

Machine Learning for Insurers and Actuaries

Arthur Charpentier

2025



Learning, with a mathematical perspective

(lecture 1)

Probabilities and Statistics

- Probability and statistics are foundational for understanding and modeling uncertainty in data.
- In machine learning, we often work with data that is uncertain or incomplete.
 - **Probability**: Quantifies uncertainty and models random events (e.g., likelihood of events, predictions).
 - **Statistics**: Analyzes data to infer properties of a population from a sample (e.g., parameter estimation, hypothesis testing).
- Key tools in ML:
 - Probability distributions.
 - Expectation and variance.
 - Maximum likelihood estimation (MLE).
 - Bayesian inference.
- A **probability distribution** describes how the values of a random variable are distributed.

Probabilities and Statistics

- Common probability distributions used in ML:
 - Gaussian Distribution
 - Bernoulli Distribution: Used for binary outcomes.
 - Poisson Distribution: Models the number of events in a fixed interval.
 - Multinomial Distribution: Generalization of the Bernoulli distribution to multiple outcomes.
- Application in Machine Learning:
 - Gaussian Naive Bayes: Assumes features follow a Gaussian distribution.
 - Generative Models: Many generative models use probability distributions to model data (e.g., Variational Autoencoders).
- **Expectation** (Mean): The expected value of a random variable, representing the central tendency.

$$\mathbb{E}[X] = \sum_x x \cdot \mathbb{P}[X = x] = \mathbf{p}^\top \mathbf{x} = \langle \vec{\mathbf{p}}, \vec{\mathbf{x}} \rangle,$$

Probabilities and Statistics

- **Variance:** Measures the spread of the distribution, i.e., how far values are from the mean.

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$$

- **Covariance:** Measures the joint variability of two random variables.

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

- **Application in Machine Learning:**

- **Gaussian Distribution / Linear Regression:** Uses covariance to understand the relationship between variables and compute coefficients.
- **Gaussian Distribution / Principal Component Analysis (PCA):** Uses the covariance matrix to identify directions of maximum variance.

- **Maximum Likelihood Estimation** is a method for estimating the parameters of a statistical model.

Probabilities and Statistics

- The likelihood function represents the probability of the observed data given the model parameters θ :

$$\mathcal{L}(\theta) = \mathbb{P}(\mathcal{D} | \theta)$$

- MLE** finds the parameter values $\hat{\theta}$ that maximize the likelihood (or the log likelihood):

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}}\{\mathcal{L}(\theta)\} \text{ or } \hat{\theta} = \underset{\theta}{\operatorname{argmax}}\{\log \mathcal{L}(\theta)\}$$

- Application in Machine Learning:**

- Gaussian Distribution
- Logistic Regression:** Uses MLE to estimate model parameters by maximizing the likelihood of the observed data.
- Gaussian Distribution
- Gaussian Mixture Models (GMM):** Uses MLE to estimate the parameters of Gaussian distributions.

Probabilities and Statistics

- **Bayesian Inference** updates the probability estimate for a hypothesis as more evidence becomes available.
- The Bayes' Theorem is central to Bayesian inference:

$$\mathbb{P}(H | \mathcal{D}) = \frac{\mathbb{P}(\mathcal{D} | H) \cdot \mathbb{P}(H)}{\mathbb{P}(\mathcal{D})}$$

where:

- **Gaussian Distribution**
- $\mathbb{P}(H | \mathcal{D})$ is the posterior (the updated belief about H after seeing data \mathcal{D}).
- **Gaussian Distribution**
- $\mathbb{P}(\mathcal{D} | H)$ is the likelihood (how likely the data is given the hypothesis).
- **Gaussian Distribution**
- $\mathbb{P}(H)$ is the prior (the belief about H before seeing the data).

- **Application in Machine Learning:**

Probabilities and Statistics

- **Naive Bayes Classifier:** Applies Bayes' Theorem to classify data.
- **Bayesian Neural Networks:** Use Bayesian methods to update the weights of the neural network.
- **Statistical Inference** is the process of drawing conclusions from data, typically based on probability models.
- **Hypothesis Testing** involves testing a null hypothesis (H_0) against an alternative hypothesis (H_1):

$$\text{p-value} = \mathbb{P}(\mathcal{D} \mid H_0)$$

If the p-value is small (e.g., < 0.05), we reject H_0 .

- **Application in Machine Learning:**
 - **Model Evaluation:** Statistical tests are used to evaluate and compare the performance of machine learning models (e.g., A/B testing).
 - **Overfitting Check:** Hypothesis testing can be used to assess whether a model is too complex and overfitting the data.

Probabilities and Statistics

- Machine Learning without Probabilistic Assumptions,
 - Many ML algorithms, such as decision trees, support vector machines (SVMs), or neural networks, can be applied without explicit probabilistic assumptions.
 - These algorithms focus on finding patterns in data through optimization, without assuming any underlying distribution.
- **Why is probability important?**
 - Probabilities help quantify uncertainty in predictions.
 - Provide mathematical guarantees about the performance of algorithms.
 - Essential for understanding generalization and overfitting
- Without probabilistic assumptions, we cannot guarantee the generalization of our models to unseen data.

Probabilities and Statistics

- **Generalization** refers to how well a model trained on a sample of data performs on new, unseen examples.
- **Probability** allows us to express and prove mathematical guarantees, such as:
 - **PAC (Probably Approximately Correct) Learning**: Ensures that a model will be approximately correct with high probability, provided certain conditions are met.
- PAC Learning formalizes the relationship between:
 - The error of the hypothesis (how far off the model's predictions are from the true distribution).
 - The confidence in the model (the probability that the model's error is small).
 - The number of training examples required to achieve these guarantees.
- PAC Learning is a formal framework that allows us to assess how well a hypothesis learned from a finite sample will generalize to unseen data.
- The key components of PAC learning:

Probabilities and Statistics

- A learning algorithm that outputs a hypothesis H .
- A target distribution \mathcal{D} over examples.
- A hypothesis H is said to be probably approximately correct if:

$$\mathbb{P}(\text{error}(H) \leq \varepsilon) \geq 1 - \delta$$

where:

- ε is the allowable error.
- δ is the probability of failure (the confidence level).

- This implies that with high probability (at least $1 - \delta$), the hypothesis h will have an error less than ε .
- Sample Complexity refers to the number of training examples needed to achieve a PAC guarantee.
- The number of examples required depends on:
 - The complexity of the hypothesis class (the set of models we consider).

Probabilities and Statistics

- The desired error ε and confidence δ .
- For a given hypothesis class H , the sample complexity can be bounded as:

$$\text{complexity}(\varepsilon, \delta) = O\left(\frac{1}{\varepsilon^2} \log\left(\frac{1}{\delta}\right)\right)$$

- This means that to guarantee with high probability that the error is less than ε , the required sample size increases as the allowed error ε decreases.
- Without probability, we cannot assess the risk of overfitting (i.e., fitting the model too closely to training data).
- Overfitting occurs when a model captures noise or irrelevant patterns in the training data.
- Probabilistic models help in controlling overfitting by:
 - Estimating uncertainty.
 - Regularizing models to avoid fitting to noise.

Probabilities and Statistics

- Using Bayesian methods to incorporate prior knowledge and regularize the model.
- **Example:** In Bayesian Neural Networks, the model's weights are treated as distributions, which allows us to quantify uncertainty in predictions and improve generalization.
- Machine Learning without Probabilistic Assumptions:
 - Many models, such as decision trees and SVMs, can be learned without making probabilistic assumptions.
 - These models often rely on optimization techniques to fit data and are effective for many practical applications.
- However, Probability Provides Mathematical Guarantees:
 - Probability gives us tools to assess generalization and make predictions about model performance on unseen data.

Probabilities and Statistics

- PAC learning formalizes the trade-off between model complexity, sample size, and generalization error.
- Probabilistic assumptions allow us to quantify uncertainty, handle noise, and ensure the reliability of the model.
- Ultimately, while deterministic models can work in practice, probabilistic reasoning is essential for understanding the theoretical guarantees of machine learning algorithms.

Probabilities and random variables

“Probability is the most important concept in modern science, especially as nobody has the slightest notion what it means,” Russell (1929), quoted in Bell (1945)

Probability and statistics rely on the concept of probability spaces, $(\Omega, \mathcal{F}, \mathbb{P})$,

- Ω (or S in some textbooks) is the sample space, the set of all possible outcomes
- \mathcal{F} a set of events on Ω , $A \in \mathcal{F}$ is an “event”
- \mathbb{P} is a function $\mathcal{F} \rightarrow [0, 1]$, called probability, satisfying some properties

e.g. $\mathbb{P}(\Omega) = 1$; for disjoint events, an additivity property: $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$; a subset property, if $A \subset B$, $\mathbb{P}(A) \leq \mathbb{P}(B)$, as in Cardano (1564) or Bernoulli (1713), or for multiple disjoint events as in Kolmogorov (1933), A_1, \dots, A_n, \dots ,

$$\mathbb{P}(A_1 \cup \dots \cup A_n \cup \dots) = \mathbb{P}(A_1) + \dots + \mathbb{P}(A_n) + \dots$$

inspired by Lebesgue (1918), etc. In this (mathematical) framework, we can finally define random variables

- X is a function $\Omega \rightarrow \mathbb{R}$ or more generally $\Omega \rightarrow \mathcal{X}$.

Probabilities and random variables

Definition 1.1: Cumulative Distribution Function

Let Y denote a (real-valued) random variable, then its cumulative distribution function is $F(y) = \mathbb{P}[Y \leq y]$, for $y \in \mathbb{R}$. We write $Y \sim F$.

Definition 1.2: Density

Absolutely continuous distributions have a density (with respect to Lebesgue measure),

$$F(y) = \int_{-\infty}^y f(z)dz \text{ or } f(y) = \frac{df(z)}{dz} \Big|_{z=y}$$

Thus, $\mathbb{P}[Y \in [a, b]] = F(b) - F(a) = \int_a^b f(z)dz$.

Probabilities and random variables

We have formal objects, mathematically well defined, but in a context of modeling does one have a univocal sense of interpretation of the result of the calculation?

See “Is the probability inherent to the event, or to our judgment?,” in Martin (2009)

There are many philosophical paradoxes when we talk about probability (and chance), e.g. I throw a coin, which falls back, out of my sight

- $\mathbb{P}(X = \text{heads}) = \mathbb{P}(X = \text{tails}) = 1/2$?
- $\mathbb{P}(X = \text{heads}) = 1$ or $\mathbb{P}(X = \text{tails}) = 1$?

Or in a legal context, “Look, the guy either did it or he didn’t do it. If he did then he is 100% guilty and if he didn’t then he is 0% guilty; so giving the chances of guilt as a probability somewhere in between makes no sense and has no place in the law,” quoted in Fenton and Neil (2018)

See also Hájek (2002) on the philosophical approach of “probability”.

Probabilities and random variables

As said by Martin (2009), “To attribute an objective meaning to the probability that an event will occur is to admit that this event is not necessary, in other words, that it is not completely determined,” [...] “If we suppose an integral and universal determinism, the probability can only receive a subjective meaning, and the probability depends on our knowledge and our ignorance”

A lot of importance is attributed to this supposedly objective probability \mathbb{P} .

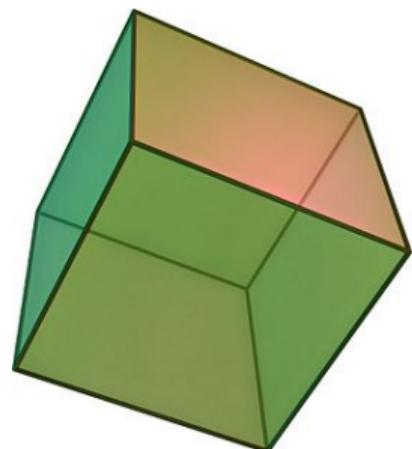
The (mathematical) probability was not born as a well defined concept within the framework of a mathematical formalism mathematical formalism, but as a tool to quantify and control situations of uncertainty, applied to the measurement of the probability of life mortality tables (for the calculation of life annuities), the calculation of the risks of error (in of error (in measurement operations), the study of the probability of testimonies and judgments, etc.

Probabilities and random variables

Cournot (1843) thus distinguished a **objective meaning** of the probability (as measure of the physical possibility of realization of a random event) and a **subjective meaning** (the probability being a judgement made on an event, this judgement being linked to the ignorance of judgment being linked to the ignorance of the conditions of the realization of the event).

Note: a probability not defined in terms of frequency can still have an objective meaning: :

There is no need to repeat throws of dice to affirm that (with a perfectly balanced die) the probability of obtaining 6 at the time of a throw is equal to $1/6$ (by symmetry of the cube)



Probabilities and random variables

But very often, the “physical” probabilities receive an objective value only posterior on the basis of the law of large numbers, the empirical frequency converge towards the probability (frequentist theory of probabilities)

$$\underbrace{\frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i \in A)}_{\text{(empirical) frequency}} \xrightarrow{\text{a.s.}} \underbrace{\mathbb{P}(X \in A)}_{\text{probability}} \text{ as } n \rightarrow \infty$$

(in some textbooks, there is a confusion between "probability" and "frequency")

Law of large numbers : $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\text{a.s.}} \mathbb{E}(X)$ as $n \rightarrow \infty$ or $\frac{1}{n} \sum_{i=1}^n X_i \approx \mathbb{E}(X)$

Probabilities and random variables

But this approach is unable to make sense of the probability of a “**single singular event**”, as noted by [von Mises \(1928, 1939\)](#).

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

See also “**Single-Case Interpretation of Probability**,” in [Reichenbach \(1971\)](#)

Probabilities and random variables

“What must be asked is whether an interpretation is adequate to account for the use of probabilities as degrees of reliability or as instruments suitable for an evaluation of predictions” [...] “The degree of probability, therefore, will be of no use after the truth about the occurrence of the event is known: probability is used as a substitute for truth so long as the truth is unknown. If the event is to happen in the future, the degree of its probability qualifies the reliability of a prediction,” Reichenbach (1971).

