

Machine Learning for Insurers and Actuaries

Arthur Charpentier

2025



Learning, with an algorithmic perspective

(lecture 2)

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)

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 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

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities) BY-NC 4.0 3 / 385

Discrimination and Insurance

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).

Definition 2.1: 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 2.1: 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

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

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

Definition 2.2: 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.

Definition 2.3: Balance Property

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

Proposition 2.2: Law of total expectations

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

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

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

Discrimination and Insurance

Homogeneous risk sharing

| | Policyholder | Insurer |
|--------------|--------------|------------|
| Loss | $E[Y]$ | $Y - E[Y]$ |
| Average loss | $E[Y]$ | 0 |
| Variance | 0 | $Var[Y]$ |

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

Discrimination and Insurance

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)

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



Proposition 2.3: 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:

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

Discrimination and Insurance

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).

Discrimination and Insurance

Proposition 2.4: 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

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

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

Discrimination and Insurance

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.

[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.

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.).

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 ,

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

Find π such that $u(\omega - \pi) = \mathbb{E}[u(\omega - Y)]$.

Price Optimization

π such that $u(\omega - \pi) = \mathbb{E}[u(\omega - Y)]$ could be seen as the willingness to pay to transfer the risk.

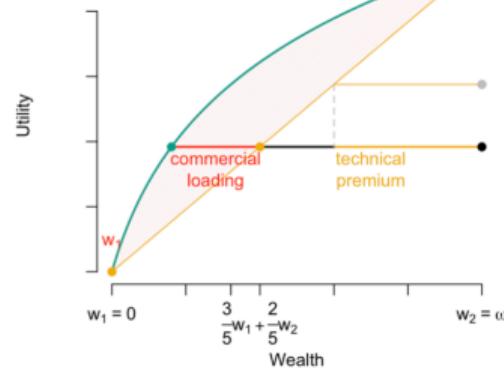
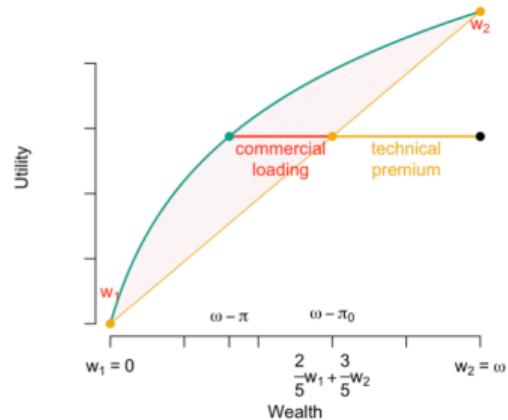
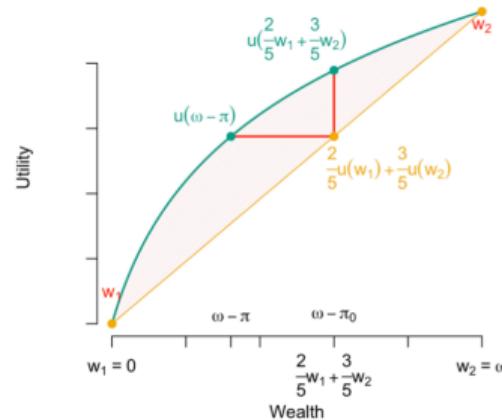
Willingness to Pay

In behavioral economics, willingness to pay (WTP) is the maximum price at or below which a consumer will definitely buy one unit of a product.[1] This corresponds to the standard economic view of a consumer reservation price. W

Price Optimization

Definition 2.4: 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 $\pi = \omega - u^{-1}(\mathbb{E}[u(\omega - Y)])$.



Price Walking

Price walking, or the loyalty penalty, is a form of price discrimination whereby longstanding, loyal customers of a service provider are charged higher prices for the same services compared to customers that have just switched to that provider. The pricing strategy is common in the insurance and telecommunications industries. It is used to acquire new customers with artificially low rates or other incentives not available to existing clients, effectively using existing customers to subsidize the prices offered to new clients. W

Ratemaking, Life Insurance

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. | | | | | | | | | | | | Table des femmes. | | | | | | | | | | | |
|-------------------|----------------|--------|----------------|--------|----------------|--------|----------------|--------|----------------|--------|----------------|-------------------|----------------|--------|----------------|--------|----------------|--------|----------------|--------|----------------|----|----|
| Années | Per- sonnes | Années | Per- sonnes | Années | Per- sonnes | Années | Per- sonnes | Années | Per- sonnes | 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 | 5 | 711 | 20 | 624 | 35 | 508 | 50 | 373 | 65 | 205 | 80 | 55 |
| 6 | 697 | 21 | 599 | 36 | 464 | 51 | 301 | 66 | 132 | 81 | 29 | 6 | 700 | 21 | 617 | 36 | 500 | 51 | 362 | 66 | 194 | 81 | 47 |
| 7 | 688 | 22 | 591 | 37 | 454 | 52 | 289 | 67 | 123 | 82 | 25 | 7 | 692 | 22 | 610 | 37 | 492 | 52 | 351 | 67 | 183 | 82 | 40 |
| 8 | 681 | 23 | 583 | 38 | 444 | 53 | 277 | 68 | 114 | 83 | 22 | 8 | 685 | 23 | 603 | 38 | 484 | 53 | 340 | 68 | 172 | 83 | 34 |
| 9 | 675 | 24 | 575 | 39 | 434 | 54 | 265 | 69 | 105 | 84 | 19 | 9 | 679 | 24 | 596 | 39 | 476 | 54 | 329 | 69 | 161 | 84 | 29 |
| 10 | 670 | 25 | 567 | 40 | 424 | 55 | 253 | 70 | 97 | 85 | 16 | 10 | 674 | 25 | 588 | 40 | 468 | 55 | 318 | 70 | 150 | 85 | 24 |
| 11 | 665 | 26 | 558 | 41 | 414 | 56 | 241 | 71 | 89 | 86 | 13 | 11 | 669 | 26 | 580 | 41 | 459 | 56 | 306 | 71 | 140 | 86 | 20 |
| 12 | 660 | 27 | 549 | 42 | 404 | 57 | 229 | 72 | 82 | 87 | 10 | 12 | 664 | 27 | 572 | 42 | 450 | 57 | 294 | 72 | 130 | 87 | 17 |
| 13 | 654 | 28 | 540 | 43 | 393 | 58 | 217 | 73 | 75 | 88 | 8 | 13 | 660 | 28 | 564 | 43 | 441 | 58 | 282 | 73 | 120 | 88 | 14 |
| 14 | 648 | 29 | 531 | 44 | 382 | 59 | 206 | 74 | 68 | 89 | 6 | 14 | 656 | 29 | 556 | 44 | 432 | 59 | 271 | 74 | 110 | 89 | 11 |
| 15 | 642 | 30 | 522 | 45 | 371 | 60 | 195 | 75 | 61 | 90 | 4 | 15 | 652 | 30 | 548 | 45 | 423 | 60 | 260 | 75 | 100 | 90 | 8 |
| 16 | 635 | 31 | 513 | 46 | 360 | 61 | 184 | 76 | 54 | 91 | 3 | 16 | 647 | 31 | 540 | 46 | 414 | 61 | 249 | 76 | 90 | 91 | 6 |
| 17 | 628 | 32 | 504 | 47 | 349 | 62 | 173 | 77 | 48 | 92 | 2 | 17 | 642 | 32 | 532 | 47 | 404 | 62 | 238 | 77 | 81 | 92 | 4 |
| 18 | 621 | 33 | 494 | 48 | 337 | 63 | 162 | 78 | 43 | 93 | 1 | 18 | 636 | 33 | 524 | 48 | 394 | 63 | 227 | 78 | 72 | 93 | 2 |
| 19 | 614 | 34 | 484 | 49 | 325 | 64 | 152 | 79 | 38 | 94 | 19 | 630 | 34 | 516 | 49 | 384 | 64 | 216 | 79 | 63 | 94 | 1 | |

Ratemaking, Life Insurance

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 | | |
|-----|-------|--------|
| x | L_x | $5p_x$ |
| 0 | 1000 | 29.0% |
| 5 | 710 | 5.6% |
| 10 | 670 | 4.2% |
| 15 | 642 | 5.5% |
| 20 | 607 | 6.6% |
| 25 | 567 | 7.9% |
| 30 | 522 | 9.2% |
| 35 | 474 | 10.5% |
| 40 | 424 | 12.5% |
| 45 | 371 | 16.6% |
| 50 | 313 | 19.2% |
| 55 | 253 | 22.9% |
| 60 | 195 | 27.2% |
| 65 | 142 | 31.7% |
| 70 | 97 | 37.1% |
| 75 | 61 | 45.9% |
| 80 | 33 | 51.5% |
| 85 | 16 | |

| women | | |
|-------|-------|--------|
| x | L_x | $5p_x$ |
| 0 | 1000 | 28.9% |
| 5 | 711 | 5.2% |
| 10 | 674 | 3.3% |
| 15 | 652 | 4.3% |
| 20 | 624 | 5.8% |
| 25 | 588 | 6.8% |
| 30 | 548 | 7.3% |
| 35 | 508 | 7.9% |
| 40 | 468 | 9.6% |
| 45 | 423 | 11.8% |
| 50 | 373 | 14.7% |
| 55 | 318 | 18.2% |
| 60 | 260 | 21.2% |
| 65 | 205 | 26.8% |
| 70 | 150 | 33.3% |
| 75 | 100 | 45.0% |
| 80 | 55 | 56.4% |
| 85 | 24 | |

freakonometrics

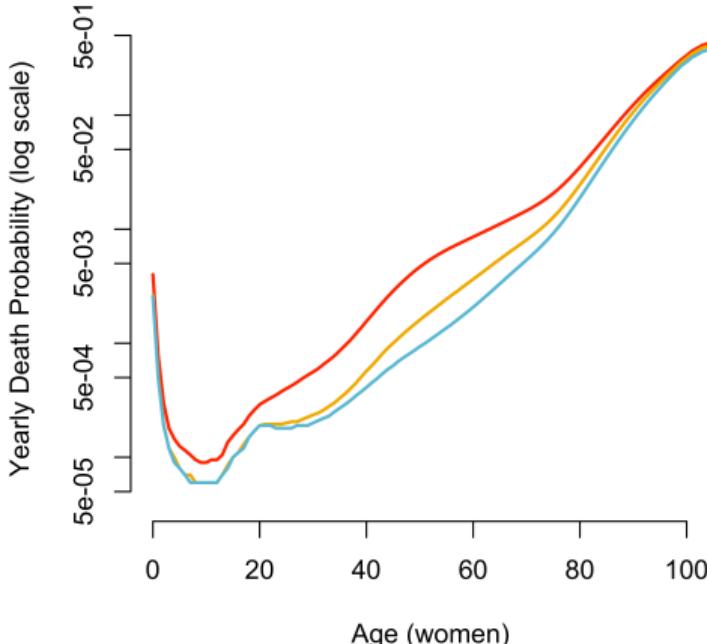
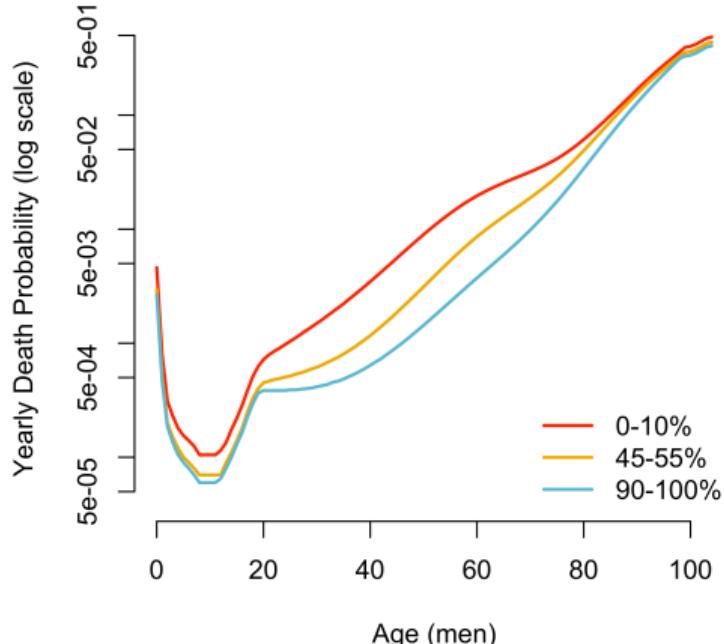
freakonometrics.hypotheses.org

Ratemaking, Life Insurance

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 | | | | women | | | |
|-----|--------|--------|---------|-------|--------|--------|---------|
| x | 0-5% | 45-50% | 95-100% | x | 0-5% | 45-50% | 95-100% |
| 0 | 100000 | 100000 | 100000 | 0 | 100000 | 100000 | 100000 |
| 10 | 99299 | 99566 | 99619 | 10 | 99385 | 99608 | 99623 |
| 20 | 99024 | 99396 | 99469 | 20 | 99227 | 99506 | 99526 |
| 30 | 97930 | 98878 | 99094 | 30 | 98814 | 99302 | 99340 |
| 40 | 95595 | 98058 | 98627 | 40 | 97893 | 98960 | 99074 |
| 50 | 90031 | 96172 | 97757 | 50 | 95021 | 97959 | 98472 |
| 60 | 77943 | 91050 | 95649 | 60 | 88786 | 95543 | 97192 |
| 70 | 59824 | 79805 | 90399 | 70 | 79037 | 90408 | 94146 |
| 80 | 38548 | 59103 | 76115 | 80 | 63224 | 79117 | 85825 |
| 90 | 13337 | 23526 | 38837 | 90 | 31190 | 45750 | 55918 |
| 100 | 530 | 1308 | 3231 | 100 | 2935 | 5433 | 8717 |

Ratemaking, Life Insurance



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

Ratemaking, Life Insurance

U.S. DECENTNIAL 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 BETWEEN TWO AGES PERIOD OF LIFE BETWEEN TWO AGES | PROPORTION DYING | OF 100,000 BORN ALIVE | | STATIONARY POPULATION | | AVERAGE REMAIN- ING LIFETIME | |
|------------------------------------------------------------------------|---------------------|-----------------------------------------------------|----------------------------------------|----------------------------|---------------------------------------------------|---------------------------------|-----|
| | | NUMBER LIVING AT BEGINNING OF AGE INTERVAL | NUMBER DYING DURING AGE INTERVAL | IN THIS AGE INTERVAL | IN THIS AND ALL SUBSEQUENT AGE INTERVALS | | |
| | | (1) | (2) | (3) | (4) | (5) | (6) |
| x to x + t | δ_x | l_x | δ'_x | l'_x | T_x | \bar{l}_x | |
| DAYS | | | | | | | |
| 0-1..... | 0.01348 | 100,000 | 1,348 | 272 | 6,411,244 | 64.11 | |
| 1-7..... | .02648 | 98,652 | 659 | 1,616 | 6,410,982 | 64.99 | |
| 7-14..... | .02442 | 97,793 | 249 | 1,211 | 6,410,753 | 64.87 | |
| 14-365..... | .01037 | 97,744 | 1,013 | 89,778 | 6,403,745 | 65.52 | |

Mortality, gender and “race”

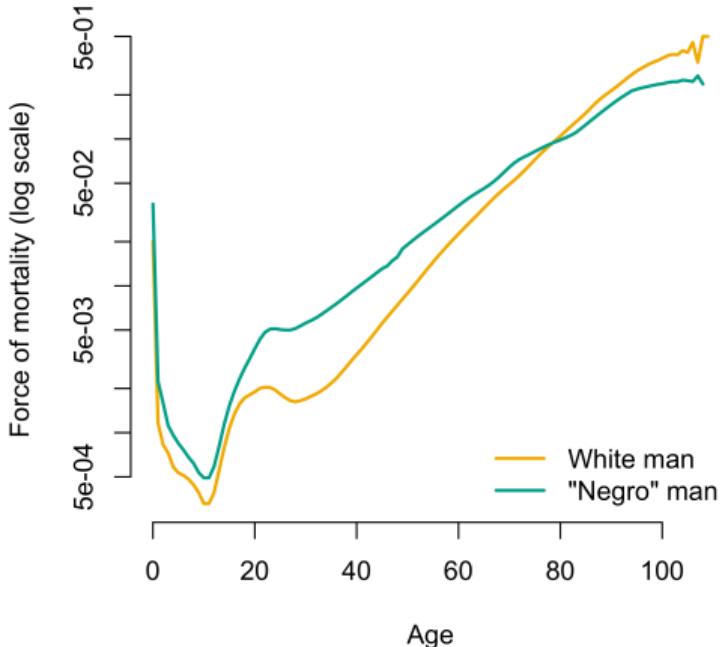
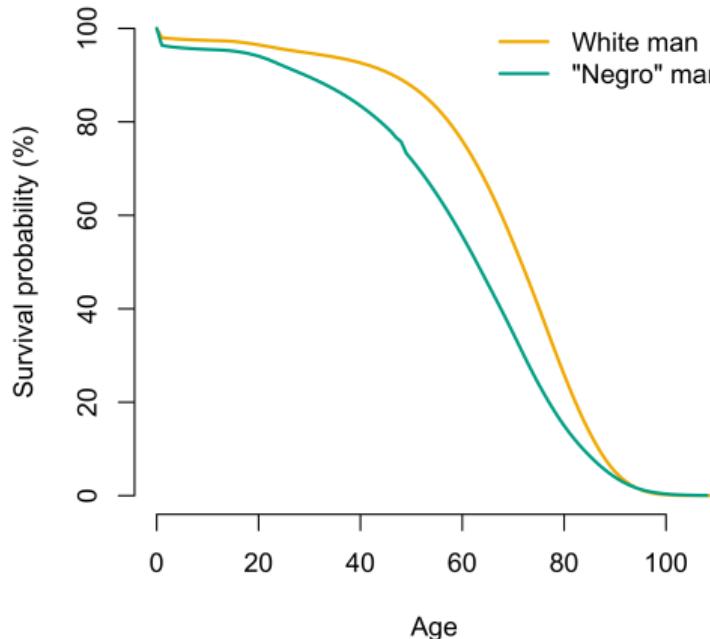
Frederick L. Hoffman

Hoffman (1896, 1918, 1931)

Ratemaking, Life Insurance

| 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% |

Ratemaking, Life Insurance



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

Ratemaking, Motor Insurance

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 \(1960\)](#), in motor insurance five risk classes can be considered,

- “*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 \(2016\)](#).

Comparative Effectiveness of Merit Rating and Class Rating

Table 1 at the end of this section shows the Canadian automobile experience¹ arranged to show what it would have looked like if there had been (1) merit rating without class rating and (2) class rating without merit rating. The premiums have been adjusted to what they would have been if all the cars had been written at 1 B rates, by use of the approximate relativities:

| Merit Rating Definition | Relativity |
|--------------------------------------------------|------------|
| A-licensed and accident free three or more years | 65 |
| X-licensed and accident free two years | 80 |
| Y-licensed and accident free one year | 90 |
| B-all others | 100 |

| Class Definitions | |
|--------------------------------------------------|-----|
| 1-pleasure, no male operator under 25 | 100 |
| 2-pleasure, non-principal male operator under 25 | 165 |
| 3-business use | 165 |
| 4-unmarried owner or principal operator under 25 | 240 |
| 5-married owner or principal operator under 25 | 165 |

freakonometrics

freakonometrics.hypotheses.org

Ratemaking, Motor Insurance

- “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)**.

Ratemaking, Motor Insurance

[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.

Ratemaking, Motor Insurance

| | CA | HI | GA | NC | NY | MA | PA | FL | TX | AL | ON | NB | NL | QC |
|--------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|---------------------------------------|-------------------------------------|-------------------------------------|---------------------------------------|-------------------------------------|---------------------------------------|-------------------------------------|-------------------------------------|
| Gender | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| Age | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> * | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> * | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| Driving experience | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Credit history | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> * | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> * | <input type="checkbox"/> | <input checked="" type="checkbox"/> * | <input type="checkbox"/> | <input checked="" type="checkbox"/> |
| Education | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Occupation | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Employment status | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Marital status | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Housing situation | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Address/ZIP code | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Insurance history | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |

CA: California, HI: Hawaii, GA: Georgia, NC: North Carolina, NY: New York, MA: Massachusetts, PA: Pennsylvania, FL: Florida, TX: Texas, AL: Alberta, ON: Ontario, NB: New-Brunswick, NL: Newfoundland-Labrador, QC: Québec.

Ordinary Least Squares

Least Squares

In regression analysis, least squares is a parameter estimation method based on minimizing the sum of the squares of the residuals (a residual being the difference between an observed value and the fitted value provided by a model) made in the results of each individual equation. \mathbf{W}

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik} + \varepsilon_i,$$

Write with a matrix form, where $p = k + 1$,

$$\underbrace{\mathbf{Y}}_{(n \times 1)} = \underbrace{\mathbf{X}}_{(n \times p)} \underbrace{\boldsymbol{\beta}}_{(p \times 1)} + \underbrace{\boldsymbol{\varepsilon}}_{(n \times 1)}$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} x_{10} & x_{11} & \dots & x_{1k} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n0} & x_{n1} & \dots & x_{nk} \end{pmatrix}, \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_k \end{pmatrix}, \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Ordinary Least Squares

\mathcal{A}_1 : Design matrix \mathbf{X} is a full rank matrix.

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

\mathcal{A}_2 : Les erreurs sont centrées, de même variance et non corrélées
 $\Leftrightarrow \mathbb{E}(\boldsymbol{\varepsilon}) = \mathbf{0}$ et $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbb{I}_n$.

$$\mathbb{E}(\hat{\boldsymbol{\beta}}) = \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}}_{=\mathbb{I}} \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \underbrace{\mathbb{E}(\boldsymbol{\varepsilon})}_{=0} = \boldsymbol{\beta}.$$

$$\text{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

$\mathcal{A}_2^{\text{gauss}}$: $\boldsymbol{\varepsilon}$ est un vecteur gaussien centré de matrice de covariance $\sigma^2 \mathbb{I}_n$, i.e.
 $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I}_n)$

Ordinary Least Squares

\mathcal{A}_3 : La matrice de design \mathbf{X} est telle que lorsque $n \rightarrow \infty$, $\frac{1}{n}(\mathbf{X}^\top \mathbf{X}) \rightarrow \mathbf{Q}$ où Q est une matrice définie positive. De plus, $h_n = \max_{1 \leq i, j \leq n} (\Pi \mathbf{x})_{ij} \rightarrow 0$ lorsque $n \rightarrow \infty$. Sous les hypothèses \mathcal{A}_1 , \mathcal{A}_2 et $\mathcal{A}_3^{\text{clt}}$, alors lorsque $n \rightarrow \infty$

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \sigma^2 \mathbf{Q}^{-1})$$

Since

$$\hat{\sigma}^{-1} \mathbf{Q}^{1/2} \sqrt{n}(\hat{\beta} - \beta) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \mathbb{I}_p)$$

ou encore

$$\hat{\sigma}^{-1}(\mathbf{X}^\top \mathbf{X})^{1/2}(\hat{\beta} - \beta) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \mathbb{I}_p)$$

soit, lorsque $n \rightarrow \infty$

$$\hat{\beta} - \beta \approx \mathcal{N}\left(0, \hat{\sigma}^2(\mathbf{X}^\top \mathbf{X})^{-1}\right).$$

Bootstrap & Tests

Consider the test of $H_0 : \beta_j = 0$,

1. compute $t_n = \frac{(\hat{\beta}_j - \beta_j)^2}{\hat{\sigma}_j^2}$

2. generate B bootstrap samples, under the null assumption H_0

3. for each bootstrap sample, compute $t_n^{(b)} = \frac{(\hat{\beta}_j^{(b)} - \hat{\beta}_j)^2}{\hat{\sigma}_j^{2(b)}}$

4. reject H_0 if $\frac{1}{B} \sum_{i=1}^B \mathbf{1}(t_n > t_n^{(b)}) < \alpha$.

Bootstrap & Tests

What does "generate B bootstrap samples, under the null assumption H_0 " mean ?

Example : $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ with $H_0 : \beta_1 = 0$.

- 2.1. Estimate the model under H_0 , i.e. $y_i = \beta_0 + \eta_i$, and save $\{\hat{\eta}_1, \dots, \hat{\eta}_n\}$
- 2.2. Define $\tilde{\eta} = \{\tilde{\eta}_1, \dots, \tilde{\eta}_n\}$ with $\tilde{\eta} = \sqrt{\frac{n}{n-1}} \hat{\eta}$
- 2.3. Draw (with replacement) residuals $\tilde{\eta}^{(b)} = \{\tilde{\eta}_1^{(b)}, \dots, \tilde{\eta}_n^{(b)}\}$
- 2.4. Set $y_i^{(b)} = \hat{\beta}_0 + \tilde{\eta}_i^{(b)}$
- 2.5. Estimate the regression model $y_i^{(b)} = \beta_0^{(b)} + \beta_1^{(b)} x_i + \varepsilon_i^{(b)}$
3. for each bootstrap sample, compute $t_n^{(b)} = \frac{(\hat{\beta}_j^{(b)} - \hat{\beta}_j)^2}{\hat{\sigma}_j^{2(b)}}$
4. reject H_0 if $\frac{1}{B} \sum_{i=1}^B \mathbf{1}(t_n > t_n^{(b)}) < \alpha$.

Regression with Weights

Suppose that instead of $\{(y_1, \mathbf{x}_1), (y_2, \mathbf{x}_2), \dots, (y_n, \mathbf{x}_n)\}$ we have $\{(y_1, \mathbf{x}_1, \omega_1), (y_2, \mathbf{x}_2, \omega_2), \dots, (y_n, \mathbf{x}_n, \omega_n)\}$. To solve

$$\sum_{i=1}^n \omega_i (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{W}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

where $\mathbf{W} = \text{diag}(\boldsymbol{\omega})$.

First order condition is $(\mathbf{X}^\top \mathbf{W} \mathbf{X}) \hat{\boldsymbol{\beta}} = \mathbf{X}^\top \mathbf{W} \mathbf{y}$.

Local Regression

Nadaraya-Watson is a simple nonparametric regression method using kernel-weighted averages

It is based on the idea that the conditional mean $\mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$ is a “local average” of Y in the neighborhood of \mathbf{x} .

Following the “law of the unconscious statistician,” in Schervish and DeGroot (2014), for some \mathbf{x} ,

$$\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] = \lim_{h \rightarrow 0} \mathbb{E}[Y | \mathbf{X} \in \mathcal{B}_h(\mathbf{x})] \text{ where } \mathcal{B}_h(\mathbf{x}) = \{\mathbf{x}' \in \mathcal{X} : \|\mathbf{x} - \mathbf{x}'\| \leq h\}$$

Local Regression

Definition 2.5: Nadaraya-Watson Estimator

Given data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, the estimate at point \mathbf{x} is:

$$\hat{m}(\mathbf{x}) = \frac{\sum_{i=1}^n K_h(\mathbf{x} - \mathbf{x}_i)y_i}{\sum_{i=1}^n K_h(\mathbf{x} - \mathbf{x}_i)}$$

where $K_h(u) = \frac{1}{h}K\left(\frac{u}{h}\right)$ for some kernel K (on \mathcal{X}) and with bandwidth h .

In the univariate case, common kernels are Gaussian, Epanechnikov, Uniform

Local Regression

It's a weighted average of the y_i 's, where weights depend on distance from x .

$$\hat{m}(x) = \underset{\beta_0 \in \mathbb{R}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n w_i(x)(y_i - \beta_0)^2 \right\} \text{ where } w_i(x) = \frac{K_h(x - x_i)}{\sum_{i=1}^n K_h(x - x_i)}$$

More generally, we can consider some local regression, also known as **LOESS** (Locally Estimated Scatterplot Smoothing) or **LOWESS**

Local Regression

Definition 2.6: LOESS Estimator

Given data $(x_1, y_1), \dots, (x_n, y_n)$, the estimate at point x is:

$$\hat{m}(x) = \hat{\beta}_0 + \hat{\beta}_1 x, \text{ where } (\hat{\beta}_0, \hat{\beta}_1) = \underset{(\beta_0, \beta_1) \in \mathbb{R}^2}{\operatorname{argmin}} \left\{ \sum_{i=1}^n w_i(x)(y_i - \beta_0 - \beta_1 x)^2 \right\}$$

where

$$w_i(x) = \frac{K_h(x - x_i)}{\sum_{i=1}^n K_h(x - x_i)}$$

where $K_h(u) = \frac{1}{h} K\left(\frac{u}{h}\right)$ for some kernel K (on \mathcal{X}) and with bandwidth h .

Local regression is very sensitive in hyper-parameter selection.

Isotonic Regression

Isotonic Regression is a type of regression that fits a non-decreasing (or non-increasing) function to data. It is used when there is a known order or monotonicity constraint in the response variable.

Definition 2.7: Isotonic Regression

Given data points $(x_1, y_1), \dots, (x_n, y_n)$ with $x_1 < x_2 < \dots < x_n$, the isotonic regression is an interpolation between $(x_1, \hat{y}_1), \dots, (x_n, \hat{y}_n)$ such that:

$$\min_{\hat{y}} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

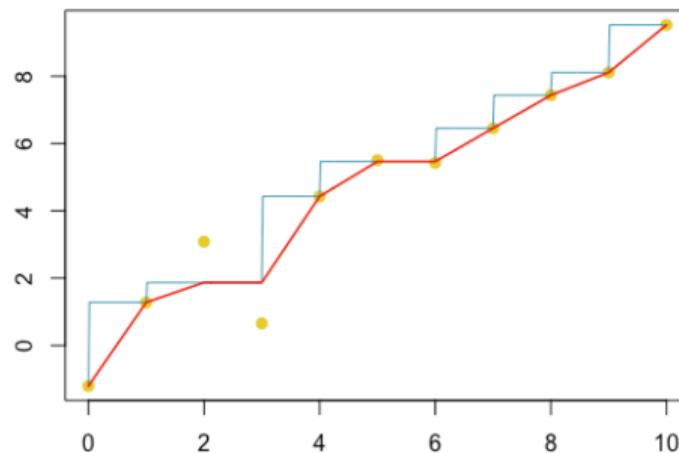
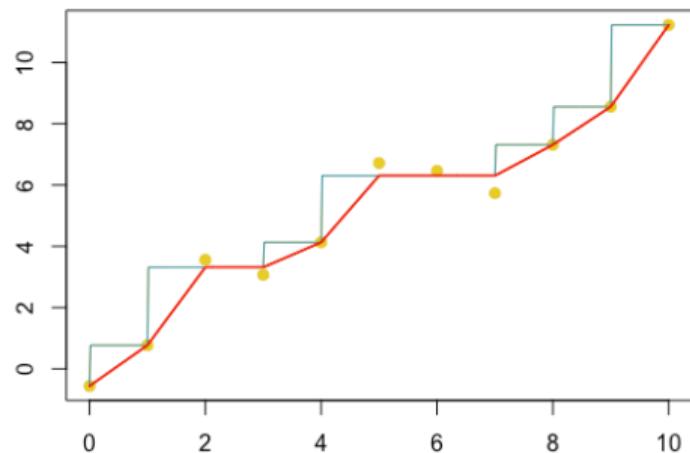
subject to:

$$\hat{y}_1 \leq \hat{y}_2 \leq \dots \leq \hat{y}_n$$

Unique solution via Pool Adjacent Violators Algorithm (PAVA), [De Leeuw et al. \(2010\)](#),

Isotonic Regression

```
1 > ir = isoreg(x=d$x, y=d$y)
2 > isofit1 = as.stepfun(ir)
3 > isofit2 = function(u) approx(ir$x, ir$vf, xout=u)$y
```



Definition 2.8: Profile likelihood

Suppose that $\theta = (\theta_1, \theta_2)$, and define

$$\mathcal{L}_p(\theta_1) = \mathcal{L}(\theta_1, \theta_2^*(\theta_1)) = \sup_{\theta_2} \{\mathcal{L}(\theta_1, \theta_2)\}$$

$\hat{\theta}_1^* \in \operatorname{argmax}\{\mathcal{L}_p(\theta_1)\}$ is called a maximum profile likelihood estimator of θ_1 .

Correlated Features, Profile Likelihood and Frisch–Waugh–Lovell

Suppose that $\beta = (\beta_1, \beta_2)$, associated to design matrix $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)$.

$$\beta_2^*(\beta_1) = (\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top (\mathbf{y} - \mathbf{X}_1 \beta_1).$$

and then, one can prove that

$$\hat{\beta}_1^* = (\mathbf{X}_1^\top (\mathbb{I} - \boldsymbol{\Pi}_{\mathbf{X}_2}) \mathbf{X}_1)^{-1} \mathbf{X}_1^\top (\mathbb{I} - \boldsymbol{\Pi}_{\mathbf{X}_2}) \mathbf{y}$$

where $\boldsymbol{\Pi}_{\mathbf{X}_2} = \mathbf{X}_2 (\mathbf{X}_2^\top \mathbf{X}_2)^{-1} \mathbf{X}_2^\top$ is the projection matrix on \mathbf{X}_2 .

Correlated Features, Profile Likelihood and Frisch–Waugh–Lovell

Under-identification is obtained when the true model is $y = \beta_0 + \mathbf{x}_1^\top \boldsymbol{\beta}_1 + \mathbf{x}_2^\top \boldsymbol{\beta}_2 + \varepsilon$, but we estimate $y = b_0 + \mathbf{x}_1^\top \mathbf{b}_1 + \eta$.

Maximum likelihood estimator for \mathbf{b}_1 is

$$\begin{aligned}\hat{\mathbf{b}}_1 &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y} \\ &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top [\mathbf{X}_{1,i} \boldsymbol{\beta}_1 + \mathbf{X}_{2,i} \boldsymbol{\beta}_2 + \varepsilon] \\ &= \boldsymbol{\beta}_1 + \underbrace{(\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \boldsymbol{\beta}_2}_{\beta_{12}} + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon}_{\nu_i}\end{aligned}$$

so that $\mathbb{E}[\hat{\mathbf{b}}_1] = \boldsymbol{\beta}_1 + \boldsymbol{\beta}_{12}$,

and the bias is null when $\mathbf{X}_1^\top \mathbf{X}_2 = \mathbf{0}$ i.e. $\mathbf{X}_1 \perp \mathbf{X}_2$.

Correlated Features, Profile Likelihood and Frisch–Waugh–Lovell

Assume that the true model is

$$\mathbf{y} = \beta_0 + \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \varepsilon$$

If $\mathbf{X}_1^\top \mathbf{X}_2 = \mathbf{0}$, $\hat{\boldsymbol{\beta}}_2$ can be estimated using

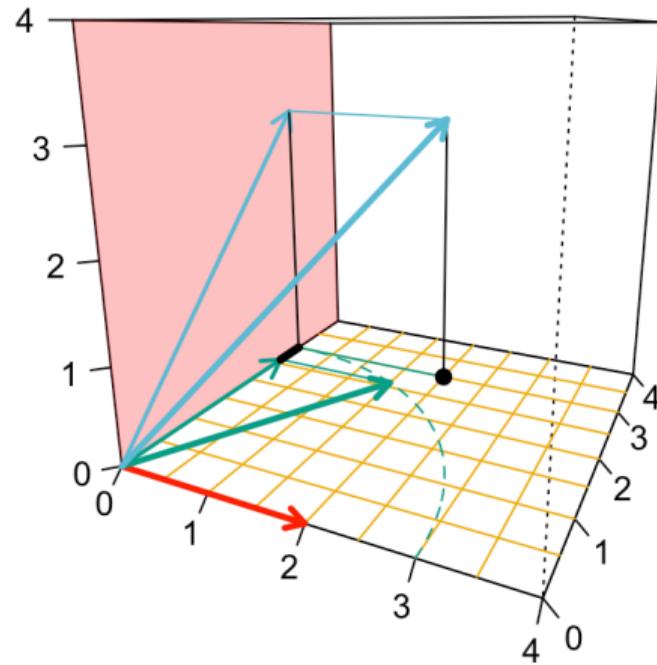
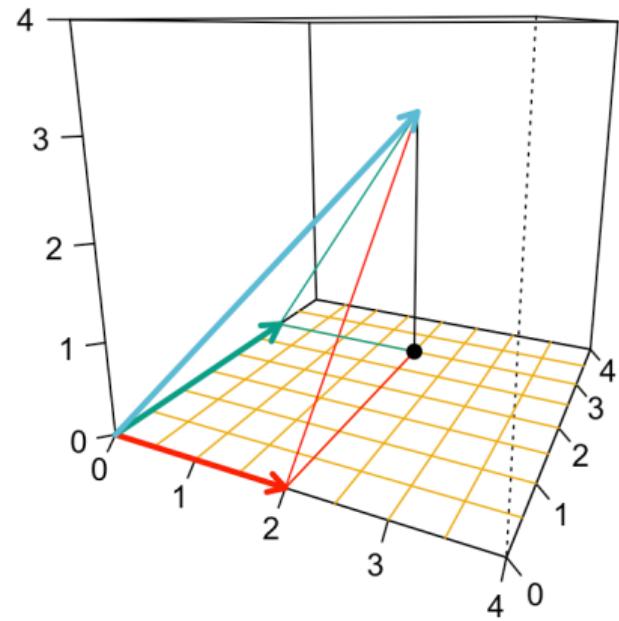
$$(b_0, \hat{\boldsymbol{\beta}}_2) = (\mathbf{X}'_2^\top \mathbf{X}'_2)^{-1} \mathbf{X}'_2^\top \mathbf{y} \text{ where } \mathbf{X}'_2 = [\mathbf{1} | \mathbf{X}_2]$$

Otherwise, let $\mathbf{y}_2^* = \Pi_{\mathcal{X}_1^\perp} \mathbf{y}$ and $\mathbf{X}_2^* = \Pi_{\mathcal{X}_1^\perp} \mathbf{X}_2$, then

$$\hat{\boldsymbol{\beta}}_2 = [\mathbf{X}_2^{*\top} \mathbf{X}_2^*]^{-1} \mathbf{X}_2^{*\top} \mathbf{y}_2^*$$

$\mathbf{X}_2^* = \mathbf{X}_2$ if $\mathbf{X}_1 \perp \mathbf{X}_2$,

Correlated Features, Profile Likelihood and Frisch–Waugh–Lovell



Binary Classifier, Logistic Regression

In statistics, A logistic model (or logit model) is a statistical model that models the log-odds of an event as a linear combination of one or more independent variables. W

$\{(y_i, \mathbf{x}_i)\}$ sample, realizations of i.i.d. vectors (Y_i, \mathbf{X}_i) where $(Y|\mathbf{X} = \mathbf{x}) \sim \mathcal{B}(p(\mathbf{x}))$,

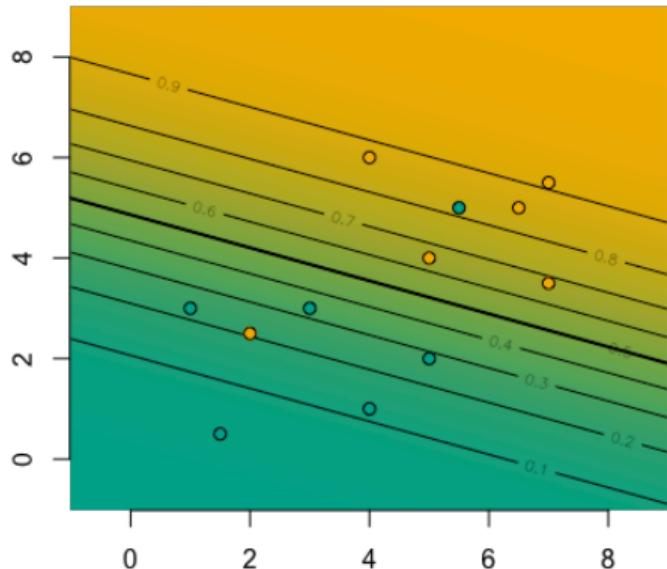
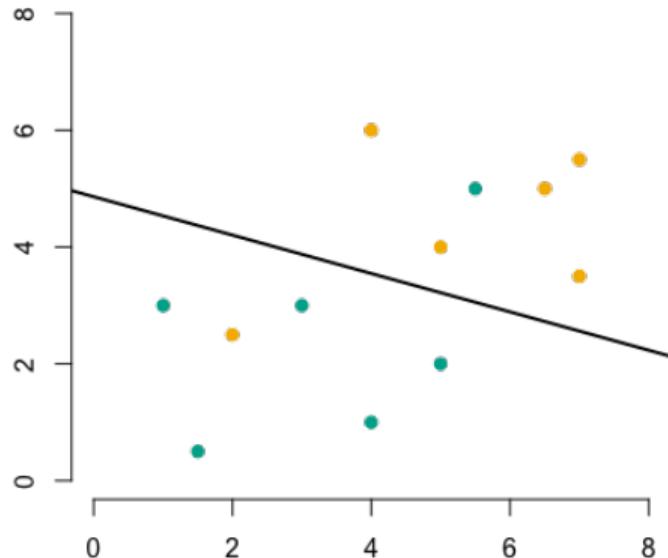
$$\text{where } \mathbb{P}[Y = 1 | \mathbf{X} = \mathbf{x}] = p(\mathbf{x}) = \frac{e^{\mathbf{x}^\top \beta}}{1 + e^{\mathbf{x}^\top \beta}} = s_{\text{logistic}}(\mathbf{x})$$

Inference: maximum likelihood

$$\log \mathcal{L}(\beta) = \sum_{i=1}^n y_i \log p_i + (1 - y_i) \log(1 - p_i), \text{ where } p_i = g^{-1}(\mathbf{x}^\top \beta) = \frac{\exp[\mathbf{x}^\top \beta]}{1 + \exp[\mathbf{x}^\top \beta]}$$

Classical first order condition, $\nabla \log \mathcal{L}(\hat{\beta}) = \mathbf{0}$ is $\mathbf{X}^\top (\mathbf{y} - \hat{\mathbf{p}}) = \mathbf{0}$ where $\hat{\mathbf{p}} = g^{-1}(\mathbf{X}\hat{\beta})$.

Binary Classifier, Logistic Regression



Binary Classifier, Linear Discriminant Analysis

$\{(y_i, \mathbf{x}_i)\}$ sample, realizations of i.i.d. vectors (Y_i, \mathbf{X}_i) where $(\mathbf{X}|Y=y) \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma})$.

$$\mathbb{P}[Y_i = y | \mathbf{X}_i = \mathbf{x}] = \frac{f(\mathbf{x}|y) \cdot \pi(y)}{f(\mathbf{x})}.$$

$$\log \left[\frac{\mathbb{P}[Y_i = 1 | \mathbf{X}_i = \mathbf{x}]}{\mathbb{P}[Y_i = 0 | \mathbf{X}_i = \mathbf{x}]} \right] = \log \left[\frac{f(\mathbf{x}|1) \cdot \pi(1)}{f(\mathbf{x}|0) \cdot \pi(0)} \right]$$

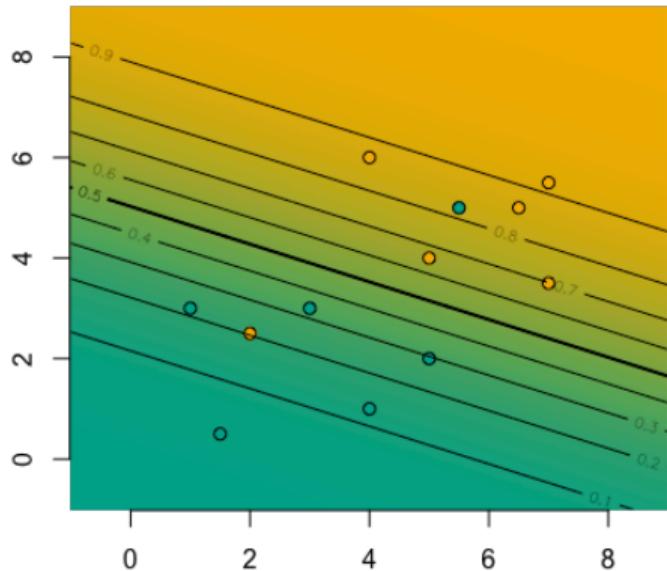
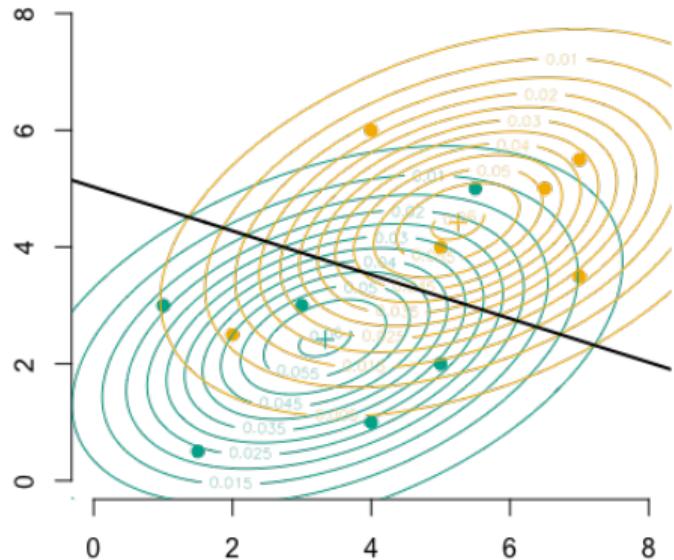
$$= \mathbf{x}^\top \underbrace{\boldsymbol{\Sigma}^{-1}[\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0]}_{\beta} + \underbrace{\left(-\frac{1}{2} [\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0]^\top \boldsymbol{\Sigma}^{-1} [\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0] + \log \frac{\pi(1)}{\pi(0)} \right)}_{\text{constant } \beta_0}$$

that is linear. And

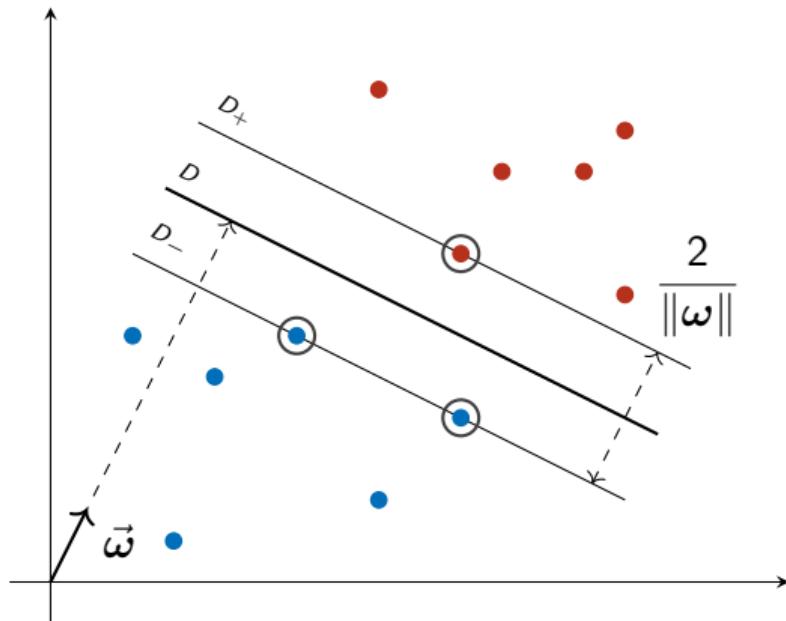
$$\mathbb{P}[Y_i = 1 | \mathbf{X}_i = \mathbf{x}] = \frac{e^{\mathbf{x}^\top \beta + \beta_0}}{1 + e^{\mathbf{x}^\top \beta + \beta_0}} = s_{\text{LDA}}(\mathbf{x})$$

where $\boldsymbol{\mu}_0$, $\boldsymbol{\mu}_1$ and $\boldsymbol{\Sigma}$ are estimated from the data (and $\boldsymbol{\Sigma}$ should be the same in the two groups).

Binary Classifier, Linear Discriminant Analysis



Binary Classifier, SVM (Support Vector Machine)



$\{(y_i, \mathbf{x}_i)\}$ collection of observations,
where $y_i \in \{-1, +1\}$.

Consider some separating hyperplane
 $D = \{\mathbf{x} : b + \mathbf{x}^\top \boldsymbol{\omega} = 0\}$.

Parameters b and $\boldsymbol{\omega}$ are solution of

$$\operatorname{argmax}_{b, \boldsymbol{\omega}} \left\{ \min_i \{ \|\mathbf{x} - \mathbf{x}_i\| : x \in \mathbb{R}^n, \text{ s.t. } b + \mathbf{x}^\top \boldsymbol{\omega} = 0 \} \right\}$$

Binary Classifier, SVM (Support Vector Machine)

$$\underset{b, \omega}{\operatorname{argmax}} \left\{ \min_i \{ \|x - x_i\| : x \in \mathbb{R}^n, \text{ s.t. } b + x^\top \omega = 0 \} \right\}$$

If $m(x) = b + x^\top \omega$, the classifier is $\hat{y} = \operatorname{sign}(m(x)) = \begin{cases} +1 & \text{if } m(x) > 0 \\ -1 & \text{if } m(x) < 0 \end{cases}$

All points are properly classified if $y_i \cdot m(x) > 0, \forall i$

Binary Classifier, SVM (Support Vector Machine)

- Perceptron problem

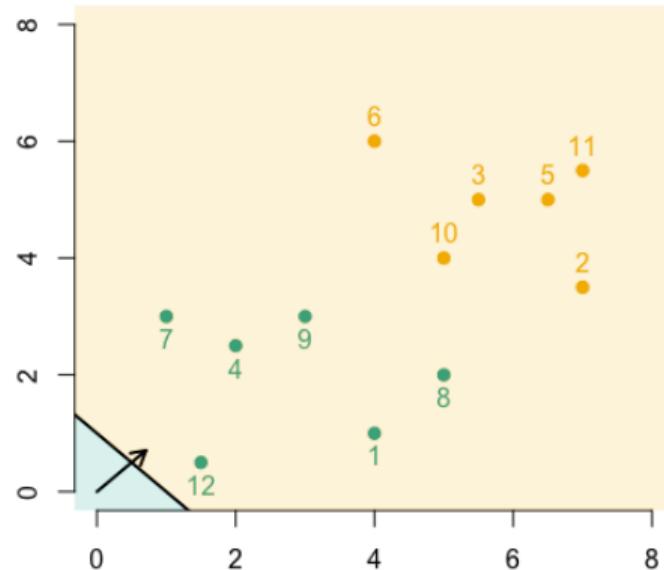
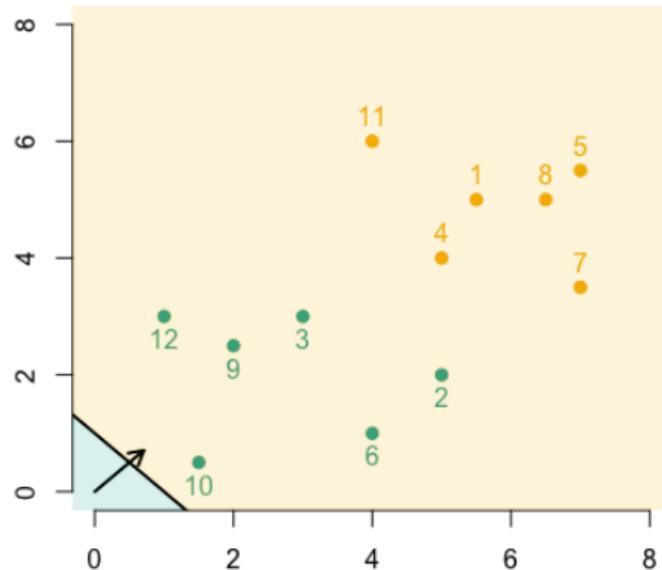
Algorithm 1: Perceptron

- 1 initialization : $\omega \leftarrow \mathbf{0}$; $\gamma \leftarrow 0.1$ (learning rate);
- 2 **for** $t=1,2,\dots$ **do**
- 3 **for** $i = 1, 2, \dots, n$ **do**
- 4 **if** i missclassified, $\omega \leftarrow \omega + \gamma \cdot \text{sign}(m(\mathbf{x}_i)) \cdot \mathbf{x}_i$
- 5 **until** all the data is correctly classified

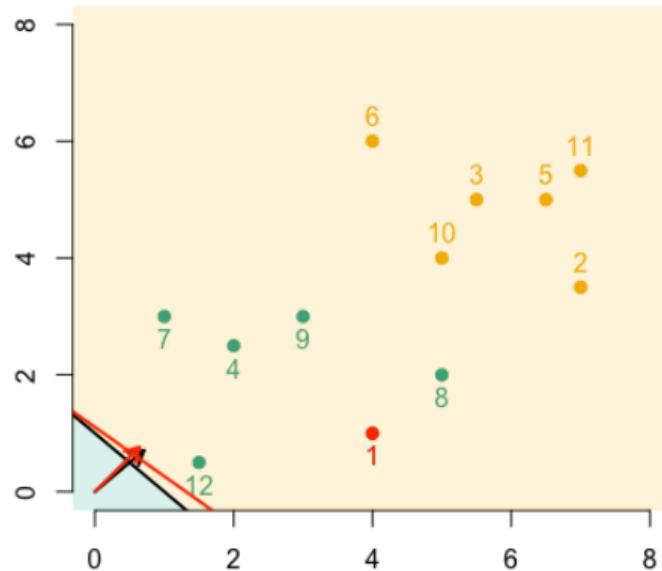
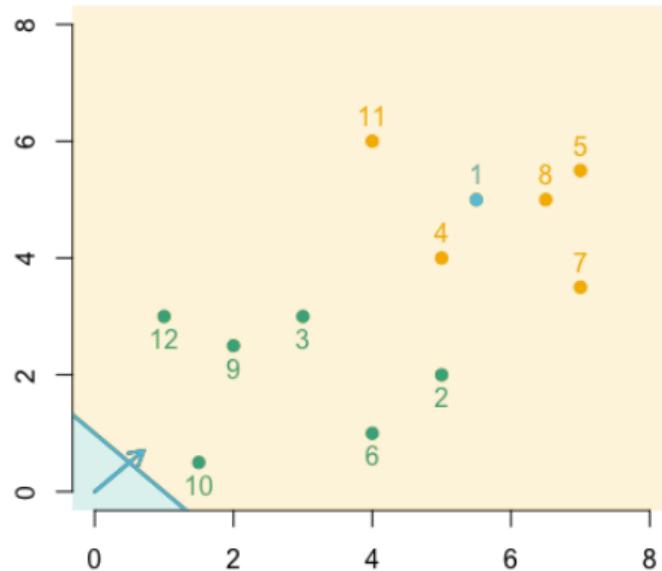
We cycle through the data points (each cycle is coined “**epoch**”)

If the data is linearly separable, then the algorithm will converge (but convergence can be slow)

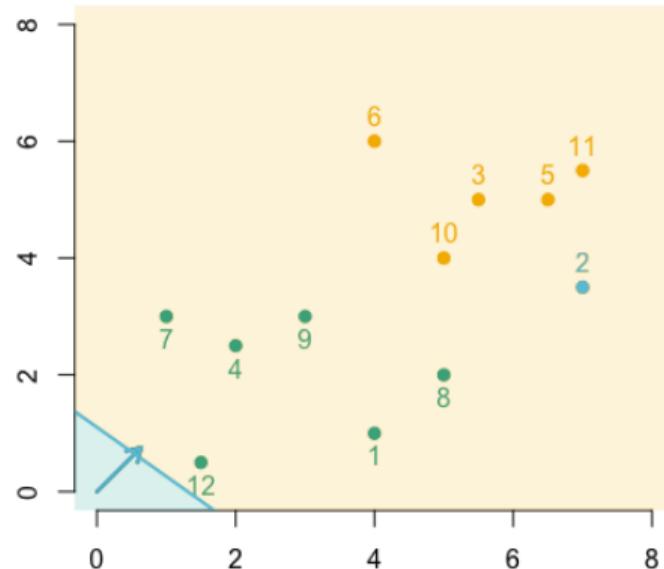
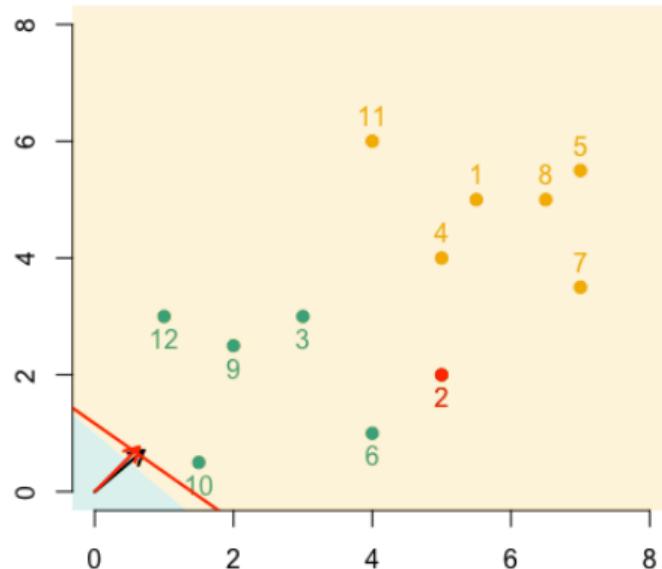
Binary Classifier, SVM (Support Vector Machine)



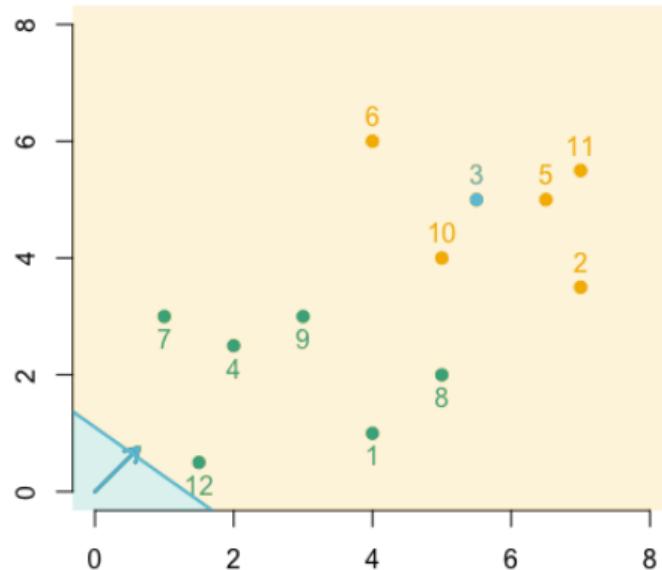
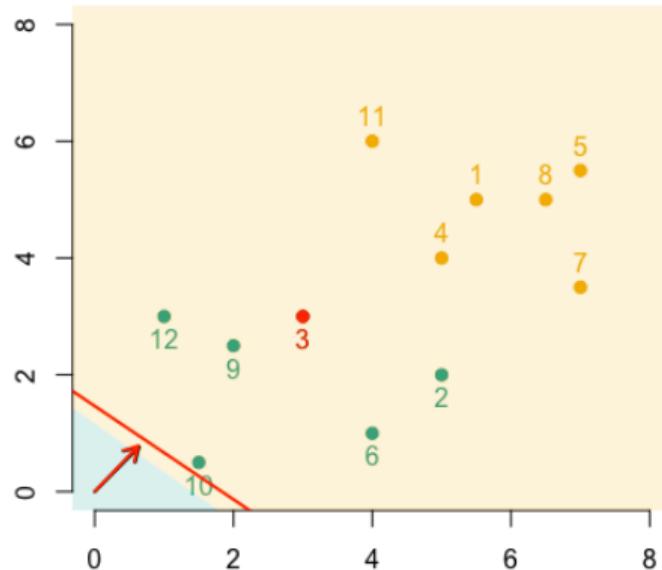
Binary Classifier, SVM (Support Vector Machine)



