Pearl, J. et al. (2009). Causal inference in statistics: An overview. *Statistics surveys*, 3:96–146.
- Pearl, J. and Mackenzie, D. (2018). *The book of why: the new science of cause and effect*. Basic books.

References

- Pearson, K. (1895). Note on regression and inheritance in the case of two parents. *Proceedings of the royal society of London*, 58(347-352):240–242.
- Perry, W. L. (2013). *Predictive policing: The role of crime forecasting in law enforcement operations*. Rand Corporation.
- Petersen, F., Mukherjee, D., Sun, Y., and Yurochkin, M. (2021). Post-processing for individual fairness. *Advances in Neural Information Processing Systems*, 34:25944–25955.
- Platt, J. et al. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3):61–74.
- Pojman, L. P. (1998). The case against affirmative action. *International Journal of Applied Philosophy*, 12(1):97–115.
- Polyanskiy, Y. and Wu, Y. (2022). *Information theory: From coding to learning*. Cambridge University Press.
- Portnoy, S. and Koenker, R. (1997). The gaussian hare and the laplacian tortoise: computability of squared-error versus absolute-error estimators. *Statistical Science*, 12(4):279–300.
- Prince, A. E. and Schwarcz, D. (2019). Proxy discrimination in the age of artificial intelligence and big data. *Iowa Law Review*, 105:1257.

References

- Prokhorov, Y. V. (1956). Convergence of random processes and limit theorems in probability theory. *Theory of Probability & Its Applications*, 1(2):157–214.
- Proschan, M. A. and Presnell, B. (1998). Expect the unexpected from conditional expectation. *The American Statistician*, 52(3):248–252.
- Puccetti, G. (2017). An algorithm to approximate the optimal expected inner product of two vectors with given marginals. *Journal of Mathematical Analysis and Applications*, 451(1):132–145.
- Puccetti, G., Rüschendorf, L., and Vanduffel, S. (2020). On the computation of wasserstein barycenters. *Journal of Multivariate Analysis*, 176:104581.
- Rao, C. R. and Mitra, S. K. (1972). Generalized inverse of a matrix and its applications. In *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Theory of Statistics*, volume 6, pages 601–621. University of California Press.
- Reichenbach, H. (1956). *The direction of time*, volume 65. University of California Press.
- Rényi, A. (1959). On measures of dependence. *Acta mathematica hungarica*, 10(3-4):441–451.
- Rényi, A. (1961). On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, volume 4, pages 547–562. University of California Press.

References

- Rhynhart, R. (2020). Mapping the legacy of structural racism in philadelphia. *Philadelphia, Office of the Controller*.
- Rizzo, M. L. and Székely, G. J. (2016). Energy distance. *wiley interdisciplinary reviews: Computational statistics*, 8(1):27–38.
- Roberts, H. V. (1968). On the meaning of the probability of rain. In *first national conference on statistical meteorology*.
- Rockafellar, R. T. (1970). *Convex Analysis*. Princeton University Press.
- Rosenbaum, P. (2018). *Observation and experiment*. Harvard University Press.
- Ross, S. M. (2014). *Introduction to probability models*. Academic press.
- Rothstein, W. G. (2003). *Public health and the risk factor: A history of an uneven medical revolution*, volume 3. Boydell & Brewer.
- Rouvroy, A., Berns, T., and Carey-Libbrecht, L. (2013). Algorithmic governmentality and prospects of emancipation. *Réseaux*, 177(1):163–196.
- Rudin, W. (1966). *Real and Complex Analysis*. McGraw-hill New York.
- Rüschendorf, L. and Uckelmann, L. (2002). On the n -coupling problem. *Journal of multivariate analysis*, 81(2):242–258.