THE THEORY OF PROBABILITY

An Inquiry into the Logical and Mathematical Foundations of the Calculus of Probability

By HANS REICHENBACH

PROFESSOR OF PHILOSOPHY IN THE UNIVERSITY OF CALIFORNIA AT LOS ANGELES

UNIVERSITY OF CALIFORNIA PRESS

BERKELEY AND LOS ANGELES • 1949

§ 71. Attempts at a Single-Case Interpretation of Probability

After the discussion of the frequency meaning of probability, the investigation must turn to linguistic forms in which the concept of probability refers to an individual event. It is on this ground that the frequency interpretation has been questioned. Some logicians have argued that such usage is based on a different concept of probability, which is not reducible to frequencies. Is the existence of two disparate concepts of probability an inescapable consequence of the usage of language?

The first interpretation of the probability of single events is the *degree of expectation* with which an event is anticipated. The feeling of expectation certainly represents a psychological factor the existence of which is indisputable; it even shows degrees of intensity corresponding to the degrees of probability. Difficulty, however, arises from the fact that the degree of expectation varies from person to person and depends on more factors than the degree of the probability of the event to which the expectation refers. Apart from the probability of an event, emotional associations will influence the feeling of expectation. If it is a desirable event, as, for instance, the passing of an examination, optimistic persons will anticipate it with too-certain expectations, whereas pessimistic persons will think of it in terms of too-uncertain expectations.

Probabilities and random variables

“If we are asked to find the probability holding for an individual future event, we must first incorporate the case in a suitable reference class. An individual thing or event may be incorporated in many reference classes, from which different probabilities will result. This ambiguity has been called the problem of the reference class” [...] “I regard the statement about the probability of the single case, not as having a meaning of its own, but as representing an elliptic mode of speech. In order to acquire meaning, the statement must be translated into a statement about a frequency in a sequence of repeated occurrences. The statement concerning the probability of the single case thus is given a fictitious meaning, constructed by a transfer of meaning from the general to the particular case,” Reichenbach (1971).

THE THEORY OF PROBABILITY

An Inquiry into the Logical and Mathematical Foundations of the Calculus of Probability

By HANS REICHENBACH

PROFESSOR OF PHILOSOPHY IN THE UNIVERSITY OF CALIFORNIA AT LOS ANGELES

UNIVERSITY OF CALIFORNIA PRESS
BERKELEY AND LOS ANGELES • 1949

§ 72. The Frequency Interpretation of the Probability of the Single Case

The analysis of meaning has suffered from too close an attachment to psychological considerations. The meaning of a sentence has been identified with the mental images associated with the utterance of the sentence. Such conception leads to meanings varying from person to person; and it will not help to find the meaning that a man would adopt if he had a clear insight into the implications of his words. Logic is interested not in what a man means but in what he should mean, that is, in the meaning that, if assumed for his words, would make his words compatible with his actions.

If we are asked to find the probability holding for an individual future event, we must first incorporate the case in a suitable reference class. An individual thing or event may be incorporated in many reference classes, from which different probabilities will result. This ambiguity has been called the problem of the reference class. Assume that a case of illness can be characterized by its inclusion in the class of cases of tuberculosis. If additional information is obtained from an X-ray, the same case may be incorporated in the class of serious cases of tuberculosis. Depending on the classification, different probabilities will result for the prospective issue of the illness.

We then proceed by considering the narrowest class for which reliable statistics can be compiled. If we are confronted by two overlapping classes, we shall choose their common class. Thus, if a man is 21 years old and has tuberculosis, we shall regard the class of persons of 21 who have tuberculosis. Classes that are known to be irrelevant for the statistical result may be disregarded. A class C is irrelevant with respect to the reference class A and the attribute class B if the transition to the common class $A \cdot C$ does not change the probability, that is, if $P(A \cdot C, B) = P(A, B)$. For instance, the class of persons having the same initials is irrelevant for the life expectation of a person.

Probabilities and random variables

For Popper (1959), probabilities correspond to physical dispositions ("propensities") inherent to the system. This propensity has a physical existence, but it is not directly observable.

The frequencies of occurrence are manifestations of these propensities. In the contrary case, it is nevertheless possible to estimate the probability of realization of the singular event, by considering this one as measured not by an "actual" frequency, but by a "potential" (or "virtual") frequency.

Finally, when an individual makes a judgment, the degree of credibility or belief that he or she gives it depends on the knowledge that the individual has (Pettigrew (2016)). depends on the knowledge that this individual has (Pettigrew (2016)). This degree of belief will be associated with a probability, which will then only have a subjective meaning.

freakonometrics

freakonometrics.hypotheses.org

Probabilities and random variables

“The probability of a diagnosis, a testimony, etc., does not measure the conformity of this judgment to reality, but the degree to which one can hypothesize this conformity. This conformity can be hypothesized”, [Martin \(2009\)](#).

This subjectivity raises concerns about their use, especially in criminal matters,

“Sometimes the ‘balance of probability’ standard is expressed mathematically as ‘50+% probability’, but this can carry with it a danger of pseudo-mathematics, as the argument in this case demonstrated. When judging whether a case for believing that an event was caused in a particular way is stronger than the case for not so believing, the process is not scientific (although it may obviously include evaluation of scientific evidence) and to express the probability of some event having happened in percentage terms is illusory, [*Nulty & Ors v Milton Keynes Borough Council*](#) cited in [Hunt and Mostyn \(2020\)](#).

See also [Jonakait \(1983\)](#), [Saini \(2011\)](#) or [Fenton et al. \(2016\)](#).

Statistics, Frequentist Approach

Strong law of large numbers (also called Kolmogorov's law), see [Loève \(1977\)](#)

Proposition 1.1: Law of Large Numbers

Consider an infinite collection of i.i.d. random variables $X, X_1, X_2, \dots, X_n, \dots$ in a probabilistic space $(\Omega, \mathcal{F}, \mathbb{P})$, then

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

Likelihood

- Fisher (1921) introduced what is today called a “likelihood interval”
- Fisher (1922) introduced the term “method of maximum likelihood”

“In 1922, I proposed the term “likelihood,” in view of the fact that, with respect to [the parameter], it is not a probability, and does not obey the laws of probability, while at the same time it bears to the problem of rational choice among the possible values of [the parameter] a relation similar to that which probability bears to the problem of predicting events in games of chance... Whereas, however, in relation to psychological judgment, likelihood has some resemblance to probability, the two concepts are wholly distinct (...) I stress this because in spite of the emphasis that I have always laid upon the difference between probability and likelihood there is still a tendency to treat likelihood as though it were a sort of probability.”

Likelihood function

A likelihood function (often simply called the likelihood) measures how well a statistical model explains observed data by calculating the probability of seeing that data under different parameter values of the model. The likelihood function, parameterized by a (possibly multivariate) parameter θ , is usually defined differently for discrete and continuous probability distributions. Given a probability density or mass function $x \mapsto f(x | \theta)$, where x is a realization of the random variable X , the likelihood function is $\theta \mapsto f(x | \theta)$, often written $\mathcal{L}(\theta | x)$. W

freakonometrics.hypotheses.org

Definition 1.3: Maximum Likelihood (1)

Let $\mathbf{y} = (y_1, \dots, y_n)$ be a sample from i.i.d. variables with distribution f_θ . The likelihood is

$$\mathcal{L}(\theta|\mathbf{y}) = \prod_{i=1}^n f_\theta(y_i)$$

And the maximum likelihood is

$$\hat{\theta}(\mathbf{y}) = \underset{\theta \in \Theta}{\operatorname{argmax}} \{\mathcal{L}(\theta|\mathbf{y})\} \text{ and } \hat{\theta}(\mathbf{Y}) = \underset{\theta \in \Theta}{\operatorname{argmax}} \{\mathcal{L}(\theta|\mathbf{Y})\}$$

Definition 1.4: Maximum Likelihood (2)

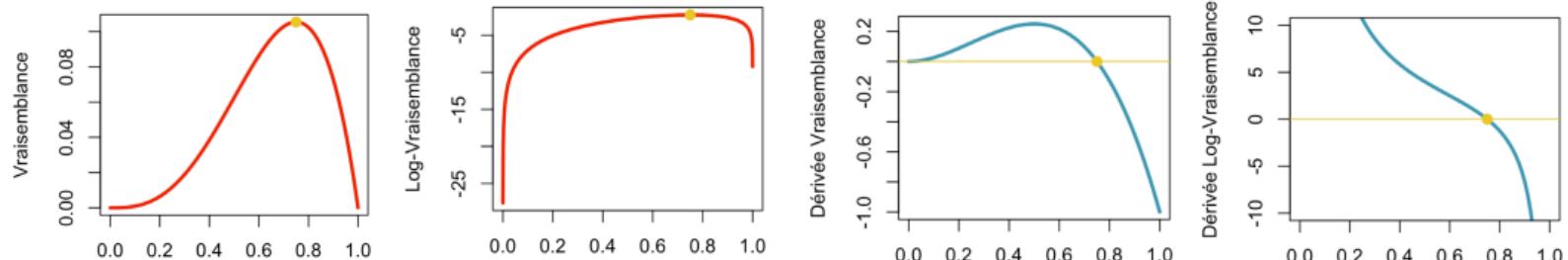
Let $\mathbf{y} = (y_1, \dots, y_n)$ be a sample from i.i.d. variables with distribution f_θ . The log-likelihood is

$$\log \mathcal{L}(\theta | \mathbf{y}) = \sum_{i=1}^n \log[f_\theta(y_i)];$$

And the maximum likelihood is

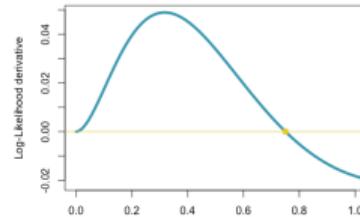
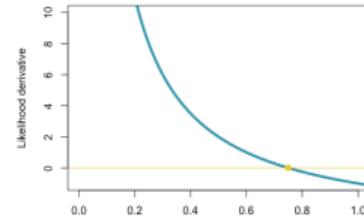
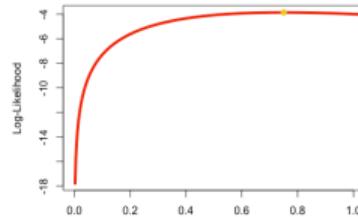
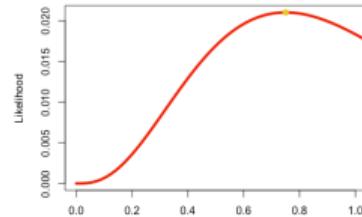
$$\hat{\theta}(\mathbf{y}) = \underset{\theta \in \Theta}{\operatorname{argmax}} \{ \log \mathcal{L}(\theta | \mathbf{y}) \} \text{ and } \hat{\theta}(\mathbf{Y}) = \underset{\theta \in \Theta}{\operatorname{argmax}} \{ \log \mathcal{L}(\theta | \mathbf{Y}) \}$$

Likelihood



- $\theta \mapsto \prod_{i=1}^n f(y_i; \theta) = \prod_{i=1}^n \theta^{y_i} (1-\theta)^{1-y_i} = \theta^3(1-\theta)$, $\mathbf{y} = \{0, 1, 1, 1\}$, $\mathbf{Y}_i \sim \mathcal{B}(\theta)$.
- $\theta \mapsto \sum_{i=1}^n \log f(y_i; \theta) = \sum_{i=1}^n y_i \log(\theta) + (1-y_i) \log(1-\theta) = 3 \log \theta + \log(1-\theta)$
- $\theta \mapsto \frac{\partial}{\partial \theta} \prod_{i=1}^n f(y_i; \theta) = \dots = 3\theta^2 - 4\theta^3$
- $\theta \mapsto \frac{\partial}{\partial \theta} \sum_{i=1}^n \log f(y_i; \theta) = \dots = \frac{3}{\theta} - \frac{1}{1-\theta}$ (score function)

Likelihood



- $\theta \mapsto \prod_{i=1}^n f(y_i; \theta) = \prod_{i=1}^n e^{-\theta} \frac{\theta^{y_i}}{y_i!} = e^{-4\theta} \theta^3$, where $\mathbf{y} = \{0, 1, 1, 1\}$, $Y_i \sim \text{Poisson}(\theta)$.
- $\theta \mapsto \sum_{i=1}^n \log f(y_i; \theta) = \sum_{i=1}^n -\theta + y_i \log \theta - \log(y_i!)$ $= -4\theta + 3 \log \theta$
- $\theta \mapsto \frac{\partial}{\partial \theta} \prod_{i=1}^n f(y_i; \theta) = \dots = -4 + \frac{3}{\theta}$
- $\theta \mapsto \frac{\partial}{\partial \theta} \sum_{i=1}^n \log f(y_i; \theta) = \dots = e^{-4\theta} (3\theta^2 - 4\theta^3)$ (score function)

Likelihood, $\mathcal{B}(p)$ Binary Sample

```
1 > y = c(1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 1, 0, 1, 0)
```

- what theory (or mathematics) tells us

$$\mathcal{L}(p; \mathbf{y}) = \prod_{i=1}^n f(y_i; p) = \prod_{i=1}^n p^{y_i} (1-p)^{1-y_i} = p^{s_n} (1-p)^{n-s_n}, \quad s_n = \sum_{i=1}^n y_i$$

$$\log \mathcal{L}(p; \mathbf{x}) = s_n \log(p) + (n - s_n) \log(1 - p)$$

$$\frac{\partial}{\partial p} \log \mathcal{L}(p; \mathbf{y}) = \frac{\partial}{\partial p} s_n \log(p) + (n - s_n) \log(1 - p) = \frac{s_n}{p} - \frac{n - s_n}{1 - p}$$

$$\left. \frac{\partial}{\partial p} \log \mathcal{L}(p; \mathbf{y}) \right|_{p=\hat{p}} = 0 \quad \text{if and only if } \frac{s_n}{\hat{p}} = \frac{n - s_n}{1 - \hat{p}}, \text{ i.e. } \hat{p} = \frac{s_n}{n} = \bar{y}$$

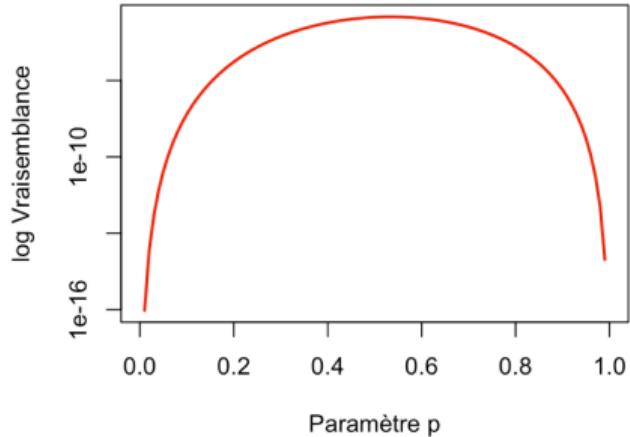
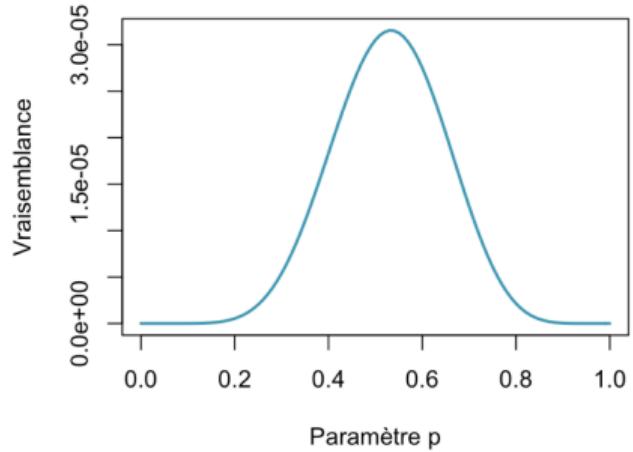
Likelihood, $\mathcal{B}(p)$ Binary Sample

- what a computer tells us ?