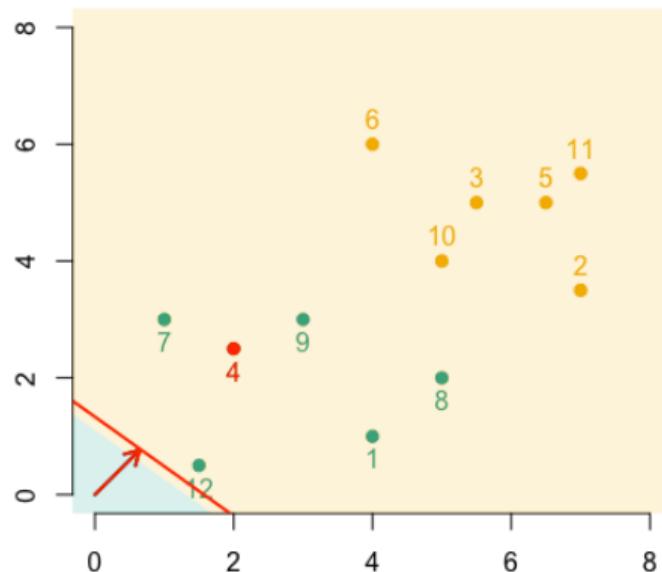
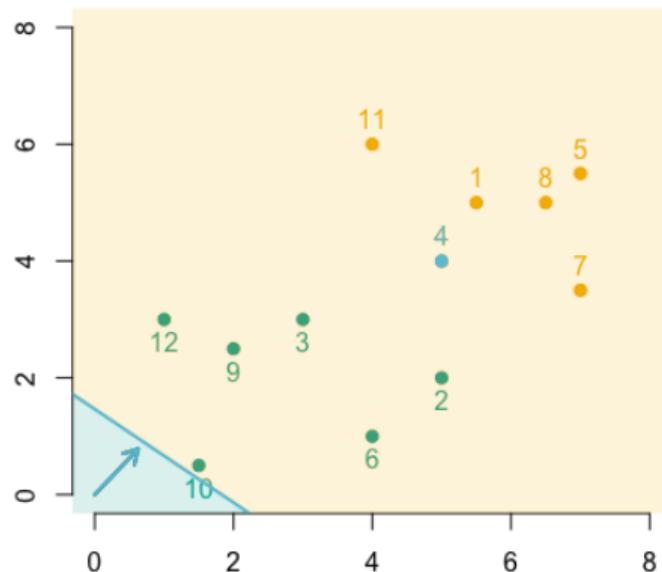
Binary Classifier, SVM (Support Vector Machine)



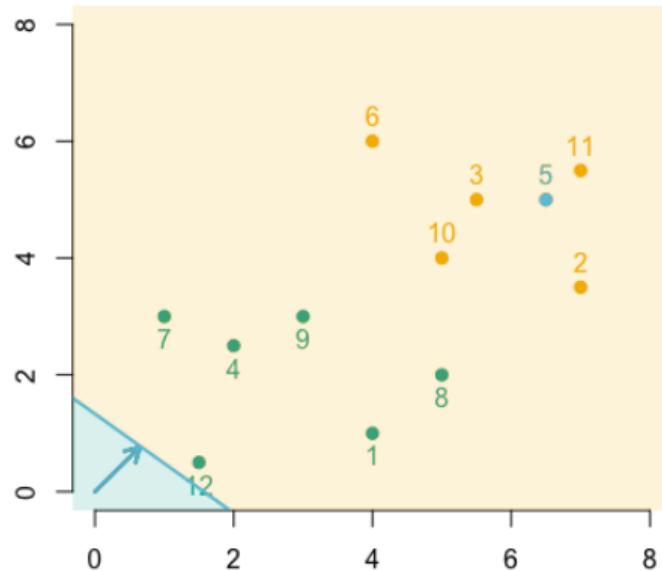
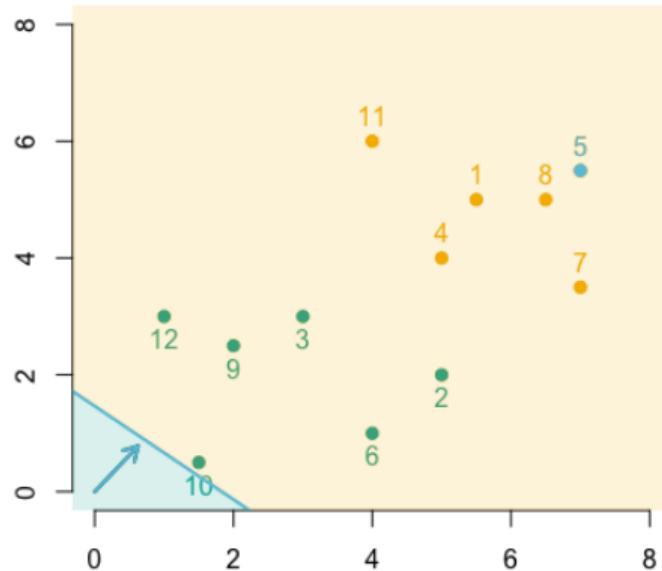
Binary Classifier, SVM (Support Vector Machine)



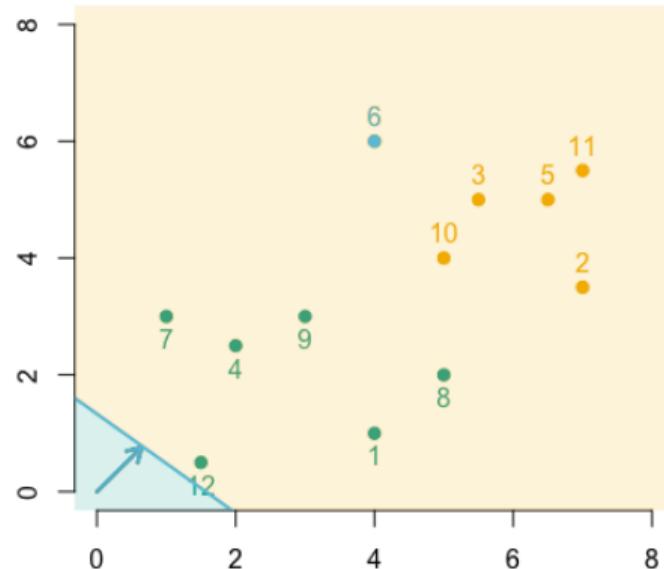
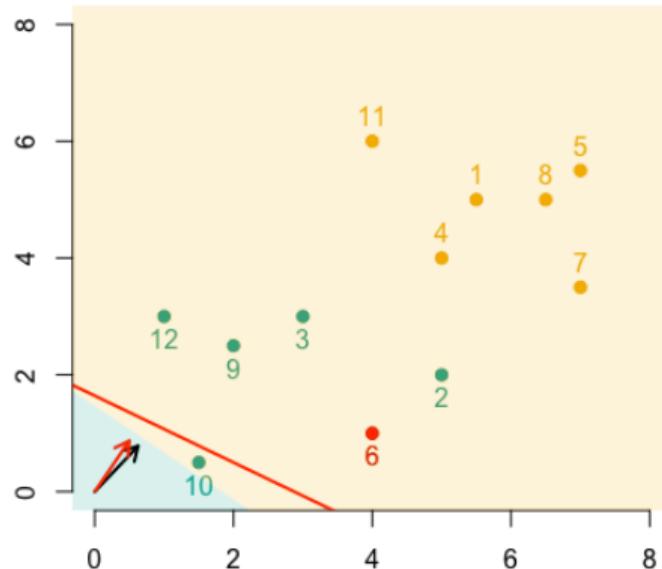
Binary Classifier, SVM (Support Vector Machine)



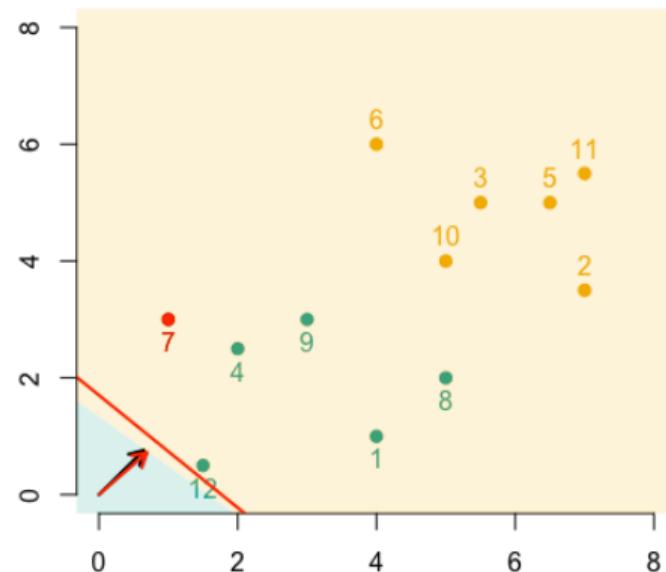
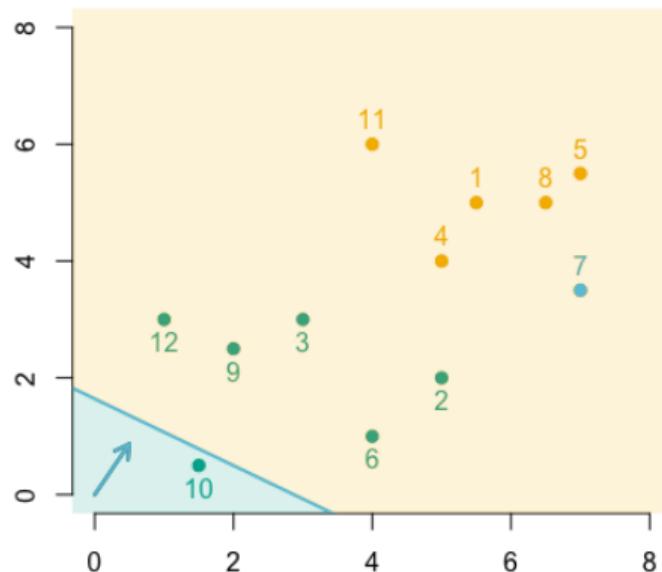
Binary Classifier, SVM (Support Vector Machine)



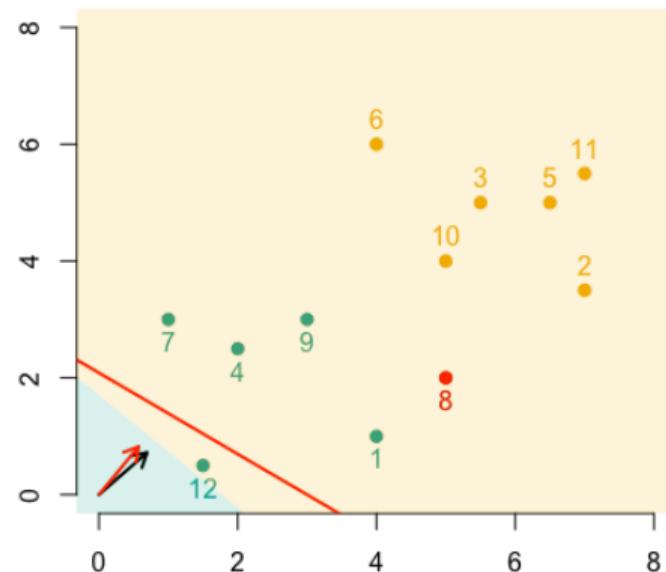
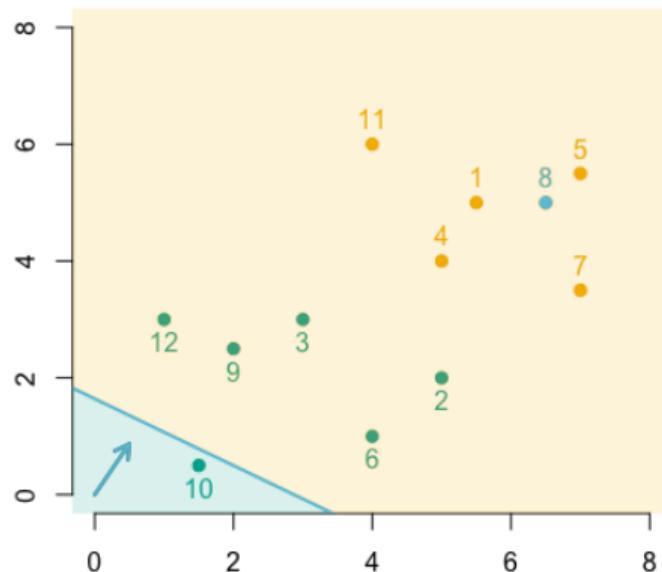
Binary Classifier, SVM (Support Vector Machine)



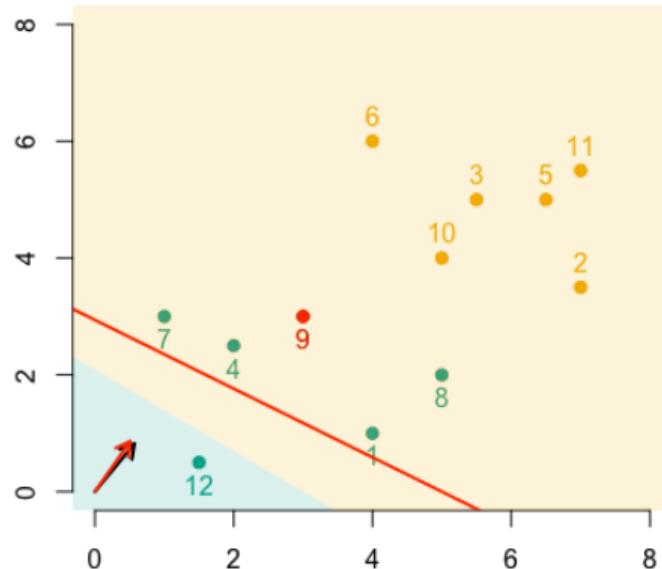
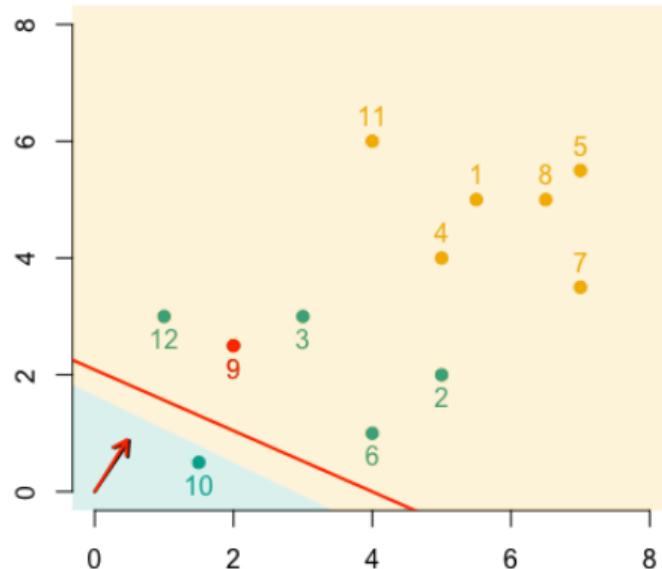
Binary Classifier, SVM (Support Vector Machine)



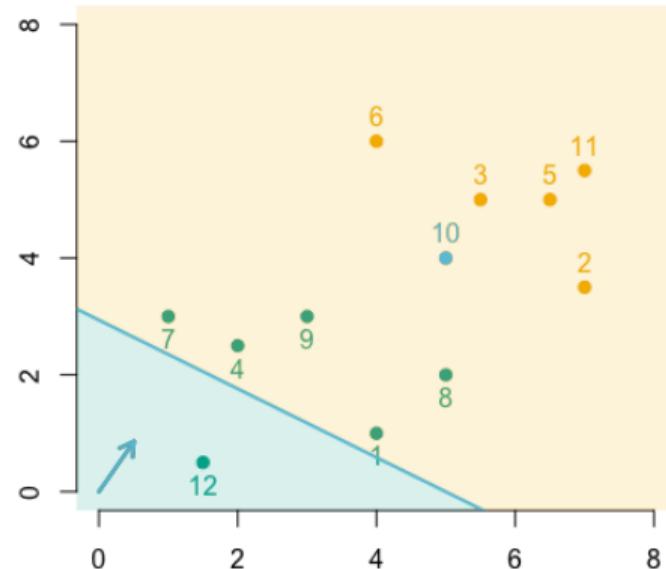
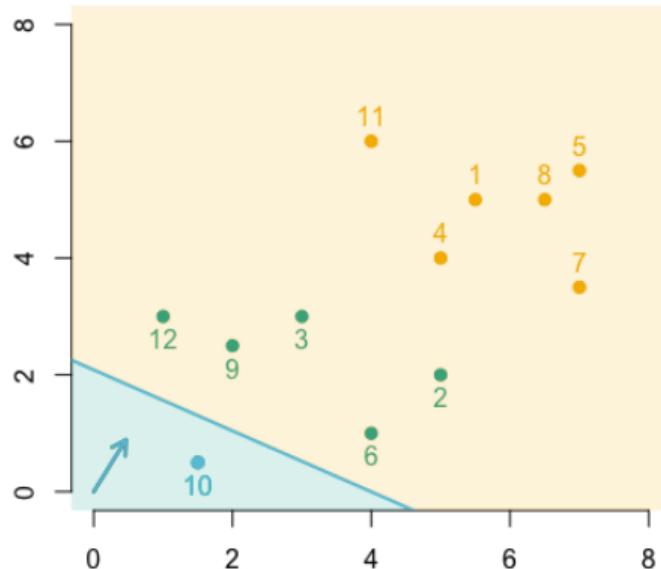
Binary Classifier, SVM (Support Vector Machine)



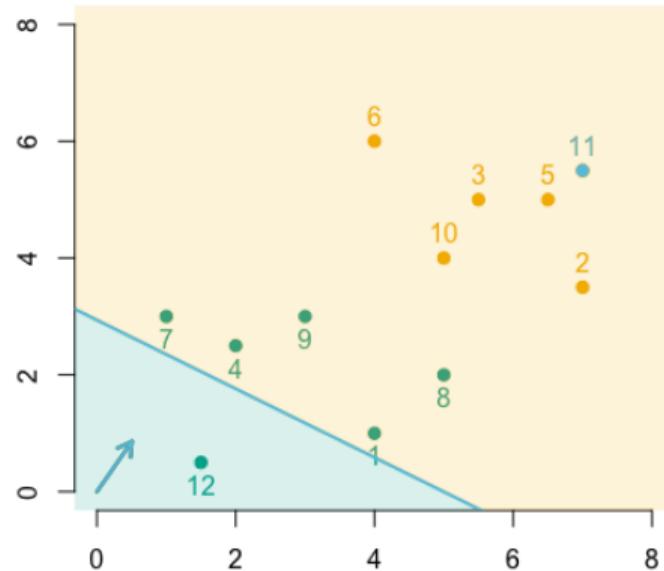
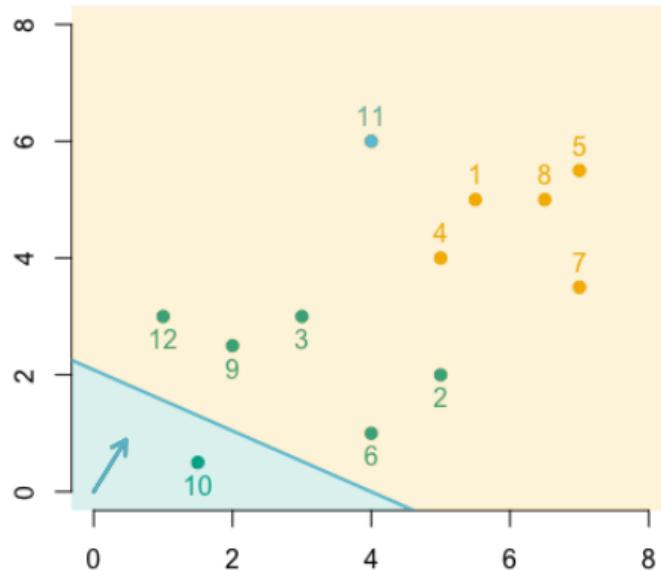
Binary Classifier, SVM (Support Vector Machine)



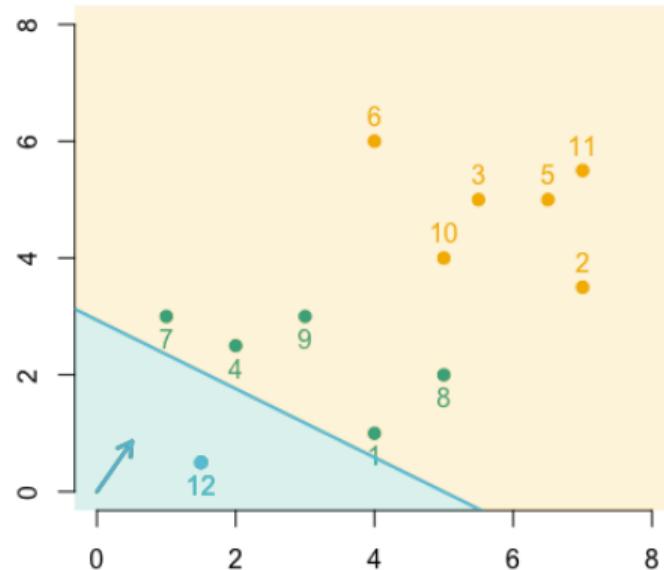
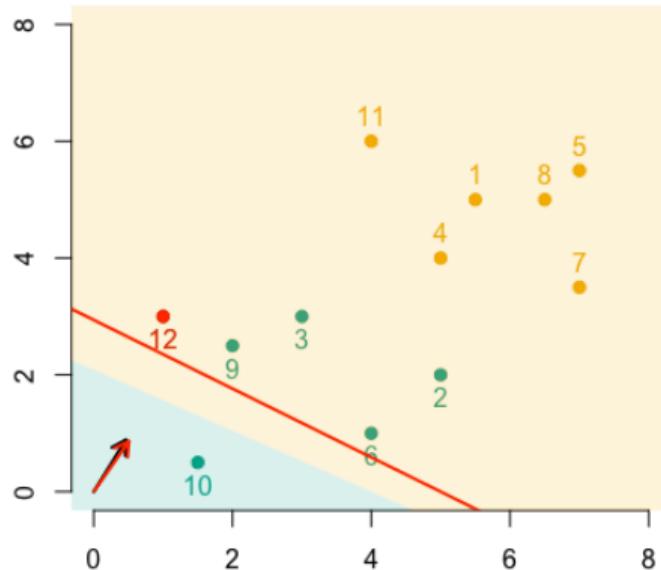
Binary Classifier, SVM (Support Vector Machine)



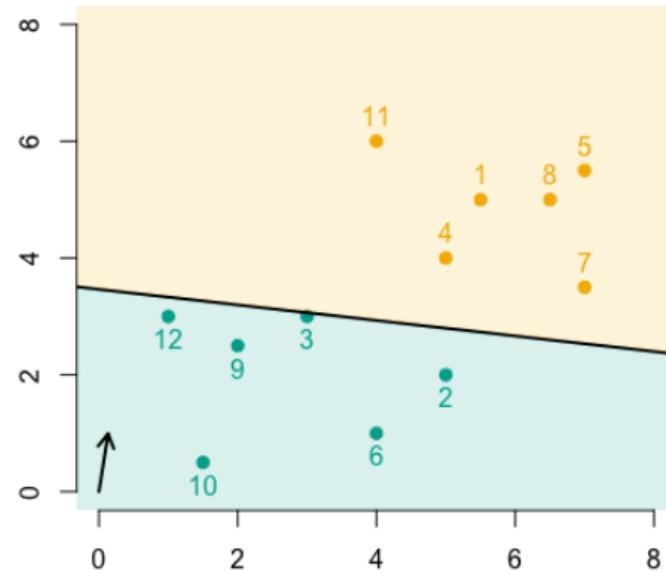
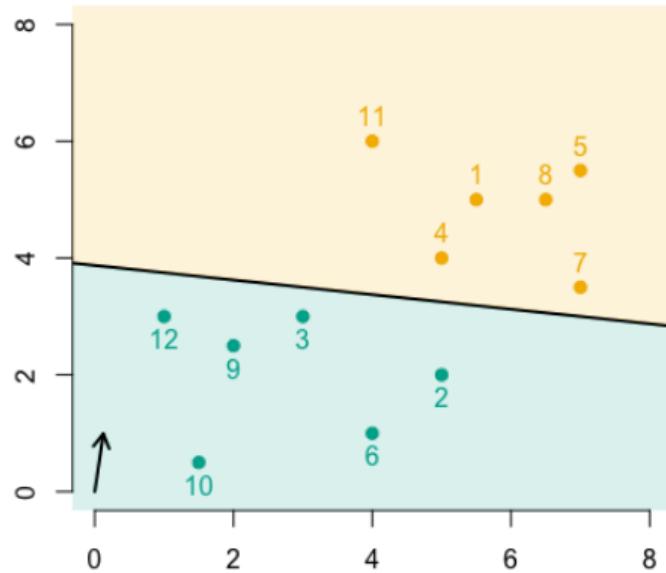
Binary Classifier, SVM (Support Vector Machine)



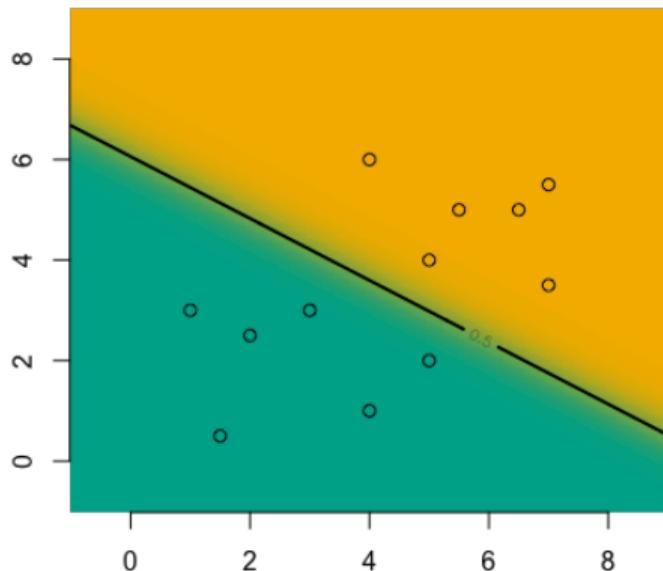
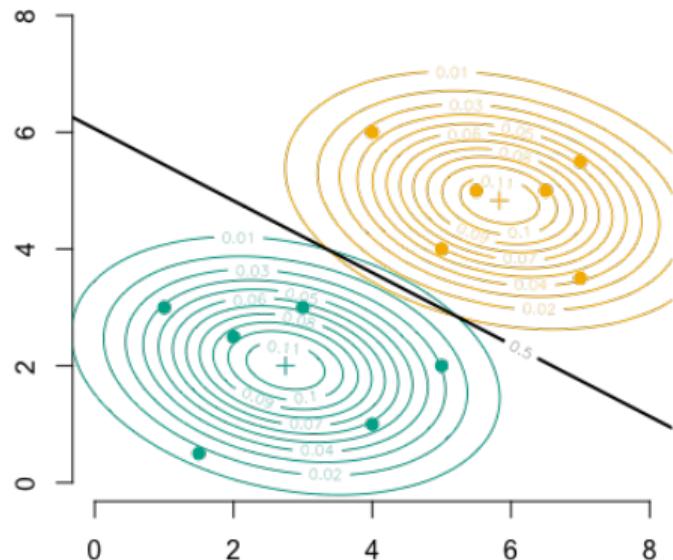
Binary Classifier, SVM (Support Vector Machine)



Binary Classifier, SVM (Support Vector Machine)



Binary Classifier, SVM (Support Vector Machine)



Binary Classifier, SVM (Support Vector Machine)

$$\operatorname{argmax}_{b,\omega} \left\{ \min_i \{ \|x - x_i\| : x \in \mathbb{R}^n, \text{ s.t. } b + x^\top \omega = 0 \} \right\}$$

- **Primal problem**

$$\operatorname{argmax}_{b,\omega} \left\{ \frac{2}{\|\omega\|}, \text{ s.t. } y \cdot (\mathbf{X}\omega + b) \geq 1 \right\}$$

This is an optimization problem subject to n linear constraints.
 $b + x^\top \omega = \pm 1$ are the **support** lines

$$\operatorname{argmin}_{b,\omega} \left\{ \|\omega\|^2, \text{ s.t. } y \cdot (\mathbf{X}\omega + b) \geq 1 \right\}$$

which is an quadratic optimization problem subject to n linear constraints

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities)

BY-NC 4.0 68 / 385

Binary Classifier, SVM (Support Vector Machine)

$$\underset{b, \omega}{\operatorname{argmax}} \left\{ \min_i \{ \|x - x_i\| : x \in \mathbb{R}^n, \text{ s.t. } b + x^\top \omega = 0 \} \right\}$$

- **Dual problem**

Consider the following “linear classifier”

$$m(x) = \sum_{i=1}^n \alpha_i y_i \cdot x_i^\top x + b = \left(\sum_{i=1}^n \alpha_i y_i x_i^\top \right) \cdot x + b$$

(see representer theorem)

Platt et al. (1999) (following Vapnik et al. (1998), Wahba et al. (1998) and Wahba (1999)) suggested “fitting a sigmoid after the SVM”,

$$\frac{e^{x^\top \omega + b}}{1 + e^{x^\top \omega + b}} = s_{\text{SVM}}(x)$$

Binary Classifier, SVM (Support Vector Machine)

$$\underset{b, \omega}{\operatorname{argmin}} \left\{ \|\omega\|^2, \text{ s.t. } \mathbf{y} \cdot (\mathbf{X}\beta + b) \geq \mathbf{1} \right\}$$

- **Soft margin problem**

$$\underset{b, \omega, \xi \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ \|\omega\|^2 + \gamma \cdot \mathbf{1}^\top \xi, \text{ s.t. } \mathbf{y} \cdot (\mathbf{X}\beta + b) \geq \mathbf{1} - \xi \right\}$$

ξ_i 's are called "**slack variables**," used when points are not "**linearly separable**".

- $\xi_i \in (0, 1)$: point is between margin and correct side of hyper-plane (margin violation)
- $\xi_i > 1$: point is misclassified

γ is a regularization parameter:

- γ small γ allows constraints to be easily ignored, i.e., large margin
- γ large γ makes constraints hard to ignore, i.e., narrow margin

Binary Classifier, SVM (Support Vector Machine)

The optimization problem was

$$\operatorname{argmin}_{b, \omega, \xi \in \mathbb{R}^n_+} \left\{ \|\omega\|^2 + \gamma \cdot \mathbf{1}^\top \xi, \text{ s.t. } \mathbf{y} \cdot (\mathbf{X}\omega + b) \geq \mathbf{1} - \xi \right\}$$

with n (linear) constraints, $y_i \cdot m(\mathbf{x}_i) \geq 1 - \xi_i$. Since $\xi_i \in \mathbb{R}_+$, saturated constraints becomes $\xi_i = \max\{0, 1 - y_i \cdot m(\mathbf{x}_i)\}$, and the optimization problem

$$\operatorname{argmin}_{b, \omega \in \mathbb{R}^p} \left\{ \|\omega\|^2 + \gamma \cdot \mathbf{1}^\top \max\{\mathbf{0}, 1 - \mathbf{y} \cdot (\mathbf{X}\omega + b)\} \right\}$$

$$\operatorname{argmin}_{b, \omega \in \mathbb{R}^p} \left\{ \underbrace{\|\omega\|^2}_{\text{regularization}} + \gamma \cdot \mathbf{1}^\top \max\{\mathbf{0}, 1 - \mathbf{y} \cdot (\mathbf{X}\omega + b)\} \right\}$$

Binary Classifier, SVM (Support Vector Machine)

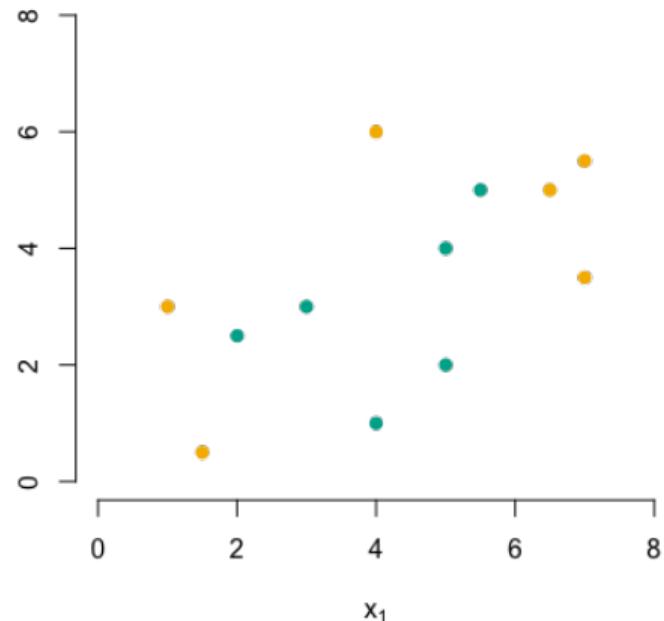
$\ell(y, \hat{y}) = \max\{0, 1 - y \cdot \hat{y}\}$ is called “**Hinge loss**”, Vito et al. (2004).

Hinge loss

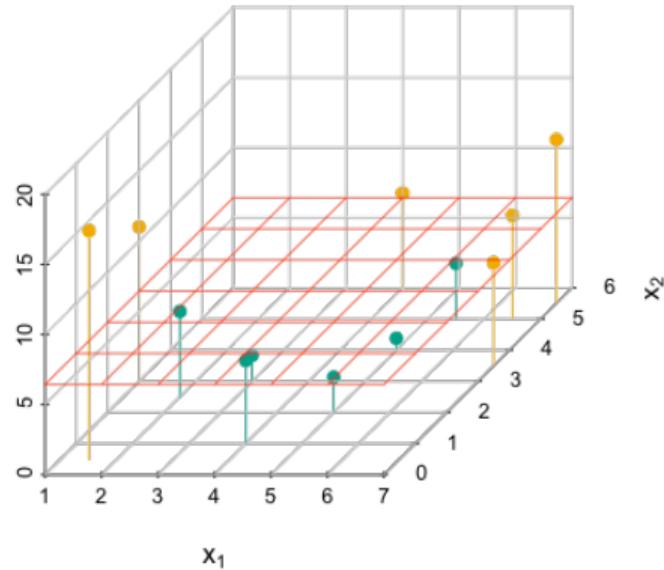
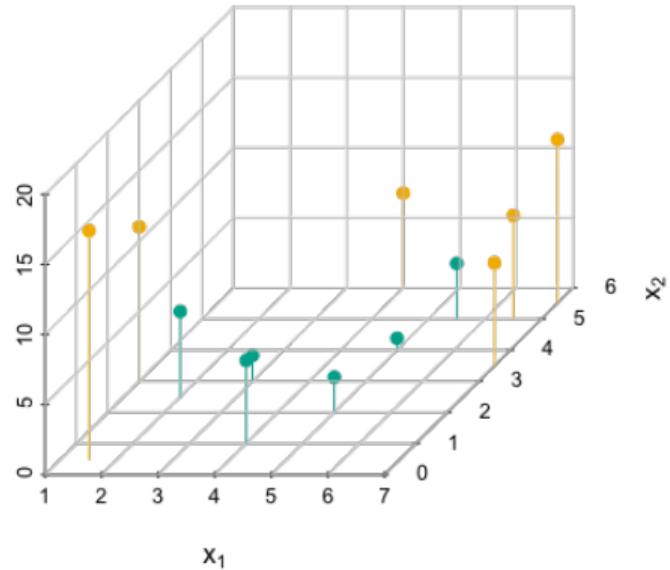
In machine learning, the hinge loss is a loss function used for training classifiers. The hinge loss is used for ‘maximum-margin’ classification, most notably for support vector machines (SVMs). For an intended output $y = \pm 1$ and a classifier score s , the hinge loss of the prediction y is defined as $\ell(y, s) = \max\{0, 1 - s \cdot y\}$.

W

Kernel trick



Kernel trick



Kernel trick

Heuristically, the dual problem was

$$m(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i \cdot \mathbf{x}_i^\top \mathbf{x} + b$$

linear kernel

$$m(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i \cdot K(\mathbf{x}_i, \mathbf{x}) + b$$

$= \varphi(\mathbf{x}_i)^\top \varphi(\mathbf{x})$

Kernel trick

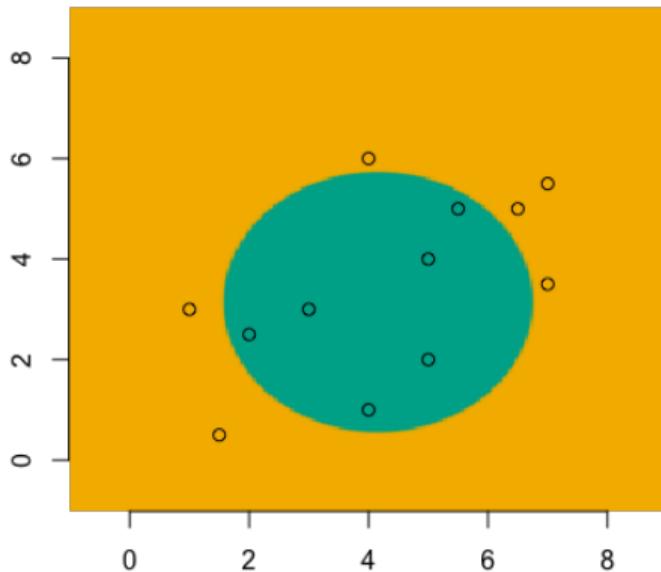
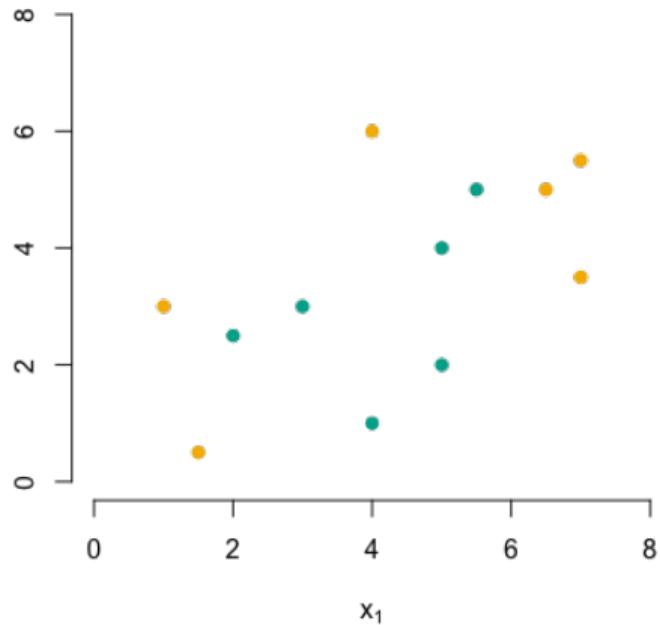
For all \mathbf{x} and \mathbf{x}' in the input space \mathcal{X} , certain functions $k(\mathbf{x}, \mathbf{x}')$ can be expressed as an inner product in another space \mathcal{V} . The function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is often referred to as a **kernel** or a kernel function. \mathbb{W}

Kernel trick

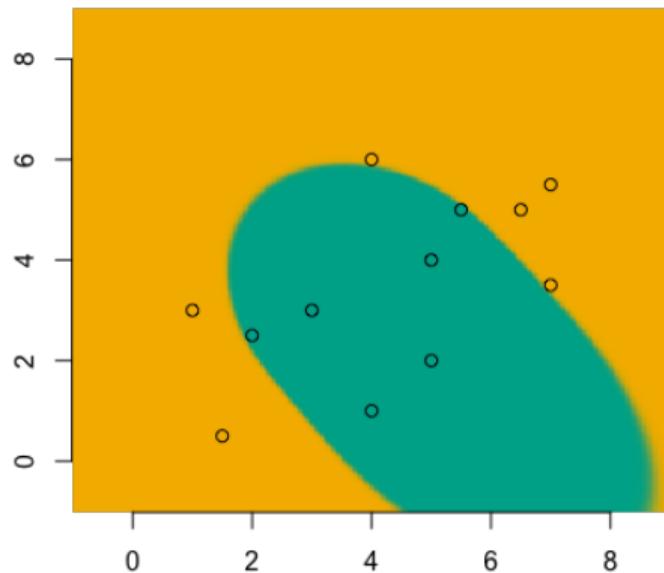
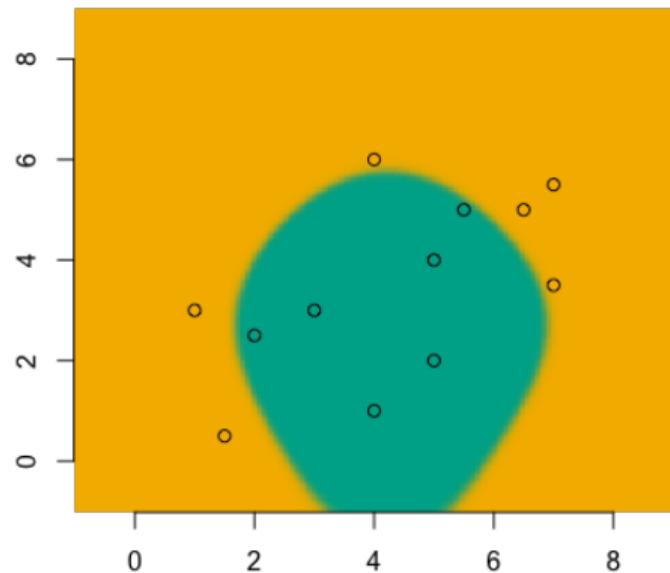
E.g. for $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$

$$\begin{cases} \text{Linear kernel : } k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}, \\ \text{Polynomial kernel : } k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\top \mathbf{y} + r)^p, r \geq 0, p \geq 1 \\ \text{Gaussian kernel (RBF) : } k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right), \sigma > 0 \end{cases}$$

Kernel trick



Kernel trick



Performance Metrics

| | | y | |
|-----------|----------|----------------|----------------|
| | | Positive | Negative |
| \hat{y} | Positive | true positive | false positive |
| | Negative | false negative | true negative |

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

$$\text{sensitivity / recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{F1 score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}}$$

freakonometrics

freakonometrics.hypotheses.org

Accuracy

Consider the case where $y \in \{0, 1\}$, and a **score** $m(\mathbf{x})$ (classically in $[0, 1]$).

E.g., for a logistic regression, $m(\mathbf{x}) = \frac{\exp[\mathbf{x}^\top \boldsymbol{\beta}]}{1 + \exp[\mathbf{x}^\top \boldsymbol{\beta}]}$.

Receiver operating characteristic

A receiver operating characteristic curve, or ROC curve, is a graphical plot that illustrates the performance of a binary classifier model (can be used for multi class classification as well) at varying threshold values. The true-positive rate is also known as sensitivity, recall or probability of detection. The false-positive rate is also known as the probability of false alarm and equals $(1 - \text{specificity})$. W

Definition 2.9: 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

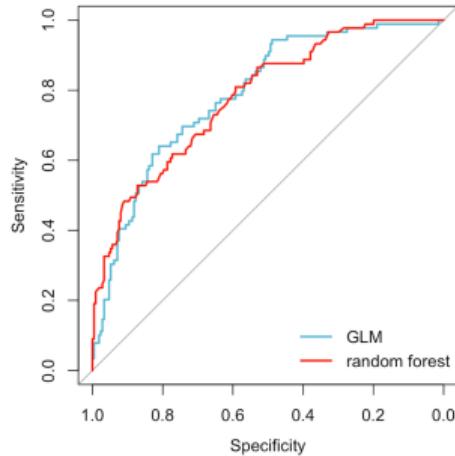
$$\begin{cases} \text{FPR}(t) = \mathbb{P}[m(\mathbf{X}) > t | Y = 0] = \mathbb{P}[m_0(\mathbf{X}) > t] \\ \text{TPR}(t) = \mathbb{P}[m(\mathbf{X}) > t | Y = 1] = \mathbb{P}[m_1(\mathbf{X}) > t]. \end{cases}$$

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)
```

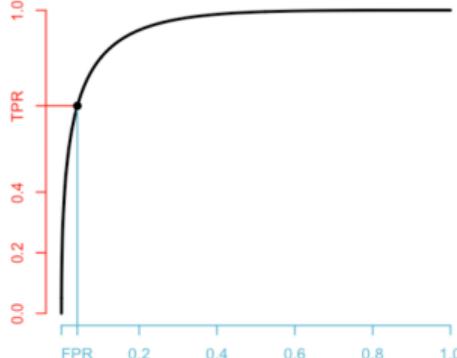
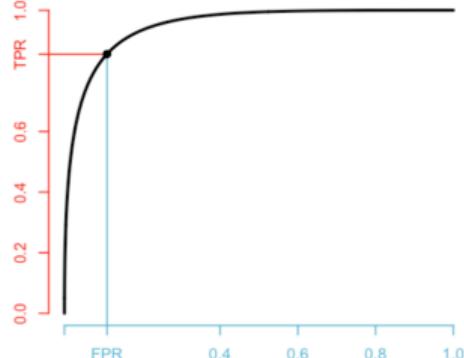
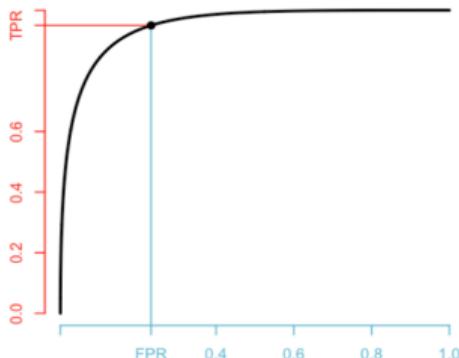
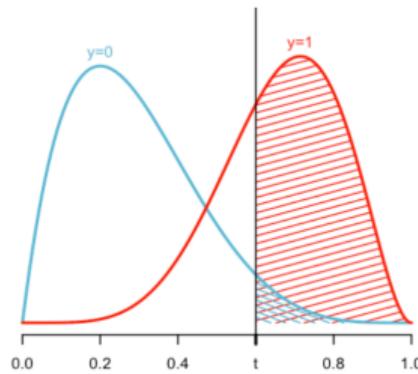
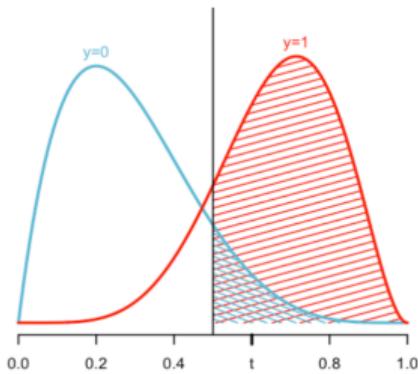
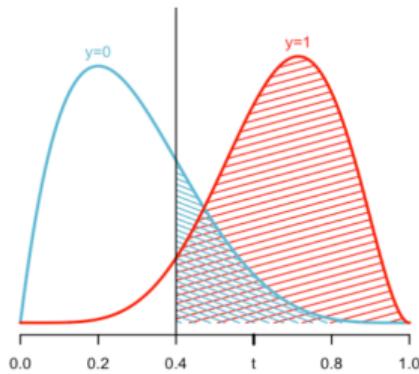


Definition 2.10: AUC, area under the ROC curve

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



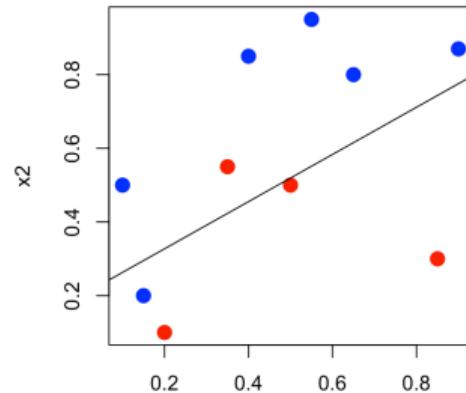
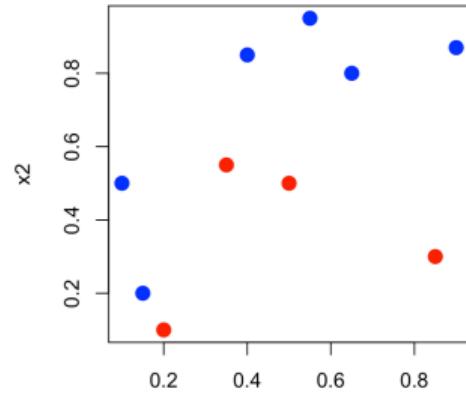
ROC Curve

```
1 > x1 = c(.4,.55,.65,.9,.1,.35,.5,.15,.2,.85)
2 > x2 = c(.85,.95,.8,.87,.5,.55,.5,.2,.1,.3)
3 > y = c(1,1,1,1,1,0,0,1,0,0)
4 > df = data.frame(x1 = x1, x2 = x2,
5                   y = as.factor(y))
6 > plot(x1,x2,col=1+y)
7 > reg = glm(y~x1+x2, data=df,
8              family=binomial(link = "logit"))
9 > b = coefficients(reg)
10 > abline(a=-b[1]/b[3],b=-b[2]/b[3])
```

$\mathbb{P}[Y=1|\mathbf{X}=\mathbf{x}] = \mathbb{P}[Y=0|\mathbf{X}=\mathbf{x}]$ si

$$e^{\mathbf{x}^\top \boldsymbol{\beta}} = 1 \text{ ou } \mathbf{x}^\top \boldsymbol{\beta} = 0$$

soit ici $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$.



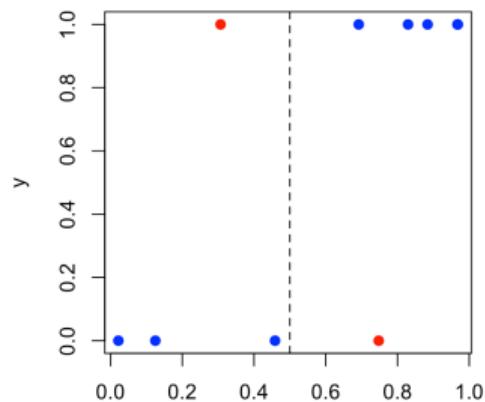
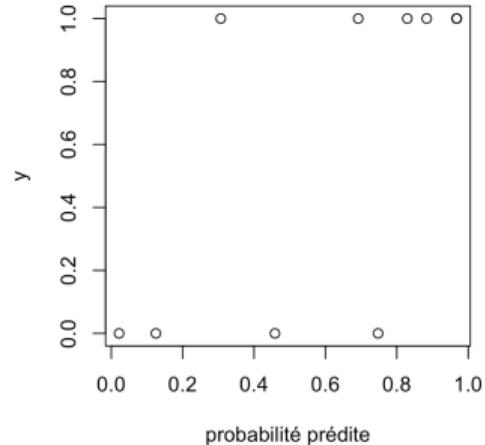
ROC Curve

Let S denote the predicted score (i.e. $\hat{m}(x_i)$)

```
1 > Y = df$y
2 > S = predict(reg, type="response")
3 > plot(S,y)
4 > threshold = .5
5 > Yhat = (S>threshold)*1
6 > plot(S,y,col=1+(y==Yhat))
7 > abline(v=threshold,lty=2)
```

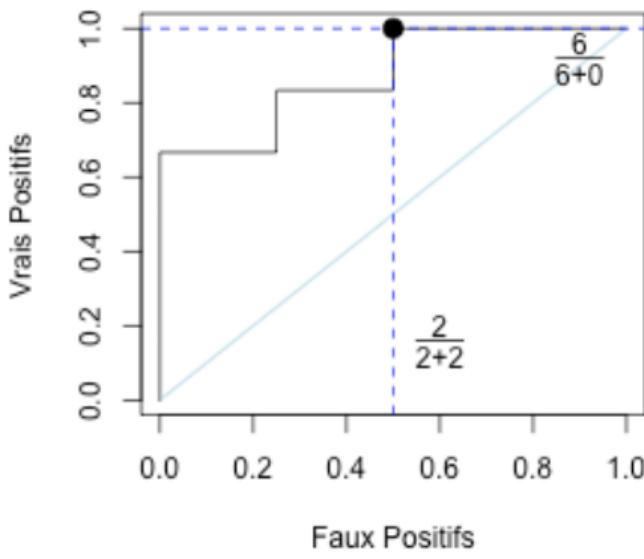
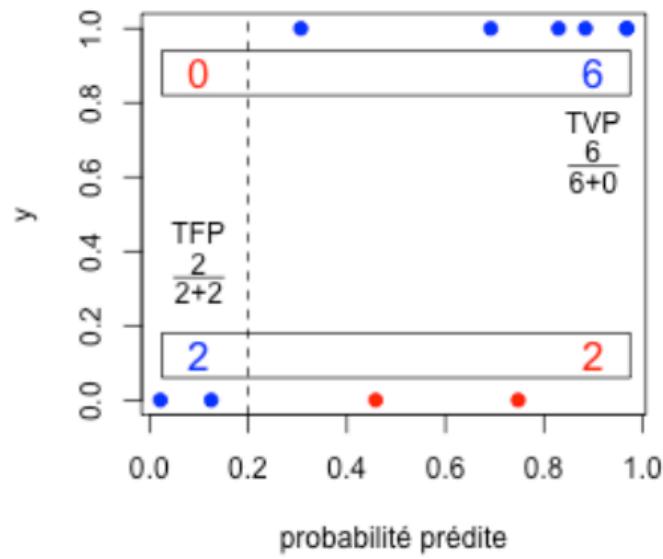
$$\hat{y}_i = \begin{cases} 1 & \text{if } \hat{m}(x_i) > \text{threshold} \\ 0 & \text{if } \hat{m}(x_i) \leq \text{threshold} \end{cases}$$