References

- Sabbagh, D. (2007). *Equality and transparency: A strategic perspective on affirmative action in American law*. Springer.
- Santambrogio, F. (2015). *Optimal transport for applied mathematicians*, volume 55. Springer.
- Sarmanov, O. (1958a). Maximum correlation coefficient (non-symmetrical case). *Doklady Akademii Nauk SSSR*, 121(1):52–55.
- Sarmanov, O. V. (1958b). The maximum correlation coefficient (symmetrical case). *Doklady Akademii Nauk SSSR*, 120(4):715–718.
- Schanze, E. (2013). Injustice by generalization: notes on the Test-Achats decision of the european court of justice. *German Law Journal*, 14(2):423–433.
- Schauer, F. (2006). *Profiles, probabilities, and stereotypes*. Harvard University Press.
- Shannon, C. E. (1948). A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423.
- Shurbet, G., Lewis, T., and Boullion, T. (1974). Quadratic matrix equations. *The Ohio Journal of Science*, 74.
- Silver, N. (2012). *The signal and the noise: Why so many predictions fail-but some don't*. Penguin.

References

- Simon, J. (1987). The emergence of a risk society-insurance, law, and the state. *Socialist Review*, (95):60–89.
- Simon, J. (1988). The ideological effects of actuarial practices. *Law & Society Review*, 22:771.
- Sinkhorn, R. (1962). On the factor spaces of the complex doubly stochastic matrices. *Notices of the American Mathematical Society*, 9:334–335.
- Sinkhorn, R. (1964). A relationship between arbitrary positive matrices and doubly stochastic matrices. *The annals of mathematical statistics*, 35(2):876–879.
- Sinkhorn, R. (1966). A relationship between arbitrary positive matrices and stochastic matrices. *Canadian Journal of Mathematics*, 18:303–306.
- Sinkhorn, R. and Knopp, P. (1967). Concerning nonnegative matrices and doubly stochastic matrices. *Pacific Journal of Mathematics*, 21(2):343–348.
- Smirnov, N. (1948). Table for estimating the goodness of fit of empirical distributions. *The annals of mathematical statistics*, 19(2):279–281.
- Spirtes, P., Glymour, C. N., and Scheines, R. (1993). *Causation, prediction, and search*. Springer Verlag.
- Squires, G. D. (2003). Racial profiling, insurance style: Insurance redlining and the uneven development of metropolitan areas. *Journal of Urban Affairs*, 25(4):391–410.

References

- Squires, G. D. and Velez, W. (1988). Insurance redlining and the process of discrimination. *The Review of Black Political Economy*, 16(3):63–75.
- Stone, D. A. (1993). The struggle for the soul of health insurance. *Journal of Health Politics, Policy and Law*, 18(2):287–317.
- Struyck, N. (1912). *Les oeuvres de Nicolas Struyck (1687-1769): qui se rapportent au calcul des chances, à la statistique général, la statistique des décès et aux rentes viagères*. Société générale néerlandaise d'assurances sur la vie et de rentes viagères.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. (2013). Intriguing properties of neural networks. *arXiv*, 1312.6199.
- Székely, G. J. (2003). E-statistics: The energy of statistical samples. *Bowling Green State University, Department of Mathematics and Statistics Technical Report*, 3(05):1–18.
- Takatsu, A. (2008). On wasserstein geometry of the space of gaussian measures. *arXiv*, 0801.2250.
- Takatsu, A. and Yokota, T. (2012). Cone structure of l2-wasserstein spaces. *Journal of Topology and Analysis*, 4(02):237–253.
- The Zebra (2022). Car insurance rating factors by state. <https://www.thezebra.com/>.
- Thorndike, R. L. (1971). Concepts of culture-fairness. *Journal of Educational Measurement*, 8(2):63–70.

References

- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288.
- Tikhonov, A. N. (1943). On the stability of inverse problems. In *Doklady Akademii Nauk*, volume 5, pages 195–198.
- Topkis, D. M. (1998). *Supermodularity and complementarity*. Princeton university press.
- Tschantz, M. C. (2022). What is proxy discrimination? In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1993–2003.
- Turner, R. (2015). The way to stop discrimination on the basis of race. *Stanford Journal of Civil Rights & Civil Liberties*, 11:45.
- Tweedie, M. C. K. (1984). An index which distinguishes between some important exponential families. *Statistics: applications and new directions (Calcutta, 1981)*, pages 579–604.
- Vallender, S. (1974). Calculation of the wasserstein distance between probability distributions on the line. *Theory of Probability & Its Applications*, 18(4):784–786.
- Van Calster, B., McLernon, D. J., Van Smeden, M., Wynants, L., and Steyerberg, E. W. (2019). Calibration: the achilles heel of predictive analytics. *BMC medicine*, 17(1):1–7.
- Van Gerven, G. (1993). Case c-109/91, Gerardus Cornelis Ten Oever v. Stichting bedrijfspensioenfonds voor het glazenwassers-en schoonmaakbedrijf. *EUR-Lex*, 61991CC0109.