One can plot the likelihood function $p \mapsto \mathcal{L}(p; \mathbf{x})$, or log-likelihood

```
1 > n = 15
2 > y = c(1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0)
3 > likelihood = function(p) prod(dbinom(y, size = 1, prob = p))
4 > vect_p = seq(0,1, by=0.01)
5 > plot(vect_p,Vectorize(likelihood)(vect_p))
```

Likelihood, $\mathcal{B}(p)$ Binary Sample



Likelihood, $\mathcal{B}(p)$ Binary Sample

- what a computer tells us ?

We can compute (numerically) the maximum of $p \mapsto \mathcal{L}(p; \mathbf{x})$

```
1 > optim(par = .5,fn = function(z) -likelihood(z))
2 $par
3 [1] 0.5333252
4
5 $value
6 [1] -3.155276e-05
```

(from the maths, we know that actually $\hat{p} = \bar{x}$)

```
1 > mean(y)
2 [1] 0.5333333
```

Likelihood, $\mathcal{B}(p)$ Binary Sample

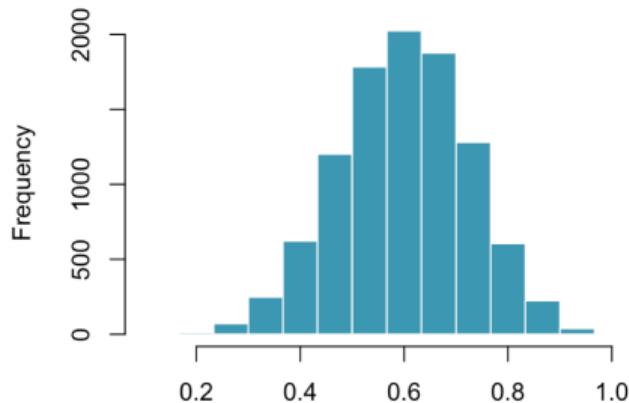
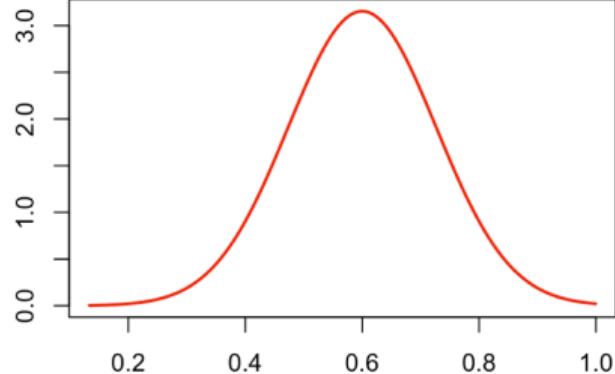
- what theory (or mathematics) tells us
since $\hat{p}(\mathbf{x}) = \bar{y}$, we can use the law of large numbers,

$$Z_n = \sqrt{n} \frac{\hat{p}(\mathbf{Y}) - p}{\sqrt{p(1-p)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$$

but since $n = 15$ this Gaussian assumption is maybe not valid...
If $p = 60\%$ the approximated distribution of $\hat{p}(\mathbf{Y})$ would be

```
1 > u = seq(2/15, 1, by=.001)
2 > plot(u, dnorm(u,.6,sqrt(.4*.6/15)))
```

Likelihood, $\mathcal{B}(p)$ Binary Sample

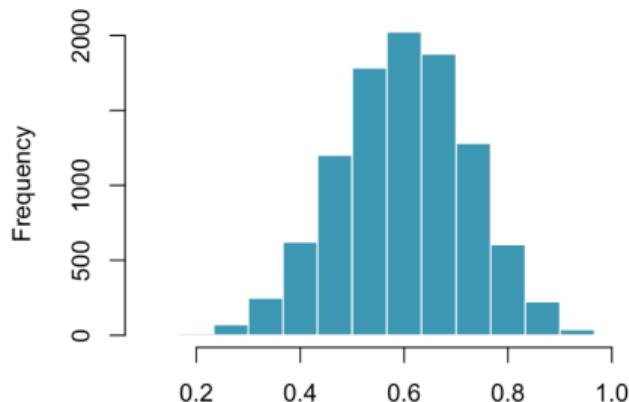
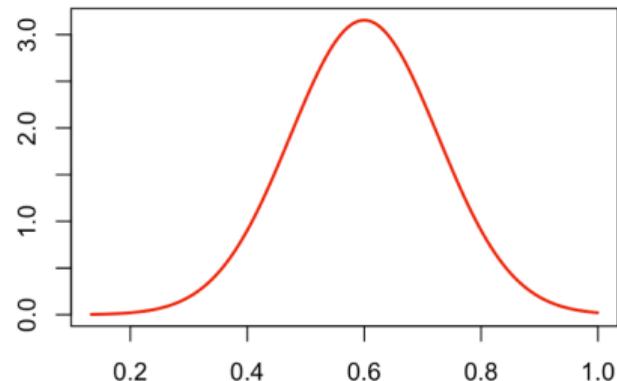


Likelihood, case $\mathcal{B}(p)$

- what a computer tells us ? (suppose that $\theta = 60\%$)

```
1 > theta=rep(NA,1e4)
2 > for(s in 1:1e4){
3 +   x=sample(0:1, size = n, prob = c(.4,.6), replace=TRUE)
4 +   neglogL = function(p) -sum(log(dbinom(y, size = 1, prob = p)))
5 +   theta[s] = optim(par = .5,fn = neglogL)$par
6 +
7 > hist(theta)
```

Likelihood, case $\mathcal{B}(p)$



Confidence Interval

This can be used to get some confidence...

Pointwise estimation: compute (simply) $\hat{\theta}(\mathbf{y})$ as a single numerical value

Definition 1.5: Confidence Interval

Let \mathbf{Y} be a random sample of i.i.d. variables with distribution f_θ . A confidence interval of level $1 - \alpha$ for the parameter θ is a (random) interval $[\hat{a}(\mathbf{Y}), \hat{b}(\mathbf{Y})]$ such that

$$\mathbb{P}\left[\theta \in [\hat{a}(\mathbf{Y}), \hat{b}(\mathbf{Y})]\right] = 1 - \alpha$$

Classically, α is 10%, 5% or 1% (or less).

The smaller α , the wider the confidence interval.

Confidence Interval, Gaussian Sample

Let $\{y_1, \dots, y_n\}$ be an i.i.d. sample with a $\mathcal{N}(\mu, \sigma_0^2)$ distribution, where σ_0^2 is assumed to be known.

$$\hat{\mu}(\mathbf{Y}) = \bar{Y}, \text{ then } \hat{\mu}(\mathbf{Y}) = \mathcal{N}\left(\mu, \frac{\sigma_0^2}{n}\right).$$

$$\text{Set } Z = \frac{\hat{\mu}(\mathbf{Y}) - \mu}{\sigma_0 / \sqrt{n}}, \quad Z \sim \mathcal{N}(0, 1).$$

The bilateral confidence interval for μ , with level $1 - \alpha$ is

$$\left[\hat{\mu}(\mathbf{Y}) - u_{\alpha/2} \frac{\sigma_0}{\sqrt{n}}, \hat{\mu}(\mathbf{Y}) + u_{\alpha/2} \frac{\sigma_0}{\sqrt{n}} \right]$$

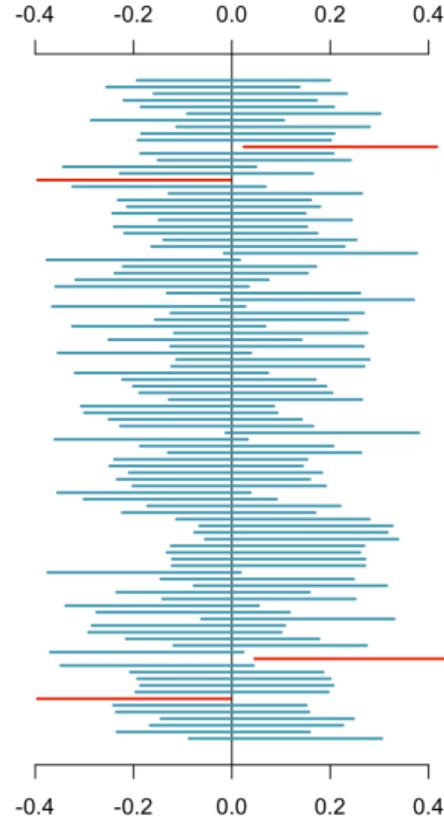
where $u_{\alpha/2} = \Phi^{-1}(1 - \alpha/2)$, where $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z \exp\left(-\frac{x^2}{2}\right) dx$.

Confidence Interval, level α

Sample $\mathcal{N}(0, 1)$ with size n ,

$$IC = \left[\hat{\mu}(\mathbf{Y}) \pm u_{\alpha/2} \frac{1}{\sqrt{n}} \right]$$

```
1 alpha = .05
2 set.seed(1)
3 n=100
4 IC = matrix(NA,100,2)
5 for(s in 1:100){
6   x = rnorm(100,0,1)
7   m = mean(x)
8   IC[s,1] = m-qnorm(1-alpha/2)*1/sqrt(n)
9   IC[s,2] = m+qnorm(1-alpha/2)*1/sqrt(n)
10 }
11 idx=which((IC[,1]<0)&(IC[,2]>0))
```

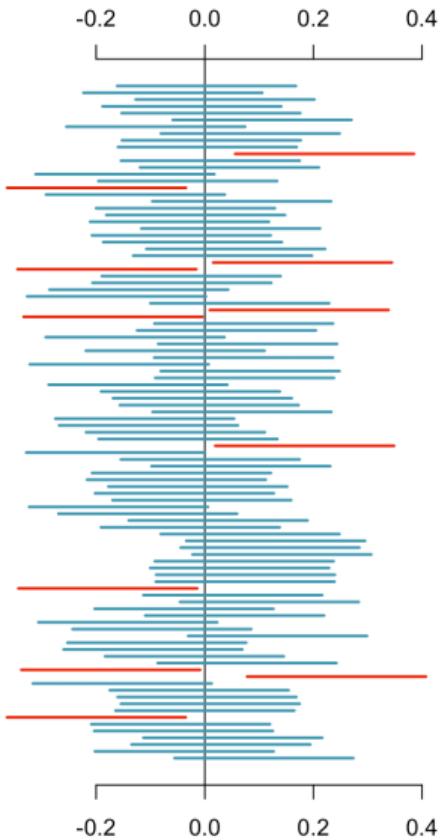


Confidence Interval, level α

Sample $\mathcal{N}(0, 1)$ with size n ,

$$IC = \left[\hat{\mu}(\mathbf{Y}) \pm u_{\alpha/2} \frac{1}{\sqrt{n}} \right]$$

```
1 alpha = .1
2 set.seed(1)
3 n=100
4 IC = matrix(NA,100,2)
5 for(s in 1:100){
6   x = rnorm(100,0,1)
7   m = mean(x)
8   IC[s,1] = m-qnorm(1-alpha/2)*1/sqrt(n)
9   IC[s,2] = m+qnorm(1-alpha/2)*1/sqrt(n)
10 }
11 idx=which((IC[,1]<0)&(IC[,2]>0))
```