```
1 > table(Yhat,Y)
2
3   Y
4 Yhat 0 1
5   0 3 1
6   1 1 5
```



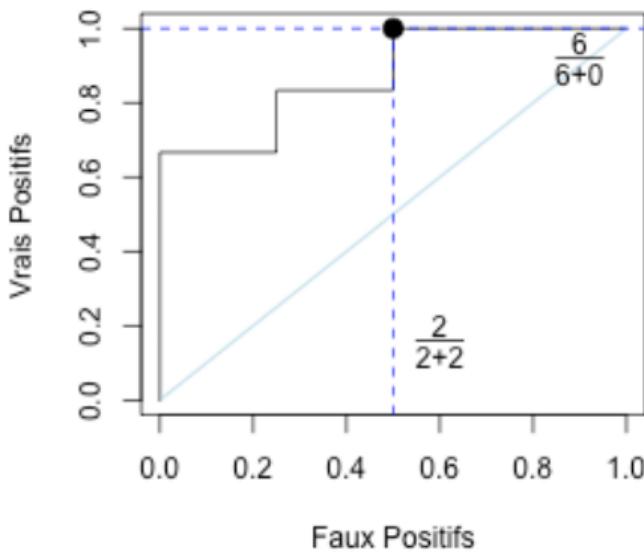
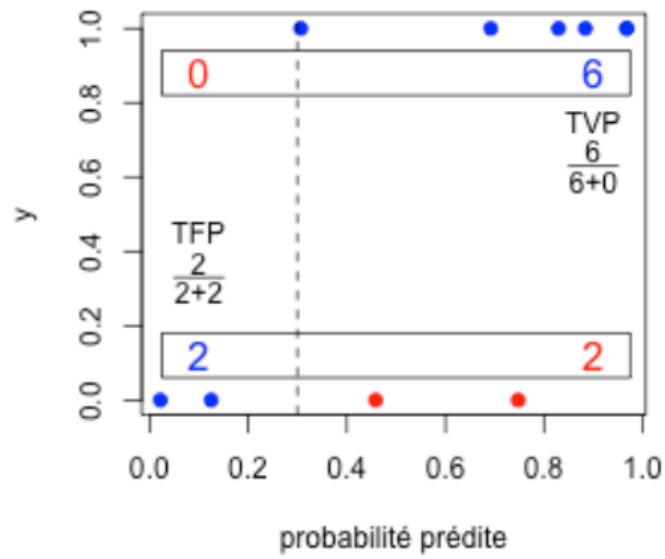
ROC Curve

$$FPR = \frac{\mathbb{P}[y = 0, \hat{y} = 1]}{\mathbb{P}[y = 0]} \text{ and } TPR = \frac{\mathbb{P}[y = 1, \hat{y} = 1]}{\mathbb{P}[y = 1]}$$



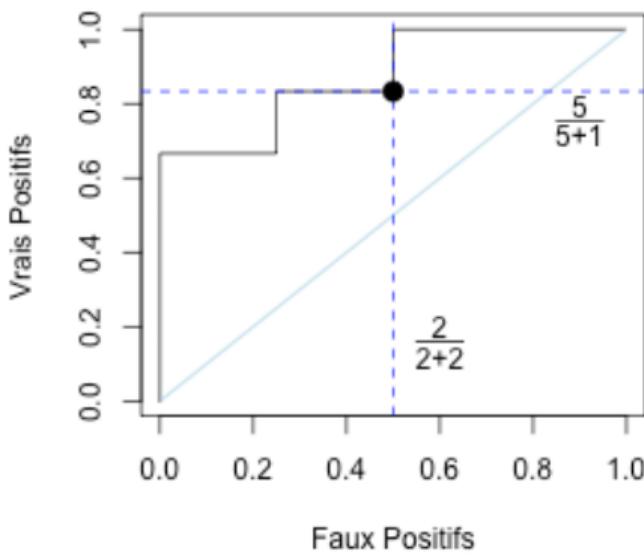
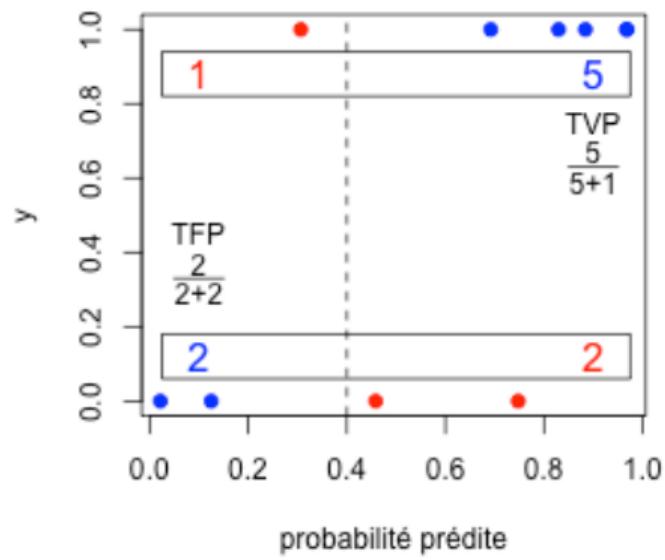
ROC Curve

$$FPR = \frac{\mathbb{P}[y = 0, \hat{y} = 1]}{\mathbb{P}[y = 0]} \text{ and } TPR = \frac{\mathbb{P}[y = 1, \hat{y} = 1]}{\mathbb{P}[y = 1]}$$



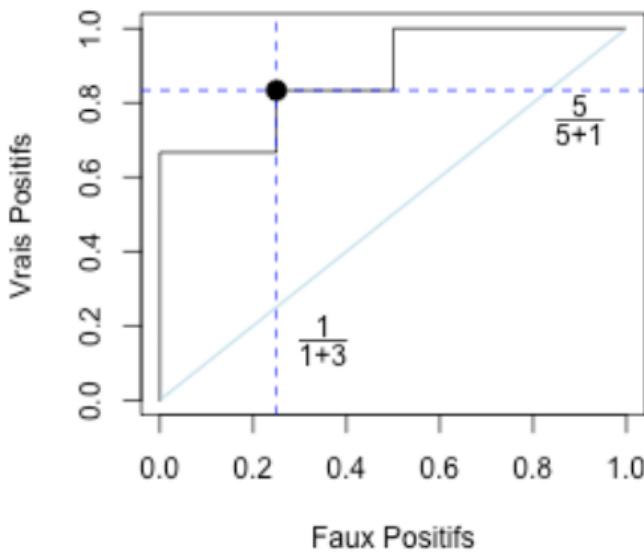
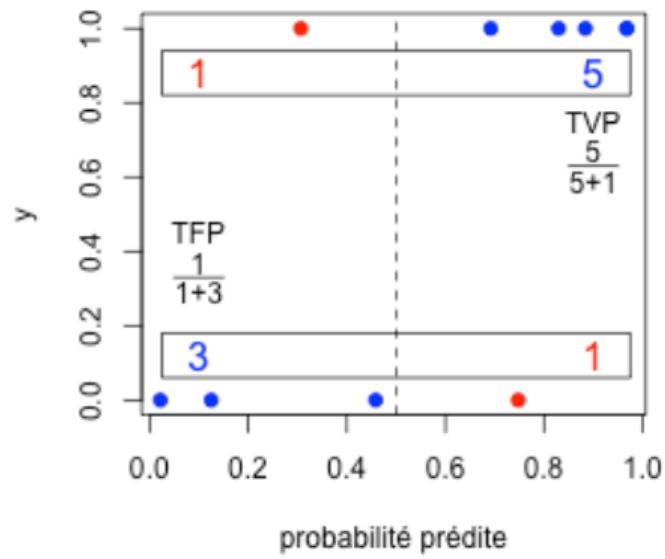
ROC Curve

$$FPR = \frac{\mathbb{P}[y = 0, \hat{y} = 1]}{\mathbb{P}[y = 0]} \text{ and } TPR = \frac{\mathbb{P}[y = 1, \hat{y} = 1]}{\mathbb{P}[y = 1]}$$



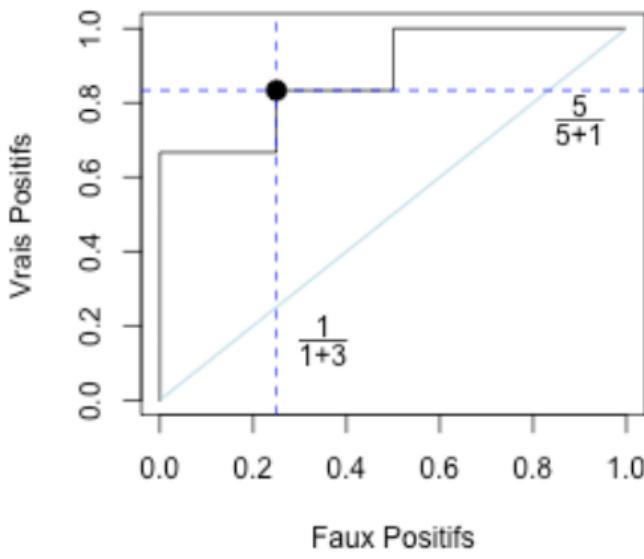
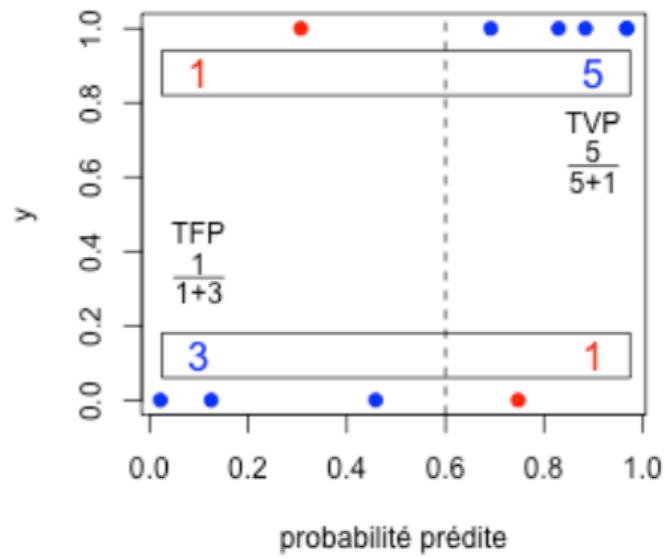
ROC Curve

$$FPR = \frac{\mathbb{P}[y = 0, \hat{y} = 1]}{\mathbb{P}[y = 0]} \text{ and } TPR = \frac{\mathbb{P}[y = 1, \hat{y} = 1]}{\mathbb{P}[y = 1]}$$



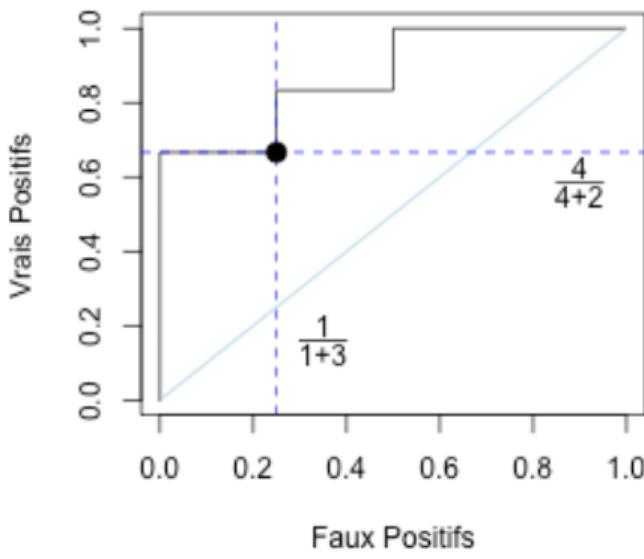
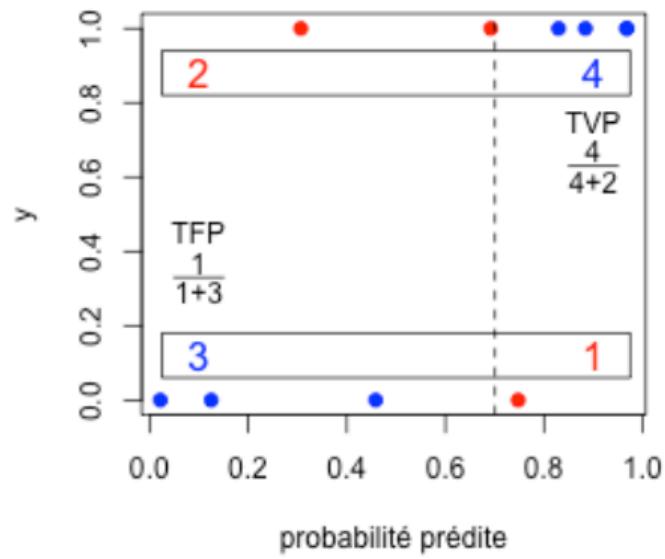
ROC Curve

$$FPR = \frac{\mathbb{P}[y = 0, \hat{y} = 1]}{\mathbb{P}[y = 0]} \text{ and } TPR = \frac{\mathbb{P}[y = 1, \hat{y} = 1]}{\mathbb{P}[y = 1]}$$



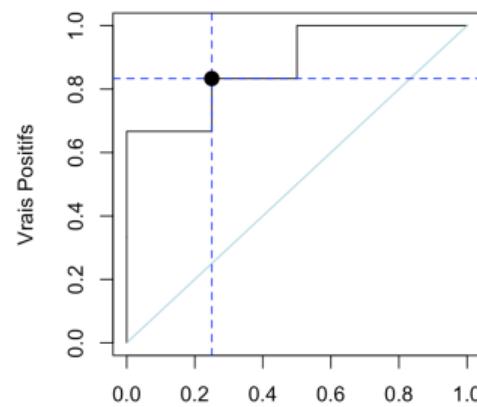
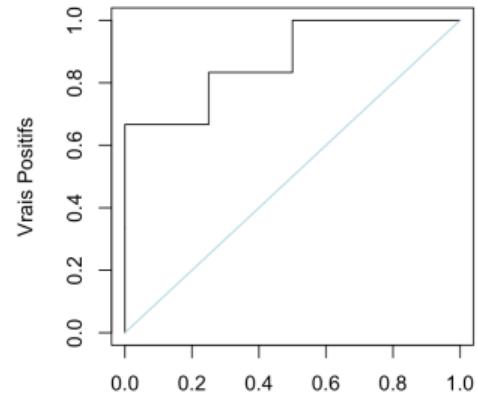
ROC Curve

$$FPR = \frac{\mathbb{P}[y = 0, \hat{y} = 1]}{\mathbb{P}[y = 0]} \text{ and } TPR = \frac{\mathbb{P}[y = 1, \hat{y} = 1]}{\mathbb{P}[y = 1]}$$



ROC Curve

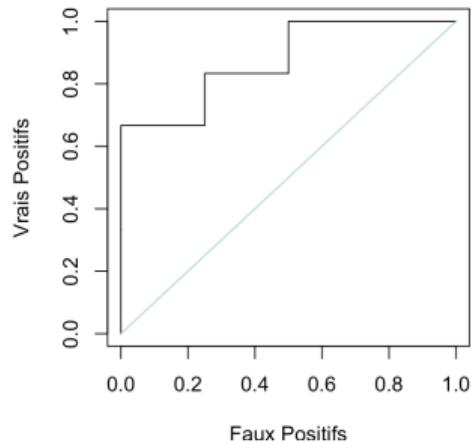
```
1 > roc.curve=function(s){  
2   Ps = (S>s)*1  
3   FP = sum((Ps==1)*(Y==0))/sum(Y==0)  
4   TP = sum((Ps==1)*(Y==1))/sum(Y==1)  
5   vect = c(FP,TP)  
6   names(vect) = c("FPR","TPR")  
7   return(vect) }  
8 > u = seq(0,1,length=251)  
9 > V = Vectorize(roc.curve)(u)  
10 > plot(t(V), type="s")  
11 > table(Yhat,Y)  
12  
13 Yhat 0 1  
14 0 3 1  
15 1 1 5  
16 > sum((Yhat)*(Y==0))/sum(Y==0)  
17 [1] 0.25  
18 > sum((Yhat==1)*(Y==1))/sum(Y==1)  
19 [1] 0.8333333
```



ROC Curve

Various R packages could be used `ROCR` (`pROC` or `plotROC`)

```
1 > library(ROCR)
2 > pred = prediction(S,Y)
3 > plot(performance(pred,"tpr","fpr"))
4 > auc.perf = performance(pred, measure = "auc")
5 > auc.perf@y.values[[1]]
6 [1] 0.875
```



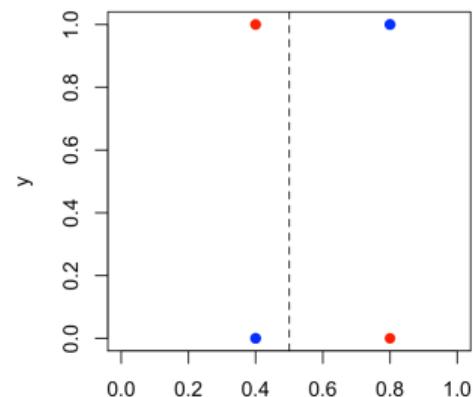
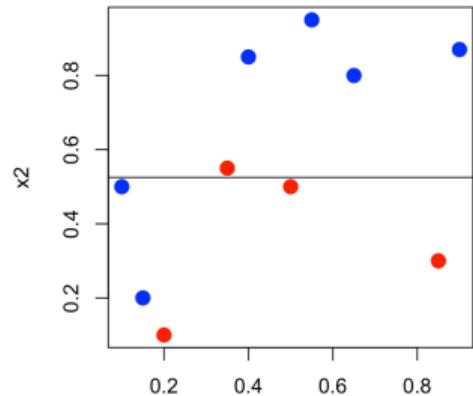
The **AUC** – Area Under the Curve – is a standard accuracy measure.

ROC Curve

If we consider some tree based model

E.g., regression of y on $\mathbf{1}_{[.525,\infty)}(x_2)$

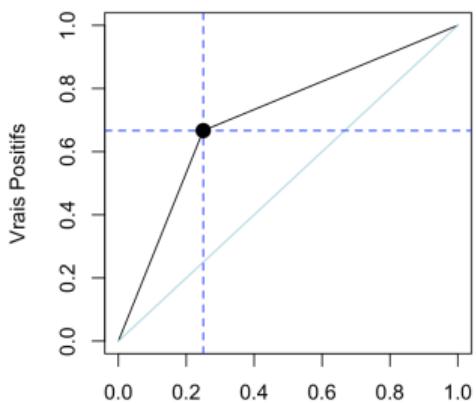
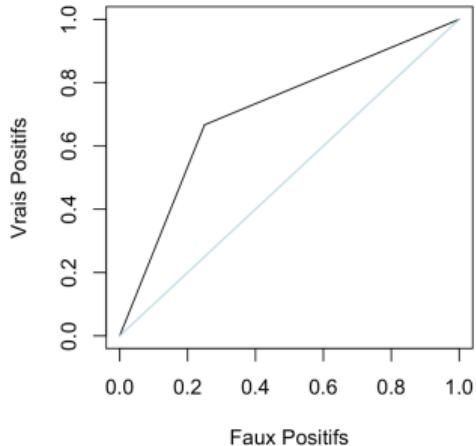
```
1 > reg = glm(y~I(x2> .525), data=df,
2                 family=binomial(link = "logit"))
3 > abline(h=.525)
4 > Y = df$y
5 > S = predict(reg,type="response")
6 > plot(S,y,xlim=0:1)
7 > threshold = .5
8 > Yhat = (S>threshold)*1
9 > table(Yhat,Y)
10
11      Y
12      Yhat 0 1
13        0 3 2
14        1 4 5
```



ROC Curve

With classes, the ROC curve is piecewise linear

```
1 > pred = prediction(S,Y)
2 > plot(performance(pred,"tpr","fpr"))
3 > table(Yhat,Y)
4 
5 Y
6 Yhat 0 1
7   0 3 2
8   1 1 4
9 > auc.perf = performance(pred, measure = "auc")
10 > auc.perf@y.values[[1]]
11 [1] 0.7083333
```



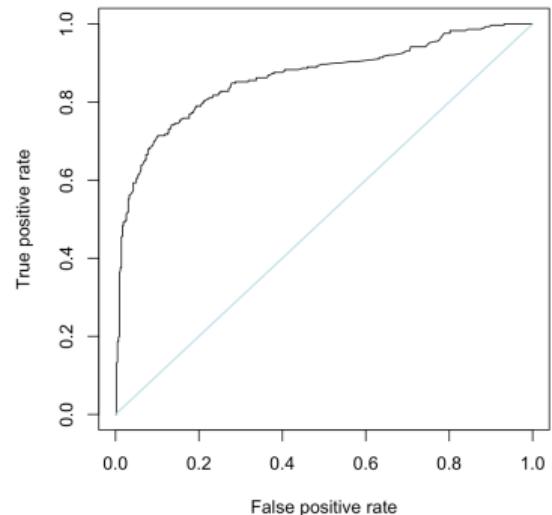
Survie des Passagers du Titanic

y survival for a Titanic passenger,

```
1 > loc = "http://freakonometrics.free.fr/titanic.RData"
2 > download.file(loc, "titanic.RData")
3 > load("titanic.RData")
4 > base = base[!is.na(base$Age),1:7]
5 > reg = glm(Survived ~ Sex+poly(Age,3)+Pclass+SibSp,
6           family = "binomial", data = base)
```

Consider a logistic regression

```
1 > library(ROCR)
2 > Y = base$Survived
3 > S = predict(reg,type="response")
4 > pred = prediction(S,Y)
5 > plot(performance(pred,"tpr","fpr"))
6 > performance(pred, measure = "auc")
7   @y.values[[1]]
8 [1] 0.8627358
```



Multi-class Classification

$$\mathbb{P}(Y = k \mid \mathbf{X} = \mathbf{x}) = \frac{e^{\beta_{k,0} + \beta_{k,1}x_1 + \beta_{k,2}x_2 + \cdots + \beta_{k,p}x_p}}{1 + \sum_{l=1}^{K-1} e^{\beta_{l,0} + \beta_{l,1}x_1 + \beta_{l,2}x_2 + \cdots + \beta_{l,p}x_p}} = \frac{e^{\mathbf{x}^\top \boldsymbol{\beta}_k}}{1 + \sum_{l=1}^{K-1} e^{\mathbf{x}^\top \boldsymbol{\beta}_l}}$$

$\forall k \in \{1, \dots, K-1\}$, while

$$\mathbb{P}(Y = K \mid \mathbf{X} = \mathbf{x}) = \frac{1}{1 + \sum_{l=1}^{K-1} e^{\beta_{l,0} + \beta_{l,1}x_1 + \beta_{l,2}x_2 + \cdots + \beta_{l,p}x_p}} = \frac{1}{1 + \sum_{l=1}^{K-1} e^{\mathbf{x}^\top \boldsymbol{\beta}_l}}$$

And in that case, the predicted class is $\hat{y} = \max_{k \in \{1, \dots, K\}} \mathbb{P}(Y = k \mid \mathbf{x})$ (majority rule).

Voronoi diagram

A Voronoi diagram is a partition of a plane into regions close to each of a given set of objects. It can be classified also as a tessellation. In the simplest case, these objects are just finitely many points in the plane. For each seed there is a corresponding region, called a Voronoi cell, consisting of all points of the plane closer to that seed than to any other. The Voronoi diagram of a set of points is dual to that set's Delaunay triangulation. \mathbf{W}

Multi-class Classification

Simplex

A k -simplex is a k -dimensional polytope that is the convex hull of its $k+1$ vertices. More formally, suppose the $k+1$ points u_0, \dots, u_k are affinely independent. Then, the simplex determined by them is the set of points

$$C(\mathbf{u}) = \left\{ \theta_0 u_0 + \cdots + \theta_k u_k \mid \sum_{i=0}^k \theta_i = 1 \text{ and } \theta_i \geq 0 \text{ for } i = 0, \dots, k \right\}.$$

The standard simplex or probability simplex is the $(k - 1)$ -dimensional simplex whose vertices are the k standard unit vectors in \mathbb{R}^k , or in other words

$$\left\{ x \in \mathbb{R}^k : x_0 + \cdots + x_{k-1} = 1, x_i \geq 0 \text{ for } i = 0, \dots, k-1 \right\}.$$

W

Multi-class Classification

Ternary plot

A ternary plot, ternary graph, triangle plot, simplex plot, or Gibbs triangle is a barycentric plot on three variables which sum to a constant.[1] It graphically depicts the ratios of the three variables as positions in an equilateral triangle. W

Viviani's theorem

Viviani's theorem, named after Vincenzo Viviani, states that the sum of the shortest distances from any interior point to the sides of an equilateral triangle equals the length of the triangle's altitude. W

Metrics related to Classifiers

Inputs, $\{y_i, \hat{y}_i\}$, where $y_i, \hat{y}_i \in \{0, 1\}$.

| | True Value y | Predicted Value \hat{y} |
|----------------------------|----------------------------------|---------------------------------------------|
| True Positive (TP) | 1 | 1 |
| True Negative (TN) | 0 | 0 |
| False Positive (FP) | 0 | 1 |
| False Negative (FN) | 1 | 0 |

| | True $y = 0$ | True $y = 1$ |
|---------------|--------------------------------|--------------------------------|
| $\hat{y} = 0$ | True Negative (TN) | False Negative (FN) |
| $\hat{y} = 1$ | False Positive (FP) | True Positive (TP) |

Metrics related to Classifiers

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Heuristics:

- Works well for balanced datasets.
- Can be misleading for imbalanced datasets.

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN}$$

Heuristics:

- Precision: Useful when false positives are costly.
- Recall: Useful when false negatives are costly.

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \tag{1}$$

Heuristics:

Metrics related to Classifiers

- Harmonic mean of precision and recall.
- Useful for imbalanced datasets.

$$\text{Balanced Accuracy} = \frac{1}{2} \left(\frac{TP}{TP+FN} + \frac{TN}{TN+FP} \right)$$

```
1 > library(caret)
2 > Y      = as.factor(base$Survived)
3 > Yhat = as.factor((predict(reg,
4   type="response") > .5)*1)
5 > confusionMatrix(Yhat, Y, mode =
6   "everything", positive="1")
7
8   Reference
9   Prediction 0 1
10  0 367 75
11  1  57 215
12
13  Accuracy : 0.8151
14  95% CI : (0.7847,
15  0.843)
```

```
1          Kappa : 0.613
2
3  Mcnemar's Test P-Value : 0.139
4
5          Sensitivity : 0.7414
6          Specificity : 0.8656
7          Pos Pred Value : 0.7904
8          Neg Pred Value : 0.8303
9          Precision : 0.7904
10         Recall : 0.7414
11         F1 : 0.7651
12         Prevalence : 0.4062
13         Detection Rate : 0.3011
14         Detection Prevalence : 0.3810
15         Balanced Accuracy : 0.8035
```

Metrics related to Classifiers

Note, for the confidence interval,

$$CI = \left[\text{Accuracy} \pm \Phi^{-1}(1 - \alpha/2) \times \sqrt{\frac{\text{Accuracy}(1 - \text{Accuracy})}{n}} \right]$$

Inputs, $\{y_i, \hat{s}(\mathbf{x})_i\}$, for some model \hat{s} ,

$$MAPE = \frac{100}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right|$$

Heuristics:

- Expressed as a percentage.
- Sensitive to small values.

Exponential Family

Definition 2.11: 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)).

Exponential Family

Pour une variable aléatoire Y de la famille exponentielle, alors

$$\mathbb{E}(Y) = b'(\theta) \text{ et } \text{Var}(Y) = b''(\theta)\varphi,$$

On définit la **fonction variance** V par

$$V(\mu) = b''([b']^{-1}(\mu)) \text{ et } \text{Var}(Y) = \varphi V(\mu)$$

Lien canonique $g = b'^{-1}$

Exponential Family

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

The **Gaussian** distribution $\mathcal{N}(\mu, \sigma^2)$

$$f(y; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2\right) = \exp\left(\frac{-y^2 + 2y\mu - \mu^2}{2\sigma^2} + \log \frac{1}{\sigma\sqrt{2\pi}}\right)$$

canonical parameter θ

$$f(y|\mu, \sigma^2) = \exp\left(\frac{y\mu - \mu^2/2}{\sigma^2} + -\frac{y^2}{2\sigma^2} \log \frac{1}{\sigma\sqrt{2\pi}}\right)$$

i.e., $\theta = \mu$, $\varphi = \sigma^2$, $b(\theta) = \theta^2/2$, canonical link $g(\mu) = \mu$ ("identity")

Famille Exponentielle

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

Exemple La loi de **Poisson** $\mathcal{P}(\lambda)$

$$f(y|\lambda) = e^{-\lambda} \frac{\lambda^y}{y!} = \exp(-\lambda + y \log \lambda - \log(y!))$$

$$f(y|\lambda) = \exp\left(\frac{y \log \lambda - \lambda}{1} - \log(y!)\right)$$

soit $\theta = \log \lambda$, $\varphi = 1$, $b(\theta) = \lambda = e^\theta$, canonical link $g(\mu) = \log(\mu)$ ("log")

Famille Exponentielle

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

Exemple La loi de **Bernoulli** $\mathcal{B}(p)$

$$f(y|\lambda) = p^y(1-p)^{1-y} = \exp(y \log p + (1-y) \log(1-p))$$

$$f(y|\lambda) = \exp\left(\frac{y[\log p - \log(1-p)] + \log(1-p)}{1}\right)$$

soit $\theta = \log \frac{p}{1-p}$, $\varphi = 1$, $b(\theta) = \log(1 + e^\theta)$

Ici $p = \frac{e^\theta}{1 + e^\theta}$, canonical link $g(\mu) = \log \frac{\mu}{1-\mu}$ ("logit")

Famille Exponentielle

Tweedie (1984) a suggéré la famille suivante

$$f(y|\mu, \varphi) = A(y, \varphi) \cdot \exp \left\{ \frac{1}{\varphi} [y\theta(\mu) - \kappa(\theta(\mu))] \right\},$$

où

$$\theta(\mu) = \begin{cases} \frac{\mu^{1-\gamma}}{1-\gamma} & \text{si } \gamma \neq 1 \\ \log \mu & \text{si } \gamma = 1 \end{cases} \quad \text{et} \quad \kappa(\theta(\mu)) = \begin{cases} \frac{\mu^{2-\gamma}}{2-\gamma} & \text{si } \gamma \neq 2 \\ \log \mu & \text{si } \gamma = 2 \end{cases}$$

La loi de Y est alors une loi Poisson composée à sauts Gamma,

$$Y \sim \mathcal{C}\mathcal{P}oi \left(\mu^{2-\gamma} \varphi(2-\gamma), \mathcal{G} \left(-\frac{2-\gamma}{\varphi(1-\gamma)}, \varphi(2-\gamma)\mu^{\gamma-1} \right) \right),$$

lorsque $\gamma \in [1, 2]$, et $V(\mu) = \mu^\gamma$.

Régression GLM

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

La log-vraisemblance est alors

$$\log \mathcal{L}(\boldsymbol{\theta}, \varphi | \mathbf{y}) = \frac{\sum_{i=1}^n y_i\theta_i - \sum_{i=1}^n b(\theta_i)}{a(\varphi)} + \sum_{i=1}^n c(y_i, \varphi).$$

On peut remplacer $a(\varphi)$ par φ/w_i , w_i est un poids individuel.

Hypothèses

- g est suffisamment régulière (au moins C^2)
- φ est supposé connu (pour l'instant)
- $n > p = k + 1$ et \mathbf{X} est de plein rang (et donc $\mathbf{X}^\top \mathbf{X}$ est définie positive)

Losses

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

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

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

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

Régression GLM

$$\log \mathcal{L}(\boldsymbol{\theta}, \varphi, \mathbf{y}) = \sum_{i=1}^n \underbrace{\left[\frac{y_i \theta_i - b(\theta_i)}{\varphi} + c(y_i, \varphi) \right]}_{\log \mathcal{L}_i}.$$

$$\nabla \log \mathcal{L} = \mathbf{X}^\top (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{X}^\top (\mathbf{y} - g^{-1}(\mathbf{X}\hat{\boldsymbol{\beta}})) = \mathbf{0}.$$

$$\frac{\partial \log \mathcal{L}_i}{\partial \beta_j} = \frac{\partial \log \mathcal{L}_i}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_j}$$

$$\frac{\partial \log \mathcal{L}_i}{\partial \theta_i} = \frac{y_i - b'(\theta_i)}{\varphi} = \frac{y_i - \mu_i}{\varphi}$$

$$\frac{\partial \theta_i}{\partial \mu_i} = \left(\frac{\partial \mu_i}{\partial \theta_i} \right)^{-1} = \frac{1}{b''(\theta_i)} = \frac{1}{V(\mu_i)} \text{ and } \frac{\partial \eta_i}{\partial \beta_j} = \frac{\partial \mathbf{x}_i^\top \boldsymbol{\beta}}{\partial \beta_j} = x_{i,j}$$

Régression GLM

Les conditions du premier ordre s'écrivent

$$\sum_{i=1}^n \frac{\partial \log \mathcal{L}_i}{\partial \beta_j} = \sum_i \frac{\partial \mu_i}{\partial \eta_i} \times \frac{y_i - \mu_i}{V(\mu_i)} x_{ij} = 0, \quad \forall j$$

Pour simplifier, notons que

$$\frac{\partial \mu_i}{\partial \eta_i} = \left(\frac{\partial \eta_i}{\partial \mu_i} \right)^{-1} = (g'(\mu_i))^{-1}$$

Les conditions du premier ordre s'écrivent

$$\mathbf{X}^\top \boldsymbol{\Omega} (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0}$$

$$\boldsymbol{\Omega} = \text{diag}((V(\mu_i)g'(\mu_i))^{-1})$$

mais on utilisera une autre écriture...

Régression GLM

Les conditions du premier ordre s'écrivent

$$\mathbf{X}^\top \mathbf{W} \boldsymbol{\Delta} (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0}$$

$$\mathbf{W} = \text{diag}((V(\mu_i)g'(\mu_i)^2)^{-1}) \text{ et } \boldsymbol{\Delta} = \text{diag}(g'(\mu_i))$$

Remarque: la matrice d'information de Fisher est $\mathbf{X}^\top \mathbf{W} \mathbf{X}$

Remarque: avec la fonction de lien canonique $\mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0}$

Régression GLM

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

La log-vraisemblance est alors

$$\log \mathcal{L}(\boldsymbol{\theta}, \varphi | \mathbf{y}) = \frac{\sum_{i=1}^n y_i\theta_i - \sum_{i=1}^n b(\theta_i)}{a(\varphi)} + \sum_{i=1}^n c(y_i, \varphi).$$

Paramètre naturel pour y_i : θ_i

Prédiction pour y_i : $\mu_i = \mathbb{E}(Y_i) = b'(\theta_i)$

Score associé à y_i : $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$

Fonction de lien : g telle que $\eta_i = g(\mu_i) = g(b'(\theta_i))$
(g doit être bijective, et croissante...)

Proposition 2.5: Statistical Properties

sous les hypothèses mentionnées auparavant

- l'estimateur du maximum de vraisemblance $\hat{\beta}$ existe et est unique,
- $\hat{\beta} \xrightarrow{P.s.} \beta$ (fortement consistant)
- lorsque $n \rightarrow \infty$, $\hat{\beta} - \beta \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \varphi \Sigma)$

où $\Sigma = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1}$ et $\mathbf{W} = \text{diag}((V(\mu_i)g'(\mu_i)^2)^{-1})$.

Proposition 2.6: Statistical Properties

sous les hypothèses mentionnées auparavant

- lorsque $n \rightarrow \infty$, $\hat{\eta} - \eta \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \varphi \mathbf{X} \boldsymbol{\Sigma} \mathbf{X}^\top)$
- lorsque $n \rightarrow \infty$, $\hat{\mu} - \mu \xrightarrow{\mathcal{L}} \mathcal{N}(\mathbf{0}, \varphi \boldsymbol{\Delta}^{-2} \mathbf{X} \boldsymbol{\Sigma} \mathbf{X}^\top)$

Déviance

La **déviance** est l'écart entre la log-vraisemblance obtenue en β , et celle obtenue avec un modèle parfait (dit saturé),

$$D = 2\varphi \times [\log \mathcal{L}(\mathbf{y}) - \log \mathcal{L}(\hat{\boldsymbol{\mu}})]$$

où $\hat{\boldsymbol{\mu}} = g^{-1}(\mathbf{X}\hat{\boldsymbol{\beta}})$. On peut aussi définir la **scaled deviance**,

$$D^* = \frac{D}{\varphi} = 2 \times [\log \mathcal{L}(\mathbf{y}) - \log \mathcal{L}(\hat{\boldsymbol{\mu}})]$$

et la **null deviance**,

$$D_0^* = 2 \times [\log \mathcal{L}(\mathbf{y}) - \log \mathcal{L}(\bar{y})]$$

$$\text{Normale : } D = \sum_{i=1}^n (y_i - \mu_i)^2, \text{ Poisson : } D = 2 \sum_{i=1}^n \left\{ y_i \ln \frac{y_i}{\mu_i} - (y_i - \mu_i) \right\}$$

$$\text{Gamma : } D = 2 \sum_{i=1}^n \left\{ -\ln \frac{y_i}{\mu_i} + \frac{y_i - \mu_i}{\mu_i} \right\}$$

Out-of-sample Risk

Consider training sample $\{(y_i, \mathbf{x}_i); i = 1, \dots, n'\}$ and validation sample $\{(y_i, \mathbf{x}_i); i = n' + 1, \dots, n\}$.

Let \hat{m}_{is} denote the model estimated on the training sample.

The empirical out-of-sample risk is

$$\hat{\mathcal{R}}_{\text{os}} = \frac{1}{n - n'} \sum_{i=n'+1}^n \ell(y_i, \hat{m}_{\text{is}}(\mathbf{x}_i))$$

Out-of-sample Risk

Consider a partition $\bigcup_{j=1}^K \mathcal{I}_j = \{1, 2, \dots, n\}$.

Let \hat{m}_{-j} denote the model estimated on the sample without set \mathcal{I}_j ,

$$\hat{m}_{-j} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \sum_{i \in \{1, \dots, n\} \setminus \mathcal{I}_j} \ell(y_i, m(\mathbf{x}_i)) \right\}$$

The empirical cross-validation risk is

$$\hat{\mathcal{R}}_{cv} = \frac{1}{K} \sum_{j=1}^K \frac{1}{n_j} \sum_{i \in \mathcal{I}_j} \ell(y_i, \hat{m}_{-j}(\mathbf{x}_i))$$

Out-of-sample Risk

Consider a bootstrap set \mathcal{I}_b of indices $\mathcal{I} = \{1, 2, \dots, n\}$.

Let \hat{m}_b denote the model estimated on the sample of points with index \mathcal{I}_b ,

$$\hat{m}_b \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \sum_{i \in \mathcal{I}_b} \ell(y_i, m(\mathbf{x}_i)) \right\}$$

The empirical out-of-bag cross-validation risk is

$$\hat{\mathcal{R}}_{\text{oob}} = \frac{1}{B} \sum_{b=1}^B \frac{1}{\text{card}(\mathcal{I} \setminus \mathcal{I}_b)} \sum_{i \in \mathcal{I} \setminus \mathcal{I}_b} \ell(y_i, \hat{m}_b(\mathbf{x}_i))$$

Out-of-sample Risk

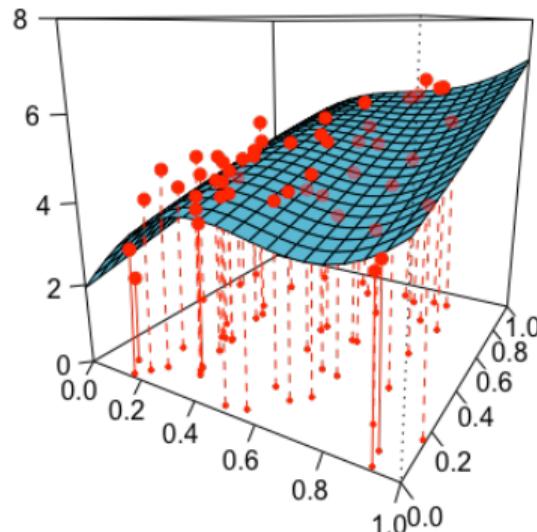
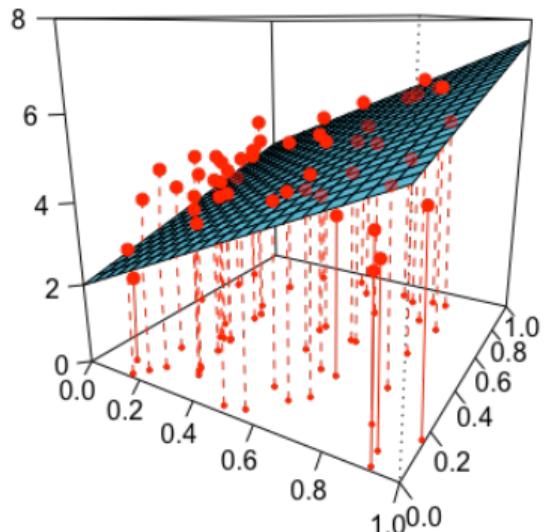
Other classical tools are Akaike's information criterion (AIC), and Schwarz's Bayesian information criterion (BIC),

$$\text{AIC} = -2 \log \mathcal{L}(\hat{\beta}) + 2\dim(\hat{\beta})$$

Preference should be given to the model that has the smallest AIC value.

$$\text{BIC} = -2 \log \mathcal{L}(\hat{\beta}) + \log(n) \cdot \dim(\hat{\beta})$$

Smoothing Splines



“Natura non facit saltus,” Gottfried Leibniz, 1704.

Smoothing Splines

Generalized Additive Models (GAMs) extend Generalized Linear Models (GLMs) by allowing non-linear relationships through smooth functions:

$$g(\mathbb{E}[Y]) = \beta_0 + f_1(X_1) + f_2(X_2) + \cdots + f_p(X_p)$$

- $g(\cdot)$ is the (known) link function,
- $f_j(X_j)$ are smooth functions modeled using splines.
- GLM: $g(\mathbb{E}[Y]) = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$
- GAM: $g(\mathbb{E}[Y]) = \beta_0 + f_1(X_1) + \cdots + f_p(X_p)$

GAM allows capturing non-linearity via smooth functions.

Splines are piecewise polynomial functions ensuring smoothness at knots.

- Basis splines (B-splines)
- Natural splines
- Thin-plate splines

Smoothing Splines

A **B-spline** is defined as:

$$f(X) = \sum_{j=1}^k \beta_j B_j(X)$$

where:

- $B_j(X)$ are basis functions,
- β_j are coefficients.

A **cubic spline** consists of piecewise cubic polynomials ensuring smoothness at knots:

$$f(x) = \sum_{j=1}^k \beta_j B_j(x)$$

subject to:

- Continuity at knots

Smoothing Splines

- Continuous first and second derivatives

A **natural spline** is a cubic spline with additional constraints to enforce linearity at boundaries.

$$f(x) = \sum_{j=1}^k \beta_j B_j(x)$$

Ensures:

- Reduced degrees of freedom
- Better extrapolation properties

GAMs estimate smooth functions by minimizing:

$$\sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int f''(x)^2 dx$$

where λ controls smoothness.

The penalty parameter λ is chosen via:

Smoothing Splines

- Cross-validation
- Generalized Cross-Validation (GCV)
- Restricted Maximum Likelihood (REML)

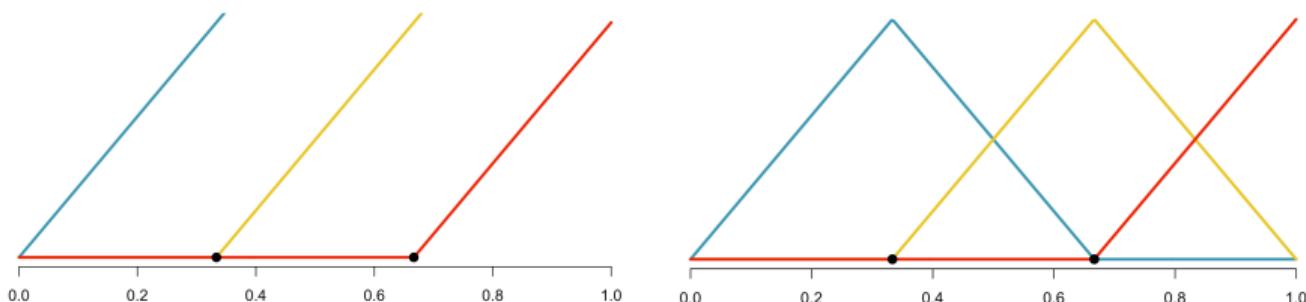
Methods for knot selection:

- Fixed knots (equally spaced or quantiles)
- Adaptive knot selection (using optimization techniques)

Controlling smoothness with penalization

Linear splines (piecewise linear continuous models) are

$$L_1(x) = 1, \ L_2(x) = x, \ L_3(x) = (x - k_1)_+, \ L_4(x) = (x - k_2)_+, \dots$$

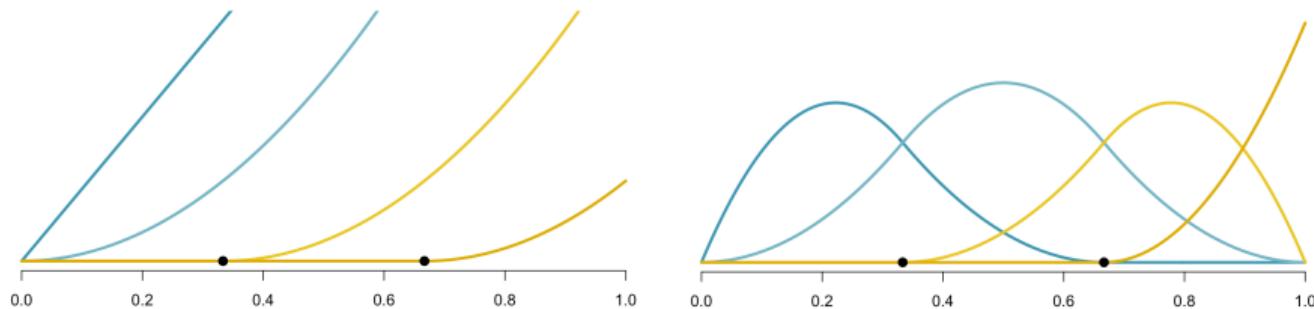


```
1 > x = sort(runif(n))
2 > X = bs(x,knots=quantile(x,p=c(1/3,2/3)),degree = 1)
3 attr(,"degree")
4 [1] 1
5 attr(,"knots")
6 33.33333% 66.66667%
7 0.3542930 0.7091861
8 attr(,"Boundary.knots")
9 [1] 0.003697588 0.989722282
```

Controlling smoothness with penalization

Quadratic splines (piecewise linear continuous models) are

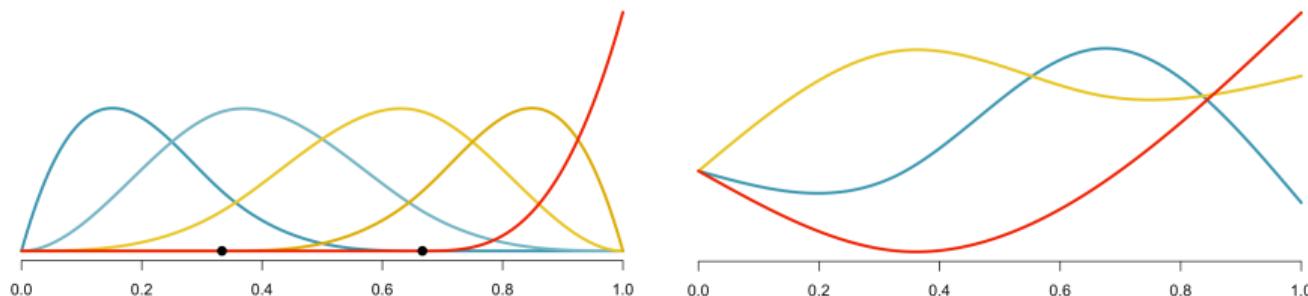
$$L_1(x) = 1, \ L_2(x) = x, \ L_3(x) = x^2, \ L_4(x) = (x - k_1)_+^2, \ \dots$$



```
1 > x = sort(runif(n))
2 > X = bs(x,knots=quantile(x,p=c(1/3,2/3)),degree = 2)
3 attr(,"degree")
4 [1] 2
5 attr(,"knots")
6 33.33333% 66.66667%
7 0.3542930 0.7091861
8 attr(,"Boundary.knots")
9 [1] 0.003697588 0.989722282
```

Controlling smoothness with penalization

Cubic splines, vs. Natural Splines



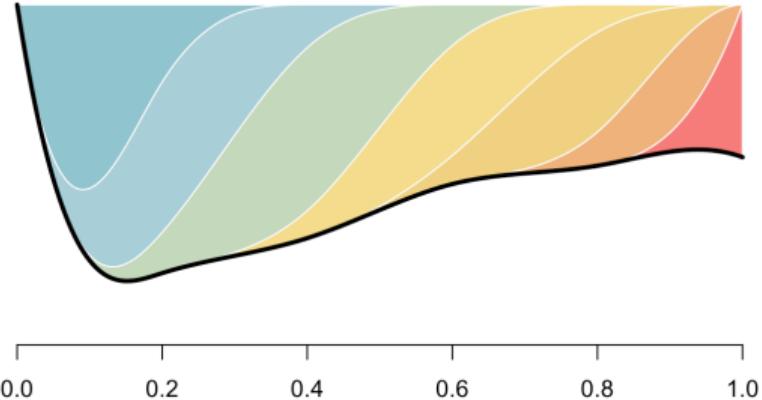
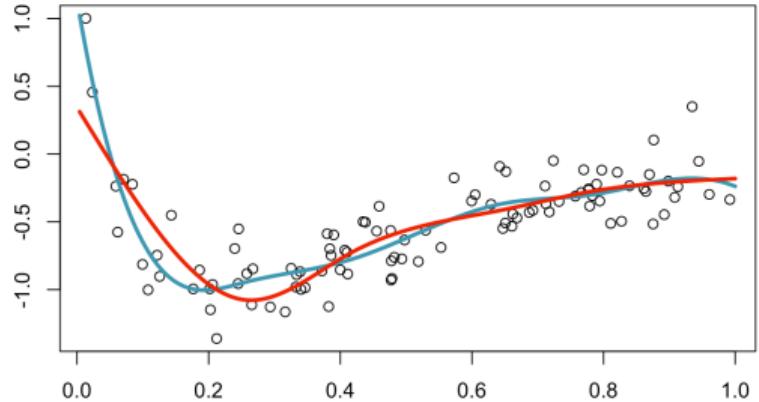
```
1 > Xb = bs(x,knots=quantile(x,p=c(1/3,2/3)),degree = 3)
2 > Xn = ns(x,knots=quantile(x,p=c(1/3,2/3)),degree = 3)
```

Polynomial models tend to be volatile at the boundaries

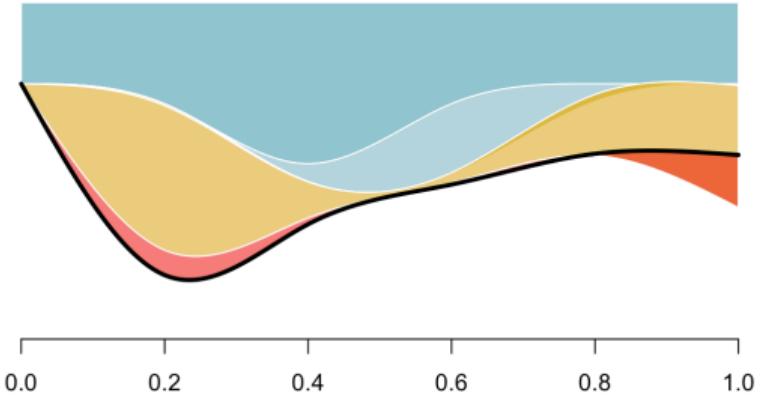
So are cubic splines

Natural cubic splines adding constraints that the function is linear beyond the boundaries of the data

Controlling smoothness with penalization



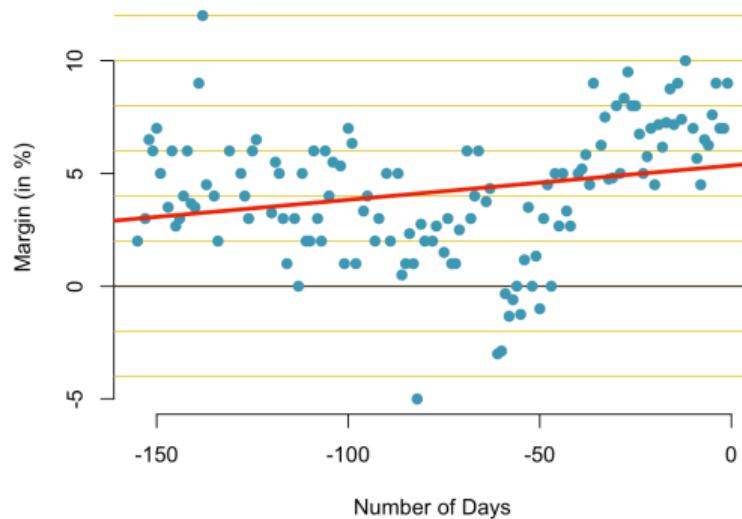
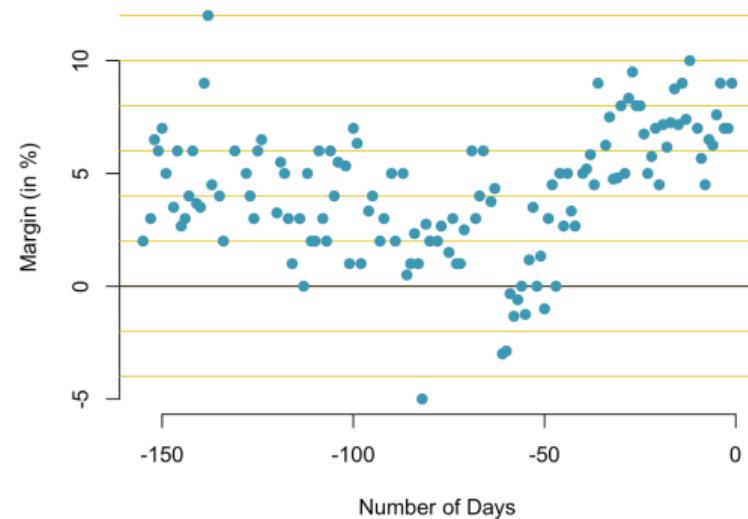
```
1 > set.seed(1)
2 > x = sort(runif(100))
3 > y = sin(log(x))+rnorm(100)/5
4 > plot(x,y)
5 > base = data.frame(x,y)
6 > q = quantile(x,p=c
+ (1/5,2/5,3/5,4/5))
7 > regb = lm(y~bs(x,knots=q))
8 > regn = lm(y~ns(x,knots=q))
```



Natura non facit saltus

Data source: <http://www.pollster.com/08USPresGEMvO-2.html>

pollsters for the popular vote between Obama and McCain (2008 US presidential election), last 150 days.

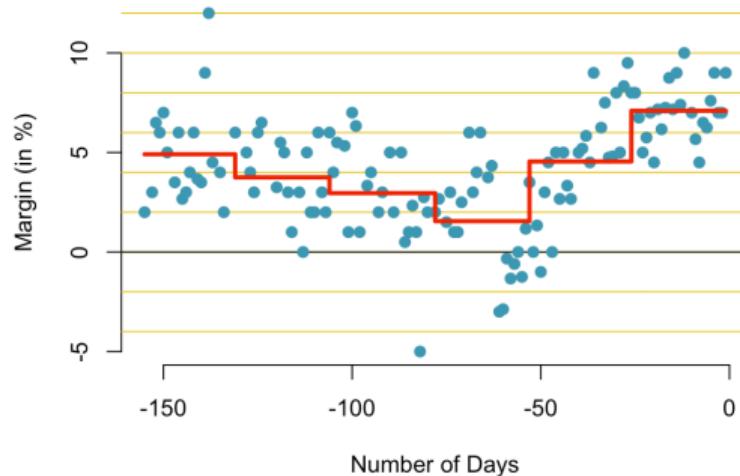
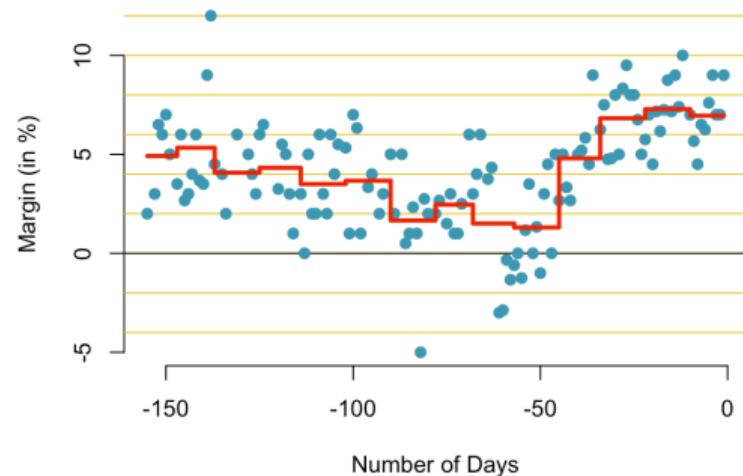


```
1 > library(dslabs)
2 > data("polls_2008")
3 > plot(polls_2008$day, polls_2008$margin*100)
```

Regressogram

From Tukey (1961), the regressogram is defined as

$$\hat{m}_a(x) = \frac{\sum_{i=1}^n \mathbf{1}(x_i \in [a_j, a_{j+1})) y_i}{\sum_{i=1}^n \mathbf{1}(x_i \in [a_j, a_{j+1}))}$$

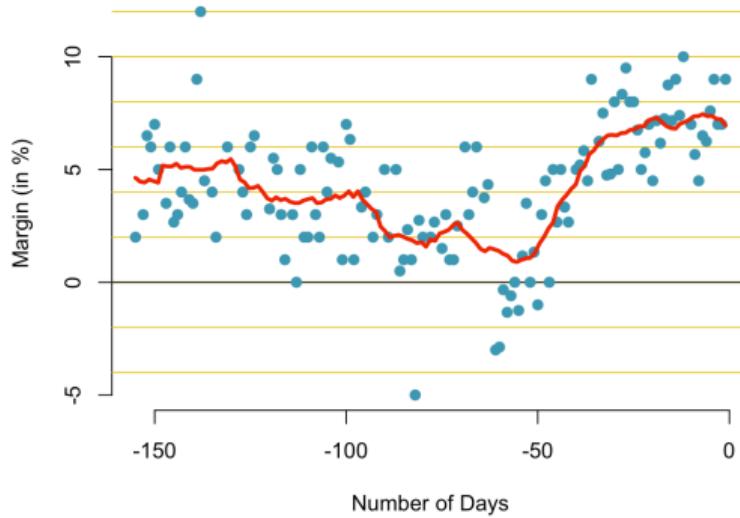
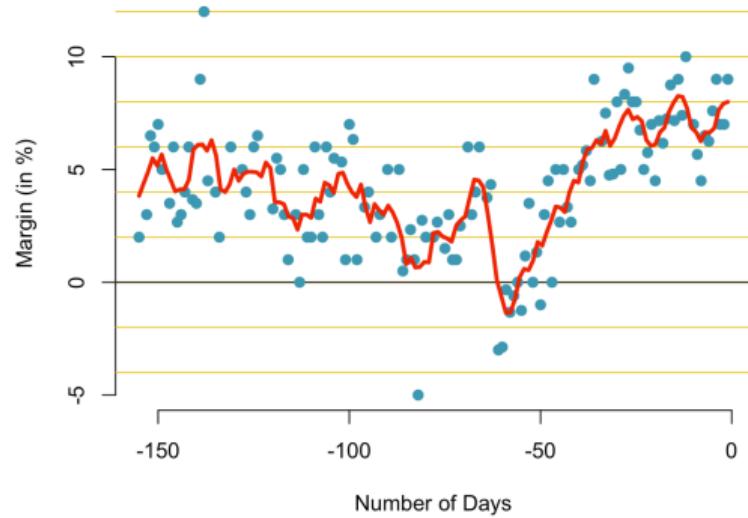


```
1 > reg=lm(margin~cut(day,seq(-160,0,length=15)),data=polls_2008)
```

Moving Regressogram

and the moving regressogram is

$$\hat{m}(x) = \frac{\sum_{i=1}^n \mathbf{1}(x_i \in [x \pm h_n]) y_i}{\sum_{i=1}^n \mathbf{1}(x_i \in [x \pm h_n])}$$



```
1 > with(polls_2008, ksmooth(day, margin, kernel = "box", bandwidth = 7))
```

with **bandwidth** h_n (size of the neighborhood around x)

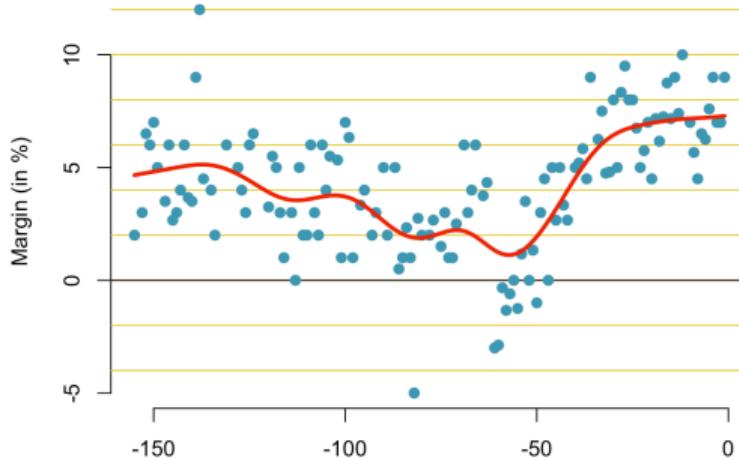
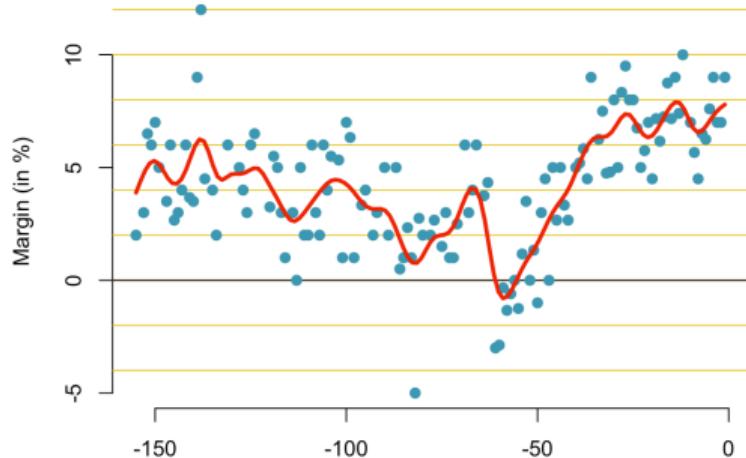
Local Regression

More generally, as moving from the histogram to kernel estimate

$$\tilde{m}(x) = \frac{\sum_{i=1}^n y_i \kappa_h(x - x_i)}{\sum_{i=1}^n \kappa_h(x - x_i)}$$

Observe that this regression estimator is a weighted average

$$\tilde{m}(x) = \sum_{i=1}^n \omega_i(x) y_i \text{ with } \omega_i(x) = \frac{\kappa_h(x - x_i)}{\sum_{i=1}^n \kappa_h(x - x_i)}$$



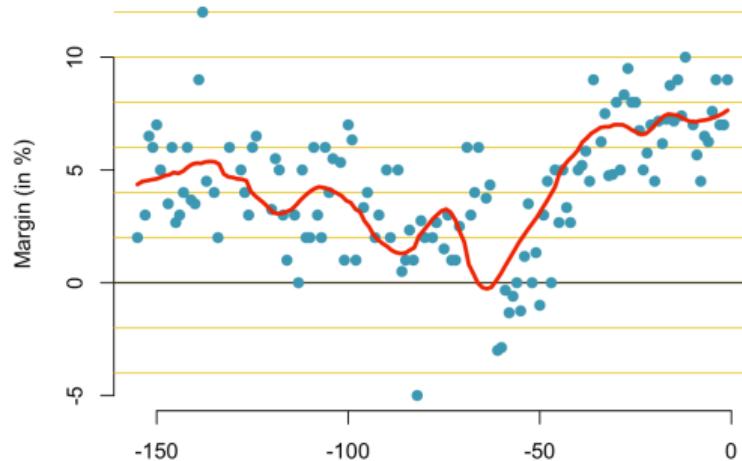
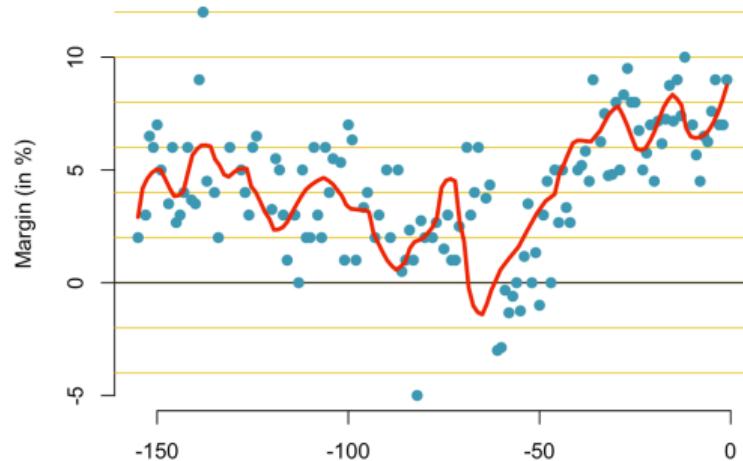
k -Nearest Neighbors

An alternative is to consider

$$\tilde{m}_k(x) = \frac{1}{n} \sum_{i=1}^n \omega_{i,k}(x) y_i$$

where $\omega_{i,k}(x) = \frac{n}{k}$ if $i \in \mathcal{I}_x^k$ with

$$\mathcal{I}_x^k = \{i : x_i \text{ one of the } k \text{ nearest observations to } x\}$$



Local Regression & k -NN