References

- Van Rijsbergen, C. (1979). Information retrieval: theory and practice. In *Proceedings of the Joint IBM/University of Newcastle upon Tyne Seminar on Data Base Systems*, volume 79.
- Vapnik, V. (1991). Principles of risk minimization for learning theory. *Advances in neural information processing systems*, 4.
- Verboven, K. (2011). Introduction: Professional collegia: Guilds or social clubs? *Ancient Society*, pages 187–195.
- Verma, S. and Rubin, J. (2018). Fairness definitions explained. In *2018 IEEE/ACM International Workshop on Software Fairness (Fairware)*, pages 1–7. IEEE.
- Villani, C. (2009). *Optimal transport: old and new*, volume 338. Springer.
- Vogel, R., Bellet, A., Clémen, S., et al. (2021). Learning fair scoring functions: Bipartite ranking under roc-based fairness constraints. In *International Conference on Artificial Intelligence and Statistics*, pages 784–792. PMLR.
- von Mises, R. (1928). *Wahrscheinlichkeit Statistik und Wahrheit*. Springer.
- von Mises, R. (1939). *Probability, statistics and truth*. Macmillan.
- von Neumann, J. (1928). Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320.

References

- Wadsworth, C., Vera, F., and Piech, C. (2018). Achieving fairness through adversarial learning: an application to recidivism prediction. *arXiv*, 1807.00199.
- Wager, S. and Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242.
- Wasserstein, L. N. (1969). Markov processes over denumerable products of spaces, describing large systems of automata. *Problemy Peredachi Informatsii*, 5(3):64–72.
- Watson, D. S., Gultchin, L., Taly, A., and Floridi, L. (2021). Local explanations via necessity and sufficiency: Unifying theory and practice. *Uncertainty in Artificial Intelligence*, pages 1382–1392.
- Wilkie, D. (1997). Mutuality and solidarity: assessing risks and sharing losses. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, 352(1357):1039–1044.
- Wilson, E. B. (1927). Probable inference, the law of succession, and statistical inference. *Journal of the American Statistical Association*, 22(158):209–212.
- Wortham, L. (1986). The economics of insurance classification: The sound of one invisible hand clapping. *Ohio State Law Journal*, 47:835.
- Wright, S. (1921). Correlation and causation. *Journal of Agricultural Research*, 20.
- Wright, S. (1934). The method of path coefficients. *The annals of mathematical statistics*, 5(3):161–215.

References

- Wu, Y., Zhang, L., Wu, X., and Tong, H. (2019). Pc-fairness: A unified framework for measuring causality-based fairness. *Advances in Neural Information Processing Systems*, 32.
- Xu, D., Yuan, S., Zhang, L., and Wu, X. (2018). Fairgan: Fairness-aware generative adversarial networks. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 570–575. IEEE.
- Xu, H., Liu, X., Li, Y., Jain, A., and Tang, J. (2021). To be robust or to be fair: Towards fairness in adversarial training. In *International conference on machine learning*, pages 11492–11501. PMLR.
- Yule, G. U. (1912). On the methods of measuring association between two attributes. *Journal of the Royal Statistical Society*, 75(6):579–652.
- Zadrozny, B. and Elkan, C. (2002). Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 694–699.
- Zafar, M. B., Valera, I., Gomez-Rodriguez, M., and Gummadi, K. P. (2017). Fairness constraints: Mechanisms for fair classification. *arXiv*, 1507.05259.
- Zafar, M. B., Valera, I., Gomez-Rodriguez, M., and Gummadi, K. P. (2019). Fairness constraints: A flexible approach for fair classification. *The Journal of Machine Learning Research*, 20(1):2737–2778.

References

- Zhang, B. H., Lemoine, B., and Mitchell, M. (2018). Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 335–340.
- Zolotarev, V. M. (1976). Metric distances in spaces of random variables and their distributions. *Mathematics of the USSR-Sbornik*, 30(3):373.