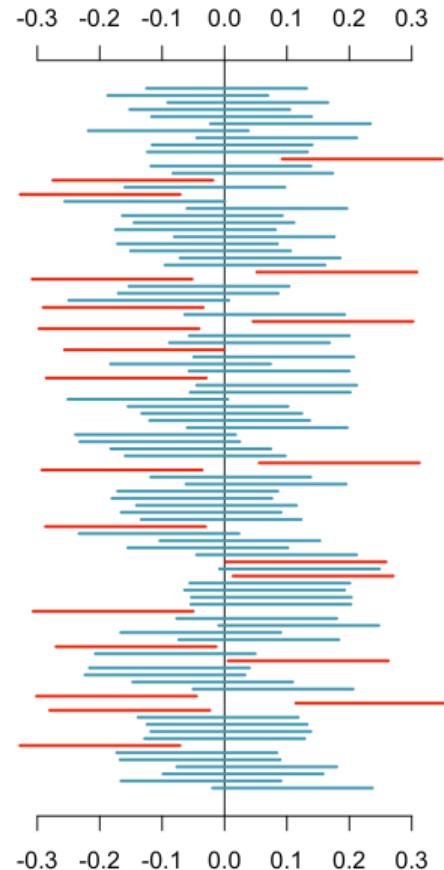


Confidence Interval, level α

Sample $\mathcal{N}(0, 1)$ with size n ,

$$IC = \left[\hat{\mu}(\mathbf{Y}) \pm u_{\alpha/2} \frac{1}{\sqrt{n}} \right]$$

```
1 alpha = .2
2 set.seed(1)
3 n=100
4 IC = matrix(NA,100,2)
5 for(s in 1:100){
6   x = rnorm(100,0,1)
7   m = mean(x)
8   IC[s,1] = m-qnorm(1-alpha/2)*1/sqrt(n)
9   IC[s,2] = m+qnorm(1-alpha/2)*1/sqrt(n)
10 }
11 idx=which((IC[,1]<0)&(IC[,2]>0))
```



Binomial distribution and Confidence

If the variables Y_1, \dots, Y_n are i.i.d. of law $\mathcal{B}(p)$, and if $\hat{p} = \frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}$,

$$\mathbb{E}[\hat{p}] = p \text{ and } \text{Var}[\hat{p}] = \frac{p(1-p)}{n}$$

More precisely, since $n\hat{p} \sim \mathcal{B}(n, p)$,

$$\mathbb{P}\left(F = \frac{k}{n}\right) = \binom{n}{k} p^k (1-p)^{n-k}, \quad \forall k = 0, 1, \dots, n.$$

If n is large enough, according to the central limit theorem

$$Z_n = \sqrt{n} \frac{\hat{p} - p}{\sqrt{p(1-p)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1)$$

Binomial distribution and Confidence

In practice, we assume the normal approximation is valid if $n \geq 30$, $np \geq 15$ and $n(1 - p) \geq 15$. Binomial model, with n large enough

Definition 1.6: Confidence Interval $\{y_1, \dots, y_n\}$, $\mathcal{B}(p)$, n large

Let $\mathbf{y} = \{y_1, \dots, y_n\}$ denote an i.i.d. sample, from $\mathcal{B}(p)$.

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

Note sometimes the confidence interval exceeds 0 or 1

```
1 > sum_y = 12
2 > n = 14
3 > sum_y/n + qnorm(c(.025,.975))*sqrt(sum_y*(n-sum_y)/n^3)
4 [1] 0.6738432 1.0404425
```

Binomial distribution and Confidence

dans un modèle binomial, avec n assez grand

Definition 1.7: Confidence Interval $\{y_1, \dots, y_n\}$, $\mathcal{B}(p)$, Wilson

Let $\mathbf{y} = \{y_1, \dots, y_n\}$ denote an i.i.d. sample, from $\mathcal{B}(p)$. A confidence interval with level α for p is

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

On obtient ces bornes en notant qu'elles correspondent aux p tels que

$$(\hat{p} - p)^2 = u_{1-\alpha/2}^2 \cdot \frac{p(1 - p)}{n} \text{ qui est l'équation de degré 2}$$

$$\left(1 + \frac{u_{1-\alpha/2}^2}{n} \right) p^2 + \left(-2\hat{p} - \frac{u_{1-\alpha/2}^2}{n} \right) p + \left(\hat{p}^2 \right) = 0 .$$

Binomial distribution and Confidence

Modèle binomial, avec n assez grand, comme

$$\frac{(\bar{X} - \bar{Y}) - (p_x - p_y)}{\sqrt{\frac{\bar{X}(1 - \bar{X})}{m} + \frac{\bar{Y}(1 - \bar{Y})}{n}}} \approx \mathcal{N}(0, 1)$$

Definition 1.8: Intervalle de confiance pour $p_x - p_y$ $\mathcal{B}(p_x)$ et $\mathcal{B}(p_y)$, Wald

Soient $\mathbf{x} = \{x_1, \dots, x_m\}$ de loi $\mathcal{B}(p_x)$ et $\mathbf{y} = \{y_1, \dots, y_n\}$ de loi $\mathcal{B}(p_y)$. Un intervalle de confiance de niveau α pour $p_x - p_y$ est

$$\left[\bar{x} - \bar{y} \pm u_{1-\alpha/2} \sqrt{\frac{\bar{x}(1 - \bar{x})}{m} + \frac{\bar{y}(1 - \bar{y})}{n}} \right]$$

Binomial distribution and Confidence

Definition 1.9: Margin of Error ($N = \infty$)

The margin of error is the precision of the result obtained given the confidence level we're willing to accept (α). The (absolute) margin of error is then equal to

$$z_{1-\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

```
1 > qnorm(.975)*sqrt(.5*(1-.5)/1000)
2 [1] 0.03098975
```

i.e. 3.1% (95% level) when $n = 1,000$ and $p \sim 50\%$

Binomial distribution and Confidence

Definition 1.10: Margin of Error ($N < \infty$)

The margin of error for a sample n drawn from a population of finite size N is written as

$$z_{1-\alpha/2} \sqrt{\frac{p(1-p)}{n}} \sqrt{\frac{N-n}{N-1}}$$

```
1 > qnorm(.975)*sqrt(.5*(1-.5)/1000)*sqrt(5000/5999)
2 [1] 0.028292
```

i.e. 2.8% (95% level) when $n = 1,000$, $p \sim 50\%$, out of $N = 6,000$.

Binomial distribution and Confidence

We can invert these formulas to determine the sampling size, based on the (absolute) margin of error

$$n = \frac{z_{1-\alpha/2}^2 \times p(1-p)}{\text{margin of error}^2}$$

when $N = \infty$, or if $N < \infty$

$$n = \frac{p(1-p) + \frac{\text{margin of error}^2}{z_{1-\alpha/2}^2}}{\frac{\text{margin of error}^2}{z_{1-\alpha/2}^2} + \frac{p(1-p)}{N}}$$

Maximum Likelihood and Confidence Intervals

- Consider a parametric model with parameter $\theta \in \Theta \subset \mathbb{R}$.
- Let $\log \mathcal{L}(\theta)$ be the log-likelihood function.
- We test the null hypothesis $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$.
- The likelihood ratio test statistic is:

$$\Lambda = -2 \left[\log \mathcal{L}(\theta_0) - \log \mathcal{L}(\hat{\theta}) \right]$$

where $\hat{\theta}$ is the MLE of θ .

- Under regularity conditions and H_0 , we have:

$$\Lambda \xrightarrow{\mathcal{L}} \chi_1^2$$

as $n \rightarrow \infty$ (Wilks' theorem, [Wilks \(1938\)](#)).

Maximum Likelihood and Confidence Intervals

- Define the $(1 - \alpha)$ confidence region as the set of values of θ for which:

$$\Lambda(\theta) = -2 [\ell(\theta) - \ell(\hat{\theta})] \leq \chi^2_{1,1-\alpha}$$

- This gives the confidence interval:

$$\left\{ \theta \in \Theta : -2 [\log \mathcal{L}(\theta) - \log \mathcal{L}(\hat{\theta})] \leq \chi^2_{1,1-\alpha} \right\}$$

- Interpretation:
 - It includes all values of θ not significantly worse than $\hat{\theta}$.
 - Useful even when the distribution of $\hat{\theta}$ is skewed or unknown.
- This interval is invariant under reparameterization.

For example

- Let y_1, \dots, y_n i.i.d. sample from $\text{Poisson}(\lambda)$

Maximum Likelihood and Confidence Intervals

- Log-likelihood function:

$$\log \mathcal{L}(\lambda) = \sum_{i=1}^n [-\lambda + X_i \log \lambda - \log(X_i!)] = -n\lambda + (\sum X_i) \log \lambda + \text{const}$$

- MLE: $\hat{\lambda} = \bar{X}$
- LRT statistic for testing $H_0 : \lambda = \lambda_0$:

$$\Lambda = -2 \left[\log \mathcal{L}(\lambda_0) - \log \mathcal{L}(\hat{\lambda}) \right] = 2n \left(\lambda_0 - \bar{y} + \bar{X} \log \frac{\bar{y}}{\lambda_0} \right)$$

- Under H_0 : $\Lambda \xrightarrow{d} \chi_1^2$
- For a confidence level $1 - \alpha$, find all λ such that:

$$\Lambda(\lambda) = 2n \left(\lambda - \bar{y} - \bar{y} \log \frac{\lambda}{\bar{y}} \right) \leq \chi_{1,1-\alpha}^2$$

Maximum Likelihood and Confidence Intervals

- This defines the confidence interval:

$$\left\{ \lambda > 0 : \Lambda(\lambda) \leq \chi^2_{1,1-\alpha} \right\}$$

(and solve this numerically).

Gamma Distribution: Maximum Likelihood

Consider some i.i.d. sample $\{y_i\}$ from $Y_i \sim \text{Gamma}(\alpha, \beta)$,

$$f(y|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} y^{\alpha-1} e^{-\beta y}$$

$$\log \mathcal{L}(\alpha, \beta | \mathbf{y}) = n\alpha \log(\beta) - n \log[\Gamma(\alpha)] + (\alpha - 1) \sum_{i=1}^n \log(y_i) - \beta \sum_{i=1}^n y_i$$

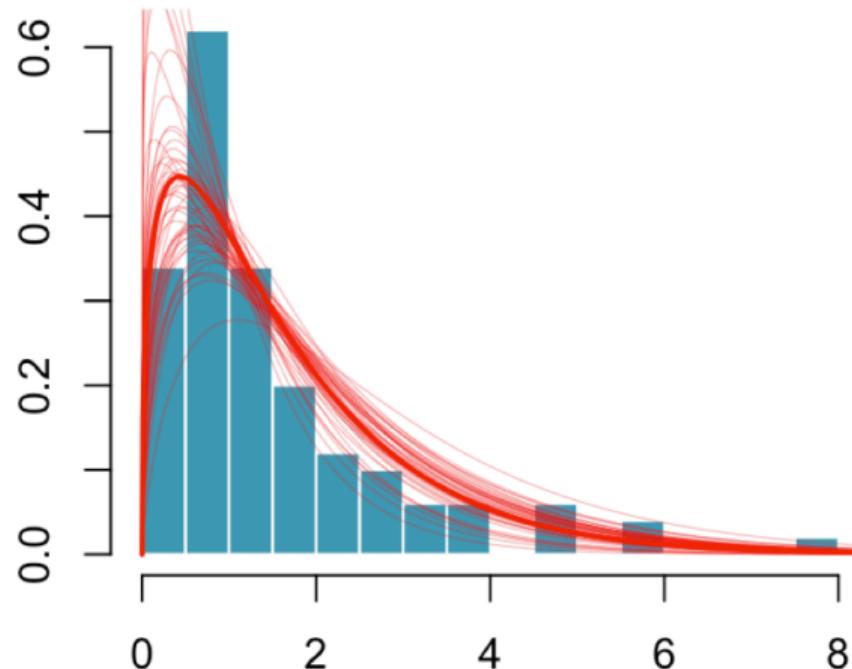
$$\nabla \log \mathcal{L}(\alpha, \beta) = \begin{pmatrix} n \log(\beta) - n[\log \Gamma(\alpha)]' + \sum_{i=1}^n \log(y_i) \\ \frac{n\alpha}{\beta} - \sum_{i=1}^n y_i \end{pmatrix}$$

$$H \log \mathcal{L}(\alpha, \beta) = \frac{1}{n(1 - \alpha[\log \Gamma(\alpha)]'')} \begin{pmatrix} \alpha & \beta \\ \beta & \beta^2 [\log(\Gamma(\alpha))]'' \end{pmatrix}$$

Gamma Distribution: Maximum Likelihood

Gamma distribution with parameter, $\theta = (\alpha, \beta)$, $\hat{\theta} = (\hat{\alpha}, \hat{\beta}) = \underset{\theta \in \mathbb{R}_+^2}{\operatorname{argmax}} \left\{ \log \mathcal{L}(\theta) \right\}$

```
1 > y = exp(rnorm(100))
2 > library(MASS)
3 > (F = fitdistr(x,"gamma"))
      shape          rate
5  1.3562877    0.8207035
6  (0.1732663) (0.1263275)
7 > log_lik = function(theta){
8 +   a = theta[1]
9 +   b = theta[2]
10 +  logL = sum(log(dgamma(y,a,b)))
11 +  return(-logL)
12 + }
13 > optim(c(1,1),log_lik)
14 $par
15 [1]  1.3558113  0.8206505
```



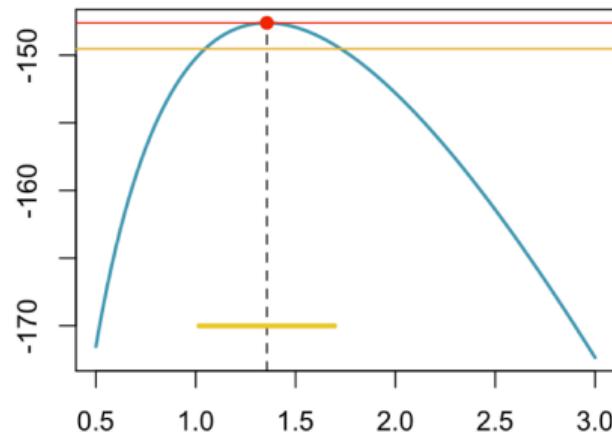
Gamma Distribution: Profile Likelihood

We can also consider if, somehow, α is the parameter of interest

$$\hat{\alpha}^* = \underset{\alpha \in \mathbb{R}_+}{\operatorname{argmax}} \left\{ \underbrace{\max_{\beta \in \mathbb{R}_+} \left\{ \log \mathcal{L}(\alpha, \beta) \right\}}_{\text{function of } \alpha} \right\}$$

Function $\alpha \mapsto \max_{\beta \in \mathbb{R}_+} \left\{ \log \mathcal{L}(\alpha, \beta) \right\}$ called **profile likelihood**

```
1 > prof_log_lik = function(a){  
2 +   b = (optim(1,function(z) -sum(log  
3 +     (dgamma(x,a,z))))$par  
4 +   return(-sum(log(dgamma(x,a,b))))  
4 + }  
5 > optim(1,prof_log_lik)  
6 $par  
7 [1] 1.356445
```



Gamma Distribution: Profile Likelihood

This can be used to get a confidence interval for the parameter of interest.