```
1 > fit = with(polls_2008, ksmooth(day, margin, kernel = "normal",
2   bandwidth = span))
2 > lines(fit$x, fit$y)
```

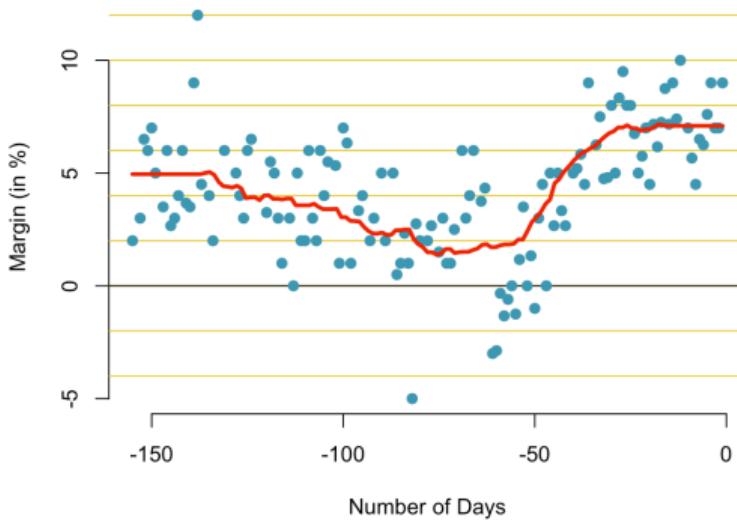
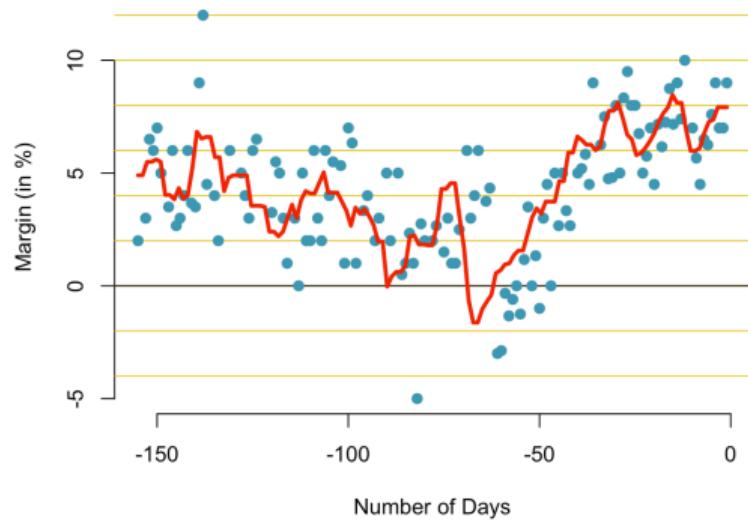
or

```
1 > library(FNN)
2 > p2=knn.reg(train = polls_2008, test = polls_2008, y = polls_2008$margin
   , k = 25)
3 > lines(polls_2008$day, p2$pred)
```

LOESS (locally weighted polynomial)

Solve

$$\tilde{m}(x) = \operatorname{argmin} \left\{ \sum_{i=1}^n \omega_i(x) (y_i - \alpha - \beta x_i)^2 \right\}, \quad \omega_i(x) = \frac{\kappa_h(x - x_i)}{\sum_{i=1}^n \kappa_h(x - x_i)}$$



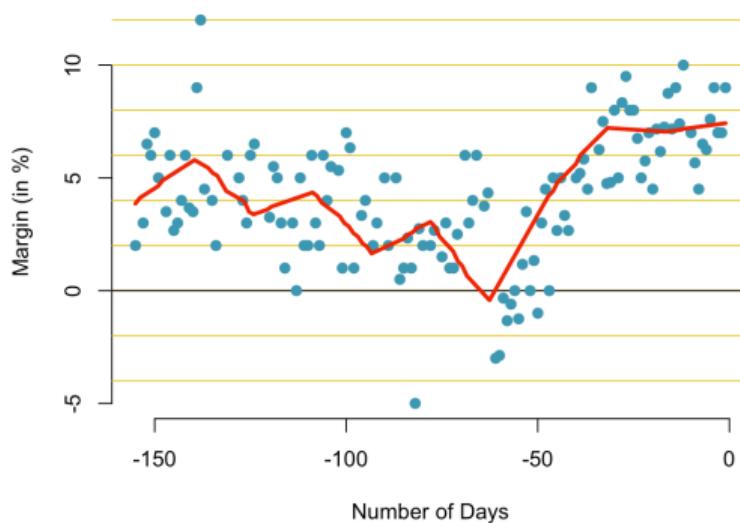
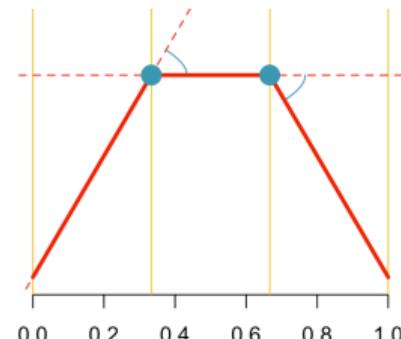
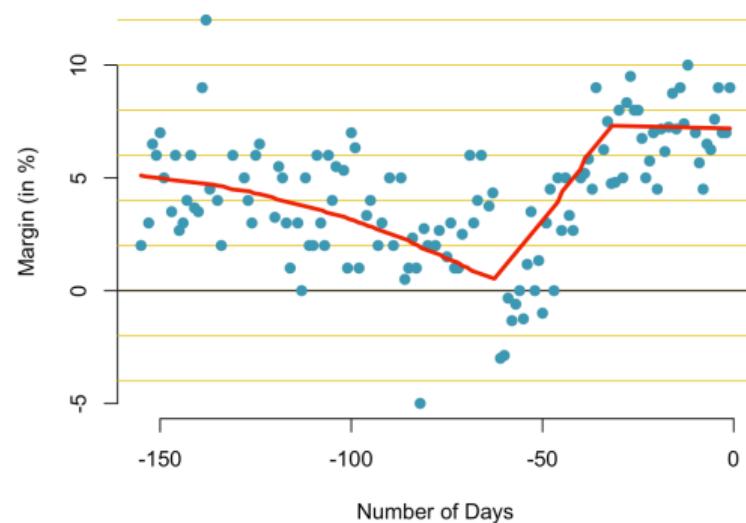
```
1 > fitL = loess(margin ~ day, degree=1, span = 7, data=polls_2008)
```

(Linear) Spline Regression

Select some knots $\{s_1, \dots, s_k\}$, then with $s_0 = 0$

$$\tilde{m}(x) = \alpha + \sum_{j=0}^k \beta_j (x - s_k)_+$$

where $(x - s)_+ = (x - s)$ if $x > s$, 0 otherwise

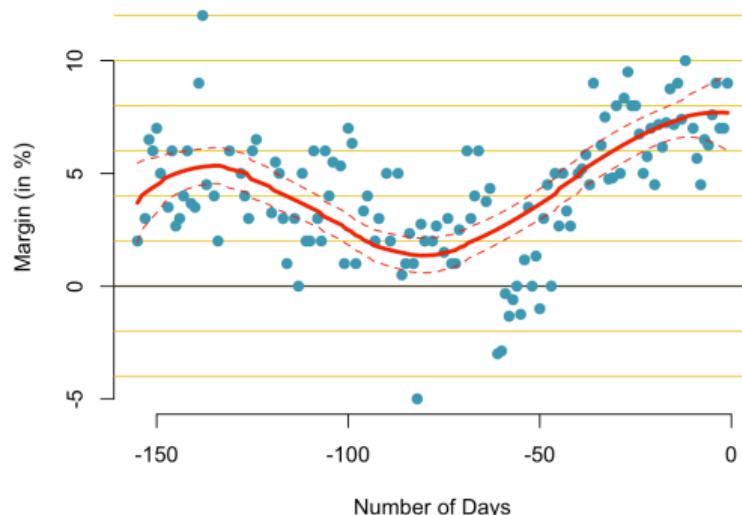
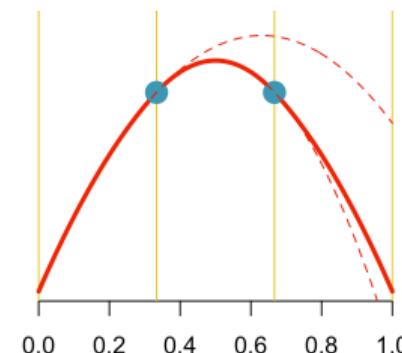
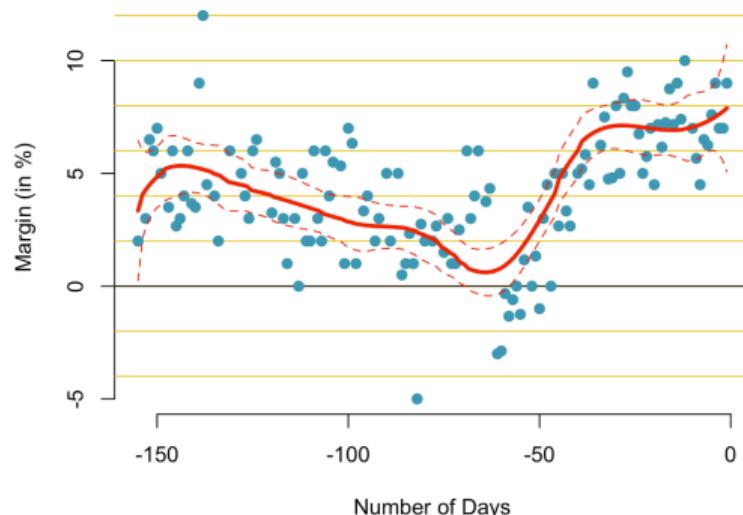


(Quadratic) Spline Regression

Select some knots $\{s_1, \dots, s_k\}$, then with $s_0 = 0$

$$\tilde{m}(x) = \alpha + \gamma x + \sum_{j=0}^k \beta_j (x - s_k)_+^2$$

where $(x - s)_+^2 = (x - s)^2$ if $x > s$, 0 otherwise



Accuracy

Likelihood-based metrics assess how well a model fits the data. These metrics are based on the likelihood function, which measures the probability of the observed data given the model parameters.

- Log-Likelihood: The natural logarithm of the likelihood function.
- Deviance: A measure of how well the model fits the data relative to a saturated model.
- AIC and BIC: Information criteria used for model selection, balancing fit and complexity.

The Log-Likelihood function is the natural logarithm of the likelihood function. It is used to estimate model parameters by maximizing the likelihood.

$$\text{Log-Likelihood} = \log \mathcal{L}(\theta|y) = \sum_{i=1}^n \log(f(y_i|\theta))$$

where

Accuracy

- y_i : Observed data points.
- θ : Model parameters.
- $f(y_i|\theta)$: “Probability” of the data given the model.

Higher log-likelihood values indicate better fitting models.

Deviance is a measure of model fit, representing the difference between the likelihood of a saturated model and the likelihood of the fitted model.

$$\text{Deviance} = -2 \log \left(\frac{\mathcal{L}(\hat{\theta})}{\mathcal{L}(\hat{\theta}_{\text{sat}})} \right)$$

where:

- $\mathcal{L}(\hat{\theta})$: Likelihood of the fitted model.
- $\mathcal{L}(\hat{\theta}_{\text{sat}})$: Likelihood of the saturated model (perfect fit).

Accuracy

Smaller deviance values indicate better fit.

AIC is used for model selection and penalizes models for their complexity. It balances goodness-of-fit with the number of parameters.

$$AIC = -2\ell(\hat{\theta}) + 2k$$

Where:

- $\ell(\hat{\theta})$: Log-Likelihood of the fitted model.
- k : Number of model parameters.

Lower AIC values indicate a better model with a good trade-off between fit and complexity.

Accuracy

BIC is similar to AIC but includes a stronger penalty for complexity. It is derived from a Bayesian perspective.

$$BIC = -2\ell(\hat{\theta}) + k \log(n)$$

Where:

- $\ell(\hat{\theta})$: Log-Likelihood of the fitted model.
- k : Number of model parameters.
- n : Sample size.

Lower BIC values suggest a better model, penalizing larger models more heavily than AIC.

Accuracy

Both AIC and BIC are used to compare models, but they differ in the severity of penalizing model complexity:

- AIC: Tends to select more complex models.
- BIC: More conservative, preferring simpler models with fewer parameters.

When comparing models, lower AIC or BIC indicates a better model. However, BIC typically selects models with fewer parameters compared to AIC.

Regularization

Regularization

In machine learning, a key challenge is enabling models to accurately predict outcomes on unseen data, not just on training data. Regularization is crucial for addressing overfitting—where a model memorizes training data details but cannot generalize to new data. The goal of regularization is to encourage models to learn the broader patterns within the data rather than memorizing it. Techniques like early stopping, ℓ_1 and ℓ_2 regularization, and dropout are designed to prevent overfitting and underfitting, thereby enhancing the model's ability to adapt to and perform well with new data, thus improving model generalization. W

With very general notations,

$$\min_{\theta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i; \boldsymbol{\theta})) + \lambda \cdot \text{Reg}(\boldsymbol{\theta}) \right\}$$

regularization loss

$\text{Reg}(\boldsymbol{\theta})$ is a **penalty function**, while $\lambda \geq 0$ is the **pregularization parameter**

Regularization

- classical LASSO “Least Absolute Shrinkage and Selection”, with squared loss (quadratic) and ℓ_1 regularization

$$\min_{\theta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - m(\mathbf{x}_i; \boldsymbol{\theta}))^2 + \lambda \cdot \sum_{j=1}^p |\theta_j| \right\}$$

- sparse solution (some weights are zero \leftrightarrow feature selection)
- “Best-subset selection regularization”, with squared loss (quadratic) and so-called “ ℓ_0 ” regularization

$$\min_{\theta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - m(\mathbf{x}_i; \boldsymbol{\theta}))^2 + \lambda \cdot \sum_{j=1}^p \mathbf{1}_{\mathbb{R} \setminus \{0\}}(\theta_j) \right\}$$

freakonometrics

freakonometrics.hypotheses.org

Regularization

- squared loss (quadratic) and ℓ_q regularization

$$\min_{\theta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - m(\mathbf{x}_i; \theta))^2 + \lambda \cdot \sum_{j=1}^p |\theta_j|^q \right\}$$

- $q \geq 1$ cost function is convex (strictly convex if $q > 1$)
- $q < 1$ cost function is not convex
- squared loss (quadratic) and elastic net regularization

$$\min_{\theta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - m(\mathbf{x}_i; \theta))^2 + \lambda \cdot \alpha \sum_{j=1}^p |\theta_j| + (1 - \alpha) \sum_{j=1}^p \theta_j^2 \right\},$$

with $\alpha \in [0, 1]$.

Regularization

- ε -insensitive loss ℓ_2 regularization

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (|y_i - m(\mathbf{x}_i; \boldsymbol{\theta})| - \varepsilon)_+ + \frac{\lambda}{2} \cdot \sum_{j=1}^p |\boldsymbol{\theta}_j|^2 \right\}$$

- quadratic programming (dual can be solved)

Suppose that variables are centered, $\bar{y} = \bar{x} = 0$

$$\min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \beta x_i)^2 \right\} : \text{(classical) ordinary least squares (OLS)}$$

$$\min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \beta x_i)^2 + \lambda \beta^2 \right\} : \text{Ridge, with OLS loss,}$$

$$\min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - \beta x_i)^2 + \lambda |\beta| \right\} : \text{Lasso, with OLS loss,}$$

Regularization

Least squares , LS = $\sum_{i=1}^n (y_i - \beta x_i)^2 = \mathbf{y}^\top \mathbf{y} - 2\mathbf{x}^\top \mathbf{y}\beta + \beta \mathbf{x}^\top \mathbf{x}\beta$

$$\frac{\partial \text{LS}}{\partial \beta} = -2 \sum_{i=1}^n x_i(y_i - \beta x_i) = -2\mathbf{x}^\top (\mathbf{y} - \beta \mathbf{x}) \text{ and } \hat{\beta}^{\text{ols}} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}} = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y}$$

Ridge, PLS = $\mathbf{y}^\top \mathbf{y} - 2\mathbf{x}^\top \mathbf{y}\beta + \beta \mathbf{x}^\top \mathbf{x}\beta + \lambda \beta^2$

$$\frac{\partial \text{PLS}}{\partial \beta} = -2\mathbf{y}^\top \mathbf{x} + 2\mathbf{x}^\top \mathbf{x}\beta + 2\beta\lambda \text{ and } \hat{\beta}_\lambda^{\text{ridge}} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x} + \lambda} = (\mathbf{x}^\top \mathbf{x} + \lambda)^{-1} \mathbf{x}^\top \mathbf{y}$$

Lasso, PLS = $\mathbf{y}^\top \mathbf{y} - 2\mathbf{x}^\top \mathbf{y}\beta + \beta \mathbf{x}^\top \mathbf{x}\beta + \lambda |\beta|$

$$\frac{\partial \text{PLS}}{\partial \beta} = -2\mathbf{y}^\top \mathbf{x} + 2\mathbf{x}^\top \mathbf{x}\hat{\beta} \pm \lambda$$

where the sign of \pm is the sign of $\hat{\beta}$.

Regularization

Let's assume that the least-squares estimator (obtained when $\lambda = 0$) is (strictly) positive, i.e. $\mathbf{x}^\top \mathbf{y} > 0$.

If λ is not too large $\widehat{\beta}_\lambda$ and $\widehat{\beta}^{\text{ols}}$ are of the same sign, and

$$-2\mathbf{y}^\top \mathbf{x} + 2\mathbf{x}^\top \mathbf{x}\widehat{\beta} + \lambda = 0, \text{ i.e. } \widehat{\beta}_\lambda^{\text{lasso}} = \frac{\mathbf{y}^\top \mathbf{x} - \lambda/2}{\mathbf{x}^\top \mathbf{x}}$$

We then increase λ until we have $\widehat{\beta}_\lambda = 0$.

By increasing it a little, $\widehat{\beta}_\lambda$ cannot become negative, and therefore

$$-2\mathbf{y}^\top \mathbf{x} + 2\mathbf{x}^\top \mathbf{x}\widehat{\beta} - \lambda = 0, \text{ i.e. } \widehat{\beta}_\lambda^{\text{lasso}} = \frac{\mathbf{y}^\top \mathbf{x} + \lambda/2}{\mathbf{x}^\top \mathbf{x}}$$

But this solution is positive ($\mathbf{y}^\top \mathbf{x} > 0$) so we must have $\widehat{\beta}_\lambda < 0$.

The only possible solution is then $\widehat{\beta}_\lambda = 0$.

Regularization

If x is centered ($\bar{x} = 0$) and reduced ($\mathbf{x}^\top \mathbf{x} = n$),

$$\hat{\beta}^{\text{ols}} = \frac{1}{n} \mathbf{x}^\top \mathbf{y}$$

$$\hat{\beta}_\lambda^{\text{ridge}} = \frac{\mathbf{x}^\top \mathbf{y}}{n + \lambda} = \frac{n}{n + \lambda} \cdot \frac{1}{n} \mathbf{x}^\top \mathbf{y} = \frac{n}{n + \lambda} \cdot \hat{\beta}^{\text{ols}}$$

$$\hat{\beta}_\lambda^{\text{lasso}} = \frac{\mathbf{x}^\top \mathbf{y} \pm \lambda/2}{n} = \frac{1}{n} \mathbf{x}^\top \mathbf{y} \pm \frac{\lambda}{2n} = \hat{\beta}^{\text{ols}} \pm \frac{\lambda}{2n}$$

Ridge Regression

- Idea: perturb the matrix \mathbf{X} to move its eigenvalues away from 0 and thus stabilize the inversion of $\mathbf{X}^\top \mathbf{X}$.
- Notation: $0 \leq \mu_1 \leq \dots \mu_p$ ordered eigenvalues of $\mathbf{X}^\top \mathbf{X}$ and let \mathbf{P} be the orthogonal matrix such that $\mathbf{X}^\top \mathbf{X} = \mathbf{P} \mathbf{D} \mathbf{P}^\top$ with $\mathbf{D} = \text{diag}(\mu_1, \dots, \mu_p)$.
- Let $\lambda \geq 0$, $\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p$ has the same eigenvectors as $\mathbf{X}^\top \mathbf{X}$ and for eigenvalues $\mu_j + \lambda$ for $j = 1, \dots, p$.
- Hoerl and Kennard (1970a): replace $\mathbf{X}^\top \mathbf{X}$ by $\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p$ in the OLS definition:

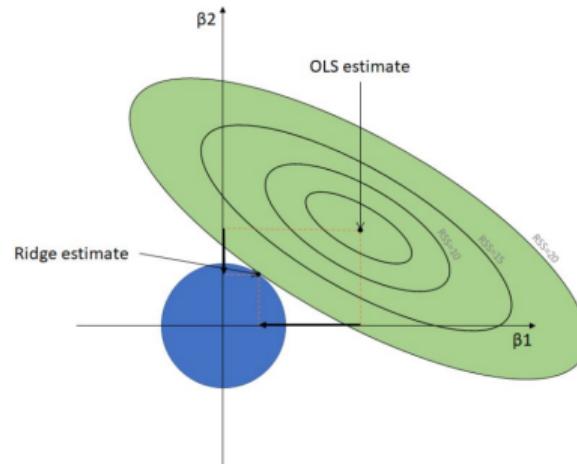
$$\widehat{\mathbf{X}}^{\text{ridge}}(\lambda) = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$$

- If $\lambda = 0$, then $\widehat{\boldsymbol{\beta}}^{\text{ridge}} = \widehat{\boldsymbol{\beta}}^{\text{mco}}$; if $\lambda \rightarrow \infty$, then $\widehat{\boldsymbol{\beta}}^{\text{ridge}} \rightarrow \mathbf{0}$.
- Fixing λ is important. What are the consequences on the bias and variance of the estimators?

Ridge Regression

- Let $\lambda \geq 0$, we can show that $\hat{\beta}^{\text{ridge}}(\lambda)$ minimizes

$$\begin{aligned}\|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda\|\beta\|^2 \\ = \sum_{i=1}^n (y_i - \mathbf{x}_i\beta)^2 + \lambda \sum_{j=1}^p \beta_j^2\end{aligned}$$



<https://rstatisticsblog.com/>

- Let $\tilde{\beta}$ be the estimator minimizing

$$\|\mathbf{y} - \mathbf{X}\beta\|^2 \text{ s.t. } \|\beta\|^2 \leq \delta.$$

$\forall \lambda > 0, \exists \delta$ so that the two solutions coincide.

Ridge Regression

- $\hat{\beta}^{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \mathbf{X}^\top \mathbf{X} \hat{\beta}$.
- $\mathbb{E}(\hat{\beta}^{\text{ridge}}) = \beta - \lambda (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \beta$.
- $\text{Var}(\hat{\beta}^{\text{ridge}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1}$.
- and the mean squared error is

$$\text{MSE}(\hat{\beta}^{\text{ridge}}) (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} (\sigma^2 (\mathbf{X}^\top \mathbf{X}) + \lambda^2 \beta \beta^\top) (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1}$$

while $\text{MSE}(\hat{\beta}) = \text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$.

Régression Ridge

If μ_j is eigenvalue of $\mathbf{X}^\top \mathbf{X}$

- Under \mathcal{A}_1 , $\text{trace}(\text{MSE}(\hat{\beta})) = \sum_{j=1}^p \frac{\sigma^2}{\mu_j};$
- while $\text{trace}(\text{MSE}(\hat{\beta}^{\text{ridge}})) = \sum_{j=1}^r \frac{\sigma^2 \mu_j + \lambda^2 (\mathbf{P}^\top \beta)_j^2}{(\mu_j + \lambda)^2}$
- Thus $\text{trace}(\text{MSE}(\hat{\beta}^{\text{ridge}})) \leq \text{trace}(\text{MSE}(\hat{\beta})) \Leftrightarrow \lambda \leq \frac{2\sigma^2}{\beta^\top \beta}.$

Shrinkage

- Let $\hat{\theta}$ be an unbiased estimator of θ , with variance σ^2 , so $\text{MSE}(\hat{\theta}) = \sigma^2$.
- Let $\lambda > 0$ et $\tilde{\theta} = \frac{\hat{\theta}}{1 + \lambda}$ alors

$$\mathbb{E}(\tilde{\theta}) = \frac{\theta}{1 + \lambda}, \quad \text{Var}(\tilde{\theta}) = \frac{\sigma^2}{(1 + \lambda)^2} \text{ and } \text{MSE}(\tilde{\theta}) = \frac{\lambda^2\theta^2 + \sigma^2}{(1 + \lambda)^2}.$$

It is possible to find λ^* such that $\text{MSE}(\tilde{\theta})$ is minimal

- $\tilde{\theta}$ is biased but with a smaller variance than $\hat{\theta}$.

Similarly, we can find λ^* such that $\lambda \mapsto \text{trace}(\text{MSE}(\hat{\beta}^{\text{ridge}}(\lambda)))$ is minimal

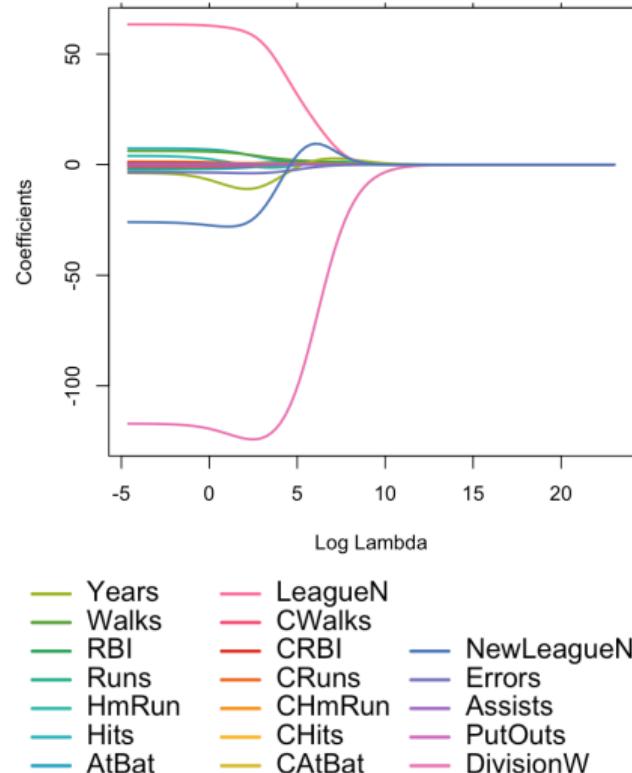
Ridge Regression

Use

```
1 > library(MASS)
2 > ?lm.ridge
```

or

```
1 > library(ISLR)
2 > library(glmnet)
3 > Hitters = na.omit(Hitters)
4 > x = model.matrix(Salary ~ ., Hitters)
   )[, -1]
5 > y = Hitters$Salary
6 > ridge_mod = glmnet(x, y, alpha =
   0)
7 > plot(ridge_mod, var="lambda")
```

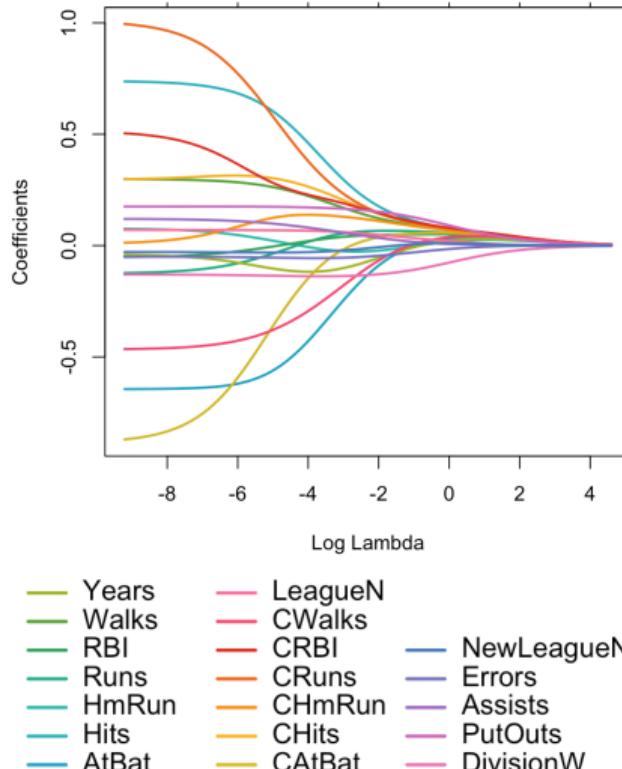


Ridge Regression

instead of the Euclidean norm $\|\beta\|_2$ use
Mahalanobis norm
i.e. center and scale x_j

$$x_j \mapsto \frac{x_j - \bar{x}_j}{s_{x_j}}$$