In a statistical context, suppose that unknown parameter can be partitioned $\theta = (\alpha, \beta)$ where α is the parameter of interest, and β is a nuisance parameter.

Consider $\{y_1, \dots, y_n\}$, a sample from distribution F_θ , so that the log-likelihood is

$$\log \mathcal{L}(\theta) = \sum_{i=1}^n \log f_\theta(y_i)$$

$\hat{\theta}^{MLE}$ is defined as $\hat{\theta}^{MLE} = \operatorname{argmax} \{\log \mathcal{L}(\theta)\}..$

Rewrite the log-likelihood as $\log \mathcal{L}(\theta) = \log \mathcal{L}_\alpha(\beta)$. Define

$$\hat{\beta}_\alpha^{pMLE} = \operatorname{argmax}_\beta \{\log \mathcal{L}_\alpha(\beta)\}$$

Gamma Distribution: Profile Likelihood

and then $\hat{\alpha}^{pMLE} = \operatorname{argmax}_{\alpha} \left\{ \log \mathcal{L}_{\alpha}(\hat{\beta}_{\alpha}^{pMLE}) \right\}$. Observe that

$$\sqrt{n}(\hat{\alpha}^{pMLE} - \alpha) \xrightarrow{\mathcal{L}} \mathcal{N}(0, [\mathbb{I}_{\alpha,\alpha} - \mathbb{I}_{\alpha,\beta} \mathbb{I}_{\beta,\beta}^{-1} \mathbb{I}_{\beta,\alpha}]^{-1})$$

The (profile) likelihood ratio test is based on

$$2 (\max \{\mathcal{L}(\alpha, \beta)\} - \max \{\mathcal{L}(\alpha_0, \beta)\})$$

If (α_0, β_0) are the true value, this difference can be written

$$2 (\max \{\mathcal{L}(\alpha, \beta)\} - \max \{\mathcal{L}(\alpha_0, \beta_0)\}) - 2 (\max \{\mathcal{L}(\alpha_0, \beta)\} - \max \{\mathcal{L}(\alpha_0, \beta_0)\})$$

Using Taylor's expansion

$$\frac{1}{\sqrt{n}} \left. \frac{\partial \mathcal{L}(\alpha, \beta)}{\partial \alpha} \right|_{(\alpha_0, \hat{\beta}_{\alpha_0})} \xrightarrow{\mathcal{L}} \mathcal{N}(0, \mathbb{I}_{\alpha_0 \alpha_0}) - \mathbb{I}_{\alpha_0 \beta_0} \mathbb{I}_{\beta_0 \beta_0}^{-1} \mathbb{I}_{\beta_0 \alpha_0}$$

and $2 (\mathcal{L}(\hat{\alpha}, \hat{\beta}) - \mathcal{L}(\alpha_0, \hat{\beta}_{\alpha_0})) \xrightarrow{\mathcal{L}} \chi^2(\dim(\alpha))$.

Gamma Distribution: Profile Likelihood

Consider some lognormal sample, and fit a Gamma distribution,

$$f(x; \alpha, \beta) = \frac{x^{\alpha-1} \beta^\alpha e^{-\beta x}}{\Gamma(\alpha)} \text{ with } x > 0 \text{ and } \theta = (\alpha, \beta).$$

```
1 > y = exp(rnorm(100))
```

Maximum-likelihood, $\hat{\theta} = \operatorname{argmax}\{\log \mathcal{L}(\theta)\}$.

```
1 > library(MASS)
2 > (F = fitdistr(x, "gamma"))
   shape      rate
  1.4214497  0.8619969
  (0.1822570) (0.1320717)
6 > F$estimate[1]+c(-1,1)*1.96*F$sd[1]
7 [1] 1.064226 1.778673
```

Gamma Distribution: Profile Likelihood

See also

```
1 > log_lik=function(theta){  
2 +   a = theta[1]  
3 +   b = theta[2]  
4 +   logL = sum(log(dgamma(y,a,b)))  
5 +   return(-logL)  
6 + }  
7 > optim(c(1,1),log_lik)  
8 $par  
9 [1] 1.4214116 0.8620311
```

We can also use [profile likelihood](#),

$$\hat{\alpha} = \operatorname{argmax}_{\beta} \left\{ \max_{\beta} \{ \log \mathcal{L}(\alpha, \beta) \} \right\} = \operatorname{argmax} \left\{ \log \mathcal{L}(\alpha, \hat{\beta}_\alpha) \right\}$$

Gamma Distribution: Profile Likelihood

```
1 > prof_log_lik = function(a){  
2 +   b = (optim(1,function(z) -sum(log(dgamma(x,a,z))))$par  
3 +   return(-sum(log(dgamma(x,a,b))))  
4 + }  
5  
6  
7 > vx = seq(.5,3,length=101)  
8 > vl = -Vectorize(prof_log_lik)(vx)  
9 > plot(vx,vl,type="l")  
10 > optim(1,prof_log_lik)  
11 $par  
12 [1] 1.421094
```

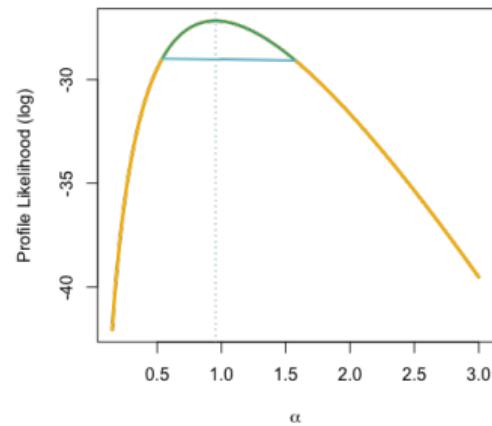
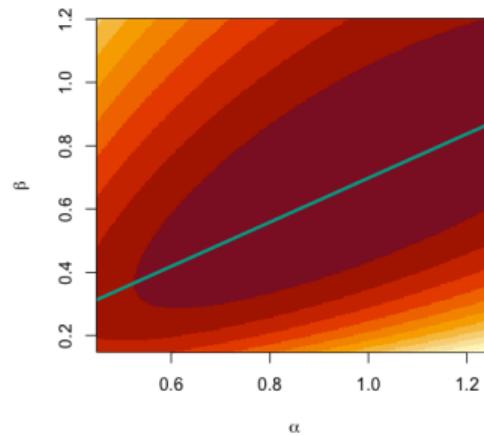
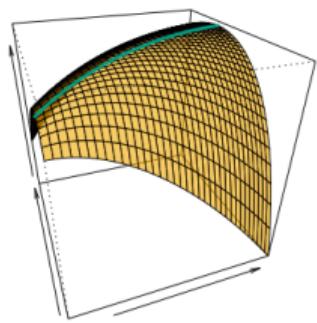
We can use the likelihood ratio test

$$2(\log \mathcal{L}_p(\hat{\alpha}) - \log \mathcal{L}_p(\alpha)) \sim \chi^2(1)$$

Gamma Distribution: Profile Likelihood

The implied 95% confidence interval is

```
1 > (b1 = uniroot(function(z) Vectorize(prof_log_lik)(z)+borne,c(.5,1.5))  
     $root)  
2 [1] 1.095726  
3 > (b2 = uniroot(function(z) Vectorize(prof_log_lik)(z)+borne,c(1.25,2.5))  
     $root)  
4 [1] 1.811809
```



Random Vectors

Let $\mathbf{X} = (X_1, \dots, X_d)$ denote a random vector in dimension d ,

- The expectation of \mathbf{X} , denoted $\mathbb{E}(\mathbf{X})$ is defined (if it exists) by the vector of dimension d , $\mathbb{E}(\mathbf{X}) = (\mathbb{E}(X_1), \dots, \mathbb{E}(X_d))^\top$.
- The covariance matrix (also called variance-covariance matrix of \mathbf{X}) is defined (if it exists) by the matrix of size (d, d) .

$$\text{Var}(\mathbf{X}) = \mathbb{E} \left((\mathbf{X} - \mathbb{E}(\mathbf{X}))(\mathbf{X} - \mathbb{E}(\mathbf{X}))^\top \right).$$

So the i, j term of this matrix represents the covariance between X_i and X_j ,

$$\text{Cov}(X_i, X_j) = \mathbb{E} [(X_i - \mathbb{E}(X_i))(X_j - \mathbb{E}(X_j))].$$

Random Vectors

Let \mathbf{X} be a random vector of dimension d , mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.
Let \mathbf{A} and \mathbf{B} be two reel matrices of size (d, p) and (d, q) and let $\mathbf{a} \in \mathbb{R}^p$ then

- $\text{Var}(\mathbf{X}) = \mathbb{E}((\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top) = \mathbb{E}(\mathbf{X}\mathbf{X}^\top) - \boldsymbol{\mu}\boldsymbol{\mu}^\top.$
- $\mathbb{E}(\mathbf{A}^\top \mathbf{X} + \mathbf{a}) = \mathbf{A}^\top \boldsymbol{\mu} + \mathbf{a}.$
- $\text{Var}(\mathbf{A}^\top \mathbf{X} + \mathbf{a}) = \mathbf{A}^\top \boldsymbol{\Sigma} \mathbf{A}.$
- $\text{Cov}(\mathbf{A}^\top \mathbf{X}, \mathbf{B}^\top \mathbf{X}) = \mathbf{A}^\top \boldsymbol{\Sigma} \mathbf{B}.$

The Gaussian Distribution

Definition 1.11: Gaussian Univariate Distribution

A Gaussian variable, with distribution $\mathcal{N}(\mu, \sigma^2)$, where $\mu \in \mathbb{R}$ and $\sigma > 0$, has density

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}\right), \quad x \in \mathbb{R}.$$

- Then $\mathbb{E}(X) = \mu$ and $\text{Var}(X) = \sigma^2$.
- Observe that if $Z \sim \mathcal{N}(0, 1)$, $X = \mu + \sigma Z \sim \mathcal{N}(\mu, \sigma^2)$.

The Gaussian Distribution

Definition 1.12: Gaussian Multivariate Distribution

The **Gaussian vector** $\mathcal{N}(\mu, \Sigma)$: $\mathbf{X} = (X_1, \dots, X_d)$ is a Gaussian vector with mean $\mathbb{E}(\mathbf{X}) = \mu \in \mathbb{R}^d$ and covariance $d \times d$ matrix $\text{Var}(\mathbf{X}) = \Sigma = \mathbb{E}((\mathbf{X} - \mu)(\mathbf{X} - \mu)^\top)$ non-degenerated (Σ is invertible) if its density is

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \sqrt{\det \Sigma}} \exp\left(-\frac{1}{2} (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu)\right), \quad \mathbf{x} \in \mathbb{R}^n,$$

see [multivariate Gaussian distribution](#)

Gaussian (multivariate) distribution

$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with density

$$f_{\mathbf{X}}(x_1, \dots, x_d) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

where $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$ and $\text{Var}(\mathbf{X}) = \boldsymbol{\Sigma}$.

Estimates are $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ and $\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$

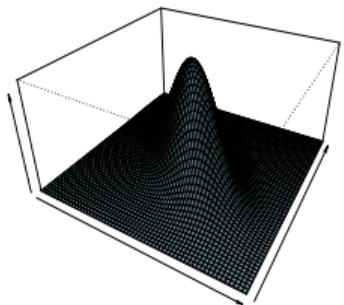
In dimension 2, $f(x, y)$ is proportional to

$$\exp\left(-\frac{1}{2(1-\rho^2)} \left[\frac{(x - \mu_X)^2}{\sigma_X^2} + \frac{(y - \mu_Y)^2}{\sigma_Y^2} - \frac{2\rho(x - \mu_X)(y - \mu_Y)}{\sigma_X \sigma_Y} \right]\right)$$

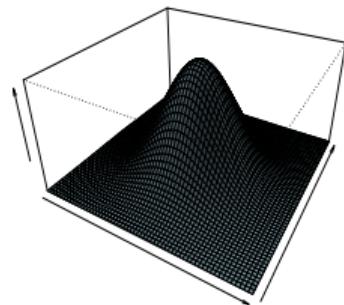
levels curves (isodensities) are ellipses.