```
1 > ys = (y-mean(y))/sd(y)
2 > xs = x
3 > for(i in 1:ncol(x)) xs[,i] = (x[,i]
   ]-mean(x[,i]))/sd(x[,i])
4 > ridge_mod_s = glmnet(xs, ys, alpha
   = 0)
5 > plot(ridge_mod_s, var="lambda")
```



LASSO Regression

- Here, we're not constraining the Euclidean $\|\boldsymbol{\beta}\| = \|\boldsymbol{\beta}\|_2$ norm (i.e. the ℓ_2 norm) but the ℓ_1 norm of the coefficients.
- The Lasso method consists in minimizing

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \text{ s.t. } \|\boldsymbol{\beta}\|_1 \leq \delta.$$

- There is no exact solution (several algorithms have been proposed - LARS, coordinate descent algorithm).
- It can be shown that the problem is equivalent to minimizing the regularized problem

$$\hat{\boldsymbol{\beta}}^{\text{lasso}}(\lambda) = \operatorname*{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1$$

($\forall \lambda > 0, \exists \delta > 0$ such that the solutions to these two problems coincide)

Ridge, LASSO and Elastic-Net

Definition 2.12: Best-subset selection (OLS), Hoerl and Kennard (1970b)

$$\hat{\beta}_\lambda^{\text{bss}} = \underset{\beta \in \mathbb{R}^k}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2 + \lambda \sum_{j=1}^k \mathbf{1}_{\beta_j \neq 0} \right\}.$$

Definition 2.13: Best-subset selection (GLM)

$$\hat{\beta}_\lambda^{\text{bss}} = \underset{\beta \in \mathbb{R}^k}{\operatorname{argmin}} \left\{ - \sum_{i=1}^n \log f(y_i | \mu_i = g^{-1}(\mathbf{x}_i^\top \beta)) + \lambda \sum_{j=1}^k \mathbf{1}_{\beta_j \neq 0} \right\}.$$

Ridge, LASSO and Elastic-Net

Definition 2.14: Ridge Estimator (OLS), Hoerl and Kennard (1970b)

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

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

Definition 2.15: Ridge Estimator (GLM)

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

Ridge, LASSO and Elastic-Net

Definition 2.16: LASSO Estimator (OLS), Tibshirani (1996)

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

Definition 2.17: LASSO Estimator (GLM)

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

Ridge, LASSO and Elastic-Net

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).

Ridge, LASSO and Elastic-Net

Suppose that we group the covariates in G groups, $\beta = (\beta_0, \beta_1, \beta_2, \dots, \beta_G)$

Definition 2.18: Group LASSO Estimator (OLS), Tibshirani (1996)

$$\hat{\beta}_\lambda^{G-\text{lasso}} = \operatorname{argmin} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2 + \sum_{g=1}^G \lambda_g \|\beta_g\|_2 \right\}.$$

Definition 2.19: Group LASSO Estimator (GLM)

$$\hat{\beta}_\lambda^{G-\text{lasso}} = \operatorname{argmin} \left\{ - \sum_{i=1}^n \log f(y_i | \mu_i = g^{-1}(\mathbf{x}_i^\top \beta)) + \sum_{g=1}^G \lambda_g \|\beta_g\|_2 \right\}.$$

Ridge, LASSO and Elastic-Net

Suppose we have an adjacency relation on the covariate components, i.e., the covariate component X_j is naturally embedded between the components X_{j-1} and X_{j+1}

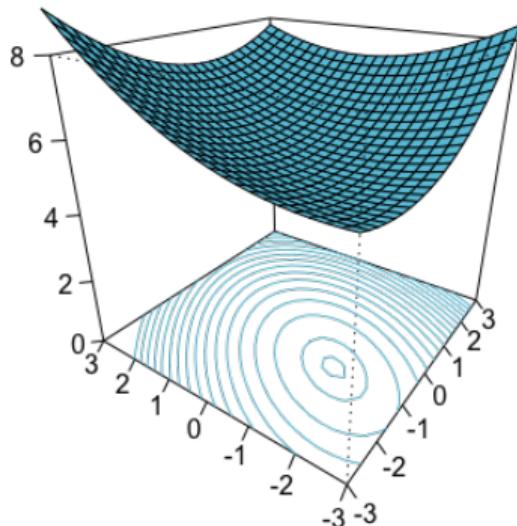
Definition 2.20: Fused LASSO Estimator (OLS), Tibshirani (1996)

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

Definition 2.21: Fused LASSO Estimator (GLM)

$$\hat{\beta}_\lambda^{F\text{-lasso}} = \operatorname{argmin} \left\{ - \sum_{i=1}^n \log f(y_i | \mu_i = g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta})) + \sum_{j=2}^k \lambda_j |\beta_j - \beta_{j-1}| \right\}.$$

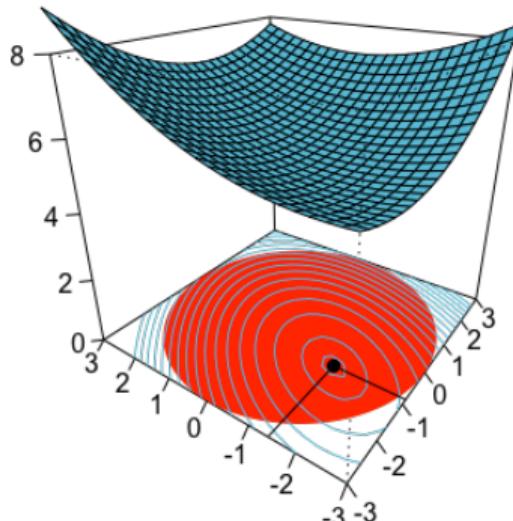
Least Square



Without any constraint, the least square problem is quadratic (and strictly convex)

$$\hat{\beta}^{\text{ols}} = \operatorname{argmin} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\}.$$

Ridge: Least Square with ℓ_2 Penalty



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

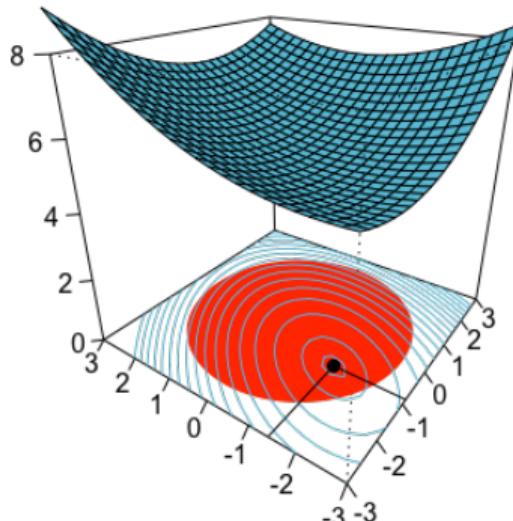
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k \beta_j^2 \leq r_\lambda$$

$$\text{or } \min_{\boldsymbol{\beta}: \boldsymbol{\beta}^\top \boldsymbol{\beta} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ small enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{ridge}} = \hat{\boldsymbol{\beta}}^{\text{ols}}$

Ridge: Least Square with ℓ_2 Penalty



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

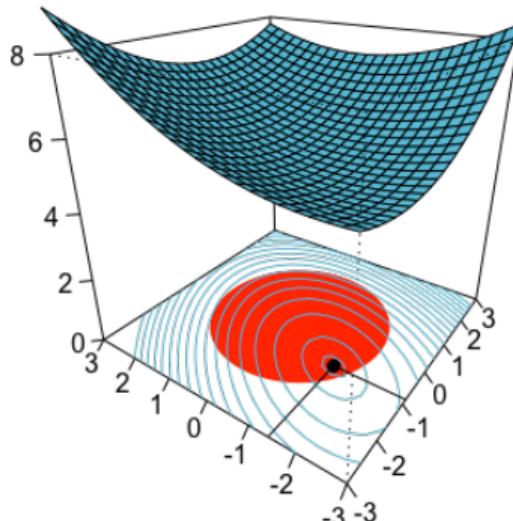
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k \beta_j^2 \leq r_\lambda$$

$$\text{or } \min_{\boldsymbol{\beta}: \boldsymbol{\beta}^\top \boldsymbol{\beta} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ small enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{ridge}} = \hat{\boldsymbol{\beta}}^{\text{ols}}$

Ridge: Least Square with ℓ_2 Penalty



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

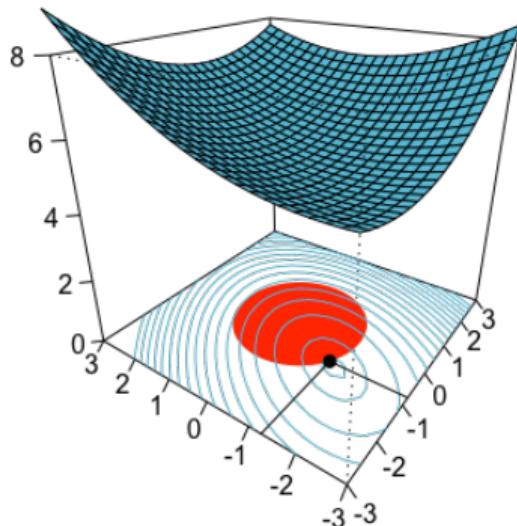
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k \beta_j^2 \leq r_\lambda$$

$$\text{or } \min_{\boldsymbol{\beta}: \boldsymbol{\beta}^\top \boldsymbol{\beta} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ small enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{ridge}} = \hat{\boldsymbol{\beta}}^{\text{ols}}$

Ridge: Least Square with ℓ_2 Penalty



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

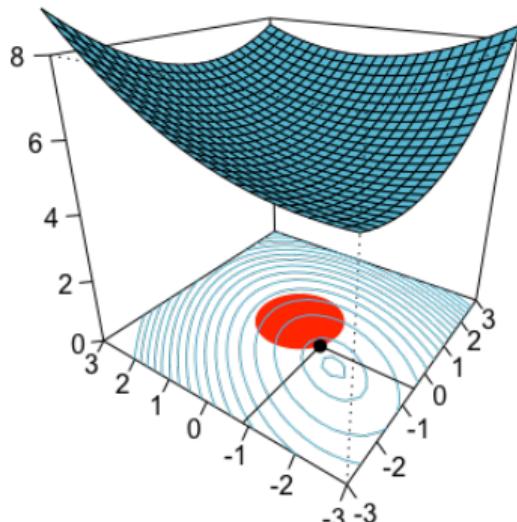
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k \beta_j^2 \leq r_\lambda$$

$$\text{or } \min_{\boldsymbol{\beta}: \boldsymbol{\beta}^\top \boldsymbol{\beta} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ large enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{ridge}} \neq \hat{\boldsymbol{\beta}}^{\text{ols}}$

Ridge: Least Square with ℓ_2 Penalty



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

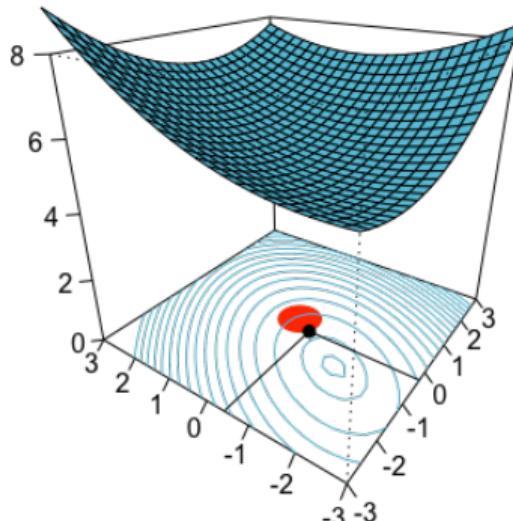
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k \beta_j^2 \leq r_\lambda$$

$$\text{or } \min_{\boldsymbol{\beta}: \boldsymbol{\beta}^\top \boldsymbol{\beta} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ large enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{ridge}} \neq \hat{\boldsymbol{\beta}}^{\text{ols}}$

Ridge: Least Square with ℓ_2 Penalty



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

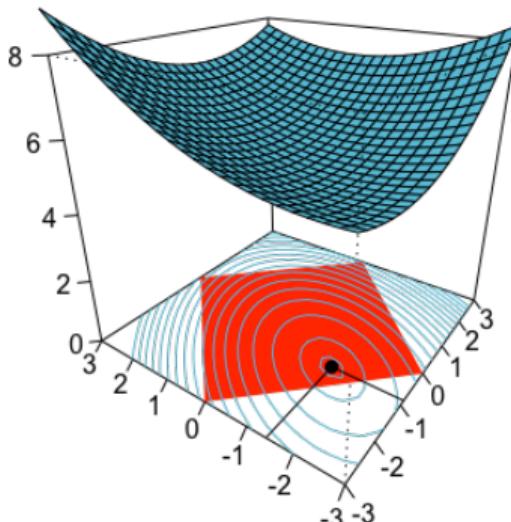
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k \beta_j^2 \leq r_\lambda$$

$$\text{or } \min_{\boldsymbol{\beta}: \boldsymbol{\beta}^\top \boldsymbol{\beta} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ large enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{ridge}} \rightarrow \mathbf{0}$

Lasso: Least Square with ℓ_1 Penalty



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

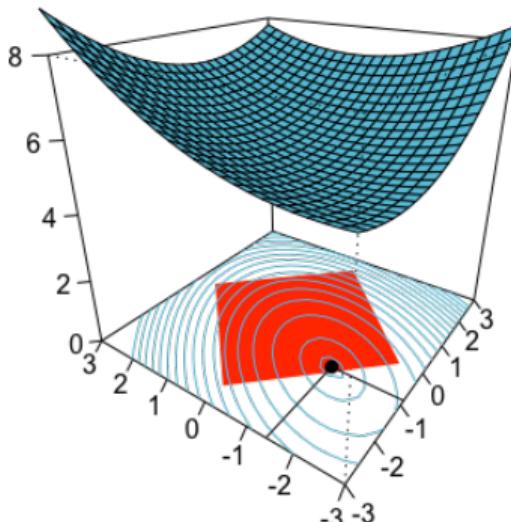
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k |\beta_j| \leq r_\lambda$$

or
$$\min_{\boldsymbol{\beta}: |\boldsymbol{\beta}|^\top \mathbf{1} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ small enough, i.e. $\hat{\beta}_\lambda^{\text{lasso}} = \hat{\beta}^{\text{ols}}$

Lasso: Least Square with ℓ_1 Penalty



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

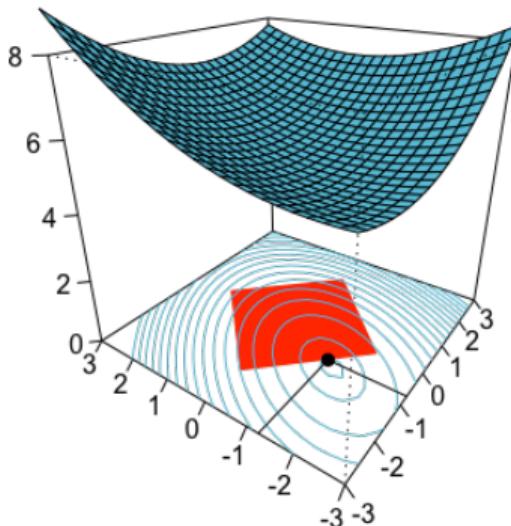
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k |\beta_j| \leq r_\lambda$$

or
$$\min_{\boldsymbol{\beta}: |\boldsymbol{\beta}|^\top \mathbf{1} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ small enough, i.e. $\hat{\beta}_\lambda^{\text{lasso}} = \hat{\beta}^{\text{ols}}$

Lasso: Least Square with ℓ_1 Penalty



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

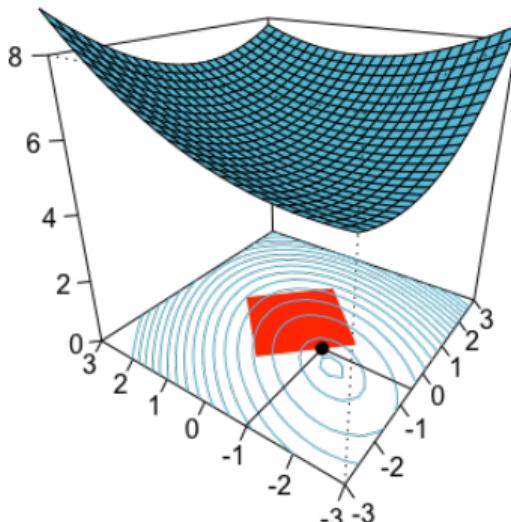
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k |\beta_j| \leq r_\lambda$$

or
$$\min_{\boldsymbol{\beta}: |\boldsymbol{\beta}|^\top \mathbf{1} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ small enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{lasso}} = \hat{\boldsymbol{\beta}}^{\text{ols}}$

Lasso: Least Square with ℓ_1 Penalty



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

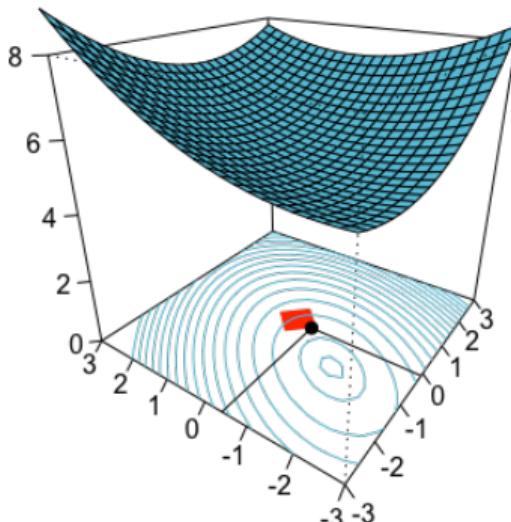
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k |\beta_j| \leq r_\lambda$$

or
$$\min_{\boldsymbol{\beta}: |\boldsymbol{\beta}|^\top \mathbf{1} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ large enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{lasso}} \neq \hat{\boldsymbol{\beta}}^{\text{ols}}$

Lasso: Least Square with ℓ_1 Penalty



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

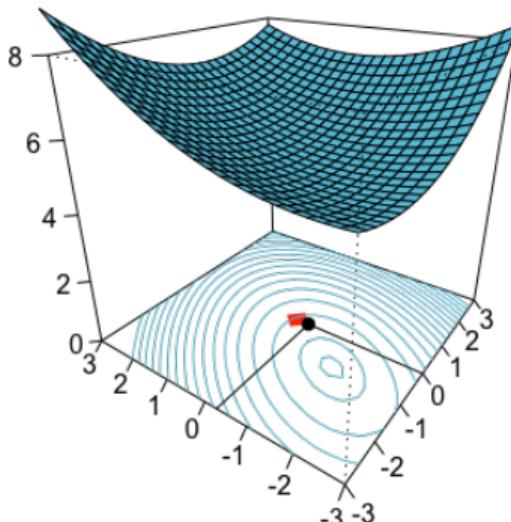
Lagrangian of constrained optimization problem

$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k |\beta_j| \leq r_\lambda$$

or
$$\min_{\boldsymbol{\beta}: |\boldsymbol{\beta}|^\top \mathbf{1} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

with λ large enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{lasso}} \rightarrow \mathbf{0}$

Lasso: Least Square with ℓ_1 Penalty



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

Lagrangian of constrained optimization problem

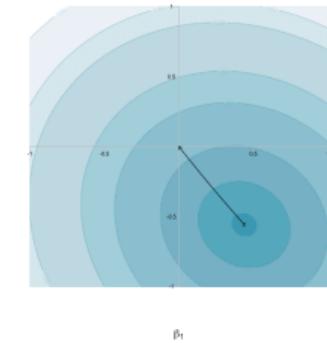
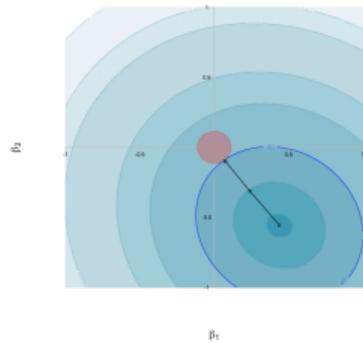
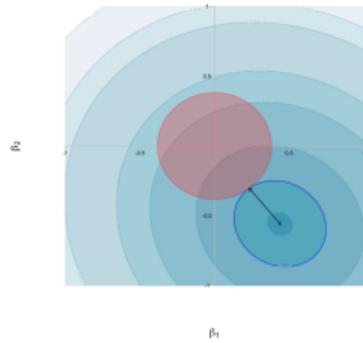
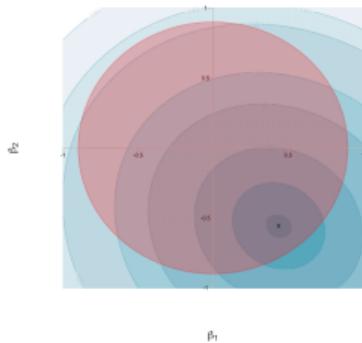
$$\min \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \right\} \text{ s.t. } \sum_{j=1}^k |\beta_j| \leq r_\lambda$$

or
$$\min_{\boldsymbol{\beta}: |\boldsymbol{\beta}|^\top \mathbf{1} \leq r_\lambda} \left\{ \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\},$$

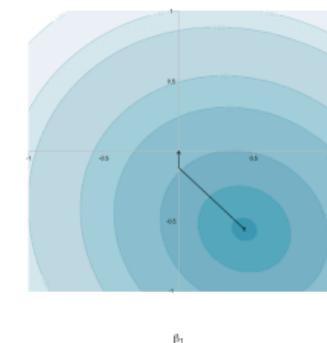
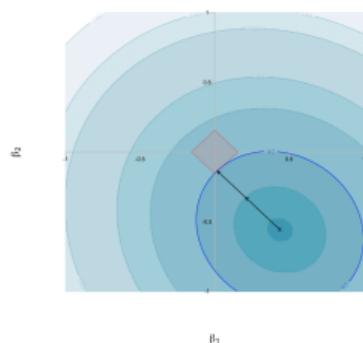
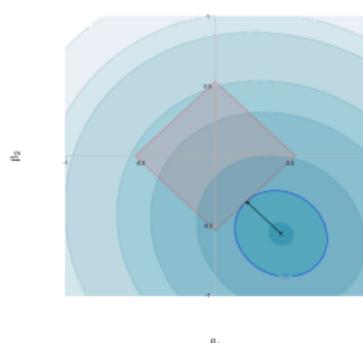
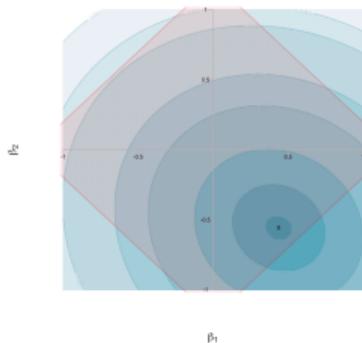
with λ large enough, i.e. $\hat{\boldsymbol{\beta}}_\lambda^{\text{lasso}} \rightarrow \mathbf{0}$

Lasso et Ridge

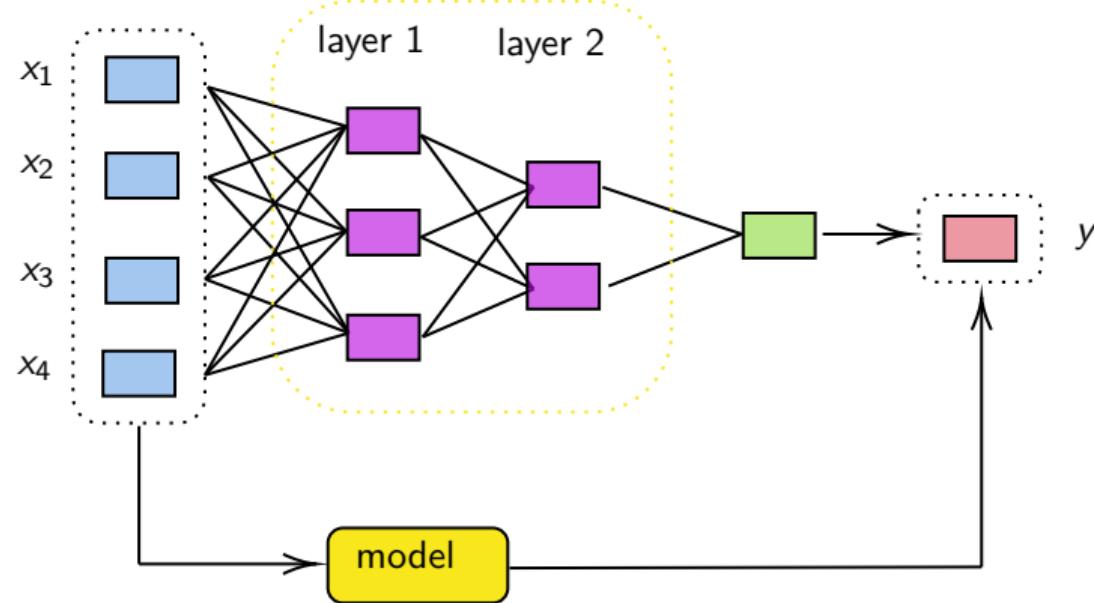
- Ridge regression



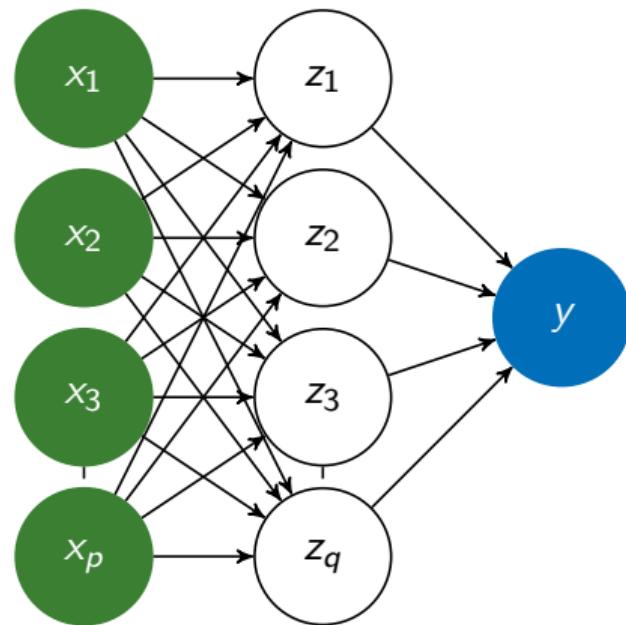
- Lasso regression



Neural Networks



Neural Networks



$$z_k = g \left(\sum_{j=1}^p b_{k,j}^{(1)} x_j \right)$$

so that $\mathbf{z}_{1 \times q} = g(\mathbf{x}_{1 \times p} \mathbf{B}_{p \times q}^{(1)})$

$$\hat{f}(\mathbf{x}) = \hat{y} = h(\mathbf{z}_{1 \times q} \mathbf{B}_{q \times 1}^{(2)})$$

...

$$\hat{f}(\mathbf{x}) = h(g(\mathbf{x}_{1 \times p} \mathbf{B}_{p \times q}^{(1)}) \mathbf{B}_{q \times 1}^{(2)})$$

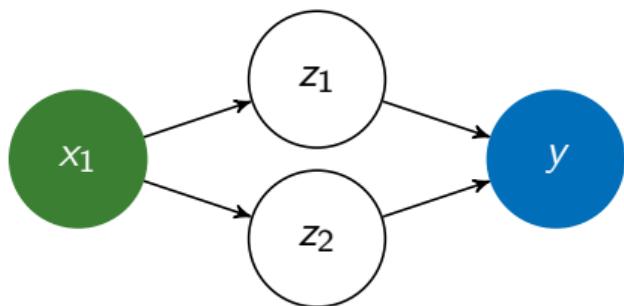
Neural Networks

Activation functions g, h ,

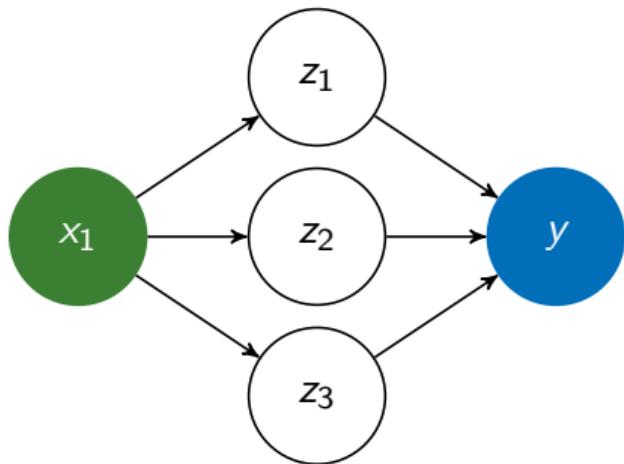
$$\begin{cases} \text{sigmoid (logistic)} & h(z) = (1 + \exp[-z])^{-1} \text{ with derivative } h' = h(1 - h) \\ \text{hyperbolic tangent} & h(z) = \tanh(z) \text{ with derivative } h' = 1 - h^2 \\ \text{ReLU (rectified linear unit)} & h(z) = z_+ = z\mathbf{1}_{z>0} \text{ with derivative } h'(z) = \mathbf{1}_{z>0} \\ \text{GELU (Gaussian error linear unit)} & h(z) = z\varphi(z) \text{ with derivative } h'(z) = \varphi(z) + z\varphi'(z) \end{cases}$$

$$\text{where } \varphi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-z^2}{2}\right), z \in \mathbb{R}.$$

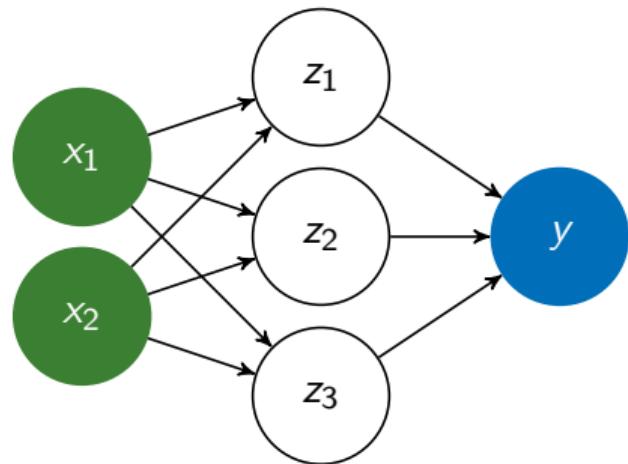
Neural Networks



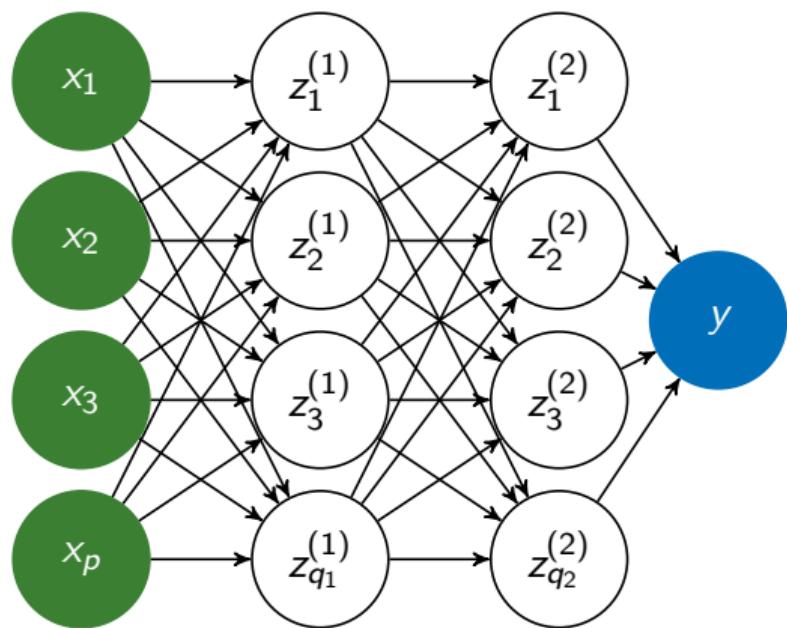
Neural Networks



Neural Networks



Neural Networks



$$z_k^{(1)} = g_1 \left(\sum_{j=1}^p b_{k,j}^{(1)} x_j \right)$$

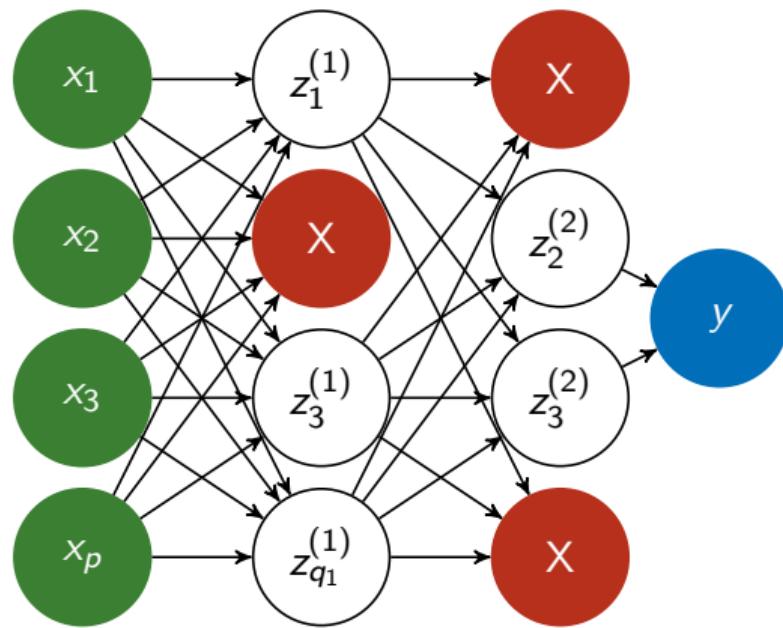
so that $\mathbf{z}_{1 \times q_1}^{(1)} = g(\mathbf{x}_{1 \times p} \mathbf{B}_{p \times q_1}^{(1)})$ and

$$z_k^{(2)} = g_2 \left(\sum_{j=1}^{q_1} b_{k,j}^{(2)} z_j^{(1)} \right)$$

so that $\mathbf{z}_{1 \times q_2}^{(2)} = g_2(\mathbf{z}_{1 \times q_1}^{(1)} \mathbf{B}_{q_1 \times q_2}^{(2)})$

$$\hat{f}(\mathbf{x}) = h(g_2(g_1(\mathbf{x}_{1 \times p} \mathbf{B}_{p \times q_1}^{(1)}) \mathbf{B}_{q_1 \times q_2}^{(2)}) \mathbf{B}_{q_2 \times 1}^{(3)})$$

Neural Networks Dropout



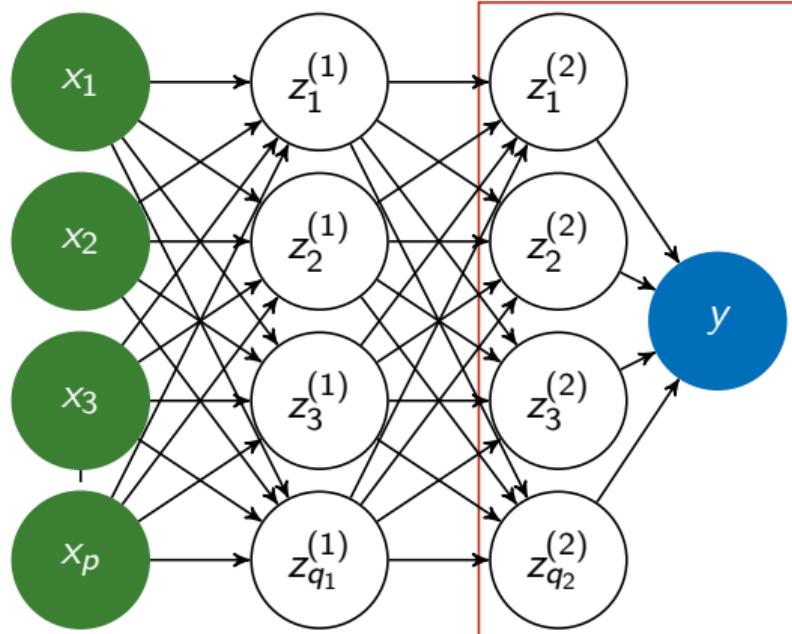
Feed-Forward Neural Network Architecture

Feedforward neural network

Feedforward refers to recognition-inference architecture of neural networks. Artificial neural network architectures are based on inputs multiplied by weights to obtain outputs (inputs-to-output): feedforward. Recurrent neural networks, or neural networks with loops allow information from later processing stages to feed back to earlier stages for sequence processing. However, at every stage of inference a feedforward multiplication remains the core, essential for backpropagation or backpropagation through time. Thus neural networks cannot contain feedback like negative feedback or positive feedback where the outputs feed back to the very same inputs and modify them, because this forms an infinite loop which is not possible to rewind in time to generate an error signal through backpropagation. W

Feed-Forward Neural Network Architecture

Feed-forward neural networks (FNNs)



feature extractor of depth 2

$$\mathbf{x} \mapsto \mathbf{z}^{(2)}(\mathbf{x}) \in \mathbb{R}^{q_2}$$

$$\mu(\mathbf{x}) = g^{-1}(\mathbf{z}^{(2)}(\mathbf{x})^\top \mathbf{b}^{(2)})$$

$\mathbf{b}^{(2)}$ is the output/readout parameter

Feed-Forward Neural Network Architecture

Universal approximation theorems are at the core of the great success of Neural Networks.

Universal approximation theorem

universal approximation theorems are theorems of the following form: Given a family of neural networks, for each function f from a certain function space, there exists a sequence of neural networks ϕ_1, ϕ_2, \dots from the family, such that $\phi_n \rightarrow f$ according to some criterion. The most popular version states that feedforward networks with non-polynomial activation functions are dense in the space of continuous functions between two Euclidean spaces, with respect to the compact convergence topology. W

freakonometrics

freakonometrics.hypotheses.org

Feed-Forward Neural Network Architecture

The main universality theorem statement says that any compactly supported continuous (regression) function can be approximated arbitrarily well by a suitable (and sufficiently large) Neural Networks.

This approximation can be w.r.t. different norms and the assumptions for such a statement to hold are comparably weak, e.g., the sigmoid activation function leads to a class of Neural Networks that are universal.

Feed-Forward Neural Network Architecture

Backpropagation

In machine learning, backpropagation^[1] is a gradient estimation method commonly used for training a neural network to compute its parameter updates. It is an efficient application of the chain rule to neural networks. Backpropagation computes the gradient of a loss function with respect to the weights of the network for a single input–output example, and does so efficiently, computing the gradient one layer at a time, iterating backward from the last layer to avoid redundant calculations of intermediate terms in the chain rule; this can be derived through dynamic programming. W

“The term back-propagation is often misunderstood as meaning the whole learning algorithm for multilayer neural networks. Backpropagation refers only to the method for computing the gradient, while other algorithms, such as stochastic gradient descent, is used to perform learning using this gradient,” Goodfellow et al. (2016)

Feed-Forward Neural Network Architecture

Convolutional neural network

A convolutional neural network (CNN) is a regularized type of feedforward neural network that learns features via filter (or kernel) optimization. This type of deep learning network has been applied to process and make predictions from many different types of data including text, images and audio. W

Recurrent neural network

Recurrent neural networks (RNNs) are a class of artificial neural networks designed for processing sequential data, such as text, speech, and time series,[1] where the order of elements is important. Unlike feedforward neural networks, which process inputs independently, RNNs utilize recurrent connections, where the output of a neuron at one time step is fed back as input to the network at the next time step. This enables RNNs to capture temporal dependencies and patterns within sequences. W

Feed-Forward Neural Network Architecture

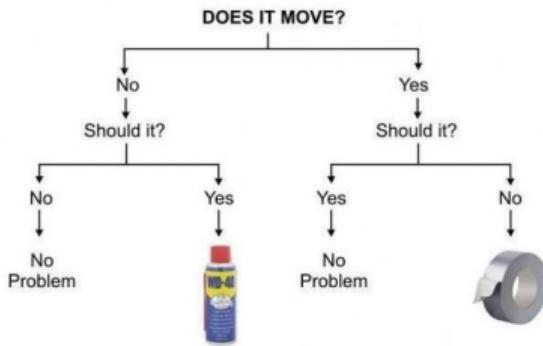
Autoencoder

An autoencoder is a type of artificial neural network used to learn efficient codings of unlabeled data (unsupervised learning). An autoencoder learns two functions: an encoding function that transforms the input data, and a decoding function that recreates the input data from the encoded representation. The autoencoder learns an efficient representation (encoding) for a set of data, typically for dimensionality reduction, to generate lower-dimensional embeddings for subsequent use by other machine learning algorithms. W

Tree Terminology

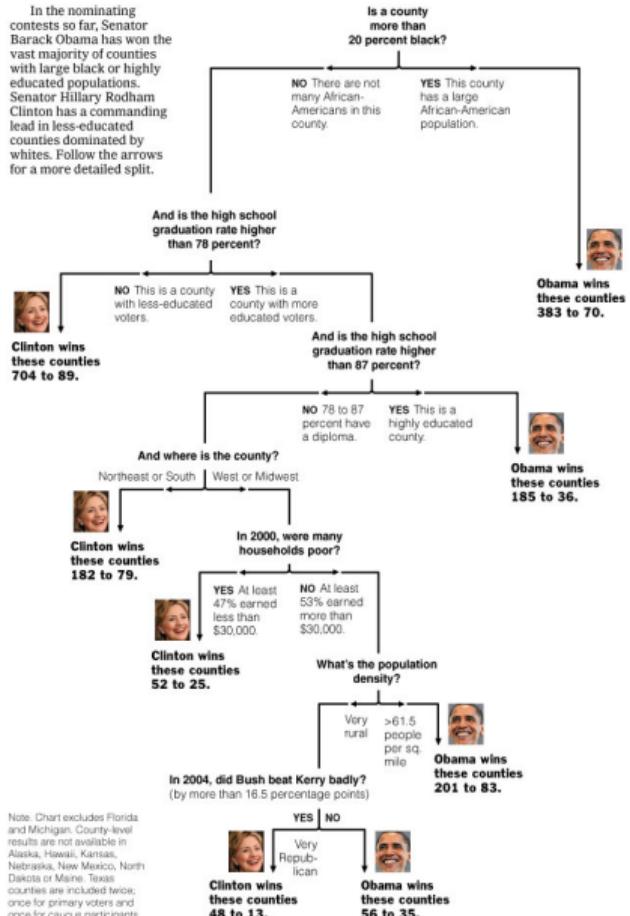
- The final regions obtained are called **terminal nodes** or, more often, **leaves** of the tree.
- The other points where the splits occur are **internal nodes** of the tree.
- The segments of the tree that connect the nodes are **branches**.

via New York Times (2008)



Decision Tree: The Obama-Clinton Divide

In the nominating contests so far, Senator Barack Obama has won the vast majority of counties with large black or highly educated populations. Senator Hillary Rodham Clinton has a commanding lead in less-educated counties dominated by whites. Follow the arrows for a more detailed split.



Note: Chart excludes Florida and Michigan. County-level results are not available in Alaska, Hawaii, Kansas, Nebraska, North Dakota or Maine. Texas counties are included twice: once for primary voters and once for caucus participants.

Decision Trees

We must therefore subdivide (partition) \mathcal{X} into separate disjoint regions,

$$\forall i, j, \ i \neq j \ (\text{i.e. } R_i \neq R_j) \Rightarrow R_i \cap R_j = \emptyset \text{ and } \bigcup_{i=1}^m R_i = \mathcal{X}.$$

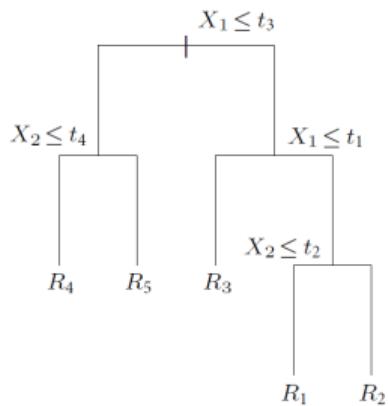
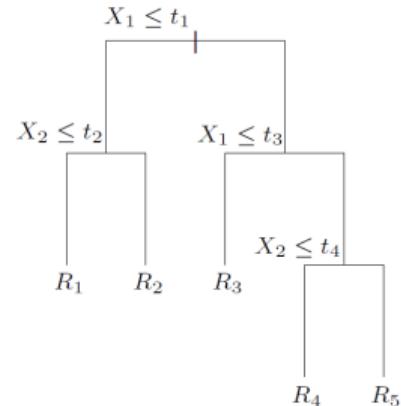
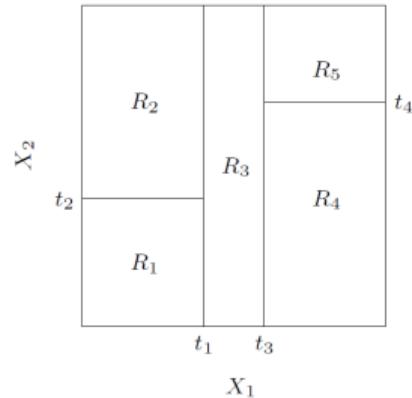
It can take the form:

$$\hat{Y}_i = \begin{cases} c_1 & \text{if } \mathbf{x}_i \in R_1 \\ c_2 & \text{if } \mathbf{x}_i \in R_2 \\ \dots \\ c_m & \text{if } \mathbf{x}_i \in R_m \end{cases} \quad \text{or for an arbitrary point } \mathbf{x}, \quad \hat{Y} = \begin{cases} c_1 & \text{if } \mathbf{x} \in R_1 \\ c_2 & \text{if } \mathbf{x} \in R_2 \\ \dots \\ c_m & \text{if } \mathbf{x} \in R_m \end{cases}$$

or can be presented as an additive model:

$$\hat{Y} = \sum_{j=1}^m c_j \mathbf{1}_{\mathbf{x} \in R_j}$$

Decision Trees



(from Friedman et al. (2001))

Objective

For a **regression** problem, we often work with the sum of squared errors (SSE) or the mean square error (MSE).

$$\text{MSE} = \sum_{j=1}^J \sum_{i: \mathbf{X}_i \in R_j} (y_i - \hat{y}_{R_j})^2,$$

where \hat{y}_{R_j} is the predicted value in the R_j region. We then look for a partition R_1, \dots, R_J that **minimizes** the SSE.

For a **classification** problem, note that

$$\text{MSE}_j = \sum_{i: \mathbf{X}_i \in R_j} (y_i - \hat{y}_{R_j})^2 = n_{0,j}(0 - \hat{y}_{R_j})^2 + n_{1,j}(1 - \hat{y}_{R_j})^2$$

Objective

$$\max\{\text{MSE}_j\} = n_{0,j} \left(\frac{n_{0,j}}{n_{0,j} + n_{1,j}} \right)^2 + n_{1,j} \left(\frac{n_{0,j}}{n_{1,j} + n_{1,j}} \right)^2 = \frac{n_{0,j}n_{1,j}}{n_{0,j} + n_{1,j}}$$

so that

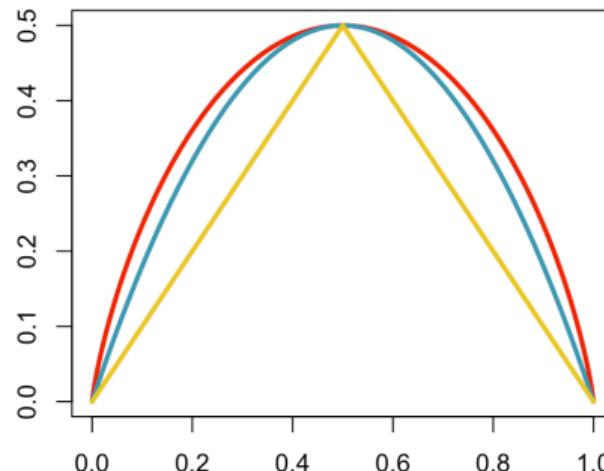
$$MSE = \sum_{j=1}^J \frac{n_{0,j}n_{1,j}}{n_{0,j} + n_{1,j}} = \sum_{j=1}^J n_j \cdot p_{1,j}(1 - p_{1,j})$$

This is **Gini impurity index**.

Objective

Formally, impurity is a function ψ , with ψ positive, symmetrical ($\psi(p) = \psi(1 - p)$), minimal in 0 (and 1).

- Misclassification, $\psi(p) = 1 - \max\{p, 1 - p\}$
- Gini, $\psi(p) = p(1 - p)$
- (cross) entropy, $\psi(p) = -p \log p - (1 - p) \log(1 - p)$



Objective

- Again, if we consider all possible partitions of the space of explanatory variables, we'll have to consider too many possibilities.
- Instead, we'll use a top-down greedy approach via recursive binary splitting.
- **Descending** algorithm: we start with all observations in the same class (root of the tree) and divide the space into smaller and smaller regions.
- **Greedy** algorithm: optimization is performed at each step without looking back or forward.

freakonometrics

freakonometrics.hypotheses.org

Recursive Binary Division

First step: we select the explanatory variable X_j and the point c such that the division into two regions

$$R_1(j, c) = \{X_1, \dots, X_J | X_j < c\} \quad R_2(j, c) = \{X_1, \dots, X_J | X_j \geq c\}$$

leads to the largest positive reduction of the objective function. If the sum of squared errors has been chosen, we then seek to minimize

$$\sum_{i: \mathbf{X}_i \in R_1(j, c)} (Y_i - \hat{Y}_{R_1(j, c)})^2 + \sum_{i: \mathbf{X}_i \in R_2(j, c)} (Y_i - \hat{Y}_{R_2(j, c)})^2.$$

(or any impurity function)

Recursive Binary Division

Next steps: in step k , we repeat the procedure described in the first step for each of the regions created in step $k - 1$, i.e., for each of the regions r , we must identify the explanatory variable X_j and the point c which will minimize the objective function after the division given by

$$\{\mathbf{X} \in r \mid X_j < c\} \text{ and } \{\mathbf{X} \in r \mid X_j \geq c\}.$$

The procedure is stopped when a certain criterion is reached.

In the absence of stopping criteria, the resulting tree will have n leaves to overfit.

Hierarchical Sequential Construction

Classification = grouping individuals into a limited number of classes

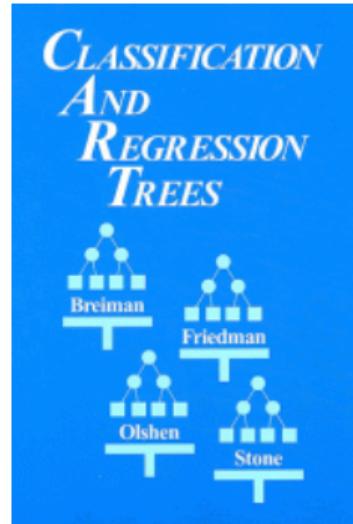
These classes are built up as we go along.

→ group together similar individuals, séparer des individus ayant des caractéristiques proches.

History: dates from the 1960s,

then Breiman et al. (1984)

tools that have become popular in machine learning



Hierarchical Sequential Construction

Downward classification :

we select the most related of the explanatory variables

to the variable to be explained Y ,

.3cm

→ gives a first division of the sample

.3cm

we repeat in each class

.3cm

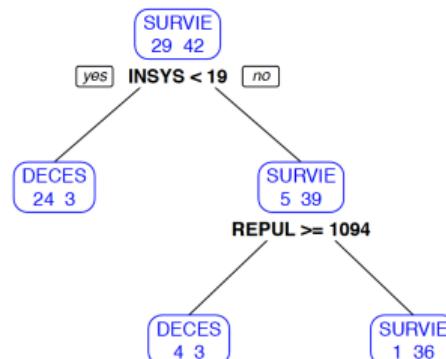
→ each class must be the most homogeneous positive, in Y .

Difference from logistic regression

.3cm

sequential use of explanatory variables

.3cm



Cutting $\mathcal{X} \subset \mathbb{R}$

$Y \in \{0, 1\}$ and $X \in \mathbb{R}$: cut according to threshold s ,

$$\begin{cases} \tilde{X} = L \text{ if } X \leq s \\ \tilde{X} = R \text{ if } X > s \end{cases}$$

| | | $\tilde{X} = L$ | $\tilde{X} = R$ | |
|---------|-----------|-----------------|-----------------|-----|
| | | $X \leq s$ | $X > s$ | |
| $Y = 0$ | $n_{L,0}$ | $n_{R,0}$ | $n_{\cdot,0}$ | |
| | $n_{L,1}$ | $n_{R,1}$ | $n_{\cdot,1}$ | |
| | | $n_{L,\cdot}$ | $n_{R,\cdot}$ | n |

Gini $\text{gini}(Y|\tilde{X}) = \sum_{x \in \{L,R\}} \frac{n_{x,\cdot}}{n} \sum_{y \in \{0,1\}} \frac{n_{x,y}}{n_{x,\cdot}} \left(1 - \frac{n_{x,y}}{n_{x,\cdot}}\right)$

Cutting $\mathcal{X} \subset \mathbb{R}$

$Y \in \{0, 1\}$ and $X \in \mathbb{R}$: we cut according to a threshold s ,

$$\begin{cases} \tilde{X} = L \text{ if } X \leq s \\ \tilde{X} = R \text{ if } X > s \end{cases}$$

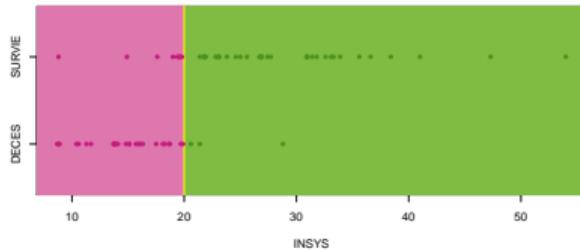
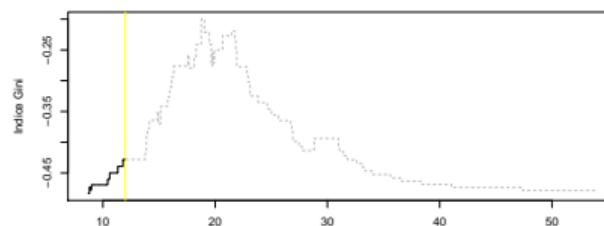
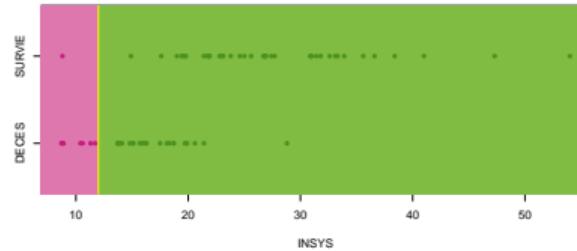
| | | $\tilde{X} = L$ | $\tilde{X} = R$ | |
|---------|-----------|-----------------|-----------------|-----|
| | | $X \leq s$ | $X > s$ | |
| $Y = 0$ | $n_{L,0}$ | $n_{R,0}$ | $n_{\cdot,0}$ | |
| $Y = 1$ | $n_{L,1}$ | $n_{R,1}$ | $n_{\cdot,1}$ | |
| | | $n_{L,\cdot}$ | $n_{R,\cdot}$ | n |

Entropy $\text{entropy}(Y|X) = - \sum_{x \in \{L,R\}} \frac{n_{x,\cdot}}{n} \sum_{y \in \{0,1\}} \frac{n_{x,y}}{n_{x,\cdot}} \log \left(\frac{n_{x,y}}{n_{x,\cdot}} \right)$

Cutting $\mathcal{X} \subset \mathbb{R}$

Gini index

$$\sum_{x \in \{A, B\}} \frac{n_{x,\cdot}}{n} \sum_{y \in \{0,1\}} \frac{n_{x,y}}{n_{x,\cdot}} \left(1 - \frac{n_{x,y}}{n_{x,\cdot}}\right)$$



Cutting $\mathcal{X} \subset \mathbb{R}$

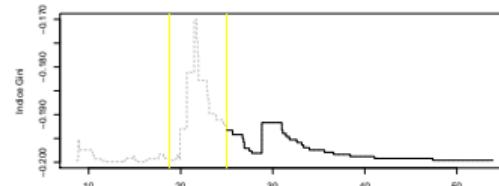
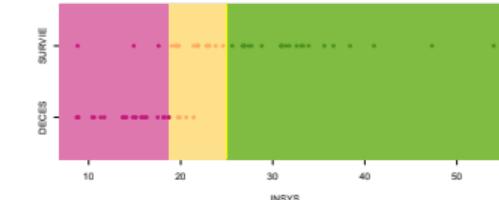
Fix threshold s , and find the second best cut

$$A = (-\infty, s_2] \quad B = (s_2, s] \quad C = (s, \infty)$$

| | $\tilde{X} = A$ $X \leq s_2$ | $\tilde{X} = B$ $X \in (s_2, s]$ | $\tilde{X} = C$ $X > s$ | |
|---------|---------------------------------|-------------------------------------|----------------------------|---------------|
| $Y = 0$ | $n_{A,0}$ | $n_{B,0}$ | $n_{C,0}$ | $n_{\cdot,0}$ |
| $Y = 1$ | $n_{A,1}$ | $n_{B,1}$ | $n_{C,1}$ | $n_{\cdot,1}$ |
| | $n_{A,\cdot}$ | $n_{B,\cdot}$ | $n_{C,\cdot}$ | n |

Using Gini index

$$\sum_{x \in \{A, B, C\}} \frac{n_{x,\cdot}}{n} \sum_{y \in \{0,1\}} \frac{n_{x,y}}{n_{x,\cdot}} \left(1 - \frac{n_{x,y}}{n_{x,\cdot}} \right)$$



Cutting $\mathcal{X} \subset \mathbb{R}$

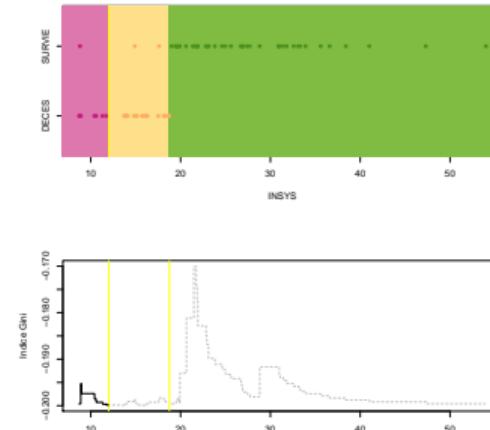
Fix threshold s , then seek for the second one

$$A = (-\infty, s] \quad B = (s, s_2] \quad C = (s_2, \infty)$$

| | $\tilde{X} = A$ $X \leq s$ | $\tilde{X} = B$ $X \in (s, s_2]$ | $\tilde{X} = C$ $X > s_2$ | |
|---------|-------------------------------|-------------------------------------|------------------------------|---------------|
| $Y = 0$ | $n_{A,0}$ | $n_{B,0}$ | $n_{C,0}$ | $n_{\cdot,0}$ |
| $Y = 1$ | $n_{A,1}$ | $n_{B,1}$ | $n_{C,1}$ | $n_{\cdot,1}$ |
| | $n_{A,\cdot}$ | $n_{B,\cdot}$ | $n_{C,\cdot}$ | n |

Gini index

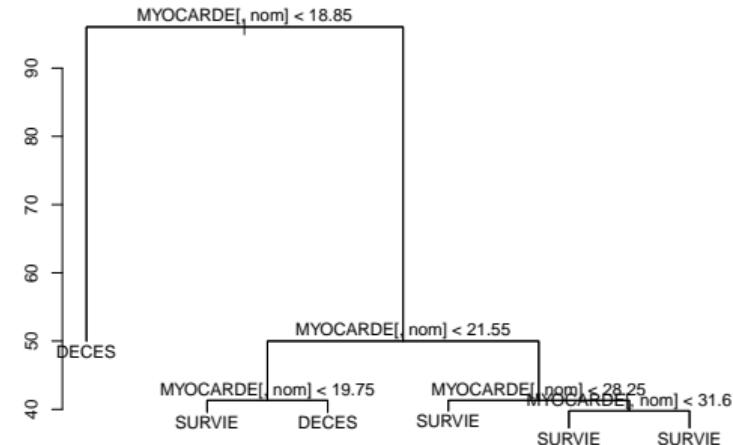
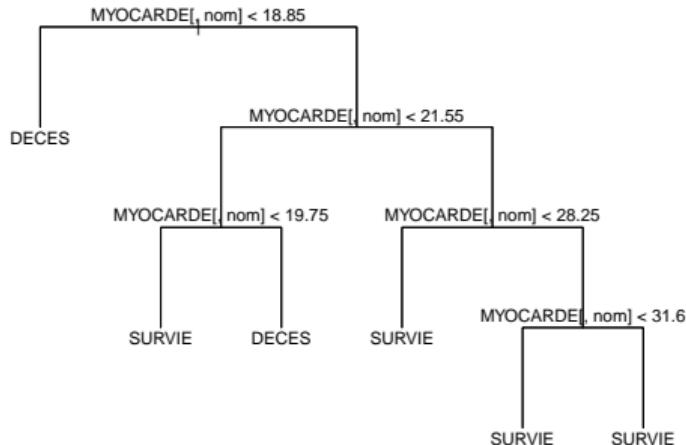
$$\sum_{x \in \{A, B, C\}} \frac{n_{x,\cdot}}{n} \sum_{y \in \{0, 1\}} \frac{n_{x,y}}{n_{x,\cdot}} \left(1 - \frac{n_{x,y}}{n_{x,\cdot}} \right)$$



Prunning Trees

Step 1 Construction of the tree by a recursive process of binary divisions

Step 2 Pruning the tree, removing branches that are too empty or not representative enough → need a pruning criterion (entropy gain)



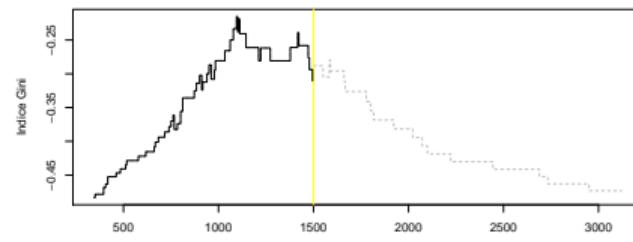
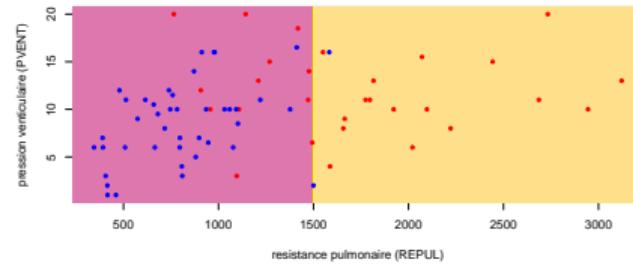
Multiple Numeric Variables

$Y \in \{0, 1\}$ and $X_1, X_2 \in \mathbb{R}$: cut X_1 based on threshold s ,

$$\begin{cases} \tilde{X} = A \text{ if } X_1 \leq s \\ \tilde{X} = B \text{ if } X_1 > s \end{cases}$$

| | $\tilde{X} = A$ | $\tilde{X} = B$ | |
|--------------|-----------------|-----------------|---------------|
| $X_1 \leq s$ | $n_{A,0}$ | $n_{B,0}$ | $n_{\cdot,0}$ |
| $X_1 > s$ | $n_{A,1}$ | $n_{B,1}$ | $n_{\cdot,1}$ |
| | $n_{A,\cdot}$ | $n_{B,\cdot}$ | n |

$$\sum_{x \in \{A, B\}} \frac{n_{x,\cdot}}{n} \sum_{y \in \{0, 1\}} \frac{n_{x,y}}{n_{x,\cdot}} \left(1 - \frac{n_{x,y}}{n_{x,\cdot}} \right)$$



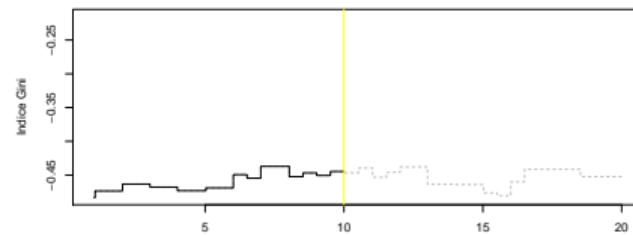
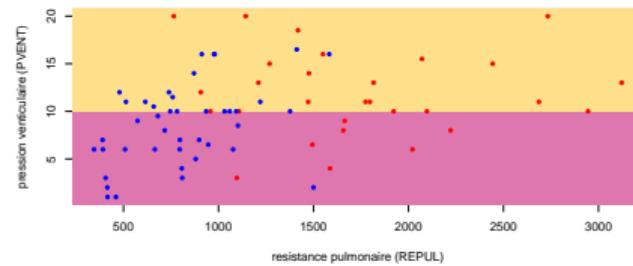
Multiple Numeric Variables

$Y \in \{0, 1\}$ and $X_1, X_2 \in \mathbb{R}$: cut X_2 based on threshold s ,

$$\begin{cases} \tilde{X} = A \text{ if } X_2 \leq s \\ \tilde{X} = B \text{ if } X_2 > s \end{cases}$$

| | $\tilde{X} = A$ | $\tilde{X} = B$ | |
|--------------|-----------------|-----------------|---------------|
| $X_2 \leq s$ | $n_{A,0}$ | $n_{B,0}$ | $n_{\cdot,0}$ |
| $X_2 > s$ | $n_{A,1}$ | $n_{B,1}$ | $n_{\cdot,1}$ |
| | $n_{A,\cdot}$ | $n_{B,\cdot}$ | n |

$$\sum_{x \in \{A, B\}} \frac{n_{x,\cdot}}{n} \sum_{y \in \{0, 1\}} \frac{n_{x,y}}{n_{x,\cdot}} \left(1 - \frac{n_{x,y}}{n_{x,\cdot}} \right)$$



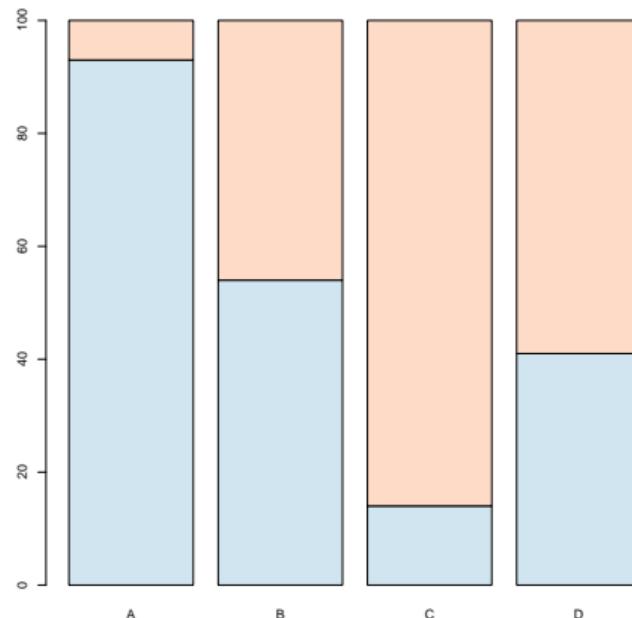
Categorical Variables

X taking values $\{a, b, c, d\}$.

Instead of making a tree by successive cuttings, we can make a tree by successive grouping.

$$\{a, b, c, d\}$$

$$\left\{ \begin{array}{l} \{(a, b), c, d\} \quad \{(a, d), b, c\} \quad \{(a, d), b, c\} \\ \{(b, c), a, d\} \quad \{\textcolor{red}{(b, d), a, c}\} \quad \{(c, d), a, b\} \\ \{(b, c, a), d\} \quad \{\textcolor{red}{(b, c, d), a}\} \quad \{(b, c), (a, d)\} \end{array} \right.$$

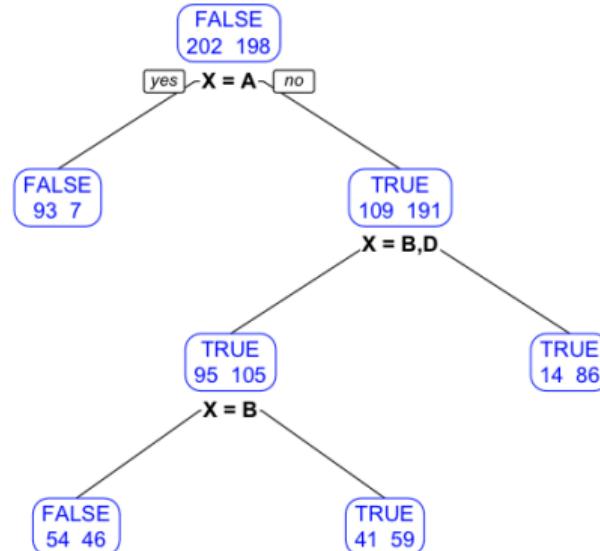


Categorical Variables

X taking values $\{a, b, c, d\}$.

Instead of making a tree by successive cuttings, we can make a tree by successive regrouping

$$\begin{array}{c} \{a, b, c, d\} \\ \left\{ \begin{array}{lll} \{(a, b), c, d\} & \{(a, d), b, c\} & \{(a, d), b, c\} \\ \{(b, c), a, d\} & \{(b, d), a, c\} & \{(c, d), a, b\} \\ \{(b, c, a), d\} & \{(b, c, d), a\} & \{(b, c), (a, d)\} \end{array} \right. \end{array}$$



Categorical Variables

X taking values $\{a, b, c, d\}$.

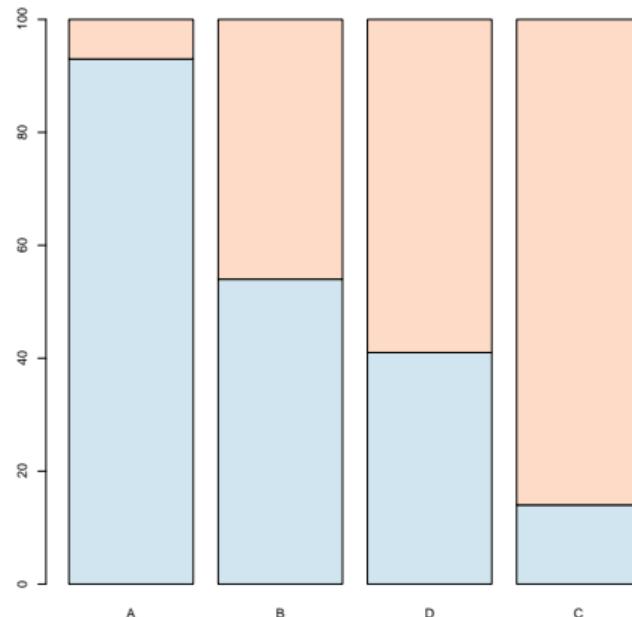
Classically, order the modalities

$$\{a, b, d, c\}$$

then find the optimum cut

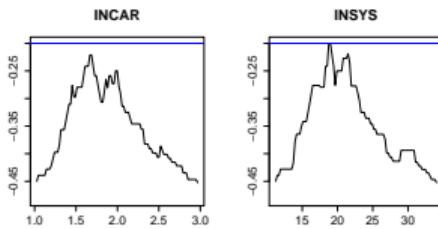
$$\{(a), (b, d, c)\} \quad \{(a, b), (d, c)\} \quad \{(a, b, d), (c)\}$$

$$\{(a), (b, d), (c)\} \quad \{(a), (b), (d, c)\}$$

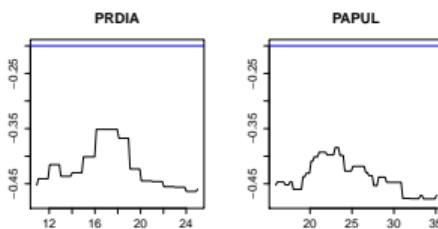


Implementation

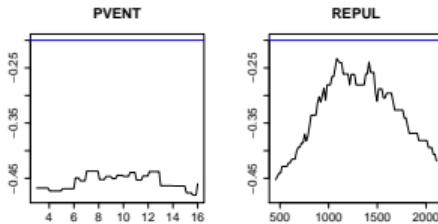
$$N_L: \{x_{i,j} \leq s\} \quad N_R: \{x_{i,j} > s\}$$



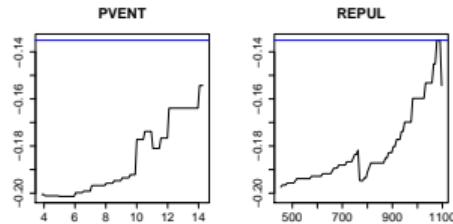
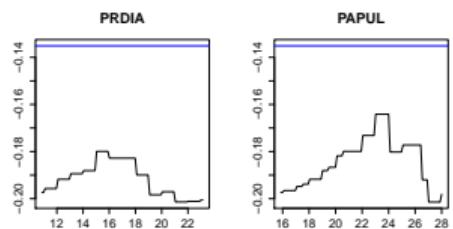
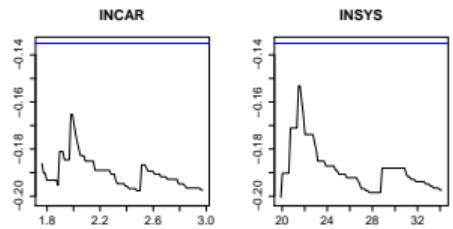
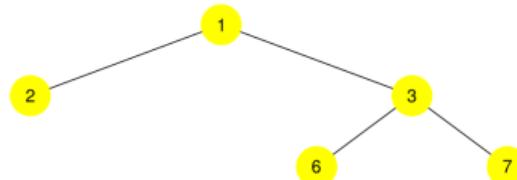
Solve $\max_{j \in \{1, \dots, k\}, s} \{\mathcal{I}(N_L, N_R)\}$



← first split



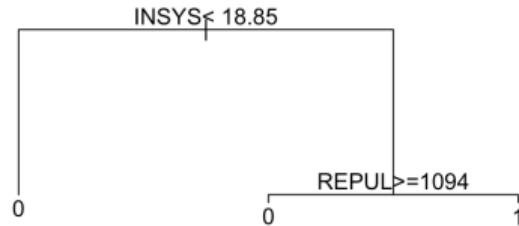
second split →



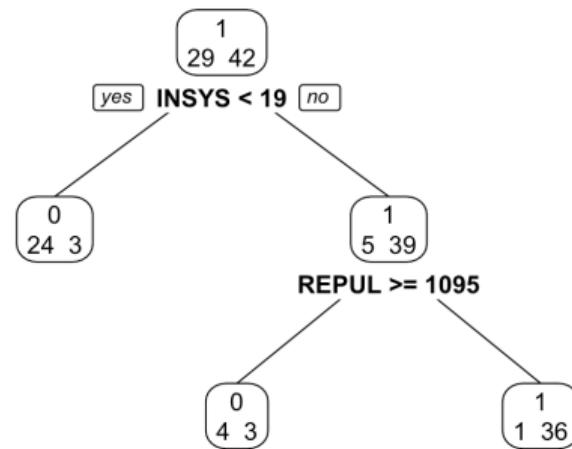
Trees with R

```
1 > loc = "http://freakonomics.free.fr/saporta.csv"
2 > myocarde = read.table(loc, header=TRUE, sep=";")
3 > levels(myocarde$PRONO) = c("0", "1")
```

```
1 > library(rpart)
2 > cart = rpart(PRONO~., data=
   myocarde)
3 > plot(cart)
4 > text(cart)
```



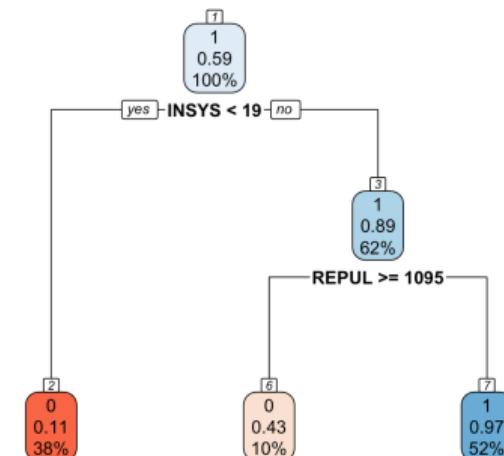
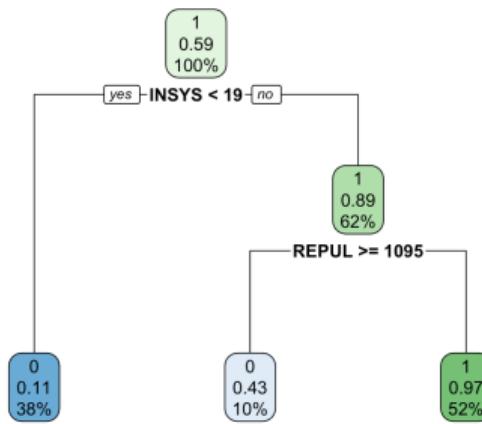
```
1 > library(rpart.plot)
2 > prp(cart, type=2, extra
      =1)
```



Trees with R

```
1 > print(cart, digits = 2)
2 1) root 71 29 1 (0.408 0.592)
3   2) INSYS< 19 27 3 0 (0.889 0.111) *
4   3) INSYS>=19 44 5 1 (0.114 0.886)
5     6) REPUL>=1094.5 7 3 0 (0.571 0.429) *
6     7) REPUL< 1094.5 37 1 1 (0.027 0.973) *
```

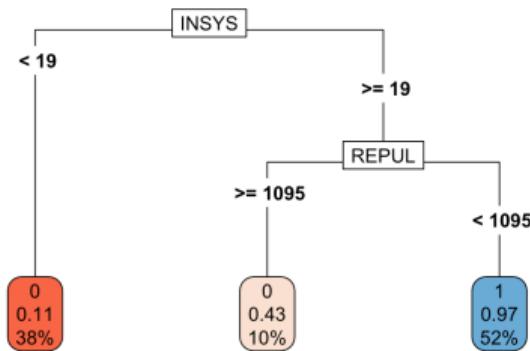
```
1 > rpart.plot(cart)
1 > rpart.plot(cart)
1 > rpart.plot(cart, box.palette="RdBu", nn=TRUE
                  )
```



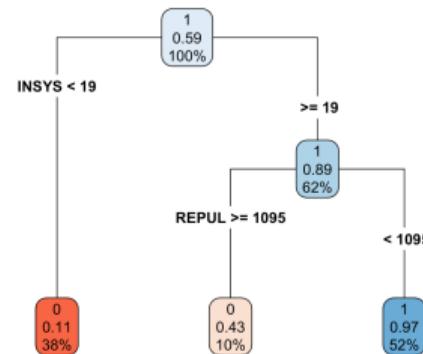
Trees with R

```
1 > rpart.rules(cart, extra = 4, cover = TRUE)
2 PRONO 0    1                                cover
3   0 [.89 .11] when INSYS < 19            38%
4   0 [.57 .43] when INSYS >= 19 & REPUL >= 1095 10%
5   1 [.03 .97] when INSYS >= 19 & REPUL < 1095 52%
```

```
1 > rpart.plot(cart, box.palandte="RdBu", type=5)
```

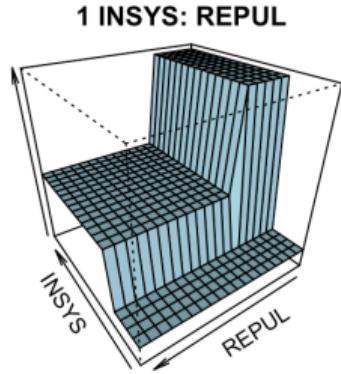
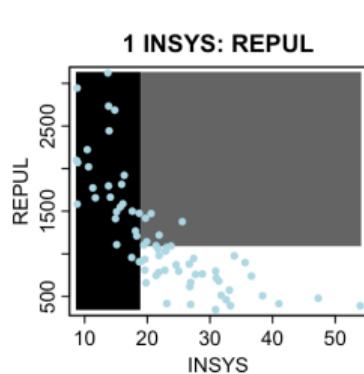


```
1 > rpart.plot(cart, box.palandte="RdBu", type=4)
```

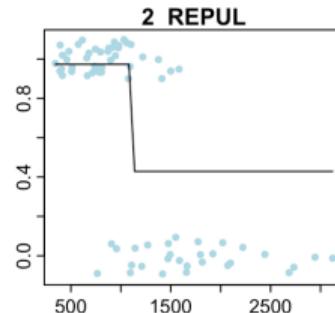
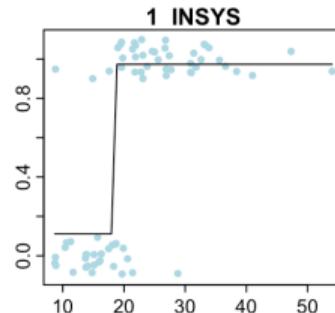


Trees with R

```
1 > library(plotmo)
2 > plotmo(cart, type = "prob",
  nresponse = "1")
```



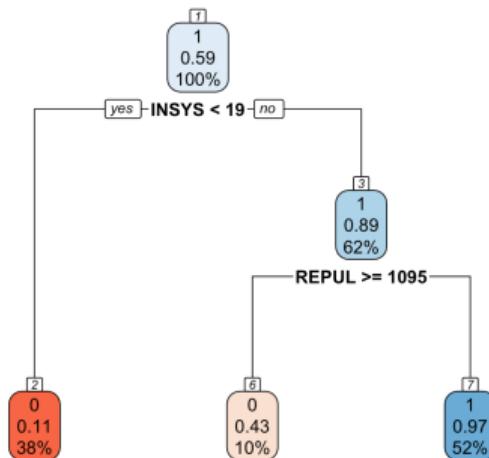
```
1 > plotmo(cart, type = "prob",
  nresponse = "1", type2 = "image",
  ngrid2 = 200)
```



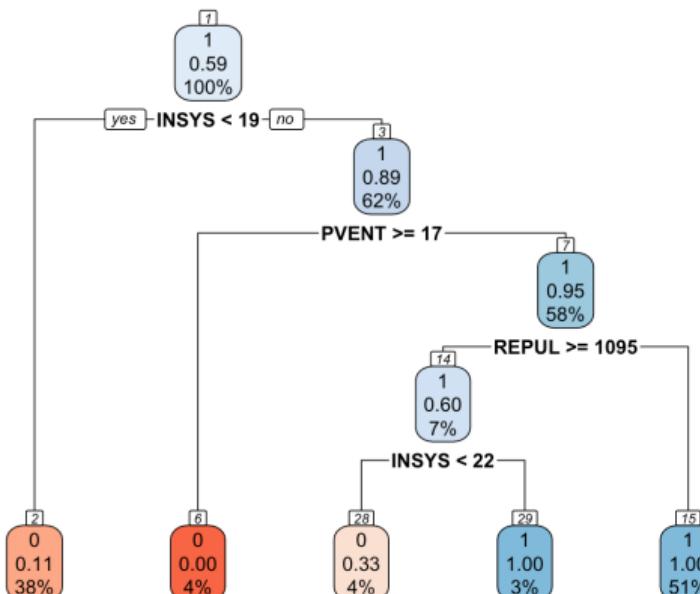
```
1 > path2root = function(node){
2   if(node == 1) node
3   else c(node, path.to.root(node %% 2))}
```

Trees with R

```
1 > cart = rpart(PRONO~.,
2   data=myocarde)
3 > rpart.plot(cart, box.
4   palandte="RdBu", nn=
5   TRUE)
```

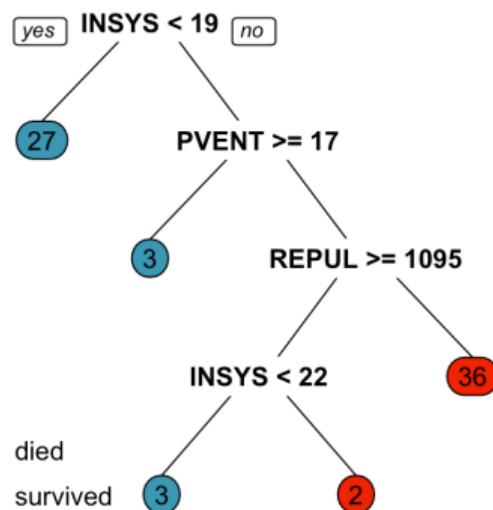


```
1 > cart = rpart(PRONO~.,
2   data=myocarde,
3   minsplit = 5,
4   cp = 0.01)
5 rpart.plot(cart, box.palandte="
6   RdBu", nn=TRUE)
```

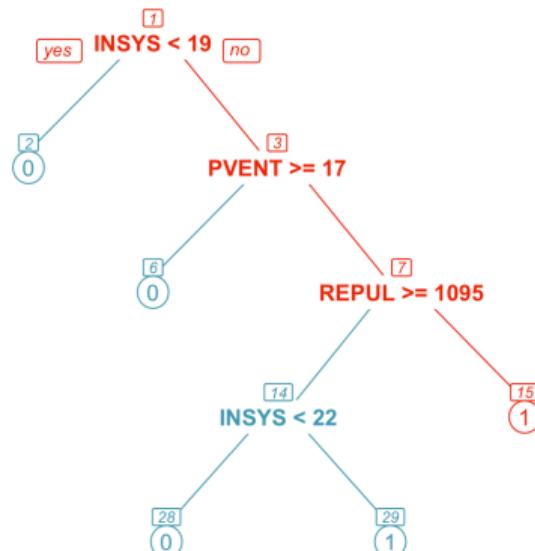


Trees with R

```
1 > boxcols =c("red","blue")
   [cart$frame$yval]
2 > prp(cart, faclen = 0,
       node.fun=only_count,
       box.col = boxcols)
```



```
1 > node = 15
2 > nodes = as.numeric(row.names(cart$frame))
3 > cols = ifelse(nodes %in%
path2root(node),"red","blue")
4 prp(cart, nn = TRUE, col = cols, branch.col =
cols, split.col = cols)
```



Stopping criteria

- Stopping criteria often include a minimum number of observations in each region.
- Using as Stopping criteria the fact of having a decrease in the objective function that exceeds a certain threshold is generally not a good idea.
- Even if an appropriate Stopping criteria is selected, there's a good chance that the tree will lead to over-fitting, i.e. it will perform (too) well on a training database but will be poor on an evaluation database.

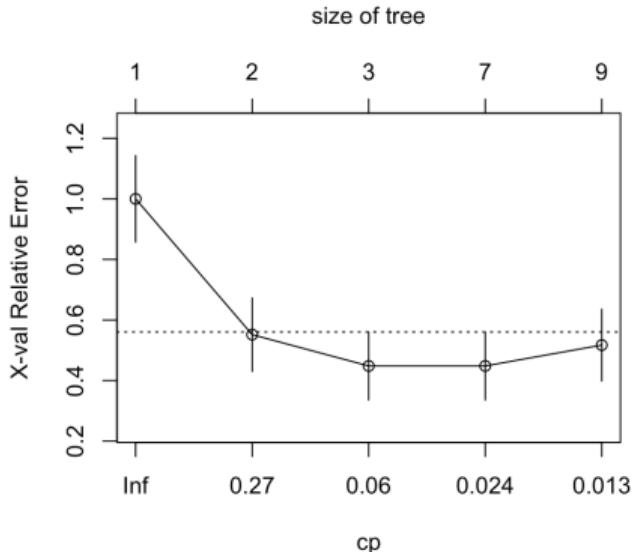
Formally, to find out if a leaf $\{N\}$ is cut into two $\{N_L, N_R\}$, we measure the variation in impurity

$$\Delta \mathcal{I}(N_L, N_R) = \mathcal{I}(N) - \mathcal{I}(N_L, N_R) = \mathcal{I}(N) - \left(\frac{n_L}{n} \mathcal{I}(N_L) + \frac{n_R}{n} \mathcal{I}(N_R) \right)$$

Stopping criteria

Stop if $\Delta\mathcal{I}(N_L, N_R)/\mathcal{I}(N)$ exceeds **cp**
(complexity parameter, par défaut 1%).