Gaussian (multivariate) distribution

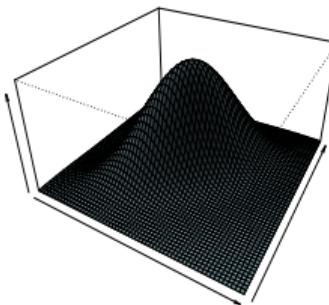
Densité du vecteur Gaussien, $r=0.7$



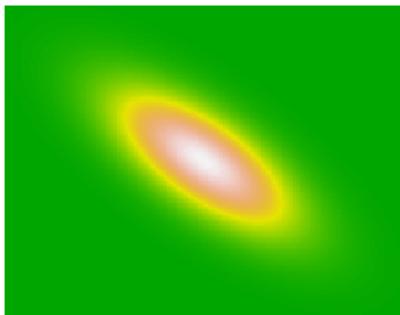
Densité du vecteur Gaussien, $r=0.0$



Densité du vecteur Gaussien, $r=-0.7$



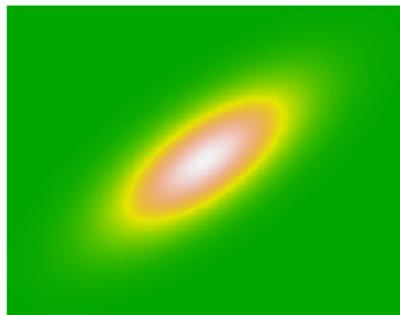
Courbes de niveau du vecteur Gaussien, $r=-0.7$



Courbes de niveau du vecteur Gaussien, $r=0.0$



Courbes de niveau du vecteur Gaussien, $r=0.7$



Quadratic Forms

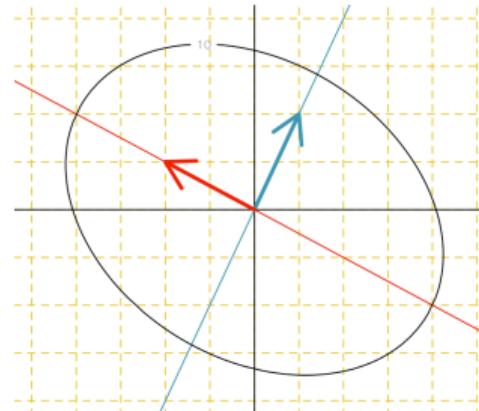
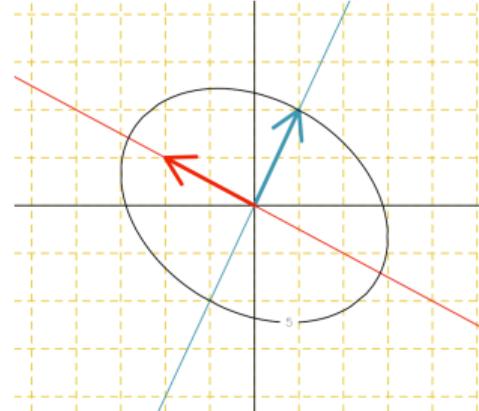
Consider $\mathbf{M} = \begin{pmatrix} a & b \\ b & c \end{pmatrix}$,
and function $\mathbf{z} \mapsto \mathbf{z}^\top \mathbf{M} \mathbf{z}$, i.e.

$$f: \begin{pmatrix} x \\ y \end{pmatrix} \mapsto \begin{pmatrix} x & y \end{pmatrix} \begin{pmatrix} a & b \\ b & c \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

or $ax^2 + 2bxy + cy^2$ is a quadratic form.

If $\mathbf{M} > 0$, points $\mathbf{z} = (x, y)$ such that $\mathbf{z}^\top \mathbf{M} \mathbf{z} = \gamma$, for some $\gamma > 0$, are on an **ellipse** (centered on $\mathbf{0}$)

Let $\lambda_1 \geq \lambda_2 > 0$ denote the eigenvalues of \mathbf{M} and \vec{v}_1 and \vec{v}_2 denote the eigenvectors.



Quadratic Forms

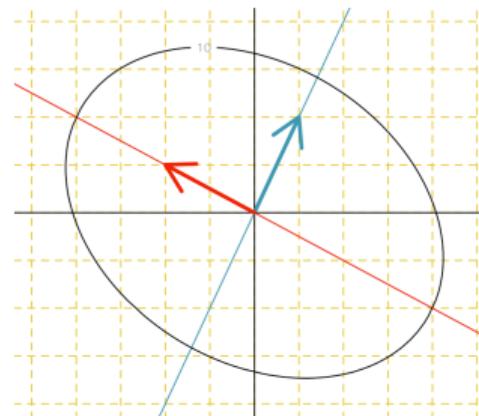
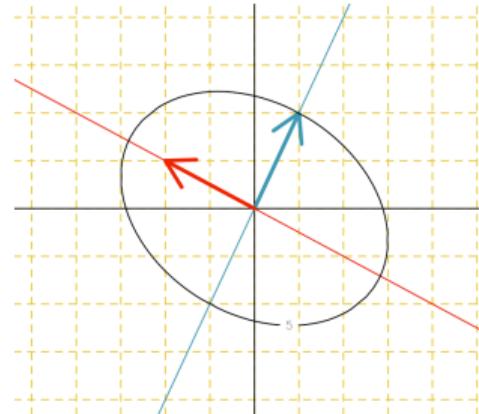
On the picture, $M = \begin{pmatrix} 0.6 & 0.2 \\ 0.2 & 0.9 \end{pmatrix}$

```
1 > M=matrix(c(.6,.2,.2,.9),2,2)
2 > eigen(M)
3 eigen() decomposition
4 $values
5 [1] 1.0 0.5
6 $vectors
7 [,1]      [,2]
8 [1,] 0.4472136 -0.8944272
9 [2,] 0.8944272  0.4472136
```

i.e. $\lambda_1 = 1$ and $\lambda_2 = 1/2$, and

$$\vec{v}_1 = \frac{1}{\sqrt{5}} \begin{pmatrix} 1 \\ 2 \end{pmatrix}, \quad \vec{v}_2 = \frac{1}{\sqrt{5}} \begin{pmatrix} -2 \\ 1 \end{pmatrix}$$

Note that $\|\vec{v}_1\| = \|\vec{v}_2\| = 1$ and $\vec{v}_1 \perp \vec{v}_2$



The Gaussian Distribution

If \mathbf{X} is a Gaussian vector, then for any i , X_i has a (univariate) Gaussian distribution, but its converse is not necessarily true.

Let $\mathbf{X} = (X_1, \dots, X_d)$ be a random vector with mean $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu}$ and with covariance matrix $\boldsymbol{\Sigma}$, if \mathbf{A} is a $k \times d$ matrix, and $\mathbf{b} \in \mathbb{R}^k$, then $\mathbf{Y} = \mathbf{AX} + \mathbf{b}$ is a Gaussian vector \mathbb{R}^k , with distribution $\mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$.

Observe that if (X_1, X_2) is a Gaussian vector, X_1 and X_2 are independent if and only if

$$\text{Cov}(X_1, X_2) = \mathbb{E}((X_1 - \mathbb{E}(X_1))(X_2 - \mathbb{E}(X_2))) = 0.$$

Beta and Dirichlet distribution

Definition 1.13: Beta distribution

X , random variable in $[0, 1]$, has distribution $\text{Beta}(\alpha, \beta)$ if its density is

$$f(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} \text{ where } B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}, \quad x \in [0, 1].$$

Proposition 1.2: Beta distribution

X_1, X_2 be two independent Gamma variables, $X_1 \sim \text{Gamma}(\alpha, \gamma)$ and $X_2 \sim \text{Gamma}(\beta, \gamma)$, then

$$Y = \frac{X_1}{X_1 + X_2} \sim \text{Beta}(\alpha, \beta).$$

Beta and Dirichlet distribution

Definition 1.14: Dirichlet distribution

\mathbf{X} , random vector in $\mathcal{S}_d \subset [0, 1]^d$, has distribution $\text{Dirichlet}(\boldsymbol{\alpha})$ if its density is

$$f(\mathbf{x}) = \frac{x_1^{\alpha_1-1} x_2^{\alpha_2-1} \cdots x_d^{\alpha_d-1}}{B(\boldsymbol{\alpha})} \text{ where } B(\boldsymbol{\alpha}) = \frac{\prod_{i=1}^d \Gamma(\alpha_i)}{\Gamma\left(\sum_{i=1}^d \alpha_i\right)}, \quad \mathbf{x} \in \mathcal{S}_d.$$

Proposition 1.3: Dirichlet distribution

X_1, X_2, \dots, X_d be d independent Gamma variables, $X_i \sim \text{Gamma}(\alpha_i, \gamma)$,

$$\mathbf{Y} = (Y_1, \dots, Y_d) = \left(\frac{X_1}{X_1 + \cdots + X_d}, \dots, \frac{X_d}{X_1 + \cdots + X_d} \right) \sim \text{Dirichlet}(\boldsymbol{\alpha}).$$

Proposition 1.4: Bayes formula

Given two events \mathcal{A} and \mathcal{B} such that $\mathbb{P}[\mathcal{B}] \neq 0$,

$$\mathbb{P}[\mathcal{A}|\mathcal{B}] = \frac{\mathbb{P}[\mathcal{B}|\mathcal{A}] \cdot \mathbb{P}[\mathcal{A}]}{\mathbb{P}[\mathcal{B}]} \propto \mathbb{P}[\mathcal{B}|\mathcal{A}] \cdot \mathbb{P}[\mathcal{A}].$$

Bayes (1763) and Laplace (1774).

Bayesian statistics ?

- Bayes formula (the “inverse problem”),
Bayes (1763), Laplace (1774)

Given two events \mathcal{A} and \mathcal{B} such that $\mathbb{P}(\mathcal{B}) \neq 0$,

$$\mathbb{P}[\mathcal{A}|\mathcal{B}] = \frac{\mathbb{P}[\mathcal{B}|\mathcal{A}] \cdot \mathbb{P}[\mathcal{A}]}{\mathbb{P}[\mathcal{B}]}.$$

“If a person has an expectation depending on the happening of an event, the probability of the event is [in the ratio] to the probability of its failure as his loss if it fails [is in the ratio] to his gain if it happens ”, Proposition 2, Bayes (1763)

“The probability of any event is the ratio between the value at which an expectation depending on the happening of the event ought to be computed, and the chance of the thing expected upon its happening ”, Bayes (1763)



Bayesian statistics ?

- Bayes formula (the “inverse problem”),
Bayes (1763), Laplace (1774)

Given two events \mathcal{A} and \mathcal{B} such that $\mathbb{P}(\mathcal{B}) \neq 0$,

$$\mathbb{P}[\mathcal{A}|\mathcal{B}] = \frac{\mathbb{P}[\mathcal{B}|\mathcal{A}] \cdot \mathbb{P}[\mathcal{A}]}{\mathbb{P}[\mathcal{B}]} \propto \mathbb{P}[\mathcal{B}|\mathcal{A}] \cdot \mathbb{P}[\mathcal{A}].$$

- subjective probabilities,
De Finetti (1937), Anscombe et al. (1963), Kahneman and Tversky (1972) Savage (1972), Jeffrey (2004)
- Non-frequentist approach of probabilities,
Neyman (1977), Bayarri and Berger (2004)
- Credibility and “experience rating”
Whitney (1918), Longley-Cook (1962), Bühlmann (1967), Klugman (1991)



Bayesian statistics ?

- Bayes formula (the “inverse problem”),
Bayes (1763), Laplace (1774)

Given two events \mathcal{A} and \mathcal{B} such that $\mathbb{P}(\mathcal{B}) \neq 0$,

$$\mathbb{P}[\mathcal{A}|\mathcal{B}] = \frac{\mathbb{P}[\mathcal{B}|\mathcal{A}] \cdot \mathbb{P}[\mathcal{A}]}{\mathbb{P}[\mathcal{B}]}.$$

- An **inverse problem** (we try to determine the causes of a phenomenon from the experimental observation of its effects)
- An **update** of beliefs (from a *prior* distribution $\mathbb{P}(\mathcal{A})$ to a *posterior* distribution $\mathbb{P}(\mathcal{A}|\mathcal{B})$)



Bayesian statistics ?

A person coughs (event B). Which hypothesis is the most credible?

$$\begin{cases} A_1 : \text{she has lung cancer} \\ A_2 : \text{she has gastroenteritis} \\ A_3 : \text{she has the flu} \end{cases}$$

With Bayes' rule $\mathbb{P}[\text{disease}|\text{symptom}] \propto \mathbb{P}[\text{symptom}|\text{disease}] \cdot \mathbb{P}[\text{disease}]$

$$\begin{cases} A_1 : \mathbb{P}[\text{disease}] \approx 0 \text{ (even if } \mathbb{P}[\text{symptom}|\text{disease}] \approx 1) \\ A_2 : \mathbb{P}[\text{symptom}|\text{disease}] \approx 0 \text{ (even if } \mathbb{P}[\text{symptom}|\text{disease}] \text{ high)} \\ A_3 : \text{two reasonable probabilities} \end{cases}$$

The practice of conditional probabilities

"Monty Hall" problem
(from *Let's make a deal*)



The practice of conditional probabilities

"Monty Hall" problem
(from *Let's make a deal*)



The practice of conditional probabilities

"Monty Hall" problem
(from *Let's make a deal*)



$$\begin{aligned}\mathbb{P}(\text{treasure behind the door}) \\ = \frac{1}{3}\end{aligned}$$

The practice of conditional probabilities

"Monty Hall" problem
(from *Let's make a deal*)



$$\begin{aligned}\mathbb{P}(\text{treasure behind the door}) \\ = \frac{1}{3}\end{aligned}$$

The practice of conditional probabilities

"Monty Hall" problem
(from *Let's make a deal*)

- strategy 1 : always switch the door
- strategy 2 : never switch the door



$\mathbb{P}(\text{strategy 2 winning})$

$= \mathbb{P}(\text{treasure behind the door choisie initialement})$

$$= \frac{1}{3}$$

(making the goat appear behind the third door does not bring no information on what's behind the first door)

The practice of conditional probabilities

"Monty Hall" problem
(from *Let's make a deal*)

- strategy 1 : always switch the door
- strategy 2 : never switch the door



$\mathbb{P}(\text{strategy 1 winning})$

$= \mathbb{P}(\text{treasure behind the other door})$

$= \mathbb{P}(\text{treasure behind the other door} | \times \text{ correct}) \cdot \mathbb{P}(\times \text{ correct})$

$+ \mathbb{P}(\text{treasure behind the other door} | \times \text{ false}) \cdot \mathbb{P}(\times \text{ false})$

$$= 0 \cdot \frac{1}{3} + 1 \cdot \frac{2}{3} = \frac{2}{3}$$

Practice of Bayesian Statistics

“Do doctors understand test results? ”, from Kremer (2014):

1 percent of adults have cancer. The vast majority of these cancers (90 percent) can be detected by a test. There is a 9 percent chance that the test will be positive in a person who does not have cancer. If the test is positive, what is the likelihood that the person actually has cancer?

- A) 9 out of 10
- B) 8 out of 10
- C) 1 out of 2
- D) 1 out of 10
- E) 1 out of 100



Practice of Bayesian Statistics

“Do doctors understand test results? ”, from Kremer (2014):

1 percent of adults have cancer. The vast majority of these cancers (90 percent) can be detected by a test. There is a 9 percent chance that the test will be positive in a person who does not have cancer. If the test is positive, what is the likelihood that the person actually has cancer?