```
1 > cart = rpart(PRONO ~ ., data =
+     myocarde,
+     minsplit=3)
2
3 > plotcp(cart)
```



Stopping criteria

```
1 > prune(cart, cp=0.06)
2 node), split, n, loss, yval, (yprob) *= terminal node
3
4 1) root 71 29 SURVIE (0.40845070 0.59154930)
5   2) INSYS< 18.85 27  3 DECES (0.888889 0.111111) *
6   3) INSYS>=18.85 44  5 SURVIE (0.113636 0.886363)
7     6) PVENT>=17.25 3  0 DECES (1.000000 0.000000) *
8     7) PVENT< 17.25 41  2 SURVIE (0.048780 0.951219) *
```

Pruning

- A better strategy is to (1) build a very large starting tree \mathcal{T}_0 and (2) cut some branches of the tree.
- We'll try to identify the $\mathcal{T}_{\text{subtree}} \subset \mathcal{T}_0$ that minimizes

$$\sum_{m=1}^{|\mathcal{T}|} \sum_{i: \mathbf{X}_i \in R_m} (Y_i - \hat{Y}_{R_m})^2 + \alpha |\mathcal{T}|,$$

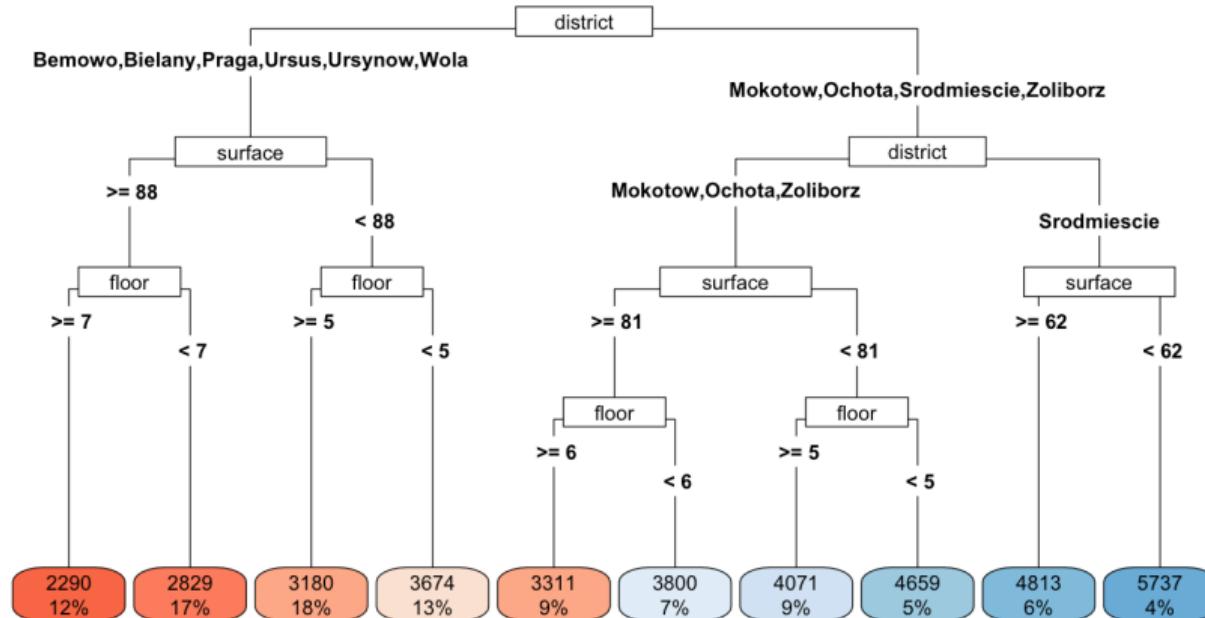
where α is a (hyper-)complexity parameter and $|\mathcal{T}|$ represents the number of leaves in the \mathcal{T} subtree.

- Similar to Ridge/Lasso regression: compromise between bias and variance, between precision and parsimony.
- The complexity parameter can be estimated by cross-validation.

Regression Trees, $y \in \mathbb{R}$

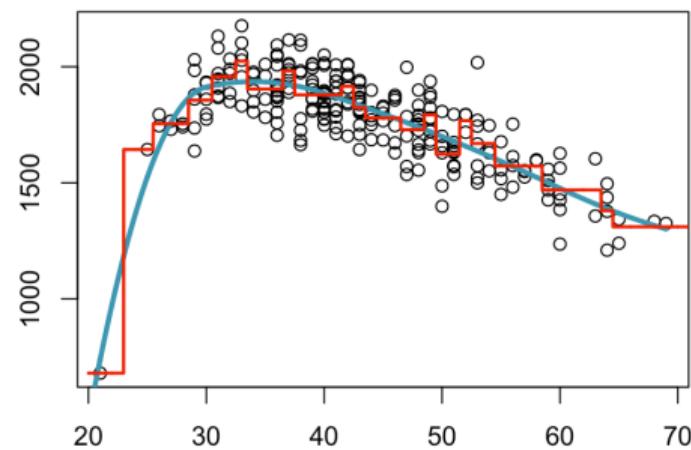
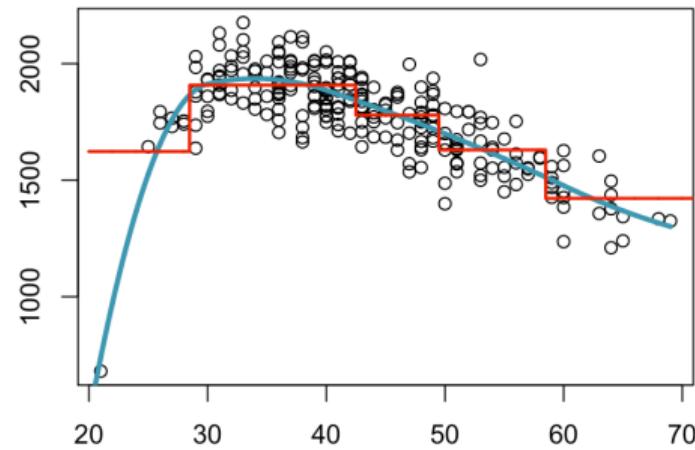
```
1 > library(DALEX)
2 > tree = rpart(m2.price~., data=apartments)
3 > rpart.plot(tree, box.palette="RdBu", type=5)
```

Regression Trees, $y \in \mathbb{R}$



Regression Trees, $y \in \mathbb{R}$

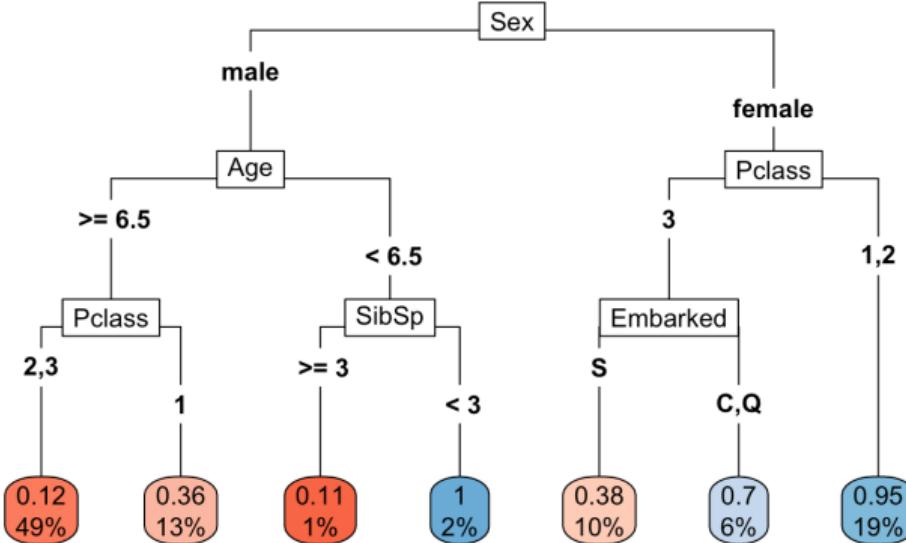
```
1 > tree = rpart(y~x, data=base)
2 > predict(tree, newdata = data.frame(x=x))
```



Classification Trees, $y \in \{0, 1\}$

```
1 > loc = "http://freakonomics.free.fr/titanic.RData"
2 > download.file(loc, "titanic.RData")
3 > load("titanic.RData")
4 > cart = rpart(Survived~, data=base[,1:7])
5 > rpart.plot(cart, box.palette="RdBu", type=5)
```

Classification Trees, $y \in \{0, 1\}$

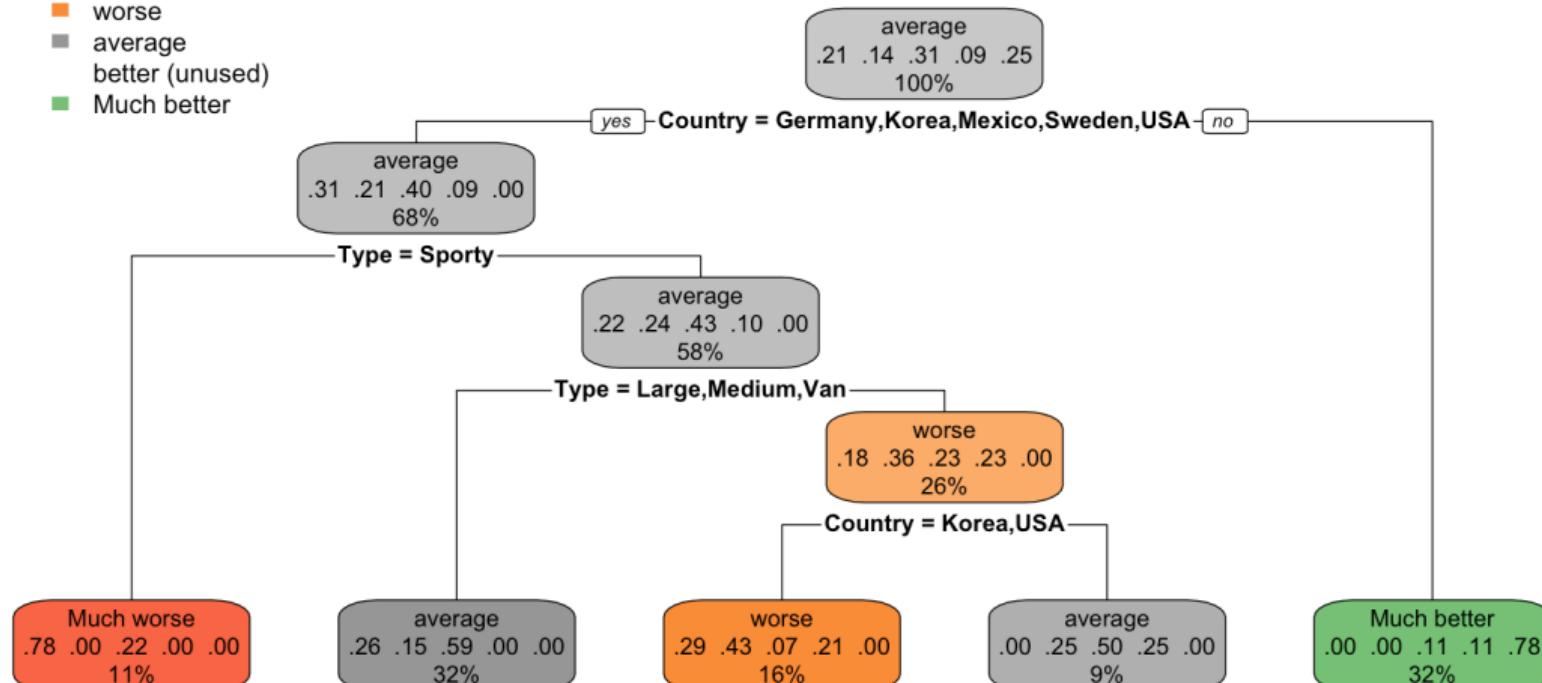


Classification Trees, $y \in \{A, B, C, D\}$

```
1 > multi.class.model <- rpart(Reliability~, data = cu.summary)
2 > rpart.plot(multi.class.model)
```

Classification Trees, $y \in \{A, B, C, D\}$

- Much worse
- worse
- average
- better (unused)
- Much better



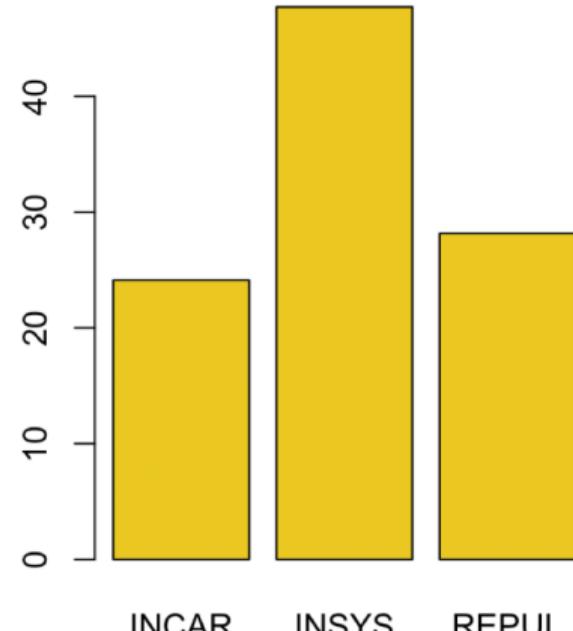
Tree Stability

```
1 > variable = rep(NA,10000)
2 > for(i in 1:10000){
3 + tree = rpart(PRONO~., data=
4 + myocarde[sample(1:71,size=47),])
5 + variable[i] = as.character(
6 + tree$frame[1,"var"])
7 + }
8 > table(variable)/100
9 INCAR INSYS REPUL
10 24.12 47.72 28.16
```

Consider

The variable used in the first node is

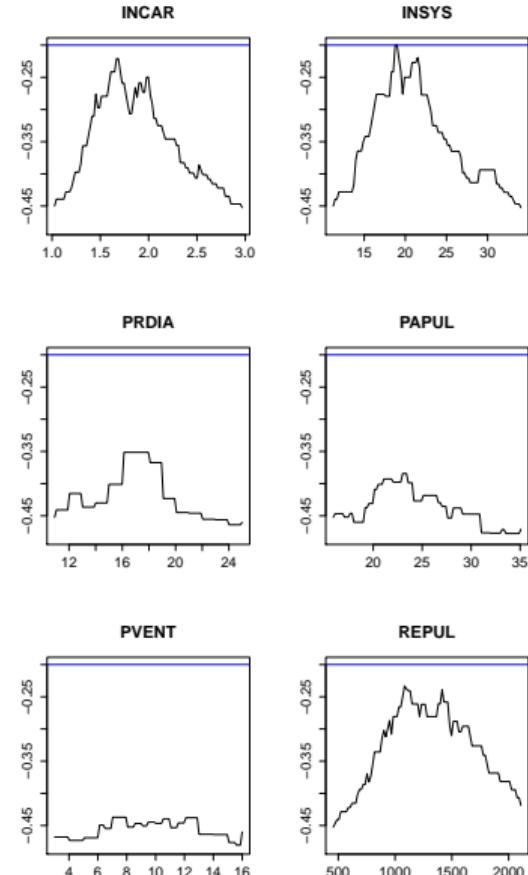
- INSYS (47.7%)
- REPUL (28.2%)
- INCAR (24.1%)



Variable Importance

We had already intuited the importance of these three variables

```
1 > cart = rpart(PRONO~., myocarde)
2 > split = summary(cart)$splits
3 > split
4       count ncat      improve
5 INCAR     71    -1  0.58621312
6 REPUL     71     1  0.55440034
7 PRDIA     71     1  0.27284114
8 PAPUL     71     1  0.20466714
```



and we can calculate the same quantities at each node

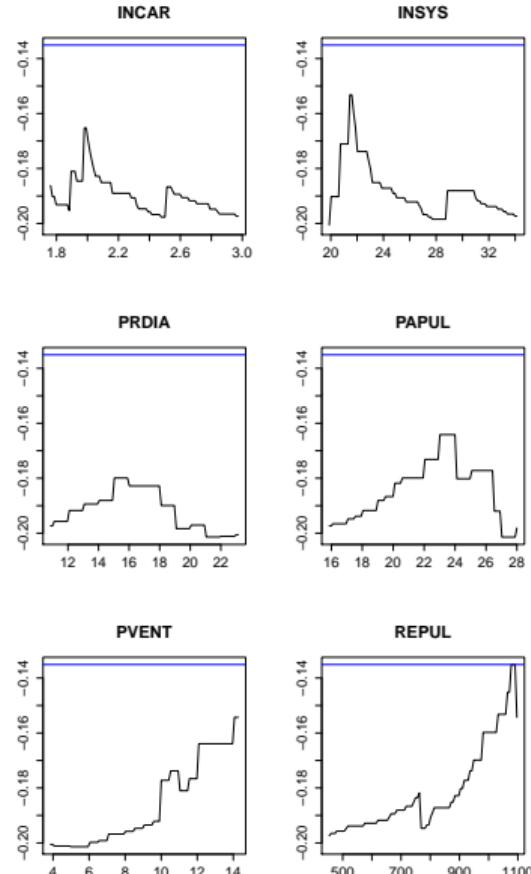
Variable Importance

for the second node

```
1 > split
      count ncat    improve
2 REPUL     27      1  0.18181818
3 PVENT     27     -1  0.10803571
4 PRDIA     27      1  0.10803571
5 PAPUL     27      1  0.10803571
7 INCAR     27      1  0.04705882
```

etc

```
1 > cart$variable.importance
2   INSYS  REPUL  INCAR  PAPUL  PRDIA
      FRCAR  PVENT
3 20.113 19.132 15.895  5.959  5.214
      3.725  0.498
```



Decision Trees

... it's good!

- No assumptions necessary.

- Intuitive visualization.

- Natural integration of high-degree interaction.

... it's bad!

- Procedure quick to overlearn.

- Procedure unstable.

Countering overfit through regularization

A standard technique is to prune trees by cutting off the least “profitable” branches (in terms of impurity)

Consider a decision tree, penalizing too many regions m .

$$\text{Regression tree } \sum_{r=1}^m \sum_{x_i \in R_r} (y_i - \hat{y}_{R_r})^2 + \lambda m$$

i.e., with “error $+\lambda m$ ”.

Countering overfit through regularization

With λ very large, we prune everything, return to a single region, and find $\hat{y} = \bar{y}$. Here, no overfit.

If λ is very small, we find the standard decision tree solution. Without stopping rules, we can form regions with a single observation such that the prediction error (on S_{train}) is 0.

Countering overfit through regularization

Although it may seem difficult to optimize the regularized objective function, there's a little trick that lets us easily re-create the equivalent of the *lasso path*.

In other words, we'll get a nested tree sequence, from a very large tree with many branches to a trunk tree with a single region.

Starting from a very large tree (say, a very small one). We'll progressively prune the branches, always choosing the one that had reduced branch diversity the least.

Countering overfit through regularization

Another perspective is to group the two regions that increase the total error the least when recombined.

Ainif, we can create a sequence of nested trees.

It could be demonstrated that this is equivalent to optimizing :

$$\sum_{r=1}^m \sum_{x_i \in R_r} (y_i - \hat{y}_{R_r})^2 + \lambda m \text{ on a value sequence of } \lambda.$$

Countering overfit through regularization

We need to choose between these two extremes, which we'll do using validation.

An interesting approach:

- Train a decision tree with many branches (β very small)
- We generate the intermediate tree sequence by pruning the branches, starting with those that reduce heterogeneity the least.
- Evaluate the performance of all trees on S_{val} .
- We choose the tree with the best performance on S_{val} .

Countering overfit through regularization

Algorithm 8.1 Building a Regression Tree

1. Use recursive binary splitting to grow a large tree on the training data, stopping only when each terminal node has fewer than some minimum number of observations.
 2. Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of α .
 3. Use K-fold cross-validation to choose α . That is, divide the training observations into K folds. For each $k = 1, \dots, K$:
 - (a) Repeat Steps 1 and 2 on all but the k th fold of the training data.
 - (b) Evaluate the mean squared prediction error on the data in the left-out k th fold, as a function of α .Average the results for each value of α , and pick α to minimize the average error.
 4. Return the subtree from Step 2 that corresponds to the chosen value of α .
-

Figure 1: Algorithme 8.1. Extrait de James et al. (2013).

Countering overfit through regularization

Although this pruning feature works well, it is rarely used in practice.
It's a slow process.

Above all, decision trees are rarely used outside the random forest and pruning is not used in forest training.

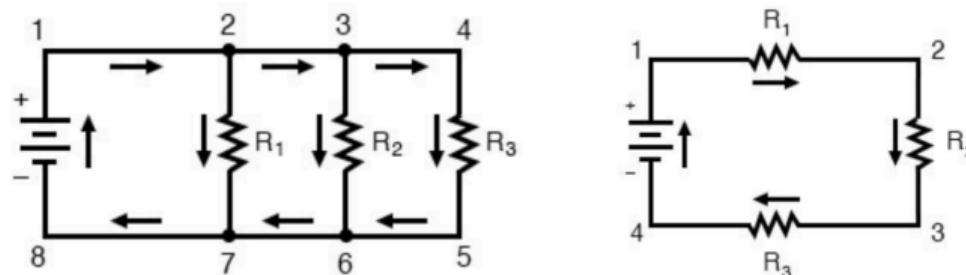
Ensemble learning

Ensemble learning

In statistics and machine learning, ensemble methods use multiple learning algorithms to obtain better predictive performance than could be obtained from any of the constituent learning algorithms alone. \mathbb{W}

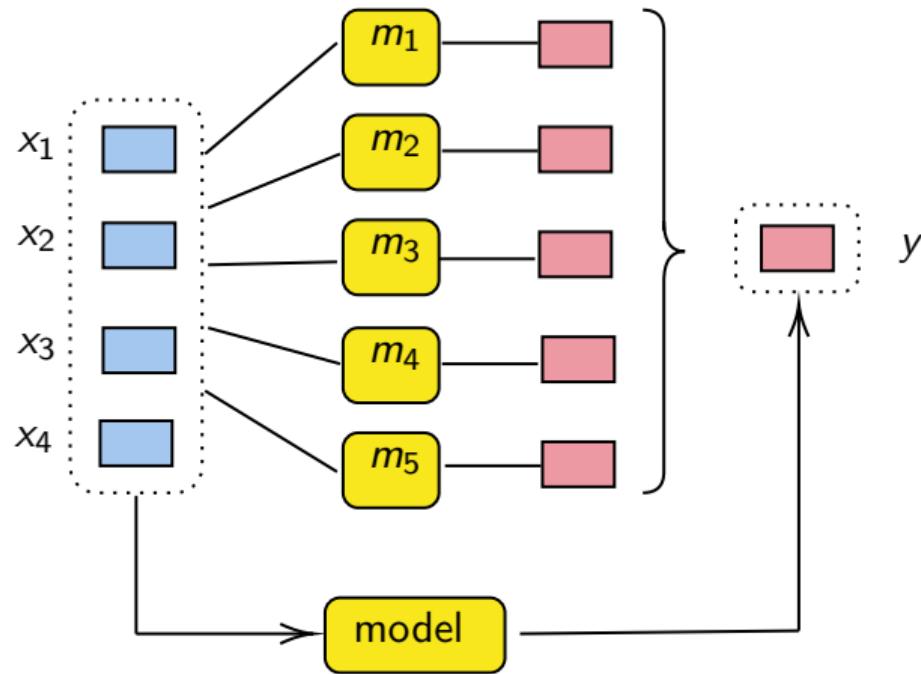
See [Mienye and Sun \(2022\)](#) for a recent survey

Ensemble learning

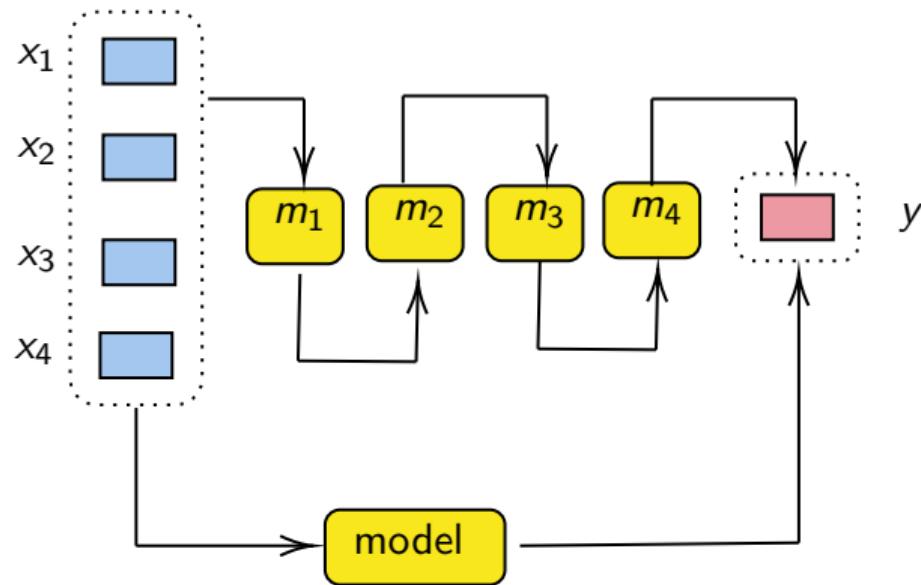


- in **parallel ensemble method**, the base learners are generated in parallel.
when **bagging**, several estimators or models are built independently and then average their predictions,
when **stacking**, we multiple classifications or regression models via a meta-classifier or a meta-regressor, trained on the same training set,
- in **sequential ensemble method**, the base learners are generated sequentially.
The overall performance can be boosted by weighing previously mislabeled examples with higher weight. This corresponds to **boosting**.

Ensemble learning (parallel learning)



Ensemble learning (series-sequential learning)



Ensemble Method

Galton (1907) guess the weight of the cow, county fair in Cornwall, England, 1906

Correct answer = 1198 lbs

787 participants, x_1, \dots, x_n

- quantile 25%: 1162 lbs
- quantile 50%: 1207 lbs
- quantile 75%: 1236 lbs
- average: 1197 lbs

Pick a single prediction x_j v.s average \bar{x} ,

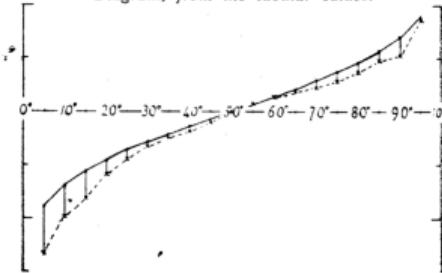
$$\mathbb{E}[(x_j - t)^2] = (\bar{x} - t)^2 + \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

where t is the truth
(ambiguity decomposition)

| Degrees of the length of Array 0-100 | Estimates in lbs. | * Centiles | | |
|--------------------------------------------|----------------------|----------------------------------------|---------------------|--------------------------------------|
| | | Observed deviates from 1207 lbs. | Normal p.e. = 37 | Excess of Observed over Normal |
| | | | - .90 | + .43 |
| 5 | 1074 | - 133 | - .90 | + .43 |
| 10 | 1109 | - 98 | - .70 | + .28 |
| 15 | 1126 | - 81 | - .57 | + .24 |
| 20 | 1148 | - 59 | - .46 | + .13 |
| $q_1 25$ | 1162 | - 45 | - .37 | + .8 |
| 30 | 1174 | - 33 | - .29 | + .4 |
| 35 | 1181 | - 26 | - .21 | + .5 |
| 40 | 1188 | - 19 | - .14 | + .5 |
| 45 | 1197 | - 10 | - .7 | + .3 |
| $m 50$ | 1207 | 0 | 0 | 0 |
| 55 | 1214 | + 7 | + .7 | 0 |
| 60 | 1219 | + 12 | + .14 | - .2 |
| 65 | 1225 | + 18 | + .21 | - .3 |
| 70 | 1230 | + 23 | + .29 | - .6 |
| $q_3 75$ | 1236 | + 29 | + .37 | - .8 |
| 80 | 1243 | + 36 | + .46 | - .10 |
| 85 | 1254 | + 47 | + .57 | - .10 |
| 90 | 1267 | + 52 | + .70 | - .18 |
| 95 | 1293 | + 86 | + .90 | - .4 |

q_1, q_3 , the first and third quartiles, stand at 25° and 75° respectively.
 m , the median or middlemost value, stands at 50°.
The dressed weight proved to be 1198 lbs.

Diagram, from the tabular values.



The continuous line is the normal curve with p.e. = 37.

The broken line is drawn from the observations.

The lines connecting them, show the differences between the observed

and the normal curve.

Ensemble Method

- **statistical interpretation:** hopefully the ensemble generalises better than a single chosen model
- **computational interpretation:** averaging can sometimes be a fast way of reaching closer to the optimal than direct optimisation
- **representational interpretation:** Averaging models of some model class can sometimes take you outside of that model class

Regression problem: averaging

Classification problem: voting (majority)

Votes

?

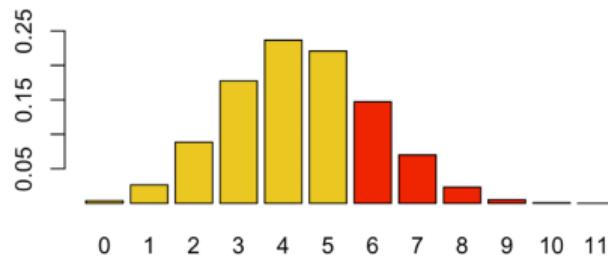
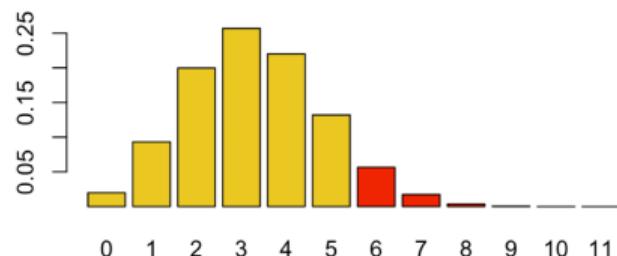
Condorcet (1785) in “Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix”

Let p denote the probability to be wrong, individually,

With n independent vote, the probability for the majority to be wrong is

$$\sum_{k \geq [(n+1)/2]} \binom{n}{k} p^k (1-p)^{n-k}$$

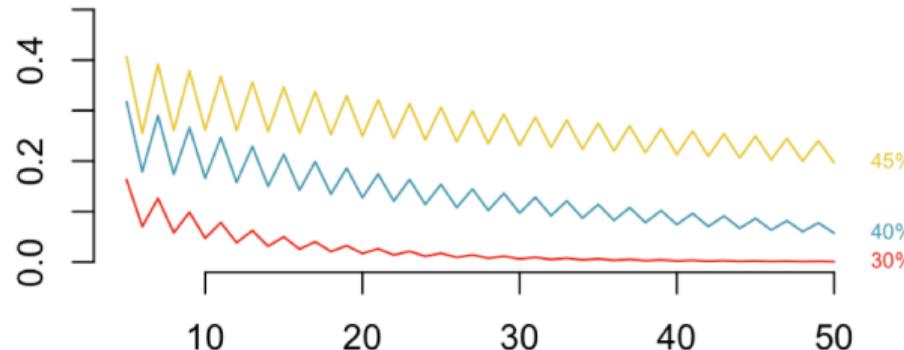
e.g. $n = 11$, $p = 30\%$ (left) and $p = 40\%$ (right)



Probability that the majority is wrong, 7.8% (left) and 24.6% (right).

Votes

Evolution of the probability that the majority is wrong, as a function of n , for various p



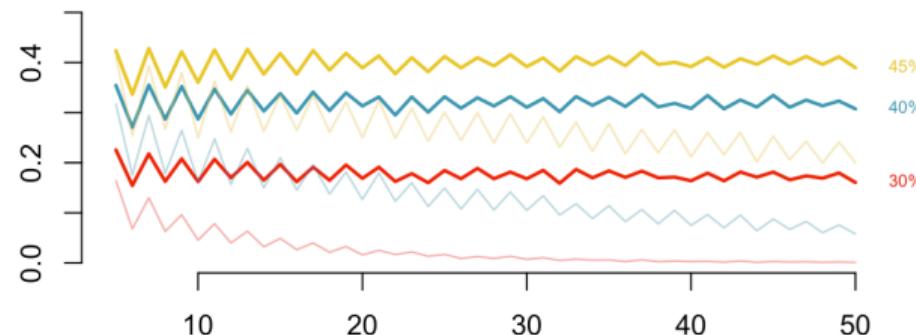
But **this is valid only if votes are independent** ! In ensemble learning, we want

- models to be reasonably* accurate
- models to be as independent as possible

* ensemble learning can turn a collection of poor learners (stumps (single-node decision trees) or linear) into a good one

Independence ?

with correlation among prediction, there is no real 'decrease'



- split training sample into m groups of n/m observations
pb, groups are small, poor predictions
- split features into m groups
pb, poor predictions if only small number of (true) predictors
- use bootstrap
pb if models are not independent enough

Overview

- Parallel ensemble methods (bagging & stacking)

- Bagging

In this method, models are generated using the same algorithm with random sub-samples of the dataset with the bootstrap sampling method to reduce the variance. In bootstrap sampling, some original examples appear more than once and some original examples are not present in the sample.

- Random Forest

In a random forest at each sample, a decision tree is used which collectively form a forest and hence, Random Forest. Rest is almost similar to a simple decision tree.

- Stacking

Stacking is an ensemble learning technique that combines multiple classifications or regression models via a meta-classifier or a meta-regressor. The base-level models are trained based on a complete training set, then the meta-model is trained on the outputs of the base level model as features.

Overview

- Sequential ensemble methods (Boosting Algorithm)
 - Gradient Boosted Decision Trees (GBDT)
uses Gradient Descent Algorithm on Decision trees to reduce loss at each level. It generalizes models by allowing optimization of an arbitrary differentiable loss function. The next model is built on the errors of the previous model
 - XG Boost
XG Boost works similar to GBDT, but it has some more features(below) which makes it more efficient. i) Regularization: Standard GBM implementation has no regularization like XGBoost, therefore it also helps to reduce overfitting. In fact, XGBoost is also known as the 'regularized boosting' technique. ii) Parallel Processing: XGBoost implements parallel processing at root nodes and becomes faster as compared to GBM. iii) High Flexibility: XGB allows for custom evaluation criteria. iv) Handling Missing Values: Can handle missing values on its own. v) Tree Pruning: XGB does

Overview

tree pruning on its own and reduces variance. vi) Built-in Cross-Validation: XGBoost allows the user to run cross-validation at each iteration of the boosting process. vii) Continue on Existing Model: User can start training an XGBoost model from its last iteration of the previous run.

- Ada Boost: each weak learner is developed as decision stumps (A stump is a tree with just a single split and two terminal nodes) that are used to classify the observations. After each classifier is trained, the classifier's weight is calculated based on its accuracy. More weight is assigned to the classifier when accuracy is more and vice-versa.
- Light GBM
It splits the tree leaf-wise with the best fit whereas other boosting algorithms split the tree depth-wise or level-wise rather than leaf-wise and tries to get a pure leaf as soon as possible. This makes Light GBM faster than other boosting algorithms and hence the word 'Light'.

Overview

- Cat Boost “CatBoost” name comes from two words “Category” and “Boosting”. CatBoost converts categorical values into numbers using various statistics on combinations of categorical features and combinations of categorical and numerical features. Hence we need not convert categorical variables to numerical like in other models.

Bayesian model averaging

Consider m models, suppose that one is the true one, but we don't know which one...
Let T denote the choice of the model, $T \in \{1, 2, \dots, m\}$,

$$\mathbb{P}[Y = 1 | \mathbf{X} = \mathbf{x}] = \sum_{t=1}^m \mathbb{P}[Y = 1, T = t | \mathbf{X} = \mathbf{x}]$$

$$\mathbb{P}[Y = 1 | \mathbf{X} = \mathbf{x}] = \sum_{t=1}^m \underbrace{\mathbb{P}[Y = 1 | \mathbf{X} = \mathbf{x}, T = t]}_{\hat{y}^{(t)}} \underbrace{\mathbb{P}[T = t | \mathbf{X} = \mathbf{x}]}_{\omega_t}$$

which is a weighted average of predictions (or posterior class prediction), where weights

$$\omega_t = \mathbb{P}[T = t | \mathbf{X} = \mathbf{x}] \propto \mathbb{P}[\mathbf{X} = \mathbf{x} | T = t]$$

are likelihoods of models.

Stacking

Estimating model likelihood can be complicated
we can learn weights using a regression
(where individual model outputs are treated as features)

```
1 url = "http://freakonometrics.free.fr/german_credit.csv"
2 credit = read.csv(url, header = TRUE, sep = ",")
3 F = c(1,2,4,5,7,8,9,10,11,12,13,15,16,17,18,19,20)
4 for(i in F) credit[,i]=as.factor(credit[,i])
```

create a training and a testing datasets

```
1 set.seed(123)
2 i_test = sample(1:nrow(credit),size=333)
3 i_train = (1:nrow(credit))[-i_test]
```

Credit scoring with ensemble techniques

- model 1 : logistic regression

```
1 GLM = glm(Creditability~., data=credit[i_train,], family= binomial)
2 pGLM = predict(GLM, newdata=credit[i_test,], type="response")
3 classGLM = c(0,1)[1+(pGLM>.5)*1]
```

- model 2 : classification tree

```
1 library(rpart)
2 CART = rpart(Creditability~., data=credit[i_train,])
3 classCART = predict(CART, newdata=credit[i_test,], type="class")
```

- model 3 : SVM (with a Gaussian kernel)

```
1 library(kernlab)
2 SVM = ksvm(Creditability~., data=credit[i_train,], kernel = "rbfdot", C
   =1)
3 classSVM =predict(SVM, newdata=credit[i_test,])
```

Credit scoring with ensemble techniques

On the correlation of those 3 models,

```
1 > df = data.frame(as.numeric(classCART)-1,
2 +                     as.numeric(classGLM),
3 +                     as.numeric(classSVM)-1)
4 > names(df)=c("CART","GLM","SVM")
5 > cor(df)
6
7      CART      GLM      SVM
8 CART 1.0000000 0.3777031 0.4888801
9 GLM  0.3777031 1.0000000 0.5528032
9 SVM  0.4888801 0.5528032 1.0000000
```

Consider a simple (majority based) voting procedure

```
1 > vote = ifelse(rowSums(df)>1.5,1,0)
```

Use AUC a classification metrics

```
1 > library(caret)
2 > auc = function(y) caret::confusionMatrix(data=as.factor(y),reference=
   credit[i_test,"Creditability"])$overall["Accuracy"]
3 > y = credit[i_test,"Creditability"]
```

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities)

BY-NC 4.0 253 / 385

Credit scoring with ensemble techniques

On the classification tree,

```
1 > table(df$CART,y)
2      0   1
3  0  42  46
4  1  52 193
5 > auc(df$CART)
6 Accuracy
7 0.7057057
```

on the GLM

```
1 > table(df$GLM,y)
2      0   1
3  0  43  43
4  1  51 196
5 > auc(df$GLM)
6 Accuracy
7 0.7177177
```

On the SVM

```
1 > table(df$SVM,y)
2      0   1
3  0  33  26
4  1  61 213
5 > auc(df$SVM)
6 Accuracy
7 0.7387387
```

on the majority votes

```
1 > table(vote,credit[i_test
2 ,y])
3 vote   0   1
4 0  37  29
5 1  57 210
6 > auc(vote)
7 Accuracy
8 0.7417417
```

Credit scoring with ensemble techniques

```
1 > df$CREDIT = as.numeric(credit[i_test,"Creditability"])-1
```

```
1 > reglm = lm(CREDIT~CART+GLM+SVM,data=df)
2 > df$STACKLM = predict(reglm)
3 > auc((df$STACKLM>.5)*1)
4 Accuracy
5 0.7627628
```

```
1 > reg = glm(CREDIT~CART+GLM+SVM,data=df,family=binomial)
2 > df$STACKGLM = predict(reg, type="response")
3 > auc((df$STACKGLM>.5)*1)
4 Accuracy
5 0.7627628
```

idea of slacking

Credit scoring with ensemble techniques

Instead of a (majority based) voting procedure, why not consider predictive probabilities?

```
1 > pGLM = pGLM
2 > pCART = predict(CART, newdata=credit[i_test,], type="prob")[,2]
3 > SVM = ksvm(Creditability~., data=credit[i_train,], kernel = "rbfdot",C
   =1, prob.model = TRUE)
4 > pSVM = predict(SVM, newdata=credit[i_test,],type="probabilities")[,2]
5 > pdf = data.frame(as.numeric(pCART),
6 +                     as.numeric(pGLM),
7 +                     as.numeric(pSVM))
8 > names(pdf)=c("CART","GLM","SVM")
9 > cor(pdf)
10          CART      GLM      SVM
11 CART  1.0000000  0.516073  0.6101085
12 GLM   0.5160730  1.000000  0.8641040
13 SVM   0.6101085  0.864104  1.0000000
```

Credit scoring with ensemble techniques

Instead of a (majority based) voting procedure, why not consider predictive probabilities?

```
1 > pGLM = pGLM
2 > pCART = predict(CART, newdata=credit[i_test,], type="prob")[,2]
3 > SVM = ksvm(Creditability~., data=credit[i_train,], kernel = "rbfdot",C
   =1, prob.model = TRUE)
4 > pSVM = predict(SVM, newdata=credit[i_test,],type="probabilities")[,2]
5 > pdf = data.frame(as.numeric(pCART),
6 +                     as.numeric(pGLM),
7 +                     as.numeric(pSVM))
8 > names(pdf)=c("CART","GLM","SVM")
9 > cor(pdf)
10          CART      GLM      SVM
11 CART  1.0000000  0.516073  0.6101085
12 GLM   0.5160730  1.000000  0.8641040
13 SVM   0.6101085  0.864104  1.0000000
```

Credit scoring with ensemble techniques

One can consider the average (probability) prediction

$$\hat{y}_i = \text{average}(\hat{y}_i^{\text{tree}}, \hat{y}_i^{\text{glm}}, \hat{y}_i^{\text{svm}}) = \frac{1}{3}(\hat{y}_i^{\text{tree}} + \hat{y}_i^{\text{glm}} + \hat{y}_i^{\text{svm}})$$

```
1 > auc((apply(pdf,1,mean)>.5)*1)
2   Accuracy
3 0.7417417
```

or

$$\hat{y}_i = \text{median}(\hat{y}_i^{\text{tree}}, \hat{y}_i^{\text{glm}}, \hat{y}_i^{\text{svm}})$$

```
1 > auc((apply(pdf,1,median)>.5)*1)
2   Accuracy
3 0.7477477
```

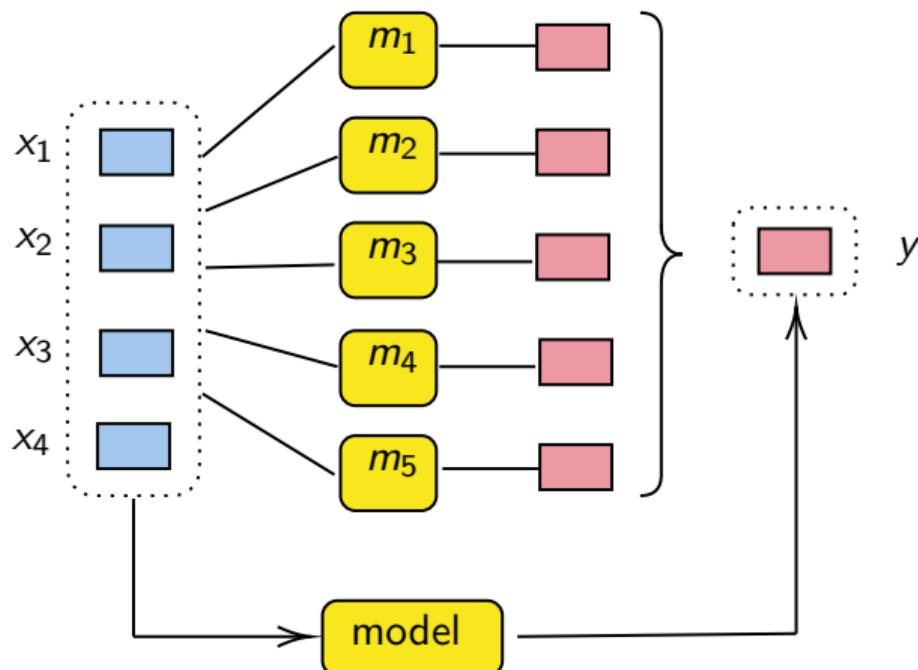
Ensemble techniques: Netflix price



2009: US\$1,000,000 **Netflix Prize**, Winner BellKor's Pragmatic Chaos (ensemble of 107 models)

“Our experience is that most efforts should be concentrated in deriving substantially different approaches, rather than refining a simple technique [...] We strongly believe that the success of an ensemble approach depends on the ability of its various predictors to expose different complementing aspects of the data. Experience shows that this is very different than optimizing the accuracy of each individual predictor”

Bagging (Bootstrap Aggregating)



Bagging (Bootstrap Aggregating)

Bootstrap aggregating, also called bagging (from bootstrap aggregating) or bootstrapping, is a machine learning (ML) ensemble meta-algorithm designed to improve the stability and accuracy of ML classification and regression algorithms. It also reduces variance and overfitting. W

$$\hat{m}_{\text{bag}}(\mathbf{x}_i) = \frac{1}{B} \sum_{j=1}^B \hat{m}_j(\mathbf{x}_i), \quad \text{where models } \hat{m}_j \text{ are trained on bootstrapped sample} \quad (2)$$

$$\hat{m}_{\text{bag}}(\mathbf{x}_i) = \operatorname{argmax}_k \sum_{j=1}^B \hat{m}_j(\mathbf{x}_i) \quad (3)$$

$$\mathbb{E}_D[\hat{m}_{\text{bag}}(\mathbf{x})] = \frac{1}{B} \sum_{j=1}^B \mathbb{E}_D[\hat{m}_j(\mathbf{x})] = \frac{1}{B} B \cdot \mathbb{E}_D[\hat{m}_j(\mathbf{x})] = \mathbb{E}_D[\hat{m}_j(\mathbf{x})]$$

Bagging (Bootstrap Aggregating)

$$\begin{aligned}\text{Var}(\hat{m}_{\text{bag}}(\mathbf{x})) &= \text{Var}\left(\frac{1}{B} \sum_{b=1}^B \hat{m}_b(\mathbf{x})\right) \\ &= \frac{1}{B^2} \text{Var}\left(\sum_{b=1}^B \hat{m}_b(\mathbf{x})\right) \\ &= \frac{1}{B^2} \sum_{b_1=1}^B \sum_{b_2=1}^B \text{Cov}(\hat{m}_{b_1}(\mathbf{x}), \hat{m}_{b_2}(\mathbf{x})) \\ &\leq \frac{1}{B^2} B^2 \text{Var}(\hat{m}_b(\mathbf{x})) = \text{Var}(\hat{m}_b(\mathbf{x}))\end{aligned}$$

because, since $b_1 \neq b_2$, we have $\text{Corr}[\hat{m}_{b_1}(\mathbf{x}), \hat{m}_{b_2}(\mathbf{x})] \leq 1$.

Bagging (Bootstrap Aggregating)

Actually, if $\text{Var}(\hat{m}_B(\mathbf{x})) = \sigma^2(\mathbf{x})$ and $\text{Corr}[\hat{m}_{b_1}(\mathbf{x}), \hat{m}_{b_2}(\mathbf{x})] = r(\mathbf{x})$,

$$\text{Var}(\hat{m}_{\text{bag}}(\mathbf{x})) = r(\mathbf{x})\sigma^2(\mathbf{x}) + \frac{1 - r(\mathbf{x})}{B}\sigma^2(\mathbf{x})$$

The greater the aggregation of different models, the lower the variability.

Tree instability makes them good candidates for aggregation.

Random Forests

Bagging benefits unstable techniques such as decision trees or neural networks. or neural networks.

In the most famous version of random forests, we push the envelope to create even more different trees.

To do this, when creating the trees, we'll randomly sample a number of predictors, say m , from which we'll select the best partitioning. from which we'll select the best partitioning.

Assuming $m = \sqrt{p}$, at each stage of tree formation (each node), we consider only a subset of the predictors.

Consequently, trees are faster to form.

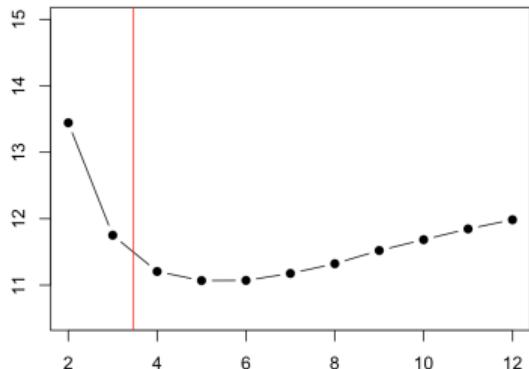
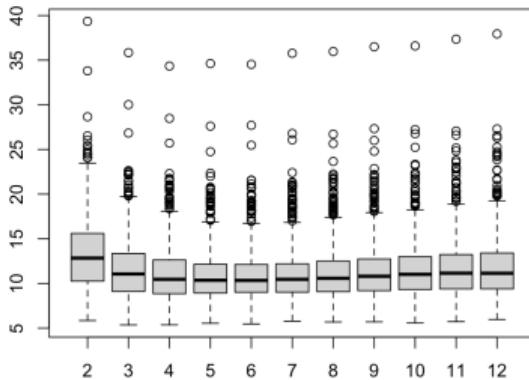
But more importantly, they're even more different, capturing all different aspects of the data.

freakonometrics

freakonometrics.hypotheses.org

Random Forests

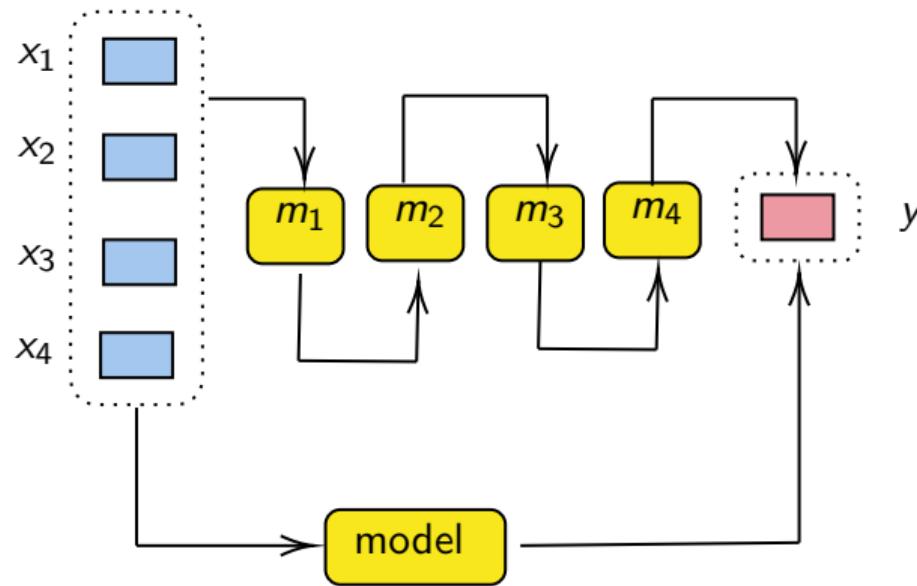
```
1 > library(randomForest)
2 > library(ISLR2)
3 > set.seed(123)
4 > n = nrow(Boston)
5 > sim = function(t){
6 +   train = sample(n, size = n*.7)
7 +   subsim = function(i){
8 +     rf.boston = randomForest(medv ~ ., data
= Boston, subset = train, mtry = i)
9 +     yhat.rf = predict(rf.boston, newdata =
Boston[-train, ])
10 +    mean((yhat.rf - Boston[-train, "medv"])
11 +      ^2)
12 +  }
13 + Vectorize(subsim)(2:12)
14 > Vectorize(sim)(1:499)
```



Boosting

The concept of boosting is based on the question posed by Kearns and Valiant (Kearns (1988); Kearns and Valiant (1989)) "Can a set of weak learners create a single strong learner?" A weak learner is defined as a classifier that is only slightly correlated with the true classification. A strong learner is a classifier that is arbitrarily well-correlated with the true classification. Robert Schapire answered the question in the affirmative in a paper published in 1990. W

Boosting



Boosting

Algorithm 2: Boosting (version 1)

- 1 initialization : k (number of trees), γ , $f_0(\mathbf{x}) = \bar{y}$;
 - 2 **for** $t = 1, 2, \dots, k$ **do**
 - 3 compute $r_{i,t} \leftarrow y_i - f_{t-1}(\mathbf{x}_i)$;
 - 4 fit a model $r_{i,t} \sim h(\mathbf{x}_i)$ for some tree h ;
 - 5 update $f_t(\cdot) = f_{t-1}(\cdot) + \gamma h(\cdot)$
-

Algorithm 3: Boosting (version 2)

- 1 initialization : k (number of trees), $f_0(\mathbf{x}) = \operatorname{argmin}_{\gamma} \sum_{i=1}^n \ell(y_i, \gamma)$;
- 2 **for** $t = 1, 2, \dots, k$ **do**
- 3 compute $r_{i,t} \leftarrow \frac{\partial \ell(y_i, \hat{y})}{\partial \hat{y}} \Big|_{\hat{y}=f_{t-1}(\mathbf{x}_i)}$;
- 4 fit a model $r_{i,t} \sim h(\mathbf{x}_i)$ for some tree h ;
- 5 update $f_t(\cdot) = f_{t-1}(\cdot) + \gamma h(\cdot)$

Algorithme AdaBoost

AdaBoost (short for Adaptive Boosting) is a statistical classification meta-algorithm formulated by Yoav Freund and Robert Schapire in 1995, who won the 2003 Gödel Prize for their work. It can be used in conjunction with many types of learning algorithm to improve performance. The output of multiple weak learners is combined into a weighted sum that represents the final output of the boosted classifier. Usually, AdaBoost is presented for binary classification, although it can be generalized to multiple classes or bounded intervals of real values. AdaBoost is adaptive in the sense that subsequent weak learners (models) are adjusted in favor of instances misclassified by previous models. In some problems, it can be less susceptible to overfitting than other learning algorithms.

W

An alternative is to see the update through weights (more weight if classification is not correct)

Algorithme AdaBoost

| At each iteration of the training process, a weight $w_{i,t}$ is assigned to each sample in the training set equal to the current error on that sample. W

On s'intéresse à un problème de classification binaire (0 ou 1) à partir d'une base de données $\mathcal{D} = \{Y_i; \mathbf{X}_i\}_{i=1,\dots,n}$.

1. Toutes les observations reçoivent un poids identique: $p_1(\mathbf{X}_i) = 1/n, i = 1, \dots, n$.

Algorithme AdaBoost (suite)

2. Pour l'étape $t = 1, \dots, T$,
 - 2a. Piger au hasard (en utilisant les poids $p_t(\mathbf{X}_i)$) une base de données d'entraînement $\mathcal{D}_t = (\mathcal{D}, p_t)$.
 - 2b. Apprendre une règle de classification r_t sur cette base de données.
 - 2c. Calculer l'**erreur apparente**

$$\epsilon_t = p_t(r_t(\mathbf{X}_i) \neq Y_i) = \sum_{i:r_t(\mathbf{X}_i) \neq Y_i} p_t(i)$$

pour \mathcal{D}_t et évaluer $\alpha_t = 0.5 \ln((1 - \epsilon_t)/\epsilon_t)$.

- 2d. Mettre à jour les poids:

$$p_{t+1}(\mathbf{X}_i) = \begin{cases} \frac{p_t(\mathbf{X}_i)}{Z_t} e^{-\alpha_t}, & r_t(\mathbf{X}_i) = Y_i \\ \frac{p_t(\mathbf{X}_i)}{Z_t} e^{+\alpha_t}, & r_t(\mathbf{X}_i) \neq Y_i, \end{cases}$$

où Z_t est une constante de normalisation assurant que $\sum_{i=1}^n p_{t+1}(\mathbf{X}_i) = 1$.

Algorithme AdaBoost (suite)

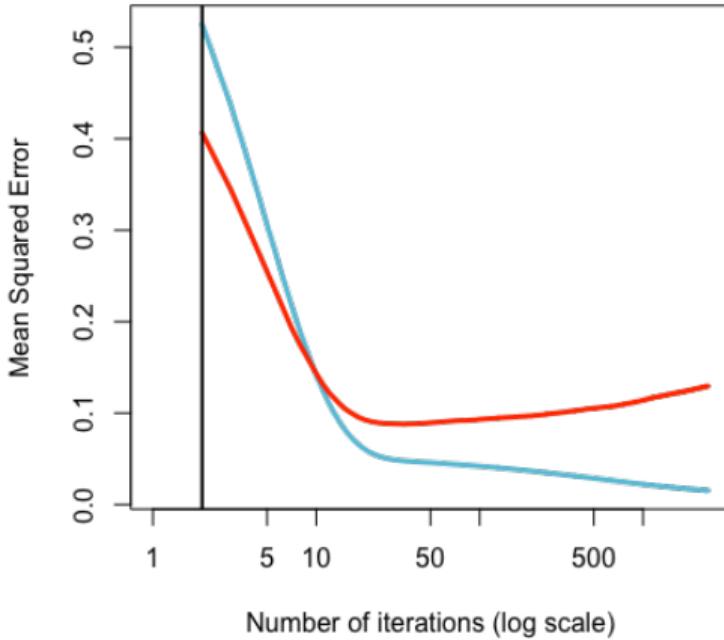
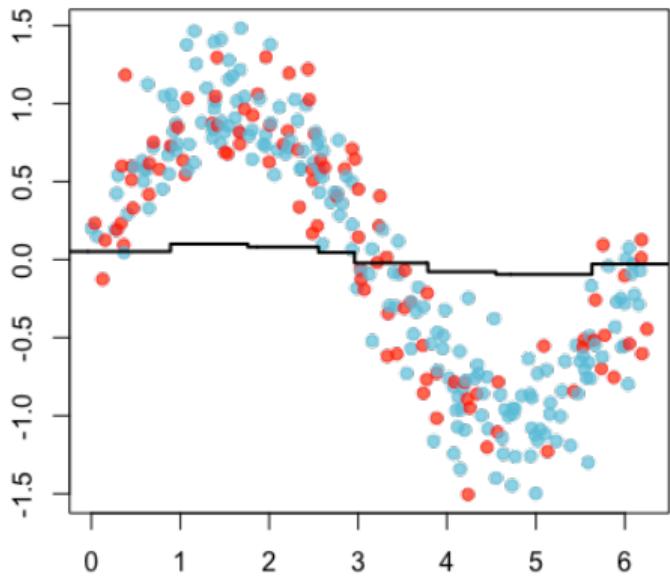
3. La classification finale faite par le modèle est alors

$$R(\mathbf{X}_i) = \begin{cases} 1, & \sum_{t=1}^T \alpha_t r_t(\mathbf{X}_i) > 0 \\ -1, & \sum_{t=1}^T \alpha_t r_t(\mathbf{X}_i) < 0. \end{cases}$$

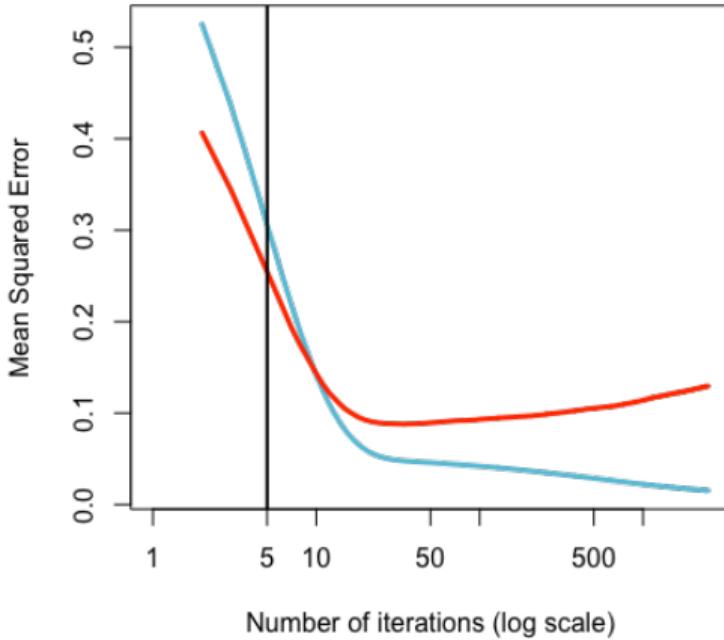
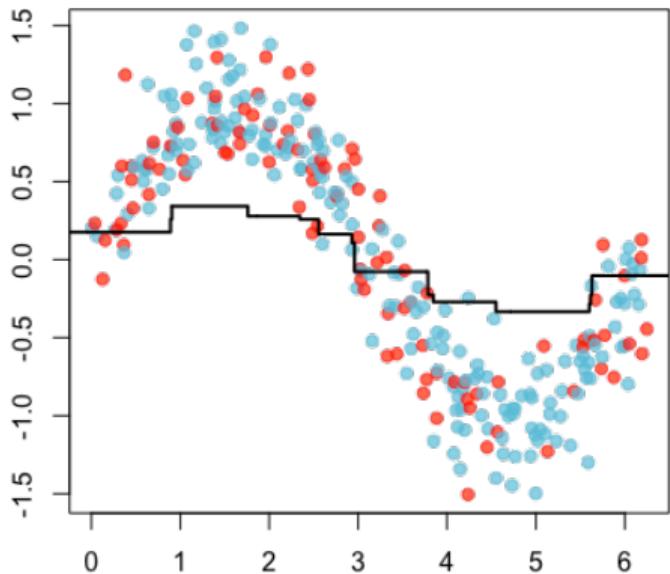
Algorithm 4: Boosting (Adaboost)

- 1 initialization : $k, \omega_i \leftarrow 1/n;$
 - 2 **for** $t = 1, 2, \dots k$ **do**
 - 3 error rate $e_t \leftarrow \frac{\sum_{i=1}^n \omega_i \mathbf{1}(y_i \neq f_{t-1}(\mathbf{x}_i))}{\sum_{i=1}^n \omega_i};$
 - 4 set $\alpha_t \leftarrow \log(1 - e_t) - \log(e_t);$
 - 5 update $\omega_i \leftarrow \omega_i e^{\alpha_t \mathbf{1}(y_i \neq f_{t-1}(\mathbf{x}_i))};$
 - 6 update $f_t(\cdot) \leftarrow f_{t-1}(\cdot) + \alpha_t h(\cdot)$
-

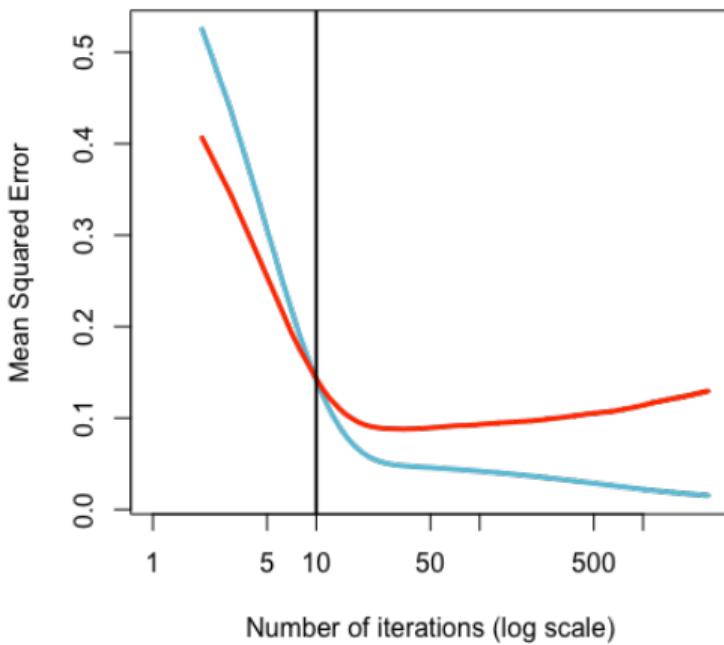
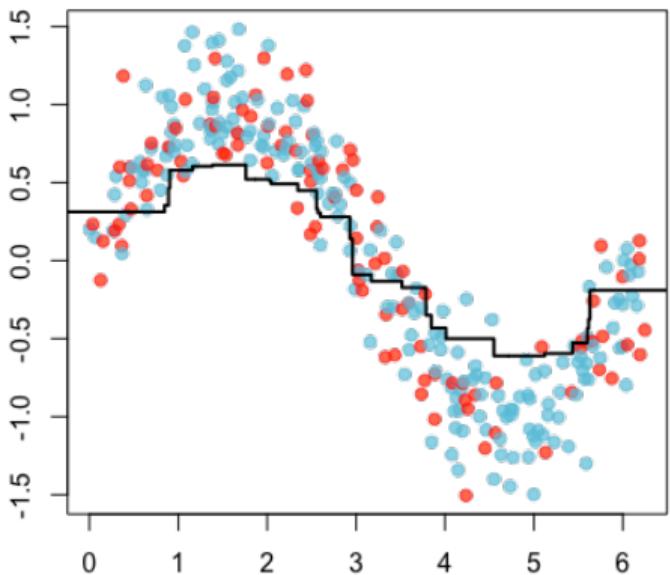
Boosting as Sequential Learning



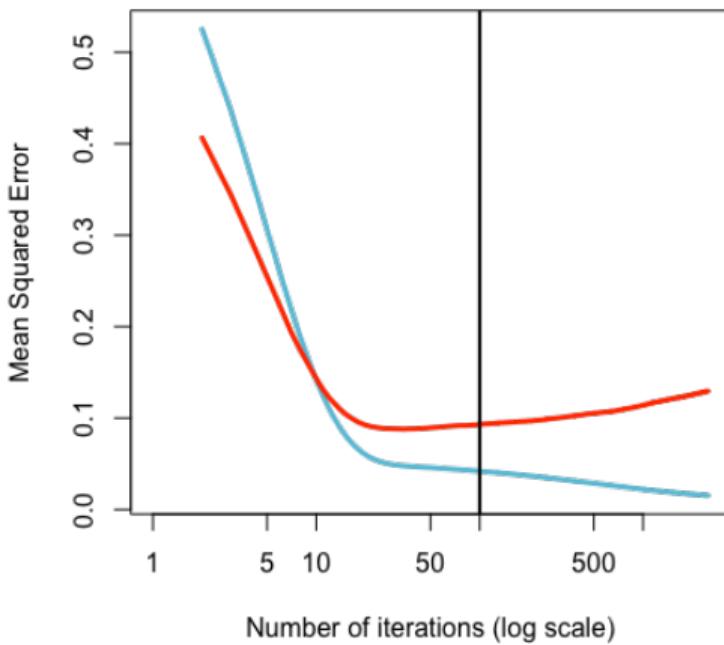
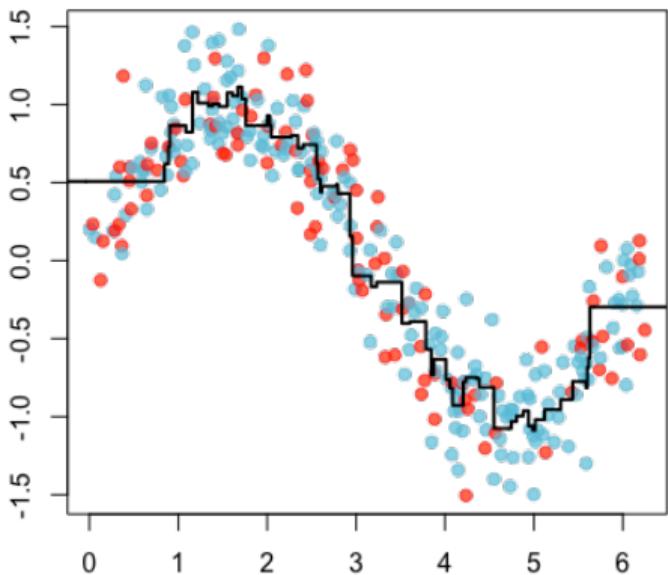
Boosting as Sequential Learning



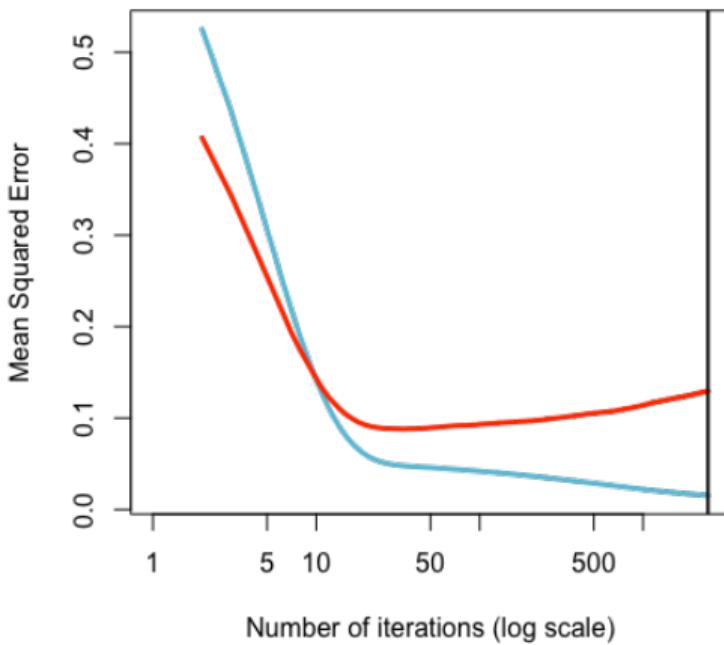
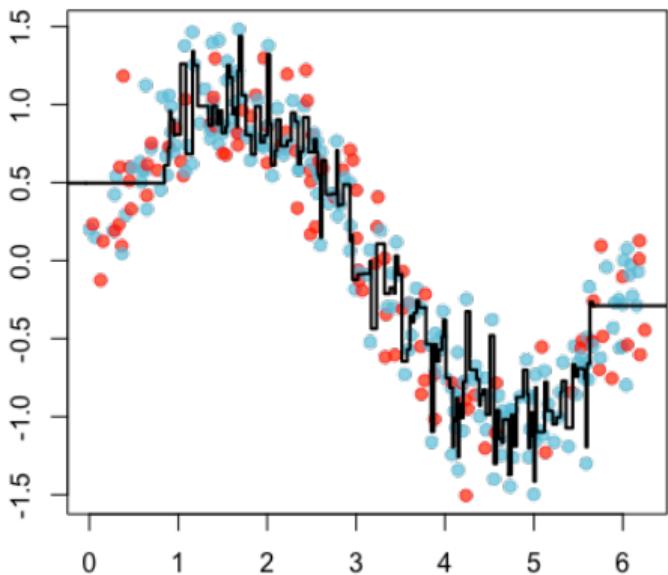
Boosting as Sequential Learning



Boosting as Sequential Learning



Boosting as Sequential Learning



Model Selection

How do we choose an appropriate hypothesis set or learning algorithm for a given problem? For a given binary hypothesis class \mathcal{H} and a function $h \in \mathcal{H}$, we have that

$$\mathcal{R}(h) - \mathcal{R}^* = \underbrace{\mathcal{R}(h) - \inf_{g \in \mathcal{H}} \{\mathcal{R}(g)\}}_{\text{estimation}} + \underbrace{\inf_{g \in \mathcal{H}} \{\mathcal{R}(g)\} - \mathcal{R}^*}_{\text{true risk}}$$

Oracles

An **oracle** is a theoretical black-box that provides access to information not directly observable, used to guide learning algorithms.

- Typically used in computational learning theory.
- Allows querying labels, concepts, or optimal actions.
- Encapsulates external knowledge or supervision.

Definition

A **labeling oracle** \mathcal{O}_L takes an input $x \in \mathcal{X}$ and returns the true label $y \in \mathcal{Y}$:

$$\mathcal{O}_L(x) = f(x)$$

where $f: \mathcal{X} \rightarrow \mathcal{Y}$ is the unknown target function.

- Used in PAC learning to model labeled data access.
- Sample complexity depends on accuracy ϵ and confidence δ .

Oracles

Definition

A **membership oracle** \mathcal{O}_M allows querying whether a hypothesis h agrees with target concept c on input x :

$$\mathcal{O}_M(x) = c(x) \in \{0, 1\}$$

- Supports efficient learning via active queries.
- Key component in the Angluin's model of exact learning.

Definition

An **equivalence oracle** \mathcal{O}_E takes a hypothesis h and returns:

- YES, if $h \equiv c$
- A counterexample x such that $h(x) \neq c(x)$, otherwise
- Used to verify hypothesis correctness.
- Forms basis of learning via hypothesis testing.
- Allows learning in $\log |C|$ steps for finite concept class C .

Model Selection

Definition

Model selection is the process of choosing the best model $\hat{h} \in \mathcal{H}$ from a hypothesis class \mathcal{H} based on observed data.

- Key step in both training and generalization.
- Goal: Minimize the true risk:

$$R(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(h(x), y)]$$

- Cannot access $R(h)$ directly — we use empirical estimates.

Empirical Risk

Given data $S = \{(x_i, y_i)\}_{i=1}^n$, define:

$$\hat{R}_S(h) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

Model Selection

- ERM selects:

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \hat{R}_S(h)$$

- Tends to overfit if \mathcal{H} is too complex.
- Bias-variance tradeoff central to model selection.

Oracle Comparison

Let $h^* = \arg \min_{h \in \mathcal{H}} R(h)$ be the **oracle model**.

A learning algorithm is **oracle-efficient** if it achieves:

$$R(\hat{h}) \leq R(h^*) + \epsilon_n$$

with high probability, where $\epsilon_n \rightarrow 0$ as $n \rightarrow \infty$.

- Bounds deviation from oracle performance.
- Central in statistical learning theory.

Model Selection

- Consider a sequence of model classes:

$$\mathcal{H}_1 \subset \mathcal{H}_2 \subset \cdots \subset \mathcal{H}_k$$

- Tradeoff between empirical risk and complexity penalty:

$$\hat{h}_j = \arg \min_{h \in \mathcal{H}_j} \hat{R}_S(h) + \lambda \cdot \text{pen}(h)$$

- Choose \hat{h} from optimal \mathcal{H}_j :

$$\hat{h} = \arg \min_j \left\{ \hat{R}_S(\hat{h}_j) + \text{pen}(\mathcal{H}_j) \right\}$$

Examples

VC dimension, Rademacher complexity, AIC, BIC.

Model Selection

- Cross-validation estimates generalization error:

$$CV_k = \frac{1}{k} \sum_{i=1}^k \hat{R}_{S^{(i)}}(h^{(i)})$$

- Select model minimizing CV_k estimate.
- Connects to oracle via goal: simulate access to $R(h)$.
- Summary:
 - Model selection balances fit vs. complexity.
 - Theoretical bounds use oracle comparisons.
 - Tools: ERM, SRM, CV, regularization.

Leave-one-Out Cross Validation

In statistics, the jackknife (jackknife cross-validation) is a cross-validation technique and, therefore, a form of resampling. It is especially useful for bias and variance estimation. The jackknife pre-dates other common resampling methods such as the bootstrap. Given a sample of size n , a jackknife estimator can be built by aggregating the parameter estimates from each subsample of size $(n - 1)$ obtained by omitting one observation. W

Leave-one-Out Cross Validation

For a linear regression, $\mathbf{y} = \mathbf{X}\beta + \epsilon$,

$$\hat{\mathbf{y}} = \mathbf{H}\mathbf{y} \text{ where } \mathbf{H} = \mathbf{X}^\top(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}$$

since $\hat{\beta} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y}$.

If we remove observation i ,

$$\hat{\beta}_{-i} = \left((\mathbf{X}\mathbf{X}^\top)^{-1} + \frac{(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{x}_i\mathbf{x}_i^\top(\mathbf{X}\mathbf{X}^\top)^{-1}}{1 - \mathbf{x}_i^\top(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{x}_i} \right) (\mathbf{X}\mathbf{y} - \mathbf{x}_i y_i)$$

and therefore $\hat{y}_{-i} = \hat{m}_{-1}(\mathbf{x}_i) = \mathbf{x}_i^\top \hat{\beta}_{-i}$ becomes

$$\begin{aligned} \hat{y}_{-i} &= \left(\mathbf{x}_i^\top(\mathbf{X}\mathbf{X}^\top)^{-1} + \frac{\mathbf{x}_i^\top(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{x}_i\mathbf{x}_i^\top(\mathbf{X}\mathbf{X}^\top)^{-1}}{1 - \mathbf{x}_i^\top(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{x}_i} \right) (\mathbf{X}\mathbf{y} - \mathbf{x}_i y_i) \\ &= \frac{1}{1 - \mathbf{x}_i^\top(\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{x}_i} \mathbf{x}_i^\top(\mathbf{X}\mathbf{X}^\top)^{-1}(\mathbf{X}\mathbf{y} - \mathbf{x}_i y_i) \end{aligned}$$

Leave-one-Out Cross Validation

$$= \frac{1}{1 - \mathbf{x}_i^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{x}_i} \left(\mathbf{x}_i^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{y} - \mathbf{x}_i^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{x}_i y_i \right),$$

i.e., there is an explicit formula.

Cross Validation

Definition 2.22: K -fold Cross Validation

Partition the index set $\mathcal{I} = \{1, 2, \dots, n\}$ into K roughly equally sized folds $\mathcal{I}_1, \dots, \mathcal{I}_K$. For $k = 1, 2, \dots, K$,

$$\hat{\mu}_{(-k)} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \sum_{i \in \mathcal{I} \setminus \mathcal{I}_k} \ell(y_i, m(\mathbf{x}_i)) \right\}$$

and

$$\hat{R}_{CV} = \frac{1}{K} \sum_{k=1}^K \frac{1}{\operatorname{card}(\mathcal{I}_k)} \sum_{i \in \mathcal{I}_k} \ell(y_i, \hat{\mu}_{(-k)}(\mathbf{x}_i))$$

bootstrap

Wrap up...

Hold-Out Cross Validation

1. Split $\{1, 2, \dots, n\}$ in \mathcal{D}_{ent} (training) and \mathcal{D}_{test} (test)
- 2 . Estimate \hat{m} on sample (y_i, \mathbf{x}_i) , $i \in \mathcal{D}_{ent}$: \hat{m}_{ent}
3. Compute $\frac{1}{|\mathcal{D}_{test}|} \sum_{i \in V} \ell(y_i, \hat{m}_{ent}(\mathbf{x}_i))$

Leave-one-Out Cross Validation (LOOCV)

1. Estimate n models : \hat{m}_{-j} on sample (y_i, \mathbf{x}_i) , $i \in \{1, \dots, n\} \setminus \{j\}$
2. Compute $\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}_{-i}(\mathbf{x}_i))$

Wrap up...

K-Fold Cross Validation

1. Split $\{1, 2, \dots, n\}$ in K groups V_1, \dots, V_K
2. Estimate K models : \hat{m}_k on sample (y_i, \mathbf{x}_i) , $i \in \{1, \dots, n\} \setminus V_k$
3. Compute $\frac{1}{K} \sum_{k=1}^K \frac{1}{|V_k|} \sum_{i \in V_k} \ell(y_i, \hat{m}_k(\mathbf{x}_i))$

Bootstrap Cross Validation

1. Generate B bootstrap samples from $\{1, 2, \dots, n\}$, I_1, \dots, I_B
2. Estimate B models : \hat{m}_b on sample (y_i, \mathbf{x}_i) , $i \in I_b$
3. Compute $\frac{1}{B} \sum_{b=1}^B \frac{1}{n - |I_b|} \sum_{i \notin I_b} \ell(y_i, \hat{m}_b(\mathbf{x}_i))$

freakonometrics

freakonometrics.hypotheses.org

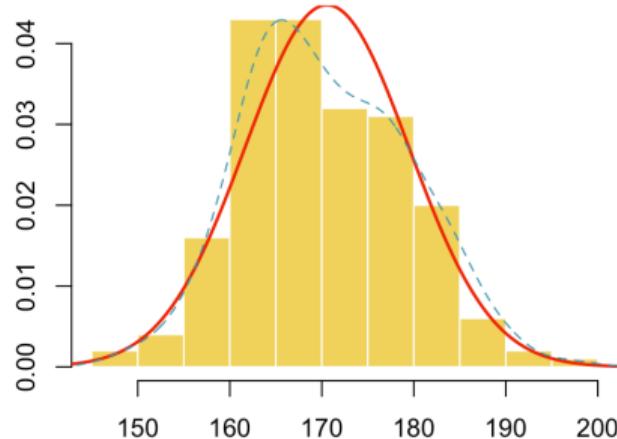
– Arthur Charpentier, April 2025 (Bermuda Financial Authorities)

© BY-NC 4.0 286 / 385

Unobserved Heterogeneity

Height of students 1. Gaussian model, $f(x) = \phi_{\bar{x}, s^2}(x)$