- A) 9 out of 10 (chosen by 50% gynecologists)
- B) out of 10
- C) 1 out of 2
- D) 1 out of 10
- E) 1 out of 100



Practice of Bayesian Statistics

1 percent of adults have cancer. The vast majority of these cancers (90 percent) can be detected by a test. There is a 9 percent chance that the test will be positive in a person who does not have cancer. If the test is positive, what is the likelihood that the person actually has cancer?

Answer: when formalizing

$$\begin{cases} \mathbb{P}[\text{cancer}] = 1\% \\ \mathbb{P}[\text{test positive}|\text{cancer}] = 90\% \\ \mathbb{P}[\text{test positive}|\text{no cancer}] = 9\% \end{cases}$$

then, using Bayes' rule

$$\mathbb{P}[\text{cancer}|\text{test positive}] = \frac{\mathbb{P}[\text{test positive}|\text{cancer}] \cdot \mathbb{P}[\text{cancer}]}{\mathbb{P}[\text{test positive}]} = \frac{90\% \times 1\%}{9\% \times 99\% +} = \frac{9}{9 + 89} \simeq \frac{1}{10}$$

valid answer is D, "1 out of 10".

Practice of Bayesian Statistics

For Gigerenzer and Hoffrage (1995), the Bayesian formulation is (too) complex.

Another presentation of the problem:

Out of 10,000 people, 100 have cancer. Of these 100, 90%, or 90, will test positive. Of the remaining 9,900, 9 percent, or 899, will test positive. Of a sample of people who test positive, what fraction actually have cancer?

Answer: 90 among (90+899), i.e. about “1 out of 10”.

Bayesianism, statistics and calculus

$$\text{posterior} = \pi(\theta|y) = \frac{\pi(\theta) \cdot \mathbb{P}(y|\theta)}{\mathbb{P}(y)} = \frac{\text{prior} \cdot \text{likelihood}}{\text{evidence}}$$

$$\text{posterior} = \pi(\theta|y) \propto \frac{\theta^{a-1}(1-\theta)^{b-1}}{B(a,b)} \cdot \binom{s}{n} \theta^s (1-\theta)^{n-s}$$

- Conjugate distributions: **Binomial - Beta**

The likelihood for binomial (Bernoulli) variables

$$\begin{cases} \mathbf{x} \mapsto f(\mathbf{x}; p) = p^s(1-p)^{n-s} \text{ where } s = \mathbf{x}^\top \mathbf{1} = x_1 + \dots + x_n \\ p \mapsto \pi(p; \mathbf{x}) = p^s(1-p)^{n-s} \text{ on } [0, 1] \text{ is a Beta distribution} \end{cases}$$

If $\begin{cases} x_i | \theta \sim \mathcal{B}(\theta) \\ \theta \sim \text{Beta}(a, b) \text{ prior} \end{cases}$ then $\theta | \mathbf{x} \sim \text{Beta}(a+s, b+n-s)$ posterior

(that can be extended to **Multinomial - Dirichlet**)

Bayesianism, statistics and calculus

- Conjugate distributions : **Poisson - Gamma**

The likelihood for Poisson variables is

$$\begin{cases} \mathbf{x} \mapsto f(\mathbf{x}; \lambda) = \frac{e^{n\lambda} \lambda^s}{x_1! \cdots x_n!} \text{ where } s = \mathbf{x}^\top \mathbf{1} = x_1 + \cdots + x_n \\ \lambda \mapsto \pi(\lambda; \mathbf{x}) \propto e^{n\lambda} \lambda^s \text{ on } \mathbb{R}_+ \text{ is a Gamma distribution} \end{cases}$$

If

$$\begin{cases} x_i | \lambda \sim \mathcal{P}(\lambda) \\ \theta \sim \text{Gamma}(a, b) \text{ a priori} \end{cases} \quad \text{then } \lambda | \mathbf{x} \sim \text{Gamma}(a + s, b + n) \text{ a posteriori}$$

Hence

$$\text{a priori } \mathbb{E}(\lambda) = \frac{a}{b} \text{ and a posteriori } \mathbb{E}(\lambda | \mathbf{x}) = \frac{a + s}{b + n}$$

intensively used in credibility theory **Bühlmann (1967)**.

Bayesianism, statistics and calculus

- Conjugate distributions : **Normal - Normal**

If variance Σ is known

$$\begin{cases} \mathbf{x}_i | \boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \end{cases} \quad \text{then } \boldsymbol{\mu} | \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$$

where
$$\begin{cases} \boldsymbol{\mu}_x = (\boldsymbol{\Sigma}_0^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1} (\boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0 + n\boldsymbol{\Sigma}^{-1}\bar{\mathbf{x}}) \\ \boldsymbol{\Sigma}_x = (\boldsymbol{\Sigma}_0^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1} \end{cases}$$

used classically in Bayesian econometrics.

Bayesianism, statistics and calculus

- Conjugate distributions : **Normal - Inverse Wishart**

If mean μ is known

$$\begin{cases} \mathbf{x}_i | \boldsymbol{\Sigma} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ \boldsymbol{\Sigma} \sim IW(\nu_0, \boldsymbol{\Psi}_0) \end{cases} \quad \text{then } \boldsymbol{\Sigma} | \mathbf{x} \sim IW(\nu_x, \boldsymbol{\Psi}_x)$$

where
$$\begin{cases} \nu_x = n + \nu \\ \boldsymbol{\Psi}_x = \boldsymbol{\Psi} + \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top \end{cases}$$

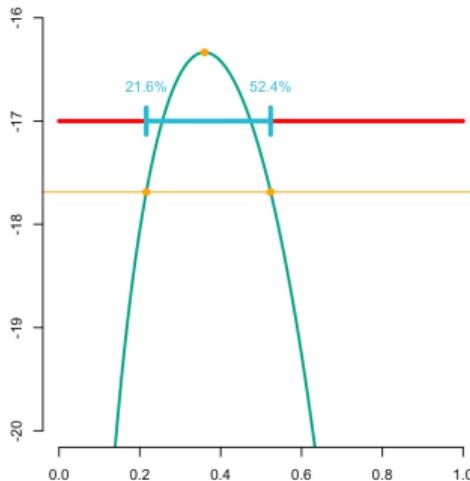
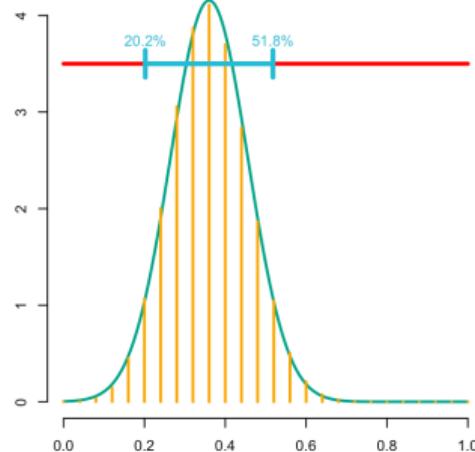
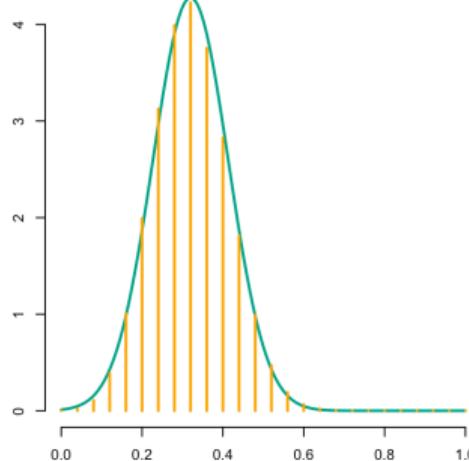
Classically used in Bayesian econometrics, for VAR models, [Adjemian and Pelgrin \(2008\)](#), or in portfolio management, [Black and Litterman \(1990, 1992\)](#) (see also [Satchell and Scowcroft \(2000\)](#) for a perspective).

Bayesianism, statistics and calculus

- Posterior distribution

Suppose $x = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0\}$, $\mathcal{B}(\theta)$

Frequentist approach, $\hat{\theta} \approx \mathcal{N}\left(\theta, \frac{\theta(1-\theta)}{n}\right)$, $\mathbb{P}\left(\theta \in [\bar{x} \pm 1.64\sqrt{\frac{\bar{x}(1-\bar{x})}{n}}]\right) \approx 90\%$

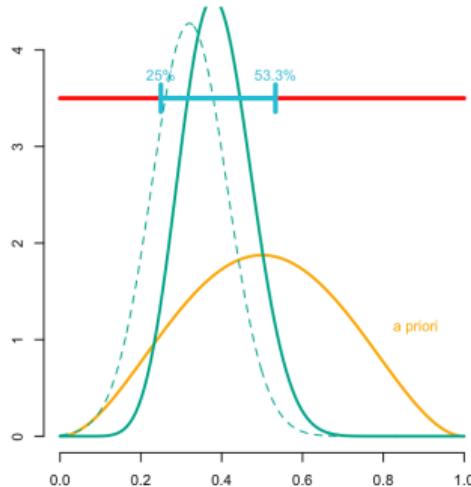
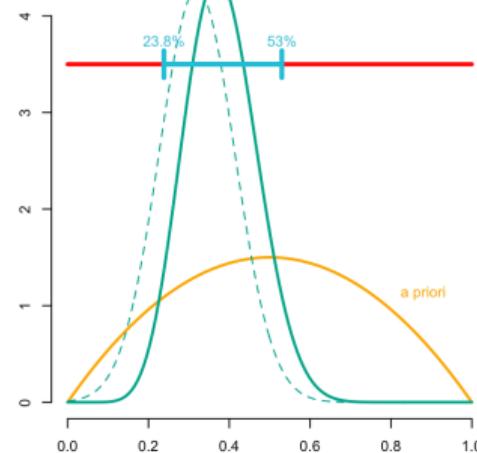
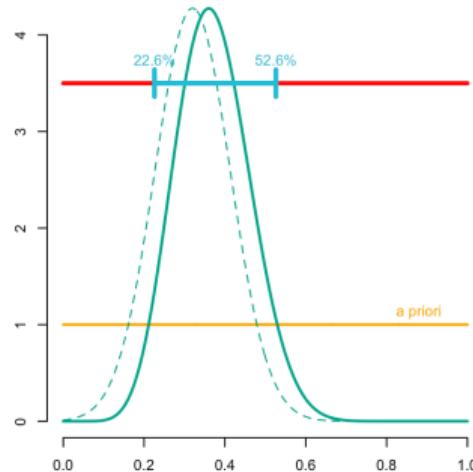


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\hat{\theta}|\mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

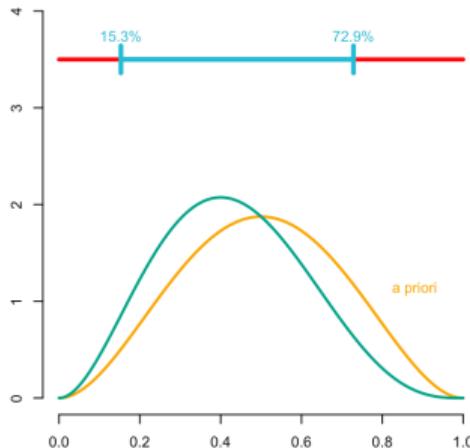
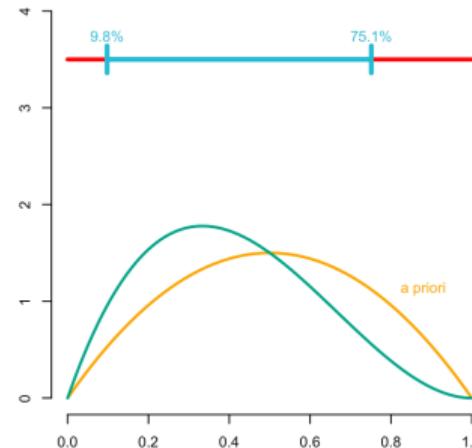
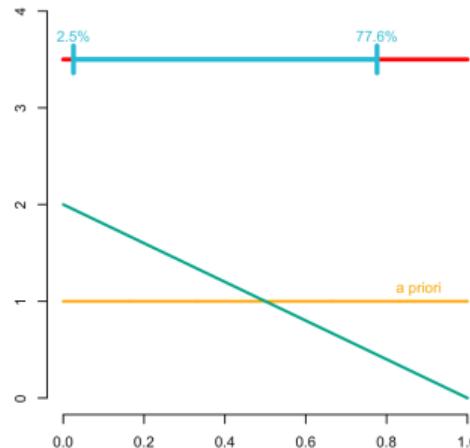


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

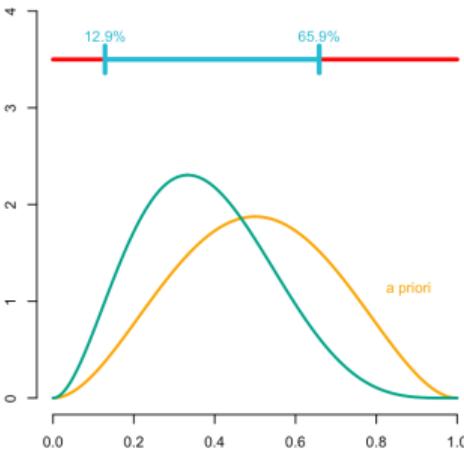
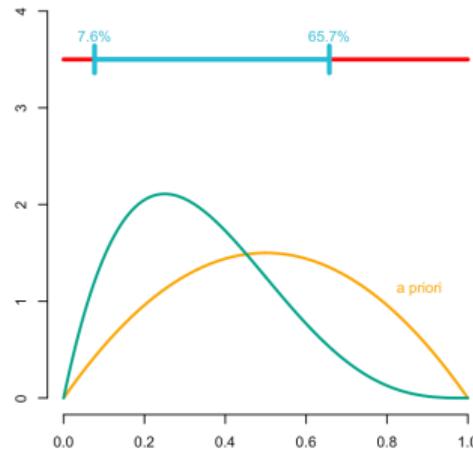
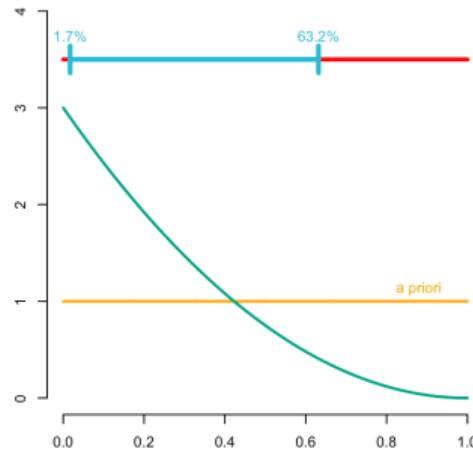