```
1 > Davis = read.table("http://freakonometrics.  
    free.fr/Davis.txt")  
2 > X = Davis$height  
3 > hist(X, proba=TRUE)  
4 > (param = fitdistr(X,"normal")$estimate)  
    mean           sd  
6 170.02000   11.97788  
7 > f1 = function(x) dnorm(x,param[1],param[2])  
8 > x = seq(100,210,by=.2)  
9 > lines(x,f1(x),lwd=2)
```



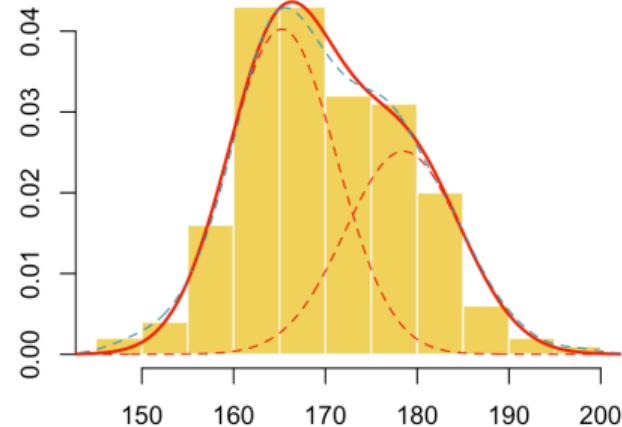
```
1 > logdf = function(x,p){  
2   p1 = p[1]  
3   m1 = p[2]; s1 = p[4]  
4   m2 = p[3]; s2 = p[5]  
5   log(p1*dnorm(x,m1,s1)+(1-p1)*dnorm(x,m2,s2))}
```

Parametric Statistical Models

Height of students 2. Mixture of 2 Gaussians,

$$f(x) = p\phi_{\mu_1, \sigma_1^2}(x) + (1 - p)\phi_{\mu_2, \sigma_2^2}(x)$$

```
1 > logL = function(parameter) -sum(
2   logdf(X,parameter))
3 > Amat = matrix(c
4   (1,-1,0,0,0,0,0,0,0,0,1,0,0,0,
5   0,0,0,0,0,1), 4, 5)
6 > bvec = c(0,-1,0,0)
7 > (param12 = constrOptim(c
8   (.5,160,180,10,10), logL, NULL,
9   ui = Amat, ci = bvec)$par)
10 [1] 0.5996263 165.2690084
11      178.4991624    5.9447675
12      6.3564746
```

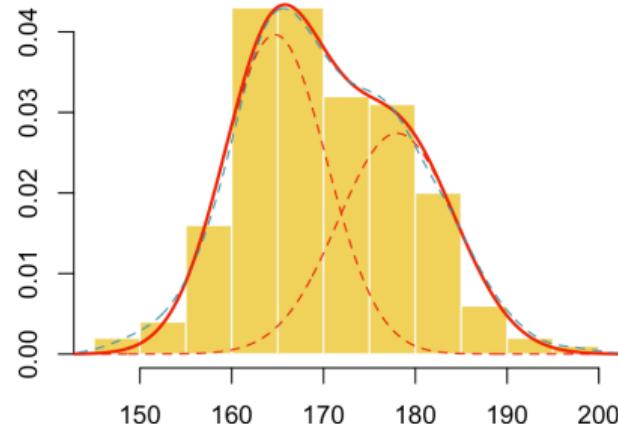


Parametric Statistical Models

Height of students 3. Conditional Gaussian

$$f(x) = \mathbb{P}(F) \cdot \phi_{\bar{x}_F, s_F^2}(x) + \mathbb{P}(M) \cdot \phi_{\bar{x}_M, s_M^2}(x)$$

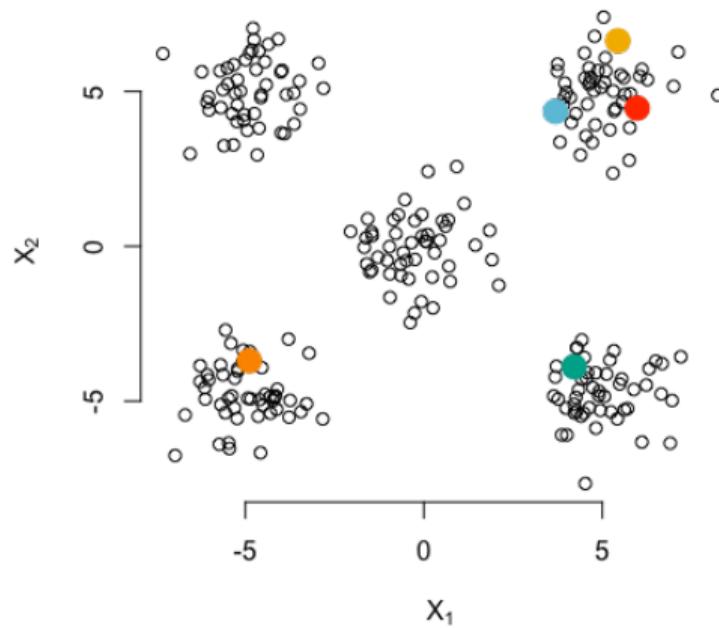
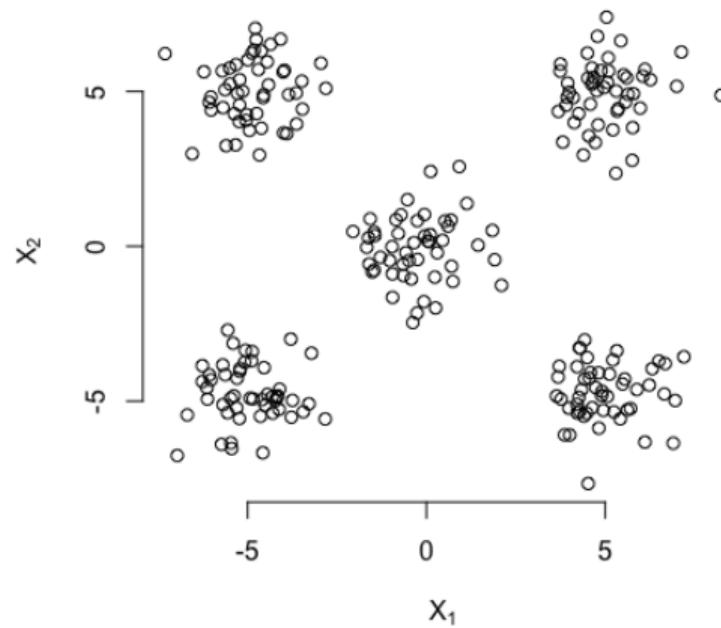
```
1 > (pM = mean(Davis$sex=="F"))
2 [1] 0.56
3 > (paramF = fitdistr(X[Davis$sex=="F",
4     ], "normal")$estimate)
      mean          sd
5 163.74107   11.59183
6 > (paramM = fitdistr(X[Davis$sex=="M",
7     ], "normal")$estimate)
      mean          sd
8 178.011364   6.404001
```



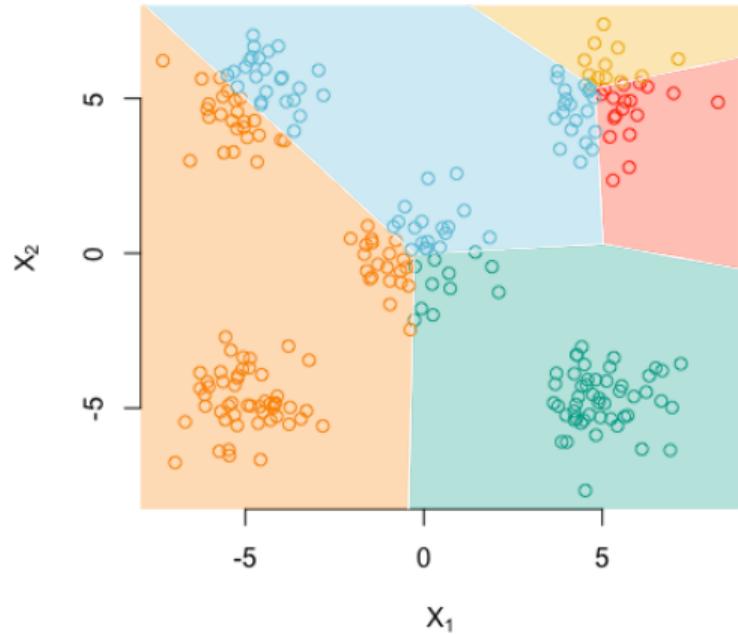
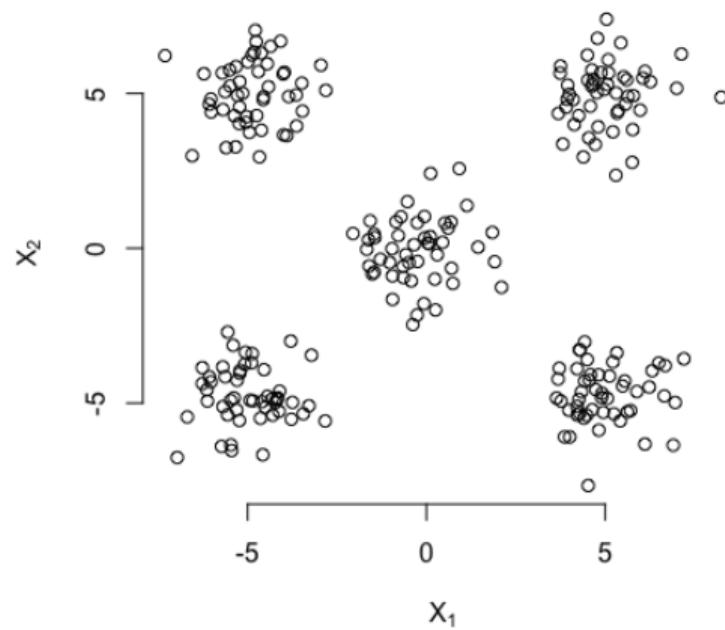
k-means

k-means clustering is a method of vector quantization, originally from signal processing, that aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean (cluster centers or cluster centroid), serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells. \mathbb{W}

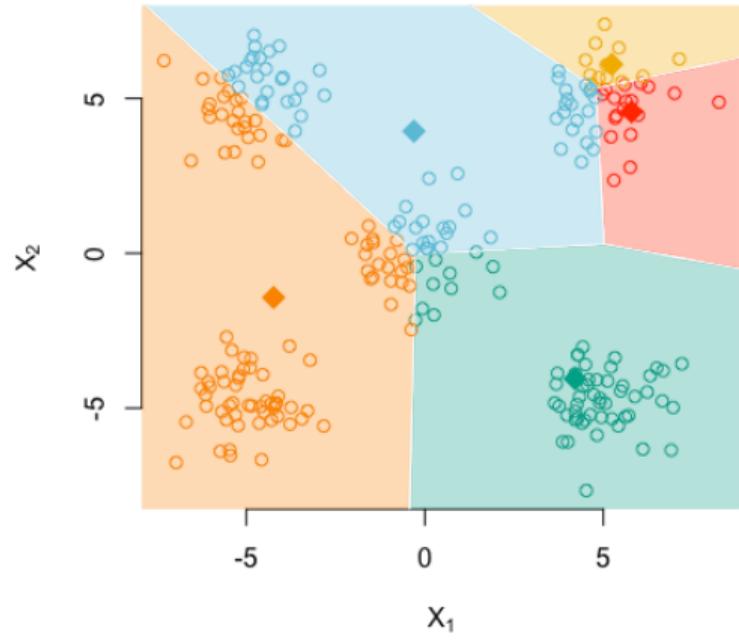
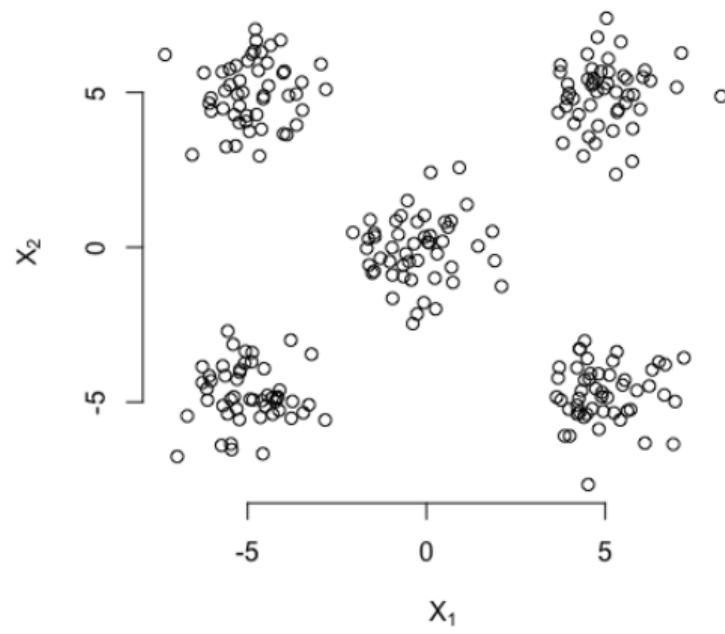
k-means



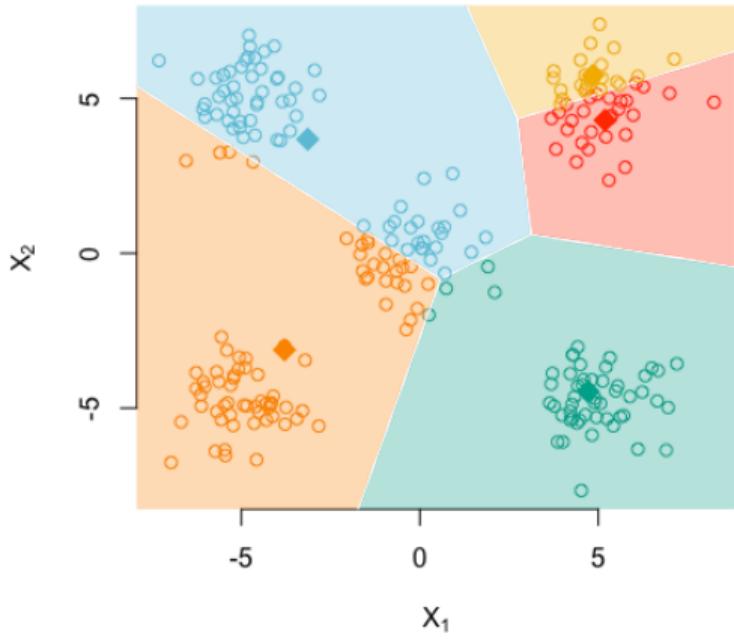
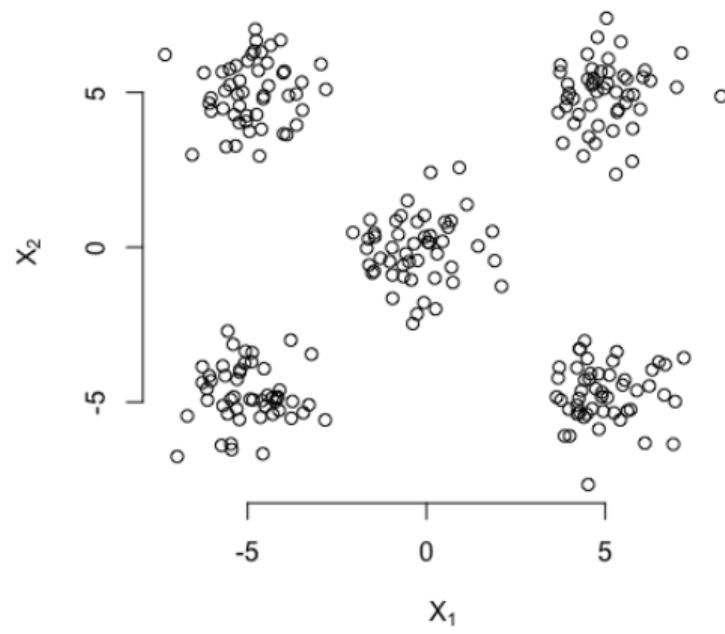
k-means



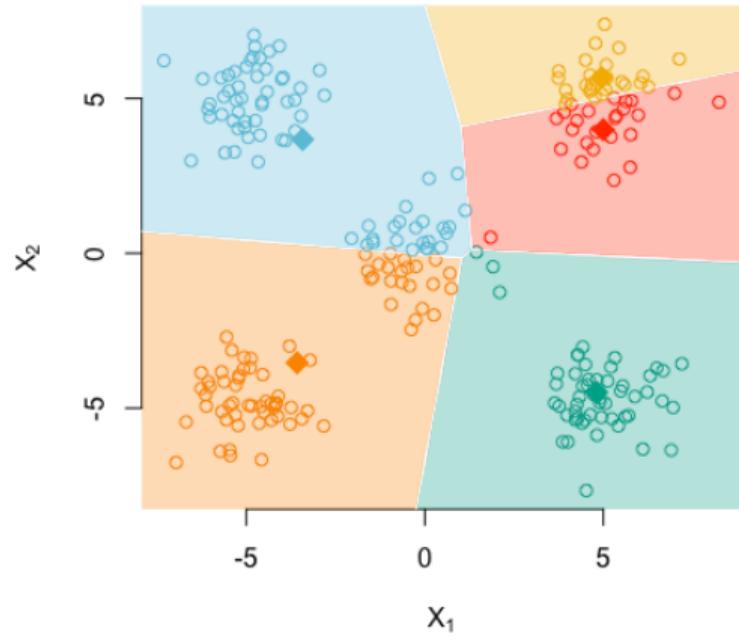
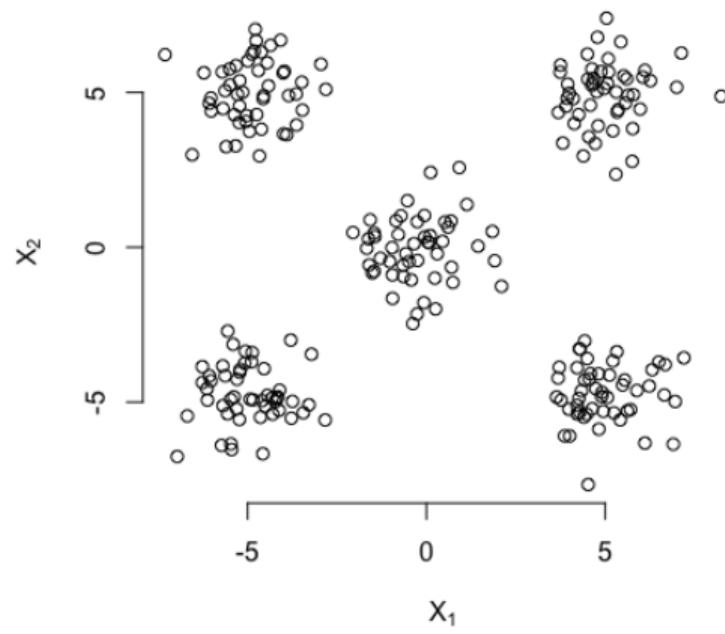
k-means



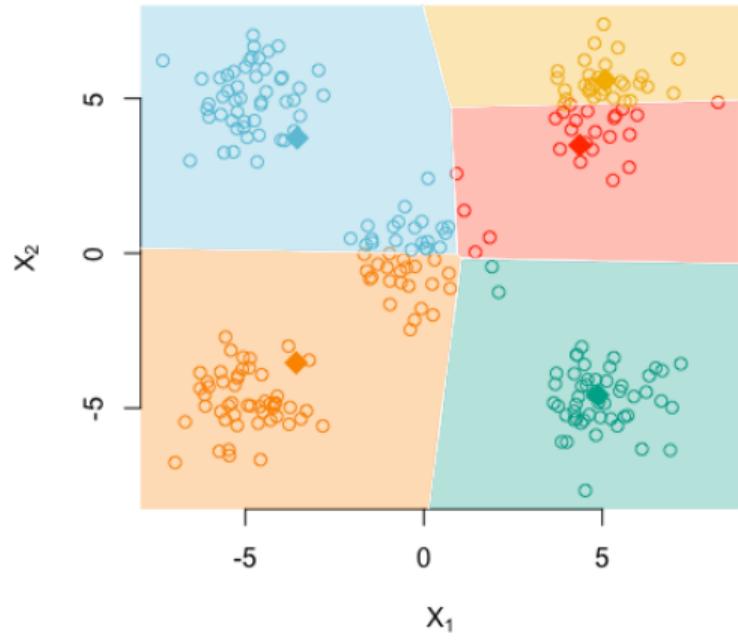
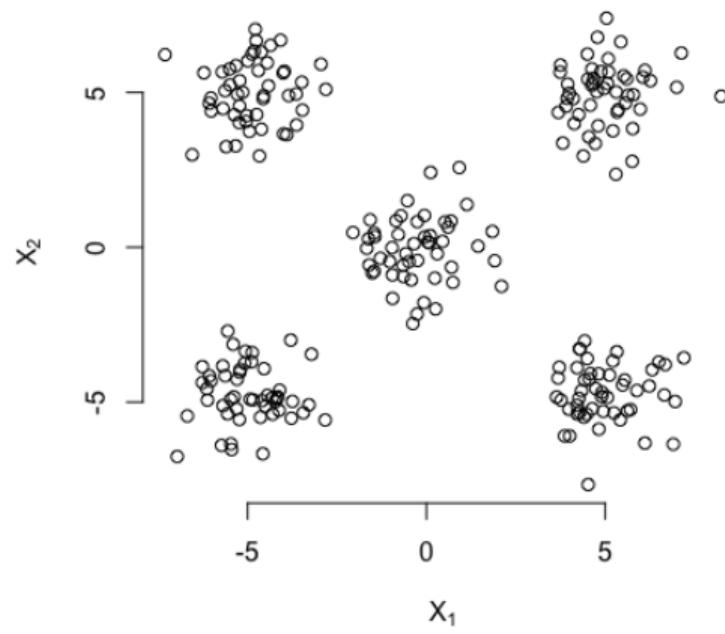
k-means



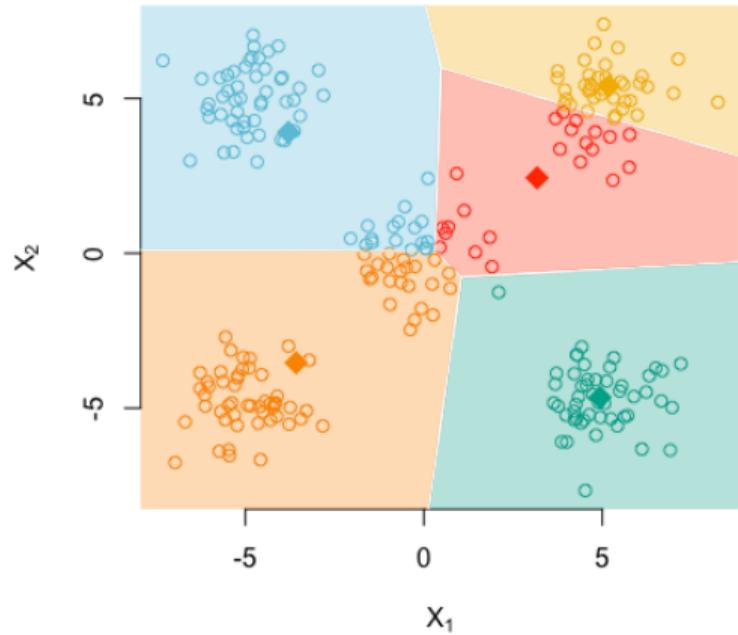
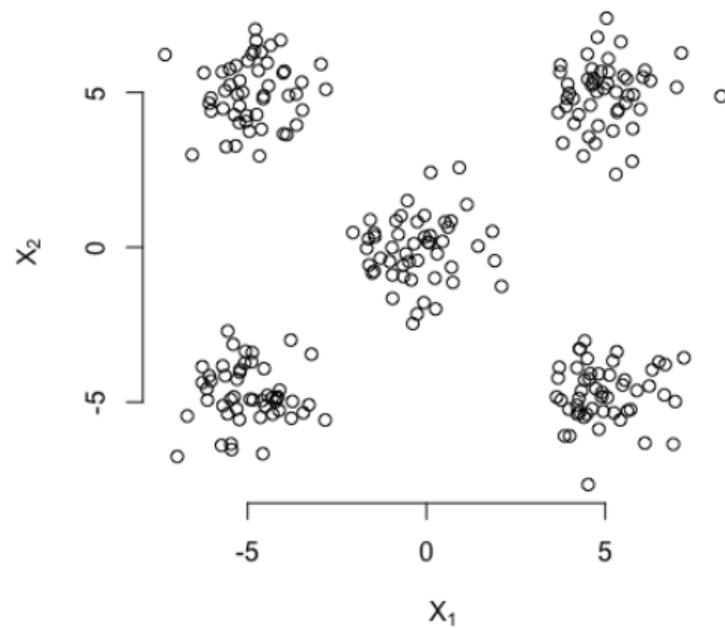
k-means



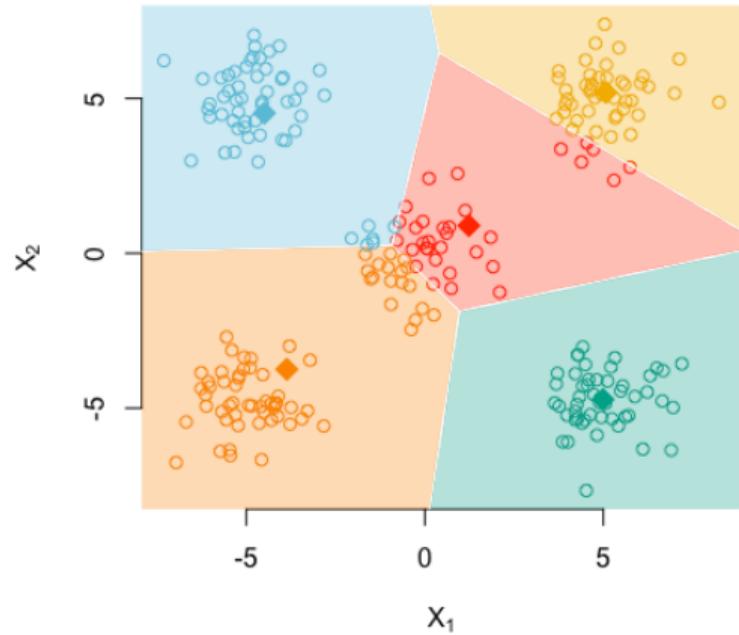
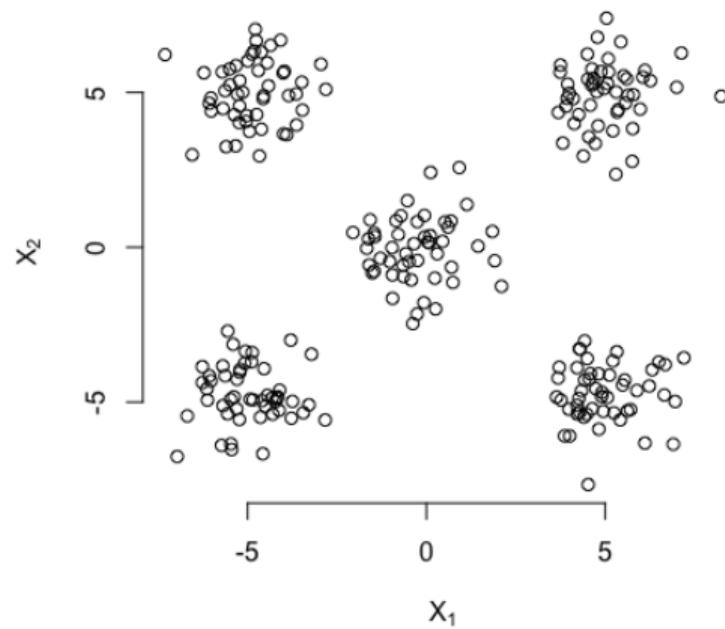
k-means



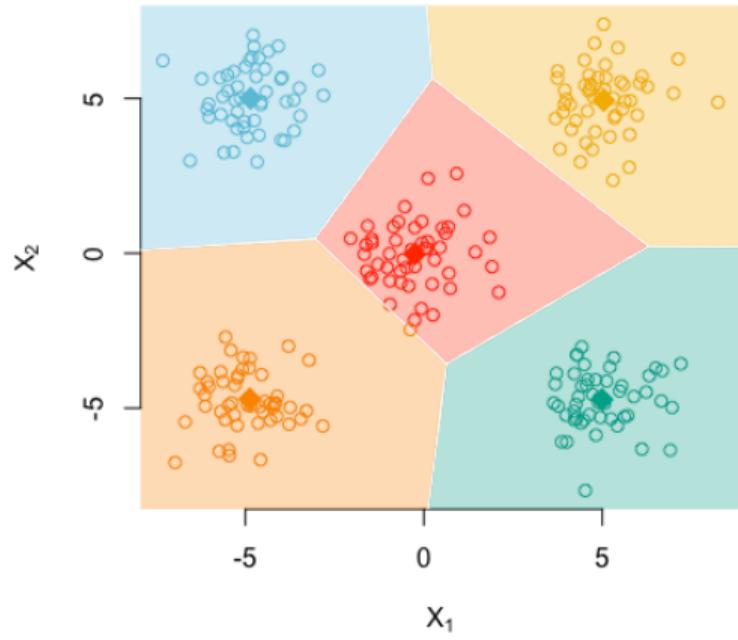
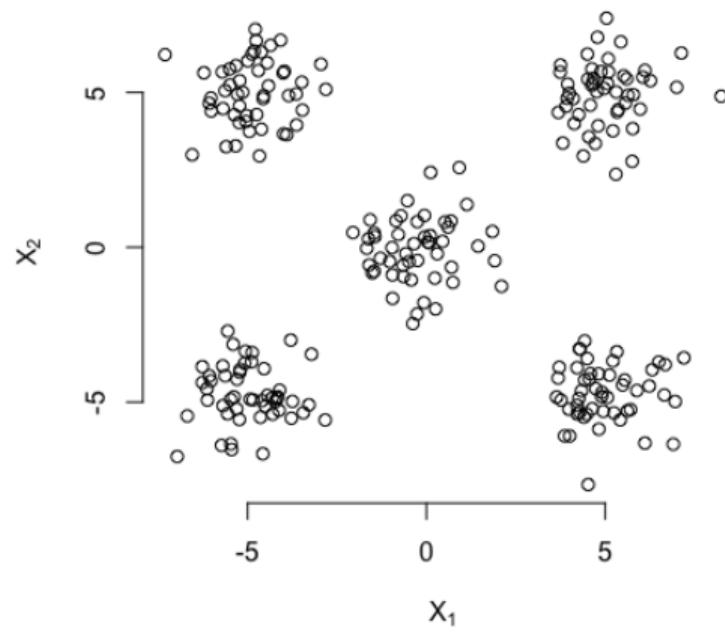
k-means



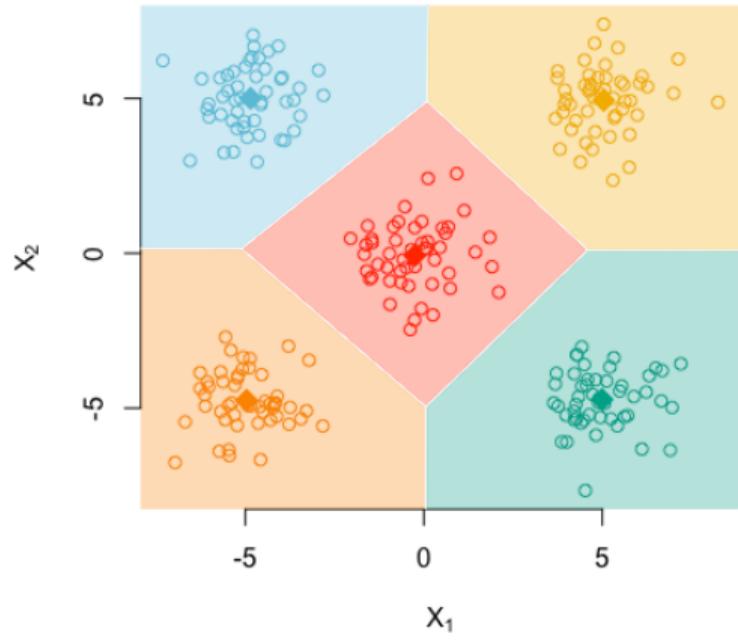
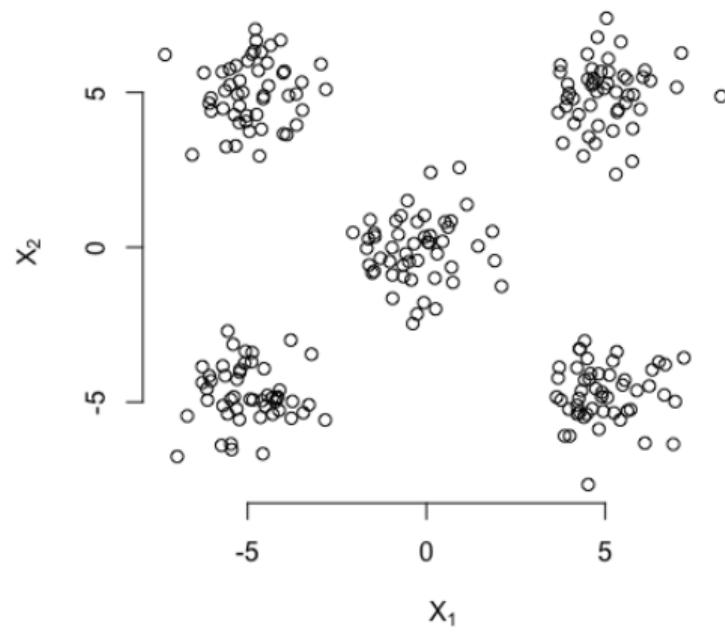
k-means



k-means



k-means



k-means

...

Given n points \mathbf{x}_i in \mathbb{R}^d , we want to find μ_1^*, \dots, μ_k^* , solutions of

$$V_{k\text{-means}} = \min \left\{ \sum_{j=1}^k \sum_{i \in S_j} \|\mathbf{x}_i - \mu_j\|^2 \right\}$$

that we can write

$$V_{k\text{-means}} = \min \left\{ \sum_{i=1}^n \min_{j=1, \dots, k} \underbrace{\|\mathbf{x}_i - \mu_j\|^2}_{\text{within inertia}} \right\}$$

Finding the optimal solution to the k -means clustering problem for n observations in d dimensions can be (exactly) solved in time $O(n^{dk+1} \log n)$... which is very long.

Given an initial set of k -means, μ_1, \dots, μ_k , alternate between the two following steps
(also called **Lloyd's algorithm**, from ?

k-means

Assignment step: we partition observations according to the Voronoi diagram generated by μ_1, \dots, μ_k 's,

$$S_i^{(t)} = \{x_j : \|x_j - \mu_i^{(t)}\|^2 \leq \|x_j - \mu_{i'}^{(t)}\|^2 \forall i' \in \{1, \dots, k\}\},$$

where each x_j is assigned to exactly one class $S_i^{(t)}$, at step t .

Update step: we compute the new centroids of the clusters,

$$\mu_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$

k-means

Key why this might work :

$$\frac{1}{\#I_j} \sum_{i,i' \in I} \sum_{j=1}^k [x_{il} - x_{i'l}]^2 = 2 \sum_{i \in I} \sum_{j=1}^k [x_{il} - \bar{x}_{jl}]^2 \text{ where } \bar{x}_{jl} = \frac{1}{\#I_j} \sum_{i \in I_j} x_{i,l}$$

The value of

$$\min_{I_1, \dots, I_K} \left\{ \sum_{j=1}^K \frac{1}{\#I_j} \sum_{i,i' \in I_j} \sum_{l=1}^k [x_{il} - x_{i'l}]^2 \right\}$$

will necessarily decrease (but no guarantee to reach the minimum, since the function is not convex).

k-means

The center of gravity of a point cloud is defined by

$$g = \sum_{i=1}^n p_i \mathbf{x}_i = \bar{\mathbf{x}}.$$

The total inertia of a point cloud is given by

$$I_T = \sum_{i=1}^n p_i d^2(\mathbf{x}_i, g),$$

where $d()$ is the Euclidean distance.

Inertie

In a classification procedure, the point cloud will be partitioned into K groups C_1, \dots, C_K whose respective weights are given by P_j such that

$$P_j = \sum_{i: \mathbf{x}_i \in C_j} p_i.$$

The center of gravity of each group is given by

$$g_j = \frac{1}{P_j} \sum_{i: \mathbf{x}_i \in C_j} p_i \mathbf{x}_i.$$

Inertia

In a classification procedure, the point cloud will be partitioned into K groups C_1, \dots, C_K whose respective weights are given by P_j such that

$$P_j = \sum_{i: \mathbf{x}_i \in C_j} p_i.$$

The center of gravity of each group is given by

$$g_j = \frac{1}{P_j} \sum_{i: \mathbf{x}_i \in C_j} p_i \mathbf{x}_i.$$

Hierarchical Clustering

The similarity between two groups can be determined in several ways.

Simple linkage (**simple linkage**): the proximity between two groups (A and B) is given by

$$\delta(A, B) = \min_{a_i \in A, b_j \in B} (d(a_i, b_j)),$$

i.e. we use the minimum of the distance between each of the point of group A and each point of group B .

Ascending Hierarchical Algorithms

Complete Linkage: the proximity between two groups (A and B) is given by

$$\delta(A, B) = \max_{a_i \in A, b_j \in B} (d(a_i, b_j)),$$

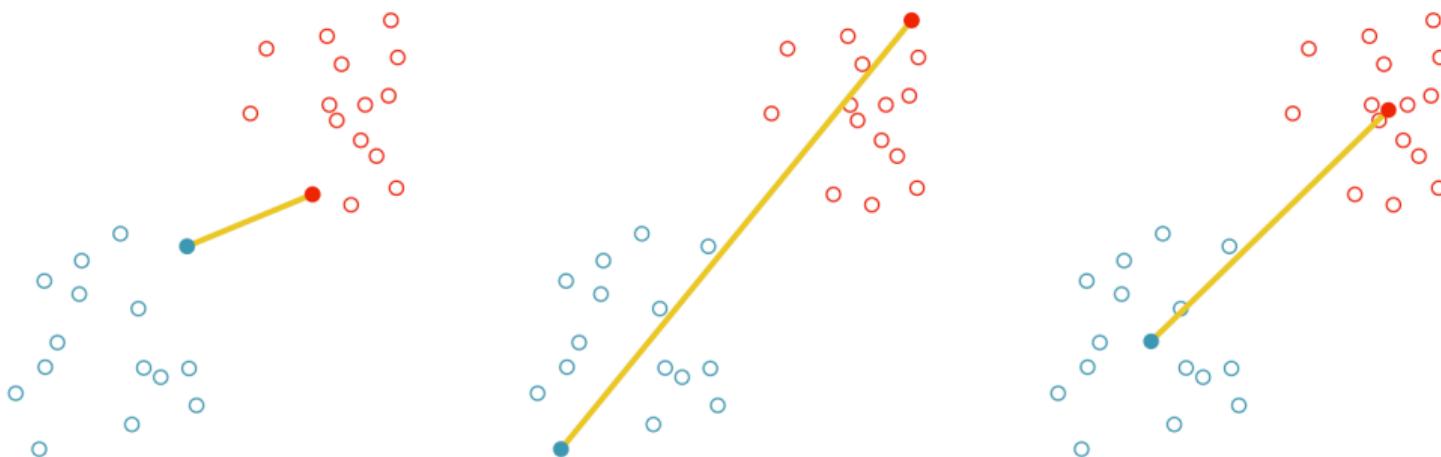
i.e. we use the maximum of the distance between each of the point of group A and each point of group B .

Ascending hierarchical algorithms

Average Linkage: the proximity between two groups groups (A and B) is given by

$$\delta(A, B) = \frac{1}{\text{card}(A)\text{card}(B)} \sum_{a_i \in A} \sum_{b_j \in B} (d(a_i, b_j)).$$

This time, we use the average of the distances between each of the points in group A and each of the points in group B .



Ascending Hierarchical Algorithms

Ward algorithm: the inertia of the point cloud is decomposed and the loss of information (loss of inertia between groups I_B) is minimized at each step. minimizes the loss of information (loss of inertia between groups I_B) at each stage: we therefore optimize the same criterion as for a dynamic reallocation classification algorithm.

At each stage, the loss of information due to the grouping of the groups A and B is given by

$$\delta(A, B) = \frac{P_A P_B}{P_A + P_B} d^2(g_A, g_B).$$

We also define the loss of information due to grouping groups $C = A \cup B$ and D by

$$\delta(C, D) = \frac{(P_A + P_D)\delta(A, D) + (P_B + P_D)\delta(B, D) - P_D\delta(A, B)}{P_A + P_B + P_D}.$$

Kaufman and Rousseeuw (2009) REPRENDRE

- Hierarchical clustering is a method of cluster analysis that builds a hierarchy of clusters.
- Two main approaches:

Spectral decomposition “★★★”

- If \mathbf{M} is a $p \times p$ matrix, then \vec{v} is an eigenvector of \mathbf{M} (associated with eigenvalue λ) if $\mathbf{M}\vec{v} = \lambda\vec{v}$.
- Every real symmetric $p \times p$ matrix \mathbf{M} is diagonalizable in an orthonormal basis.
- That is, $\mathbf{M} = \mathbf{V}\Lambda\mathbf{V}^\top$ where $\Lambda = \text{diag}(\lambda)$, and $\mathbf{V} = [\vec{v}_1, \dots, \vec{v}_n]$ with $\mathbf{V}^{-1} = \mathbf{V}^\top$.
- If \mathbf{M} is invertible ($\lambda_j \neq 0$), then $\mathbf{M}^{-1} = \mathbf{V}\Lambda^{-1}\mathbf{V}^\top$ where $\Lambda = \text{diag}(\lambda)$, and $\mathbf{V} = [\vec{v}_1, \dots, \vec{v}_n]$.

Principal Component Analysis

Principal Component Analysis (PCA) is a dimensionality reduction technique. It identifies new axes (principal components) that maximize the variance in the data.

Given a data matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ (centered), PCA finds orthonormal vectors $\vec{v}_1, \dots, \vec{v}_k$ such that:

$$\vec{v}_1 = \arg \max_{\|\vec{v}\|=1} \text{Var}(\mathbf{X}\vec{v})$$

The sample covariance matrix is:

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X}$$

PCA seeks the eigenvectors of \mathbf{S} .

Principal Component Analysis

Eigen decomposition:

$$\mathbf{S} = \mathbf{V}\Lambda\mathbf{V}^\top$$

where \mathbf{V} contains the principal components (eigenvectors).

Let \mathbf{V}_k be the matrix with the first k eigenvectors.

The projection of \mathbf{X} onto the k -dimensional subspace:

$$\mathbf{Z} = \mathbf{X}\mathbf{V}_k$$

\mathbf{Z} is the low-dimensional representation of the data.

The eigenvalues $\lambda_1, \dots, \lambda_p$ represent the variance explained by each principal component.

The proportion of variance explained by the i -th component:

$$\frac{\lambda_i}{\sum_{j=1}^p \lambda_j}$$

Principal Component Analysis

Often use a scree plot to determine how many components to keep.

PCA rotates the coordinate axes to align with directions of maximum variance.

The first principal component defines the direction of highest variance.

Each subsequent component is orthogonal to the previous ones.

PCA is often viewed from two perspectives:

Individuals: the rows of the data matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$

Variables: the columns of \mathbf{X}

Let \mathbf{X} be centered (and possibly scaled), then:

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X} \quad (\text{covariance of variables})$$

$$\mathbf{G} = \frac{1}{p} \mathbf{X} \mathbf{X}^\top \quad (\text{Gram matrix of individuals})$$

These matrices are linked through the **duality** of eigen decompositions.

Principal Component Analysis

Eigen decomposition of the covariance matrix:

$$\mathbf{S} = \mathbf{V}\Lambda\mathbf{V}^T, \quad \mathbf{V}^T\mathbf{V} = \mathbf{I}$$

Principal component scores (individuals in PC space):

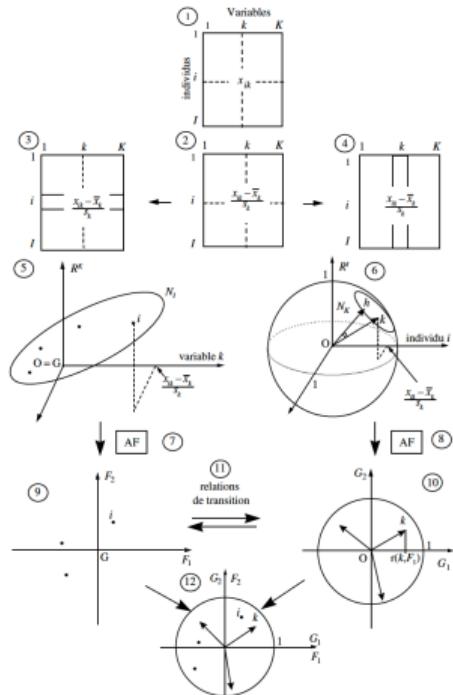
$$\mathbf{Z} = \mathbf{X}\mathbf{V}$$

Dual decomposition (of the Gram matrix):

$$\mathbf{G} = \mathbf{Z}\mathbf{Z}^T = \mathbf{X}\mathbf{V}\Lambda\mathbf{V}^T\mathbf{X}^T$$

Individuals and variables are linked via the **same eigenvalues**.

A **biplot** represents both individuals and variables in the same low-dimensional space.



Principal Component Analysis

Let \mathbf{X} be standardized. The loading vectors (columns of \mathbf{V}) can be plotted together with:

$$\text{Scores (individuals)} : \mathbf{Z} = \mathbf{X}\mathbf{V}$$

$$\text{Loadings (variables)} : \mathbf{V}\sqrt{\Lambda}$$

The angles between variable vectors approximate correlations.

The projection of an individual onto a variable axis estimates its value on that variable.

Partial Least Squares (PLS) regression finds components that explain variance in both \mathbf{X} and \mathbf{Y} .

It is useful when predictors are highly collinear or $p > n$.

PLS finds latent vectors $\mathbf{t}_i = \mathbf{X}\vec{w}_i$ such that:

$$\text{Cov}^2(\mathbf{X}\vec{w}_i, \mathbf{Y}\vec{c}_i) \quad \text{is maximized}$$

Principal Component Analysis

Iteratively extracts components from \mathbf{X} and \mathbf{Y} .

PCA: maximizes variance in \mathbf{X} only.

PLS: maximizes covariance between \mathbf{X} and \mathbf{Y} .

PLS components are better suited for prediction tasks.

Autoencoders

Definition 2.23: Dissimilarity

$L : \mathbb{R}^q \times \mathbb{R}^q \rightarrow \mathbb{R}_+$ such that $L(\mathbf{x}, \mathbf{x}') = 0$ if and only if $\mathbf{x} = \mathbf{x}'$.

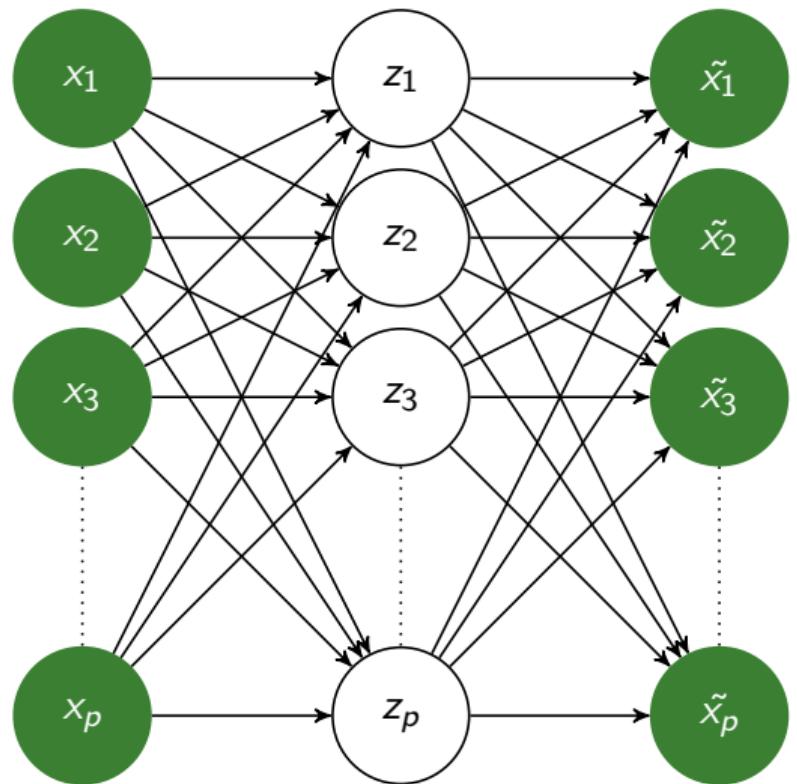
Definition 2.24: Auto-encoder

Let $p < q$. An auto-encoder is a pair (Φ, Ψ) of mappings, called **encoder** and **decoder** respectively,

$$\Phi : \mathbb{R}^q \rightarrow \mathbb{R}^p \text{ and } \Psi : \mathbb{R}^p \rightarrow \mathbb{R}^q$$

such that their composition $\Psi \circ \Phi$ has a small reconstruction error w.r.t. the chosen dissimilarity function L , i.e. for all \mathbf{x} , $L(\mathbf{x}, \Psi \circ \Phi(\mathbf{x}))$ is small.

Autoencoders



$$\mathbf{X}^\top \mathbf{X} = \mathbf{V} \Lambda \mathbf{V}^\top$$

or

$$(\mathbf{V} \mathbf{X})^\top (\mathbf{X} \mathbf{V}) = \mathbf{V}^\top \mathbf{V} \Lambda \mathbf{V}^\top \mathbf{V} = \Lambda$$

Spectral decomposition

- Si \mathbf{M} est une $p \times p$ matrice, \vec{v} est un vecteur propre de \mathbf{M} (associé à λ) si $\mathbf{M}\vec{v} = \lambda\vec{v}$
- Si \mathbf{M} peut être diagonalisée, $\mathbf{M} = \mathbf{V}\Lambda\mathbf{V}^\top$ où $\Lambda = \text{diag}(\lambda)$, et $\mathbf{V} = [\vec{v}_1, \dots, \vec{v}_n]$.
- Si \mathbf{M} est inversible, $\mathbf{M}^{-1} = \mathbf{V}\Lambda^{-1}\mathbf{V}^\top$ où $\Lambda = \text{diag}(\lambda)$, et $\mathbf{V} = [\vec{v}_1, \dots, \vec{v}_n]$.
- **Eckart–Young–Mirsky** théorème ([Eckart and Young \(1936\)](#), [Mirsky \(1960\)](#)):

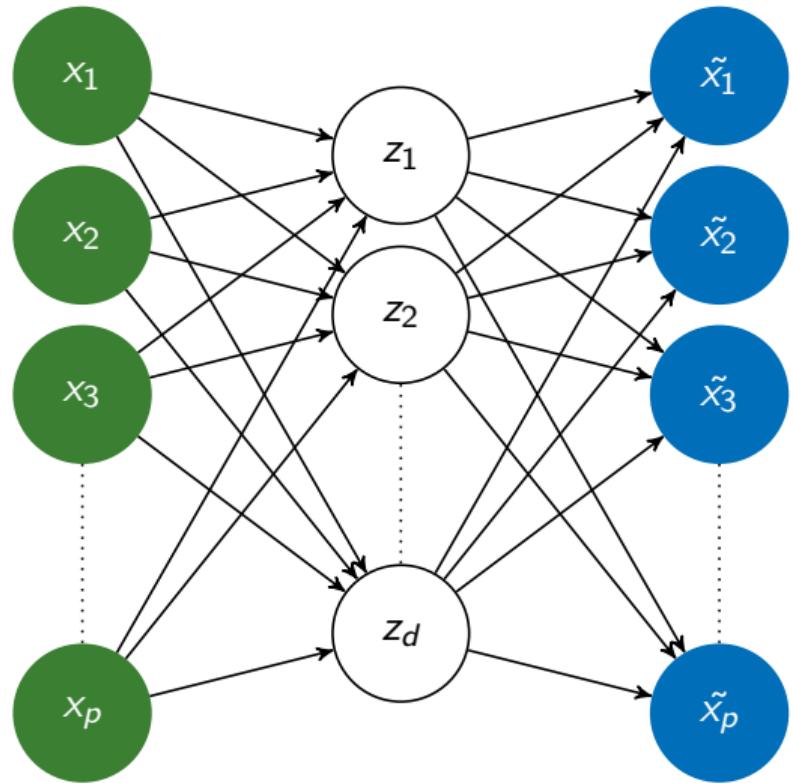
$$\mathbf{M}^* \in \underset{\mathbf{H}:p \times p}{\operatorname{argmin}} \{ \|\mathbf{M} - \mathbf{H}\|_F \text{ s.t. } \text{rank}(\mathbf{H}) \leq d \}$$

Si $\mathbf{M} = \mathbf{V}\Lambda\mathbf{V}^\top = [\mathbf{V}_d \quad \mathbf{V}_{p-d}] \begin{bmatrix} \Lambda_r & \mathbf{0} \\ \mathbf{0} & \Lambda_{n-r} \end{bmatrix} [\mathbf{V}_r \quad \mathbf{V}_{p-d}]^\top$, $\mathbf{M}^* = \mathbf{V}_d \Lambda_d \mathbf{V}_d^\top$

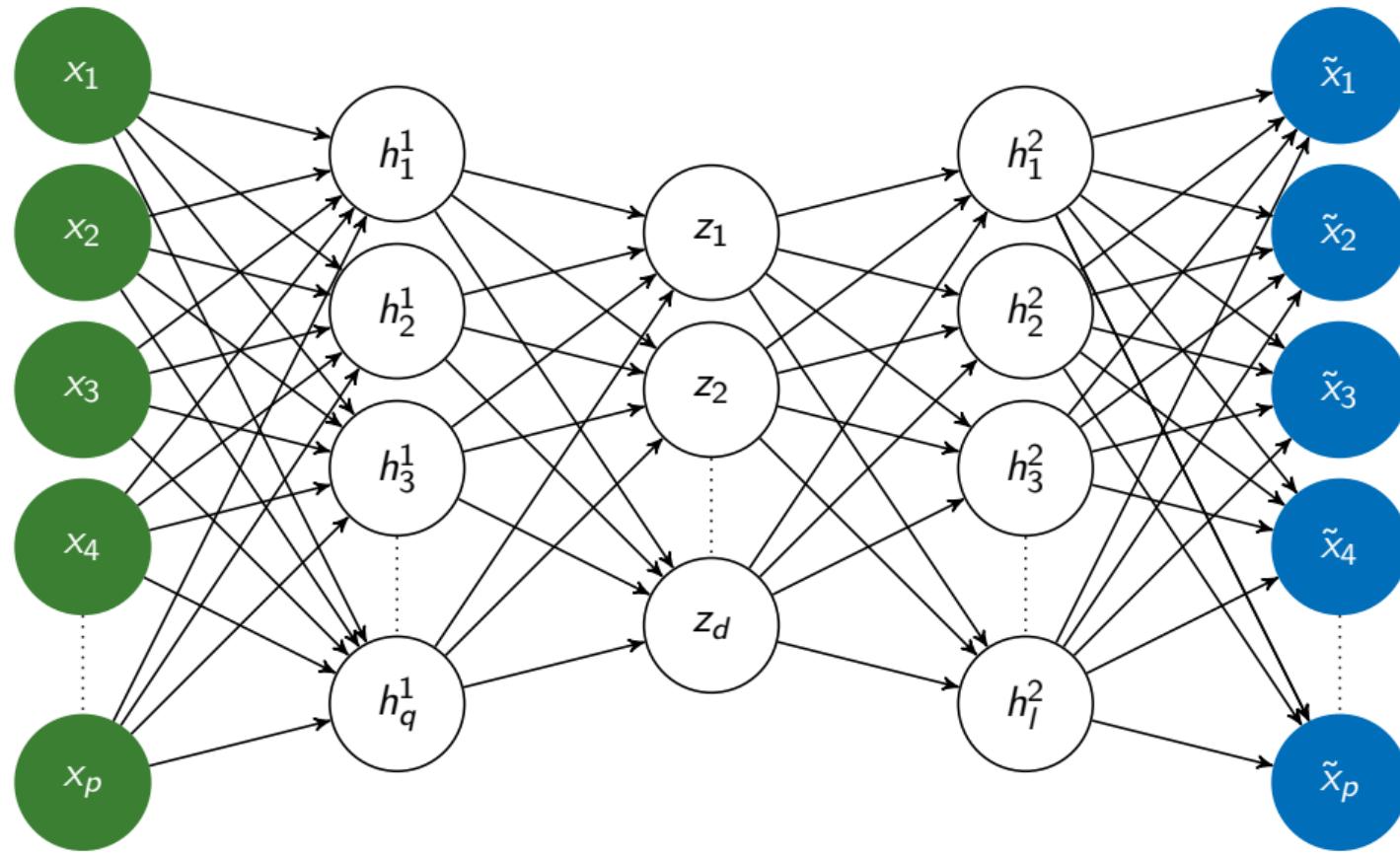
(ce qui sera unique si $\lambda_d > \lambda_{d+1}$)

- La approximation est $\|\mathbf{M} - \mathbf{M}^*\|_F = \sqrt{\lambda_{d+1}^2 + \dots + \lambda_p^2}$

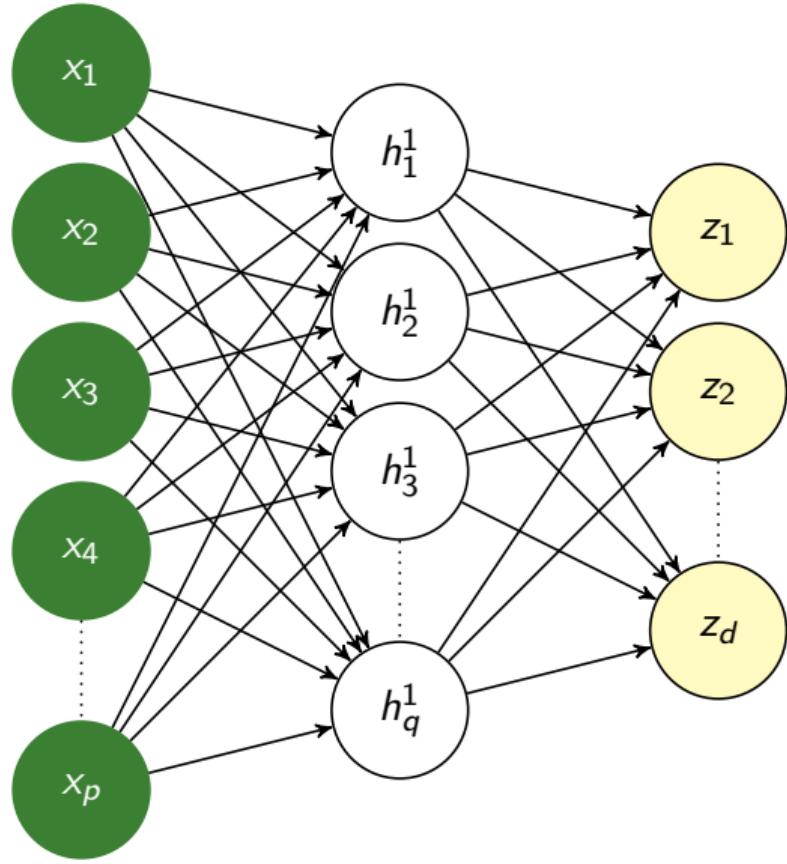
Principal Component Analysis



Autoencoders



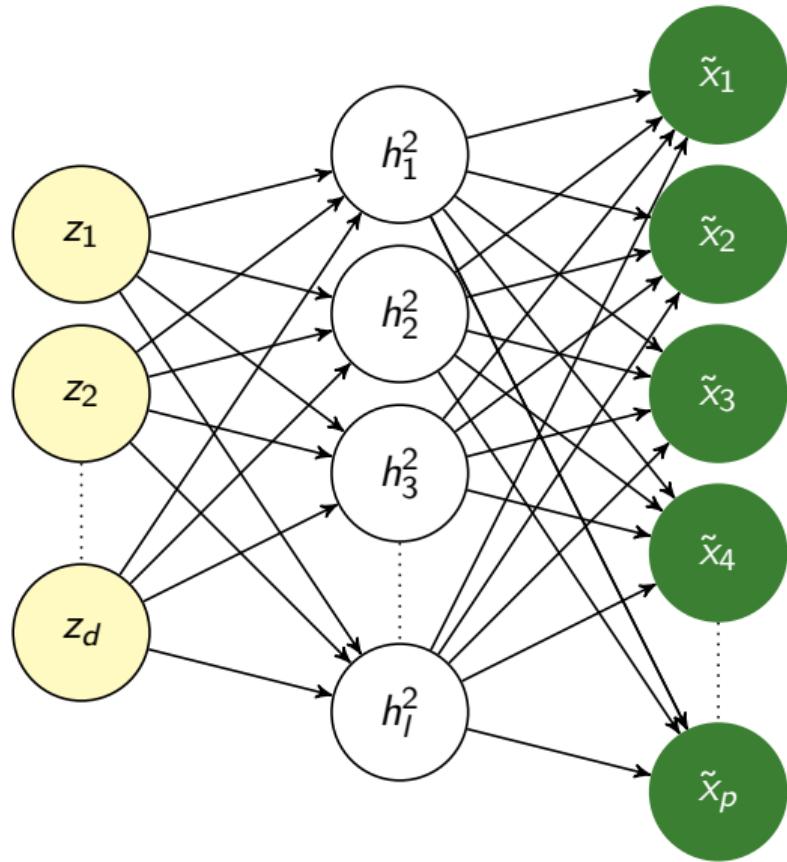
Encoder



Let $f : \mathcal{X} \rightarrow \mathcal{Z}$ be an encoding function
then, for example:

$$\mathbf{z} = f(\mathbf{x}) = \phi_2(\phi_1(\mathbf{x}_{1 \times p} \mathbf{B}_{p \times q}^{(1)} \mathbf{B}_{q \times d}^{(2)})$$

Decoder



Let $g : \mathcal{Z} \rightarrow \mathcal{X}$ be a neural network-type
then, for example:

$$\tilde{x} = g(\mathbf{z}) = \phi_4(\phi_3(\mathbf{z}_{1 \times d} \mathbf{B}_{d \times l}^{(3)}) \mathbf{B}_{l \times p}^{(4)})$$

Autoencoders

The error function of a **linear autoencoder** (PCA) is

$$\|\mathbf{X} - \tilde{\mathbf{X}}\|^2 = \|\mathbf{X} - \mathbf{P}^\top \mathbf{P}\mathbf{X}\|^2 = \sum_{i=1}^n (\mathbf{P}^\top \mathbf{P}\mathbf{x}_i - \mathbf{x}_i)^\top (\mathbf{P}^\top \mathbf{P}\mathbf{x}_i - \mathbf{x}_i)$$

The error function of a **nonlinear autoencoder** is

$$\|\mathbf{X} - \tilde{\mathbf{X}}\|^2 = \|\mathbf{X} - g \circ f(\mathbf{X})\|^2 = \sum_{i=1}^n (g \circ f(\mathbf{x}_i) - \mathbf{x}_i)^\top (g \circ f(\mathbf{x}_i) - \mathbf{x}_i)$$

Autoencoders

Encoder Φ_p / Decoder Ψ_p

$$\Psi_p \circ \Phi_p : \mathbf{X} \mapsto \sum_{j=1}^p \langle \mathbf{X}, \vec{\mathbf{v}}_j \rangle \vec{\mathbf{v}}_j$$

with reconstruction error of

$$\mathbf{X} - \Psi_p \circ \Phi_p(\mathbf{X}) = \sum_{j=p+1}^q \langle \mathbf{x}, \vec{\mathbf{v}}_j \rangle \vec{\mathbf{v}}_j$$

Generative Modeling

Definition

A **generative model** learns the joint probability distribution $p(x, y)$ of the data and labels, and can generate new samples from it:

$$p(x) = \int p(x|z)p(z)dz$$

where $p(x|z)$ is the likelihood, and $p(z)$ is the prior distribution of hidden variables.

The goal is to learn how data is generated.

Examples include Gaussian Mixture Models (GMMs) and Variational Autoencoders (VAEs).

Can be used for tasks such as image generation, text synthesis, and anomaly detection.

Generative Modeling

Definition

A **Gaussian Mixture Model** (GMM) is a specific case of a mixture model where each component is a Gaussian distribution:

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \sigma_k^2)$$

where:

$\mathcal{N}(x|\mu_k, \sigma_k^2)$ is a Gaussian distribution with mean μ_k and variance σ_k^2 .

π_k are the mixture coefficients, satisfying $\sum_k \pi_k = 1$.

GMMs are widely used in clustering and density estimation tasks.

The Expectation-Maximization (EM) algorithm is typically used for training GMMs.

Generative Modeling

Generative Process

Generative AI methods aim to model the data distribution $p(x)$ and generate new samples x' . For example, in a GMM:

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \sigma_k^2)$$

The goal is to sample new data $x' \sim p(x)$.

Generative models are foundational for advanced AI techniques, such as GANs and VAEs.

These models aim to capture data distributions in a way that allows realistic data generation.

Example: Generative models can create images, music, or even text based on learned patterns.

Generative Modeling

Generative Adversarial Networks (GANs): Two neural networks (generator and discriminator) compete to generate realistic data.

Variational Autoencoders (VAEs): Use probabilistic graphical models to learn latent variable representations.

Both GANs and VAEs use latent variables to model complex data distributions and generate new samples.

Generative Process in VAEs

A VAE learns a probabilistic mapping:

$$p(x) = \int p(x|z)p(z)dz$$

where $p(x|z)$ is the likelihood, and $p(z)$ is the prior on the latent variable z .

These models are used for unsupervised learning, data generation, and anomaly detection.

Variational Autoencoders

Kingma and Welling (2013)

- Variational Autoencoders (VAEs) are generative models that learn latent representations of data.
- Unlike standard autoencoders, VAEs impose a probabilistic structure on the latent space.
- VAEs use variational inference to approximate intractable posterior distributions.

Gaussian Latent Variable Model

- Encoder (inference/recognition model):

$$q(\mathbf{z} \mid \mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}_x(\mathbf{x}), \boldsymbol{\Sigma}_x(\mathbf{x}))$$

- Decoder (generative model):

$$p(\mathbf{x} \mid \mathbf{z}) \sim \mathcal{N}(\boldsymbol{\mu}_z(\mathbf{z}), \boldsymbol{\Sigma}_z(\mathbf{z}))$$

Variational Autoencoders

- the decoder is the generative component: once trained, it can take random samples $\mathbf{z} \sim \pi$ from the latent space and produce synthetic data points by generating new data $\mathbf{x}' \sim p(\cdot | \mathbf{z})$ that resembles the original dataset.

Latent Variable Model

- Assume observed data $\mathbf{x} \in \mathbb{R}^d$ is generated from a latent variable $\mathbf{z} \in \mathbb{R}^k$, i.e., there exists an (unknown) deterministic mapping $\mathbf{x} \mapsto \mathbf{z}$
- The generative process:

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})\pi(\mathbf{z})d\mathbf{z} \quad (10)$$

- The posterior distribution $p(\mathbf{z}|\mathbf{x})$ is typically intractable.

Variational Inference

- Instead of directly computing $p(\mathbf{z}|\mathbf{x})$, approximate it using $q(\mathbf{z}|\mathbf{x})$.

freakonometrics

freakonometrics.hypotheses.org

Variational Autoencoders

- Use Kullback-Leibler (KL) divergence to measure approximation quality:

$$D_{KL}(q(z|x)||p(z|x)) \quad (11)$$

Evidence Lower Bound (ELBO)

- The marginal log-likelihood can be rewritten as:

$$\log p(x) = \mathbb{E}_{q(z|x)} \left[\log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x)||p(z|x)) \quad (12)$$

- The ELBO is defined as:

$$\mathcal{L}(x) = \mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{KL}(q(z|x)||p(z)) \quad (13)$$

- Maximizing ELBO minimizes $D_{KL}(q(z|x)||p(z|x))$.

Reparameterization Trick

Variational Autoencoders

- Directly sampling from $q(z|x)$ is not differentiable.
- Introduce a reparameterization:

$$z = \mu(x) + \sigma(x) \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I) \quad (14)$$

- This allows gradients to flow through μ and σ during optimization.

VAE Architecture

- Encoder: $q(z|x) = \mathcal{N}(z; \mu(x), \Sigma(x))$
- Decoder: $p(x|z) \sim \mathcal{N}(x; f(z), \sigma^2 I)$
- Loss function:

$$\mathcal{L} = \mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{KL}(q(z|x)||p(z)) \quad (15)$$

KL Divergence Computation

- Assume $q(z|x) = \mathcal{N}(z; \mu, \sigma^2 I)$ and $p(z) = \mathcal{N}(0, I)$.

Variational Autoencoders

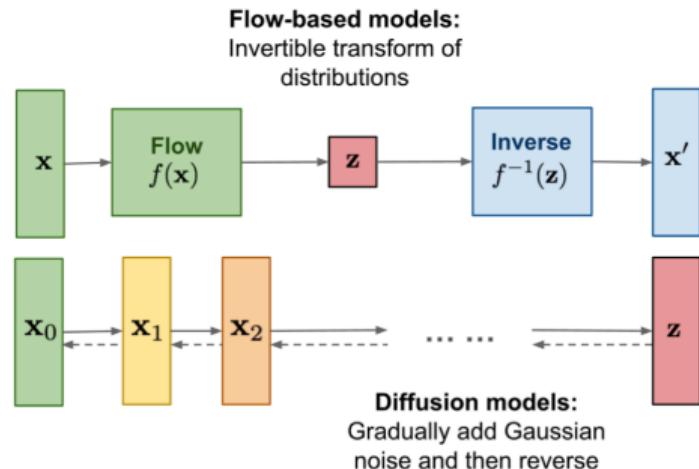
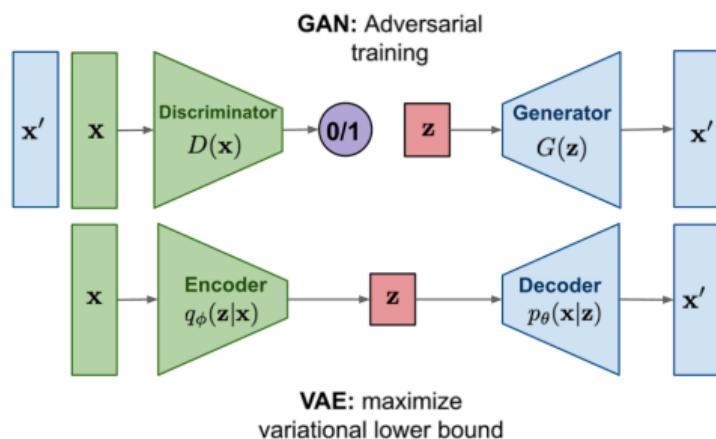
- Then,

$$D_{KL}(q(z|x)||p(z)) = \frac{1}{2} \sum_{i=1}^k \left(\sigma_i^2 + \mu_i^2 - 1 - \log \sigma_i^2 \right) \quad (16)$$

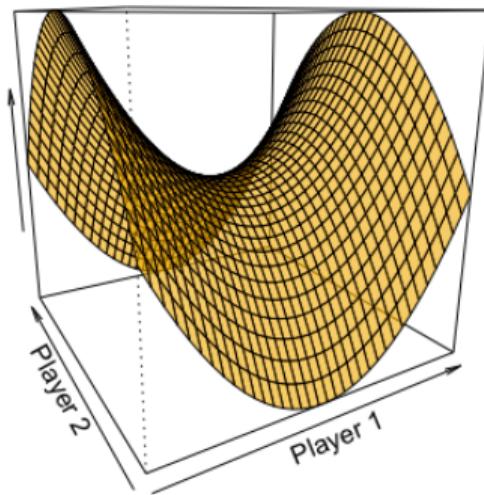
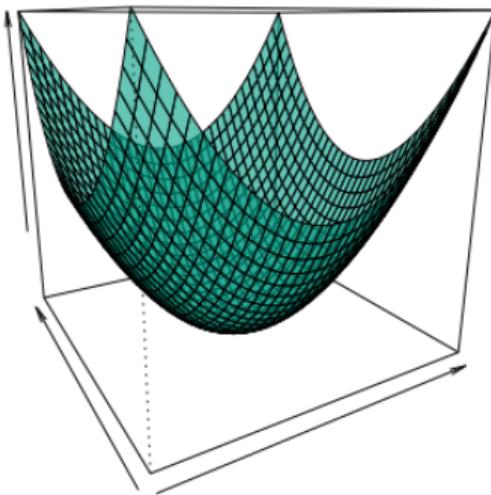
Sampling and Generation

- After training, generate new samples by drawing $z \sim \mathcal{N}(0, I)$ and decoding $x \sim p(x|z)$.
- Ensures meaningful interpolation in latent space.

Overview



Adversarial learning



Adversarial learning

Definition

Adversarial learning focuses on the study of models' vulnerabilities when exposed to perturbations in the input data.

Adversarial examples are intentionally crafted inputs that cause machine learning models to make incorrect predictions.

Adversarial learning explores how to defend against these attacks or how to use adversarial examples for training.

Adversarial learning

Definition 2.25: Adversarial example

An **adversarial example** \mathbf{x}' is an input to a model that is close to an original input \mathbf{x} , but causes the model to different predictions, i.e., given a model $m : \mathcal{X} \rightarrow \mathcal{Y}$, a distance metric d on \mathcal{Y} and a norm $\|\cdot\|$ on \mathcal{X} , and thresholds $\varepsilon, \alpha > 0$, an input $\mathbf{x}' \in \mathcal{X}$ is called an **adversarial example** for an original input $\mathbf{x} \in \mathcal{X}$ if:

$$\begin{cases} \|\mathbf{x} - \mathbf{x}'\| < \varepsilon & \text{small input perturbation} \\ d(m(\mathbf{x}), m(\mathbf{x}')) > \alpha & \text{large model output change} \end{cases}$$

Adversarial examples are typically generated by small perturbations.

These perturbations are often imperceptible to humans but lead to model errors.

Example: Slightly altering pixels in an image to mislead an image classifier.

Adversarial learning

Types of Attacks

Adversarial attacks are methods used to generate adversarial examples. Common attacks include:

Fast Gradient Sign Method (FGSM): Perturbation is calculated by:

$$\mathbf{x}' = \mathbf{x} + \varepsilon \cdot \text{sign}(\nabla_{\mathbf{x}} \ell(m_{\theta}(\mathbf{x}), y))$$

where $\ell(m_{\theta}(\mathbf{x}), y)$ is the loss function, and ε is the perturbation magnitude.

Projected Gradient Descent (PGD): Iterative version of FGSM, where the adversarial example is refined over multiple steps.

These attacks are designed to optimize the perturbation such that the model misbehaves.

The goal is to maximize the loss function with respect to the input x .

Adversarial learning

Definition

Adversarial training is a technique to defend against adversarial attacks by including adversarial examples in the training process.

The model is trained on both original and adversarial examples.

This forces the model to learn to be robust to perturbations:

$$\mathcal{L}_{adv}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\max_{x' \in \mathcal{A}} \ell(m_\theta(x'), y) \right]$$

where \mathcal{A} represents the adversarial set.

Example: Adding adversarial examples generated by FGSM to the training data.

The model learns to be less sensitive to small perturbations.

Adversarial Training

The adversarial training process includes adversarial examples in the training data to improve the model's robustness.

Adversarial learning

During adversarial training, adversarial examples are generated and used to update the model.

The model learns to be more robust to these perturbations.

Adversarial learning addresses the vulnerability of machine learning models to small, carefully crafted perturbations.

Adversarial examples can be used to improve model robustness through adversarial training.

Applications:

Security: Protecting AI systems from adversarial attacks (e.g., image classification, speech recognition).

Robustness testing: Evaluating the stability of models in real-world environments.

Generative models: Using adversarial learning principles in GANs to train robust models.

Adversarial learning

Adversarial attacks generate input perturbations that exploit the model's weaknesses.

These attacks can be used to fool models into making wrong predictions.

Adversarial learning

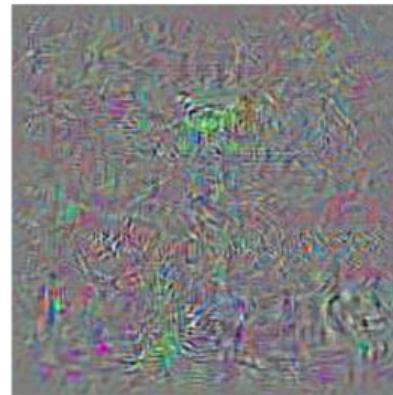
Adversarial learning has to do with robustness of learning algorithm, Szegedy et al. (2013), “are neural network stables?”.

“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

Adversarial learning

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...

Adversarial learning

$$\min_{\theta} \left\{ \max_{\varepsilon \in \mathcal{E}} \left\{ \mathbb{E}_{(\mathbf{X}, Y) \sim \mathbb{P}} [\ell(m_{\theta}(\mathbf{X} + \varepsilon), Y)] \right\} \right\}$$

training a robust classifier
creating an adversarial example

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 2.7: 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

$$\mathbf{y}_*^\top A \mathbf{x}_* = \max_{\mathbf{x} \in \mathcal{S}_m} \min_{\mathbf{y} \in \mathcal{S}_n} \{\mathbf{y}^\top A \mathbf{x}\} = \min_{\mathbf{y} \in \mathcal{S}_n} \max_{\mathbf{x} \in \mathcal{S}_m} \{\mathbf{y}^\top A \mathbf{x}\}.$$

Adversarial learning

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

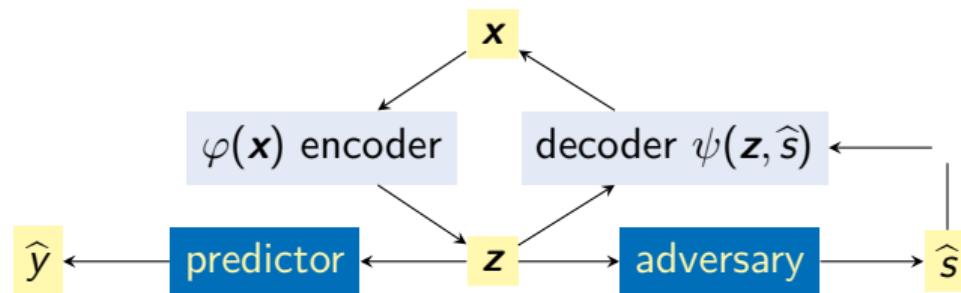
Adversarial learning

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)

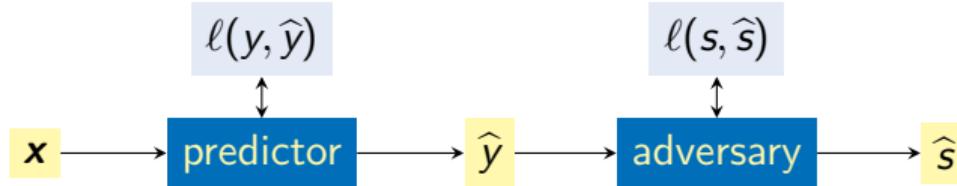
Learning adversarially fair and transferable representations, Madras et al. (2018)



Adversarially learning fair representations, Beutel et al. (2017)

Fair Adversarial Debiasing Approach, Zhang et al. (2018)

Adversarial learning



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

Sequential learning

Classical supervised machine learning, the learner

- observes training data with labels,
- builds a program to minimize the training error
- controls the error of new data if they are similar to the training data

but

- the environment may evolve over time
- the data may be available sequentially

Sequential learning

At each time step $t = 1, \dots, T$

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t \in \Theta$ (compact decision/parameter set);
- the environment chooses a loss function $\ell_t : \Theta \rightarrow [0, 1]$;
- the player suffers loss $\ell_t(\theta_t)$ and observes
 - the losses of every actions: $\ell_t(\theta)$, for all $\theta \in \Theta$: full-information feedback
 - the loss of the chosen action only: $\ell_t(\theta_t)$: bandit feedback.

Goal. Minimize the cumulative loss $\sum_{t=1}^T \ell_t()$

Multi-Armed Bandit

In probability theory and machine learning, the multi-armed bandit problem (sometimes called the K -armed bandit problem) is a problem in which a decision maker iteratively selects one of multiple fixed choices (i.e., arms or actions) when the properties of each choice are only partially known at the time of allocation, and may become better understood as time passes. A fundamental aspect of bandit problems is that choosing an arm does not affect the properties of the arm or other arms. The multi-armed bandit problem is a classic reinforcement learning problem that exemplifies the exploration–exploitation tradeoff dilemma.

W

Exploration–exploitation dilemma

The exploration–exploitation dilemma, also known as the explore–exploit tradeoff, is a fundamental concept in decision-making. It is depicted as the balancing act between two opposing strategies. Exploitation involves choosing the best option based on current knowledge of the system (which may be incomplete or misleading), while exploration involves trying out new options that may lead to better outcomes in the future at the expense of an exploitation opportunity. Finding the optimal balance between these two strategies is a crucial challenge in many decision-making problems whose goal is to maximize long-term benefits

W

Multi-armed bandit

Consider K possible actions (arms)

Each arm k is associated an unknown probability distribution with mean μ_k
sequentially pick an arm k_t and get reward $X_{k,t}$, with mean $\mu_{k,t}$

Goal: maximize the expected cumulative reward $\mathbb{E} \left[\sum_{t=1}^T X_{k,t} \right]$

Exploration vs Exploitation trade-off.

Historical motivation, [Thompson \(1933\)](#): clinical trials, for each patient t in a clinical study

- choose a treatment k_t
- observe response to the treatment X_{kt} Goal: maximize the number of patient healed (or find the best treatment)

Successful because of many applications coming from Internet: recommender systems, online advertisements,

Multi-armed bandit

At each time step $t = 1, \dots, T$

- the player observes a context $x_t \in \mathcal{X}$ (optional step)
- the player chooses an action $\theta_t = k_t \in \Theta = \{1, \dots, K\}$ (compact decision/parameter set);
- the environment chooses a loss function (by sampling the arms);
- the player suffers loss and observes
 - the losses of every actions: full-information feedback
 - the loss of the chosen action only: bandit feedback.

The goal of the player is to minimize his cumulative loss:

$$r_T := \sum_{t=1}^T \ell_t(\theta_t)$$

REPRENDRE

Definition 2.26: Regret

The regret of the player with respect to a fixed parameter $\theta^* \in \Theta$ after T time steps is

$$r_T(\theta^*) := \sum_{t=1}^T \ell_t(\theta_t) - \sum_{t=1}^T \ell_t(\theta^*)$$

The regret (or uniform regret) is defined as $r_T^* = \sup_{\theta^* \in \Theta} \{r_T(\theta^*)\}$

Multi-armed bandit

Instead of a deterministic choice, one can consider a probability distribution,

$$\theta_t \in \mathcal{S}_K = \left\{ \mathbf{p} \in [0, 1]^K : \sum_{k=1}^K p_k = 1 \right\}$$

How to choose the weights At round t the player needs to choose a weight vector $\theta_t \in \mathcal{S}_K$. How to choose the weights? The player should

- give more weight to actions that performed well in the past.
- not give all the weight to the current best action, otherwise it would not work

The exponentially weighted average forecaster (EWA) also called Hedge performs this trade-off by choosing a weight that decreases exponentially fast with the past errors.

Multi-armed bandit

Algorithm 5: Exponentially Weighted Average forecaster

- 1 initialization : $\eta > 0$ and $\mathbf{p}_1 \propto (1, 1, \dots, 1)$;
 - 2 **for** $t=1,2,\dots,T$ **do**
 - 3 select \mathbf{p}_1 ; incur loss $\ell_t(\mathbf{p}_t) = \mathbf{p}_t^\top \mathbf{g}_t$ and observe $\mathbf{g}_t \in [-1, +1]^K$; update for all k ;
 - 4 $p_{t+1}(k) \leftarrow \frac{\exp[-\eta \sum_{i=1}^t g_i(k)]}{\sum_{j=1}^K \exp[-\eta \sum_{i=1}^t g_i(j)]}$
-

In that case, it is possible to get an upper bound for regret

Multi-armed bandit

Proposition 2.8: Regret bound for EWA

Given T , for all sequences of loss vectors $\mathbf{g}_1, \dots, \mathbf{g}_T \in [1, +1]^K$, EWA achieves the bound

$$r_T := \sum_{t=1}^T \ell_t(\mathbf{p}_t) - \min_{\mathbf{p} \in \mathcal{S}_K} \sum_{t=1}^T \ell_t(\mathbf{p}) \leq \eta \sum_{t=1}^T \sum_{k=1}^K p_t(k) g_t(k)^2 + \frac{\log(K)}{\eta}$$

since $\ell_t(\mathbf{p}) = \sum_{k=1}^K p(k) g_t(k) = \mathbf{p}^\top \mathbf{g}_t$. If $\eta = \sqrt{\frac{\log(K)}{T}}$, EWA satisfies

$$r_T \leq 2\sqrt{T \log(K)}$$

This regret bound is optimal, see [Cesa-Bianchi and Lugosi \(2006\)](#)

Reinforcement Learning

Reinforcement learning (RL) is an interdisciplinary area of machine learning and optimal control concerned with how an intelligent agent should take actions in a dynamic environment in order to maximize a reward signal. Reinforcement learning is one of the three basic machine learning paradigms, alongside supervised learning and unsupervised learning. W

Bayesianism as a learning process

Markov Decision Process

Markov decision process (MDP), also called a stochastic dynamic program or stochastic control problem, is a model for sequential decision making when outcomes are uncertain. W

Reinforcement learning

Reinforcement learning is an interdisciplinary area of machine learning and optimal control that has, as main objective, finding an approximately optimal policy for MDPs where transition probabilities and rewards are unknown. W

Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

Thompson sampling

Thompson sampling is a heuristic for choosing actions that address the exploration-exploitation dilemma in the multi-armed bandit problem. It consists of choosing the action that maximizes the expected reward with respect to a randomly drawn belief. W

See [Russo et al. \(2018\)](#)

Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

$$\alpha_0 = 1, \beta_0 = 1 \quad 0.1331$$

$$\alpha_0 = 1, \beta_0 = 1 \quad 0.1282$$



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

1 $\alpha_1 = 2, \beta_1 = 1$ 0.7369

$\alpha_1 = 1, \beta_1 = 1$ 0.8081



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

1 $\alpha_2 = 2$, $\beta_2 = 1$ 0.8835

0 $\alpha_2 = 1$, $\beta_2 = 2$ 0.3092



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

1 1 $\alpha_3 = 3, \beta_3 = 1$ 0.9407

0 0 $\alpha_3 = 1, \beta_3 = 2$ 0.8079



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad 1 \quad \alpha_4 = 3, \beta_4 = 2 \quad 0.6529$$

$$0 \quad \alpha_4 = 1, \beta_4 = 2 \quad 0.5452$$



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad 1 \quad 0 \quad 1 \quad \alpha_5 = 4, \beta_5 = 2 \quad 0.5835$$

$$0 \quad \alpha_5 = 1, \beta_5 = 2 \quad 0.8632$$



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad 1 \ 0 \ 1 \quad \alpha_6 = 4, \beta_6 = 2 \quad 0.7858$$

$$0 \quad \boxed{1} \quad \alpha_6 = 2, \beta_6 = 2 \quad 0.4509$$



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

1 1 0 1 1 $\alpha_7 = 5, \beta_7 = 2$ 0.8871

0 1 $\alpha_7 = 2, \beta_7 = 2$ 0.1643



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 1 \quad \alpha_8 = 6, \beta_8 = 2 \quad 0.8052$$

$$0 \quad \quad \quad 1 \quad \quad \alpha_8 = 2, \beta_8 = 2 \quad 0.9383$$



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad 1 \ 0 \ 1 \quad 1 \ 1 \quad \alpha_9 = 6, \beta_9 = 2 \quad 0.5769$$

$$0 \quad 1 \quad 0 \quad \alpha_9 = 2, \beta_9 = 3 \quad 0.6047$$



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad 1 \ 0 \ 1 \quad 1 \ 1 \quad \alpha_{10} = 6, \beta_{10} = 2 \quad 0.692$$

$$0 \quad 1 \quad 0 \ 0 \quad \alpha_{10} = 2, \beta_{10} = 4 \quad 0.5244$$



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|-------------------|------------------|--------|
| 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | $\alpha_{15} = 9$ | $\beta_{15} = 4$ | 0.5462 |
| 0 | | 1 | 0 | 0 | | | | $\alpha_{15} = 2$ | $\beta_{15} = 4$ | 0.2837 |



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|--------------------|------------------|--------|
| 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | $\alpha_{20} = 11$ | $\beta_{20} = 6$ | 0.5201 |
| 0 | | 1 | | 0 | 0 | | 1 | | | | | $\alpha_{20} = 3$ | $\beta_{20} = 4$ | 0.2459 |



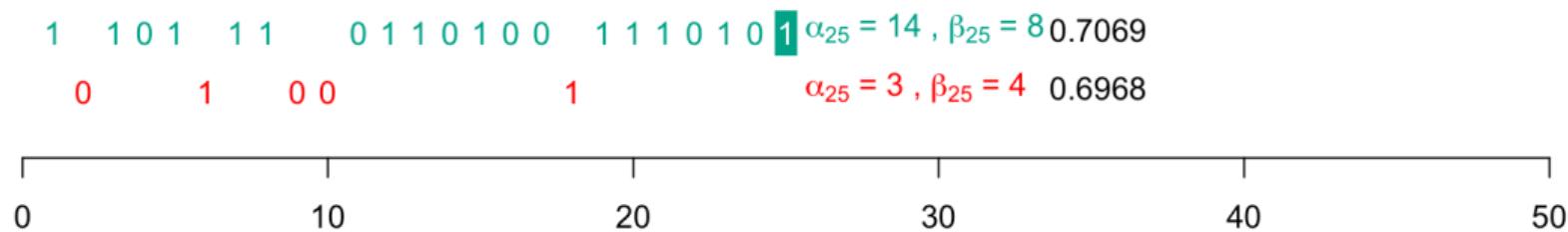
Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1,\dots,K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
 - ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$



Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$

$$1 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 1 \ 1 \ 1 \ 0 \ 1 \ 0 \ 1 \ 0 \ 1 \ 1 \ \boxed{1} \ \alpha_{30} = 17, \beta_{30} = 9 \ 0.7365$$

$$0 \ 1 \ 0 \ 0 \ 1 \ \alpha_{30} = 3, \beta_{30} = 5 \ 0.0937$$



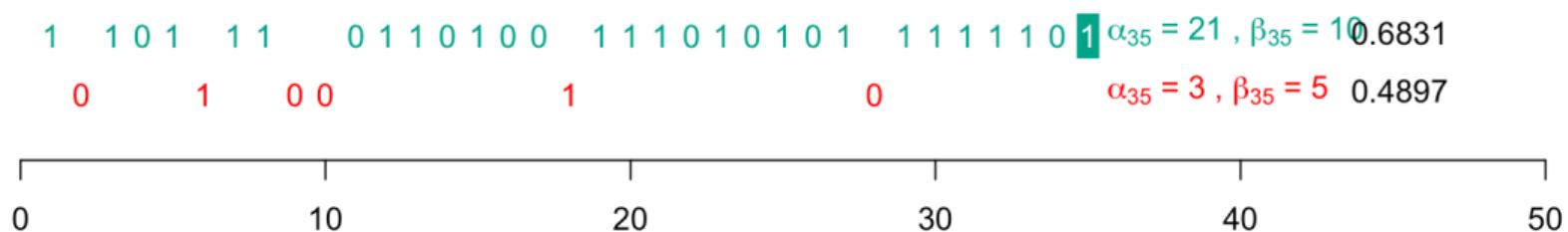
Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among K alternatives, that yield $\mathbf{X} = (X_1, \dots, X_K)$, $X_k \sim \mathcal{B}(\theta_k)$. Assume (prior) $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$. At time t , draw K Beta variables (independents) $B_k \sim \text{Beta}(\alpha_k, \beta_k)$, and select $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$.

Consider updating $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$,

- ▶ simulated data, i.i.d., $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d., $X_2 \sim \mathcal{B}(24\%)$



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



$$\alpha_0 = 1, \beta_0 = 1 \quad 0.4267$$

$$\alpha_0 = 1, \beta_0 = 1 \quad 0.8151$$



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



$$\alpha_1 = 1, \beta_1 = 1 \quad 0.4473$$

$$1 \quad \alpha_1 = 2, \beta_1 = 1 \quad 0.6376$$



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



$$\alpha_3 = 1, \beta_3 = 1 \quad 0.5552$$

$$1 \ 0 \ 0 \quad \alpha_3 = 2, \beta_3 = 3 \quad 0.8841$$



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



$$\alpha_4 = 1, \beta_4 = 1 \quad 0.1631$$

$$1 \ 0 \ 0 \ 1 \quad \alpha_4 = 3, \beta_4 = 3 \quad 0.5958$$



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

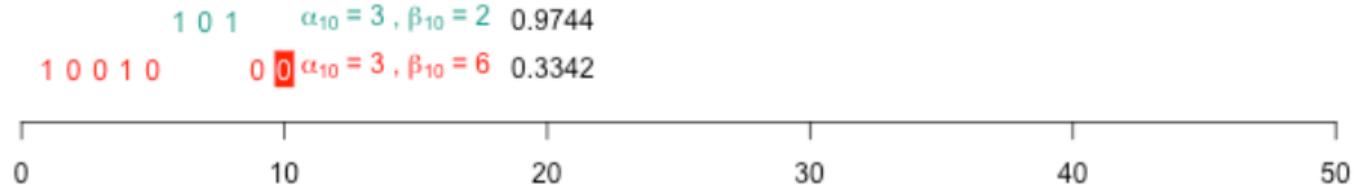
- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

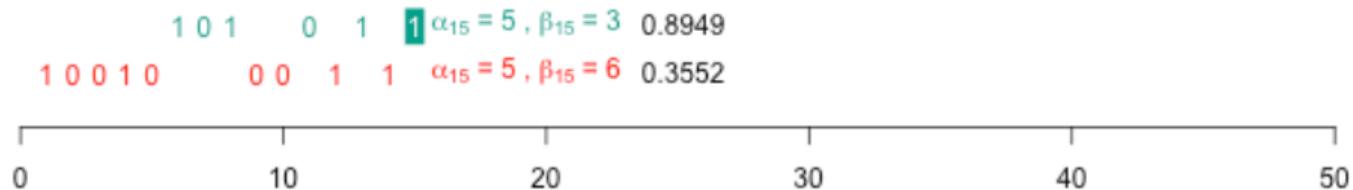
- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

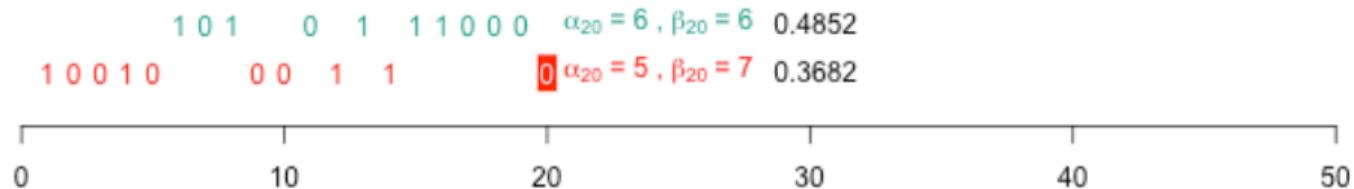
- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

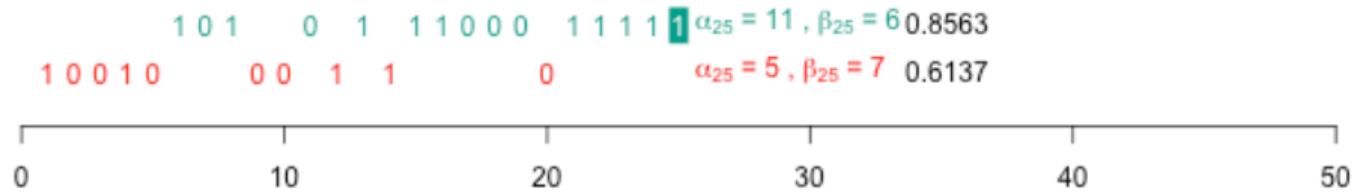
- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



Bayesianism as a learning process

We can use that approach in the context of Monty Hall

- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



References

- Avraham, R. (2017). Discrimination and insurance. In Lippert-Rasmussen, K., editor, *Handbook of the Ethics of Discrimination*, pages 335–347. Routledge.
- Bailey, R. A. and Simon, L. J. (1960). Two studies in automobile insurance ratemaking. *ASTIN Bulletin: The Journal of the IAA*, 1(4):192–217.
- Ben-Tal, A., El Ghaoui, L., and Nemirovski, A. (2009). *Robust optimization*, volume 28. Princeton university press.
- Bengio, Y., Goodfellow, I., and Courville, A. (2017). *Deep learning*, volume 1. MIT press.
- Beutel, A., Chen, J., Zhao, Z., and Chi, E. H. (2017). Data decisions and theoretical implications when adversarially learning fair representations. *arXiv*, 1707.00075.
- Blanpain, N. (2018). L'espérance de vie par niveau de vie-méthode et principaux résultats. *INSEE Document de Travail*, F1801.
- Box, G. E., Luceño, A., and del Carmen Paniagua-Quinones, M. (2011). *Statistical control by monitoring and adjustment*, volume 700. John Wiley & Sons.
- Breiman, L., Friedman, J., Olshen, R., and Stone, C. (1984). *Classification and regression trees*. Taylor & Francis.

freakonometrics

freakonometrics.hypotheses.org

References

- 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.
- Cesa-Bianchi, N. and Lugosi, G. (2006). *Prediction, learning, and games*. Cambridge university press.
- Charpentier, A. (2014). Mesures de risque. In Drosbeke, J.-J. and Saporta, G., editors, *Approches statistiques du risque*. Éditions Technip.
- Condorcet, J. A. N. (1785). *Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix*. Imprimerie royale.
- Danskin, J. M. (1967). *The theory of max-min and its application to weapons allocation problems*. Springer.
- De Leeuw, J., Hornik, K., and Mair, P. (2010). Isotone optimization in r: pool-adjacent-violators algorithm (pava) and active set methods. *Journal of statistical software*, 32:1–24.
- 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.

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities)

BY-NC 4.0 385 / 385

References

- Denuit, M. and Charpentier, A. (2004). *Mathématiques de l'assurance non-vie: Tome I Principes fondamentaux de théorie du risque*. Economica.
- Eckart, C. and Young, G. (1936). The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218.
- 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.
- Friedman, J., Hastie, T., Tibshirani, R., et al. (2001). *The elements of statistical learning*, volume 1. Springer series in statistics New York.
- Galton, F. (1907). Vox populi (the wisdom of crowds). *Nature*, 75(7):450–451.
- Gandy, O. H. (2016). *Coming to terms with chance: Engaging rational discrimination and cumulative disadvantage*. Routledge.
- Glenn, B. J. (2003). Postmodernism: the basis of insurance. *Risk Management and Insurance Review*, 6(2):131–143.
- Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. (2016). *Deep learning*. MIT Press.
- Harcourt, B. E. (2015). Risk as a proxy for race: The dangers of risk assessment. *Federal Sentencing Reporter*, 27(4):237–243.

References

- Hoerl, A. E. and Kennard, R. W. (1970a). Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics*, 12(1):55.
- Hoerl, A. E. and Kennard, R. W. (1970b). 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.
- 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.
- James, G., Witten, D., Hastie, T., Tibshirani, R., et al. (2013). *An introduction to statistical learning*, volume 112. Springer.
- Kaufman, L. and Rousseeuw, P. J. (2009). *Finding groups in data: an introduction to cluster analysis*. John Wiley & Sons.
- Kearns, M. (1988). Thoughts on hypothesis boosting. *Unpublished manuscript*, 45:105.
- Kearns, M. and Valiant, L. (1989). Cryptographic limitations on learning boolean formulae and finite automata. *Proceedings of the twenty-first annual ACM symposium on Theory of computing*, 41(1):67–95.

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities) BY-NC 4.0 385 / 385

References

- Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. *arXiv*, 1312.6114.
- 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.
- 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.
- McCullagh, P. and Nelder, J. (1989). *Generalized linear models*. Chapman & Hall.
- Mienye, I. D. and Sun, Y. (2022). A survey of ensemble learning: Concepts, algorithms, applications, and prospects. *Ieee Access*, 10:99129–99149.
- Mirsky, L. (1960). Symmetric gauge functions and unitarily invariant norms. *The quarterly journal of mathematics*, 11(1):50–59.
- 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*.
- Neumann, J. v. and Morgenstern, O. (1947). *Theory of games and economic behavior*. Princeton University Press.
- 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.

References

- Russo, D. J., Van Roy, B., Kazerouni, A., Osband, I., Wen, Z., et al. (2018). A tutorial on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96.
- Schervish, M. J. and DeGroot, M. H. (2014). *Probability and statistics*, volume 563. Pearson Education London, UK:.
- 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.
- Struyck, N. (1912). *Les oeuvres de Nicolas Struyck (1687-1769): qui se rapportent au calcul des chances, à la statistique général, à 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.
- Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294.
- Thompson, W. R. (1935). On the theory of apportionment. *American Journal of Mathematics*, 57(2):450–456.

References

- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288.
- Tukey, J. W. (1961). Curves as parameters, and touch estimation. In Neyman, J., editor, *Proceedings of the 4th Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, pages 681–694. University of California Press, University of California Press.
- Tweedie, M. C. K. (1984). An index which distinguishes between some important exponential families. *Statistics: applications and new directions (Calcutta, 1981)*, pages 579–604.
- Vapnik, V. N., Vapnik, V., et al. (1998). Statistical learning theory.
- Vito, E. D., Rosasco, L., Caponnetto, A., Piana, M., and Verri, A. (2004). Some properties of regularized kernel methods. *Journal of Machine Learning Research*, 5(Oct):1363–1390.
- von Neumann, J. (1928). Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320.
- Wadsworth, C., Vera, F., and Piech, C. (2018). Achieving fairness through adversarial learning: an application to recidivism prediction. *arXiv*, 1807.00199.
- Wahba, G. (1999). Support vector machines, reproducing kernel hilbert spaces and the randomized gacv. *Advances in Kernel Methods-Support Vector Learning*, 6:69–87.
- Wahba, G., Lin, X., Gao, F., Xiang, D., Klein, R., and Klein, B. (1998). The bias-variance tradeoff and the randomized gacv. *Advances in Neural Information Processing Systems*, 11.

References

- Wortham, L. (1986). The economics of insurance classification: The sound of one invisible hand clapping. *Ohio State Law Journal*, 47:835.
- 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.
- Zhang, B. H., Lemoine, B., and Mitchell, M. (2018). Mitigating unwanted biases with adversarial learning. *arXiv*, 1801.07593:335–340.