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $x = \{0, 1\}$, $\mathcal{B}(\theta)$

Bayesian approach , $\widehat{\theta}|x \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

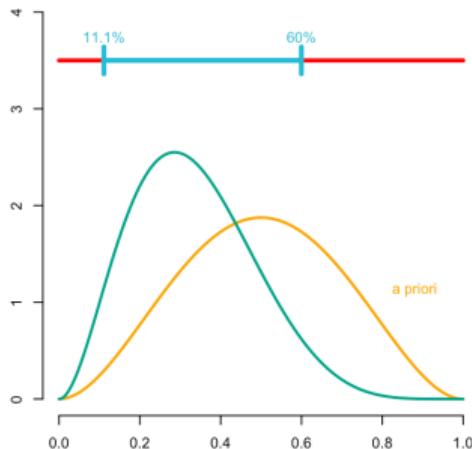
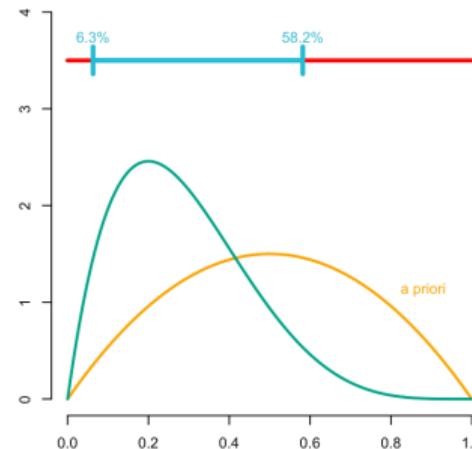
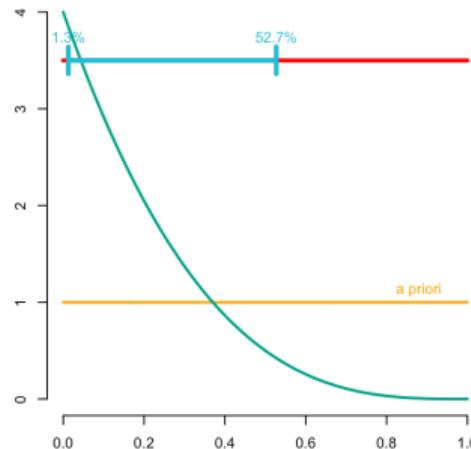


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach , $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

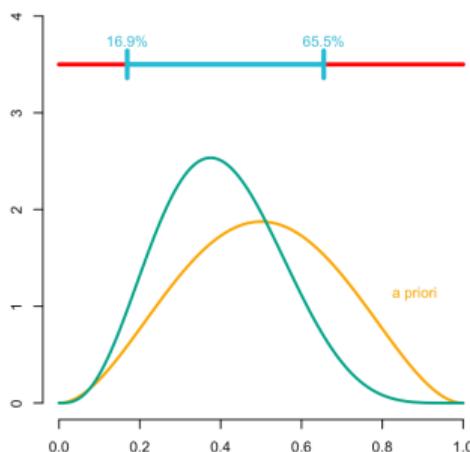
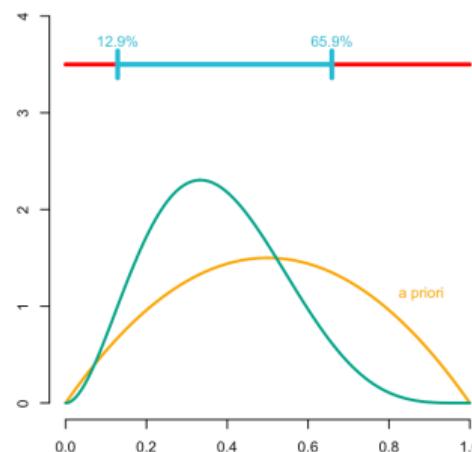
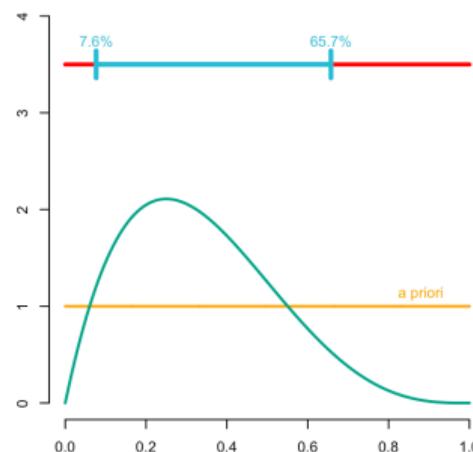


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1\}$, $\mathcal{B}(\theta)$

Bayesian approach , $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

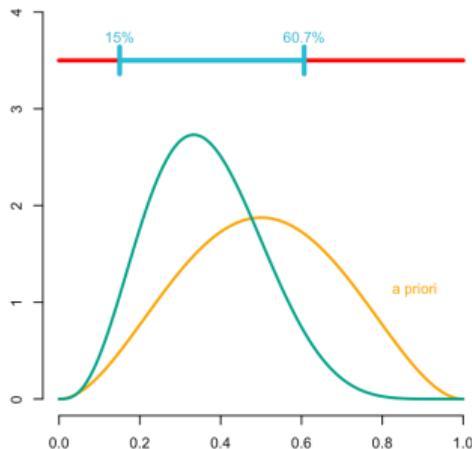
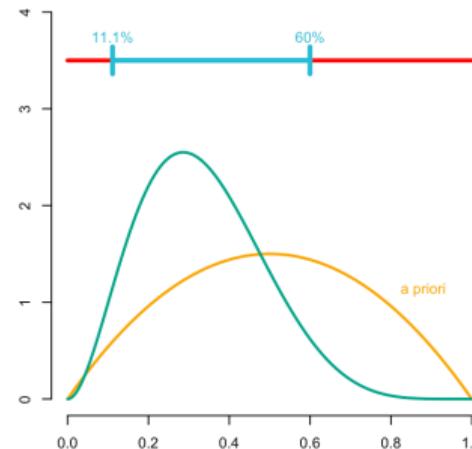
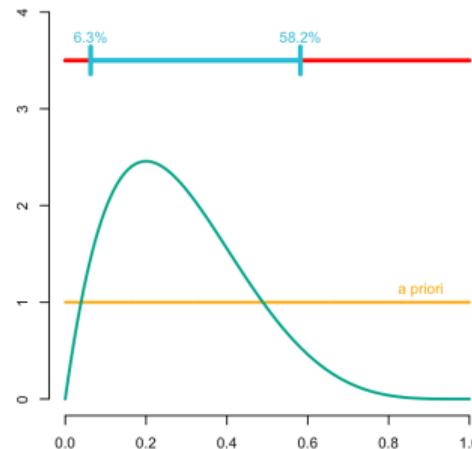


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach , $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

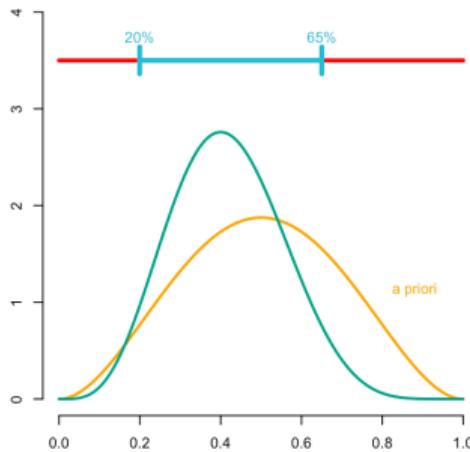
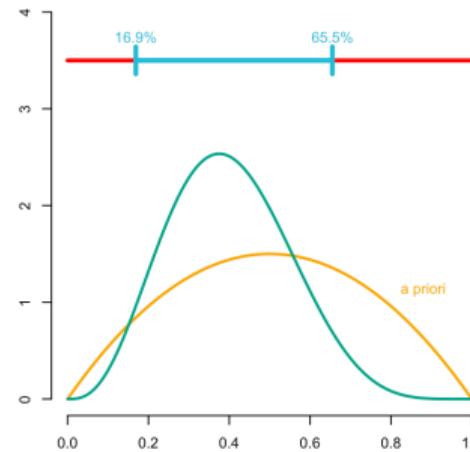
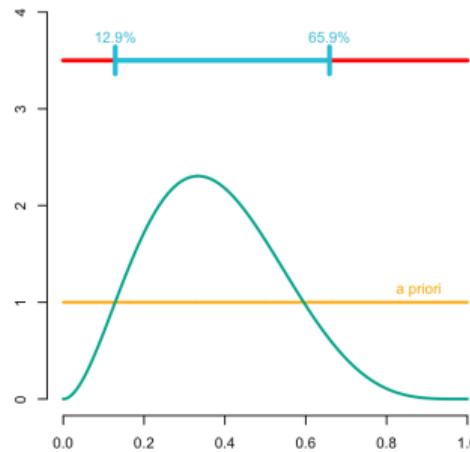


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

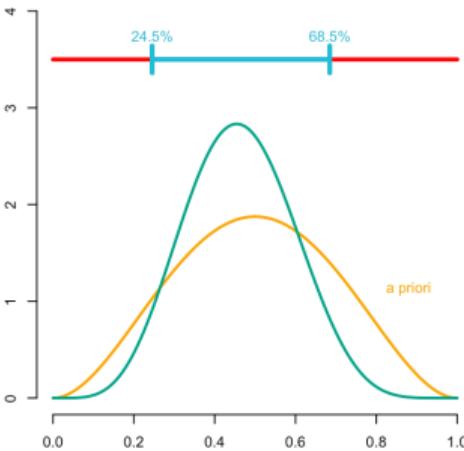
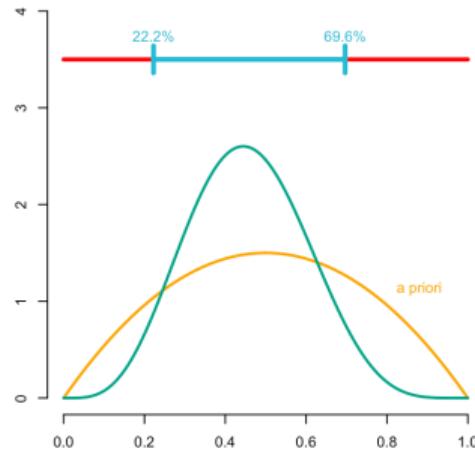
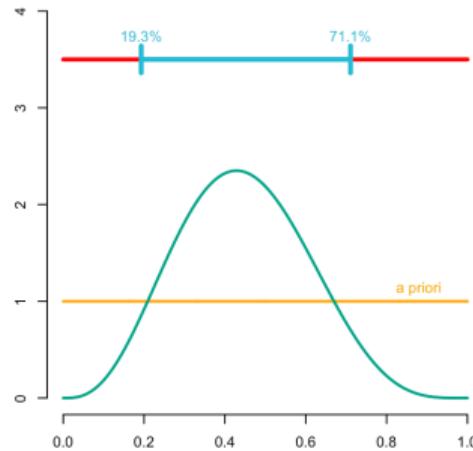


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

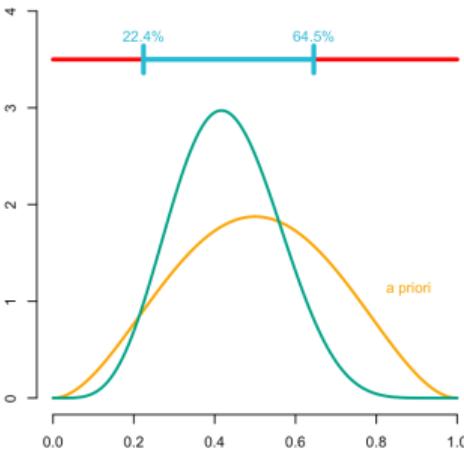
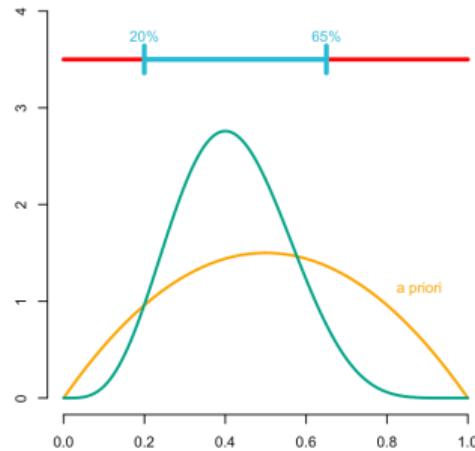
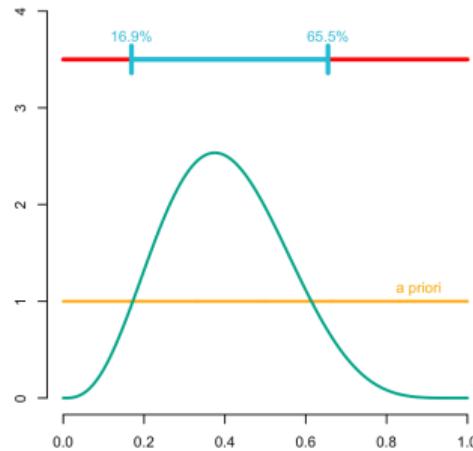


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\hat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

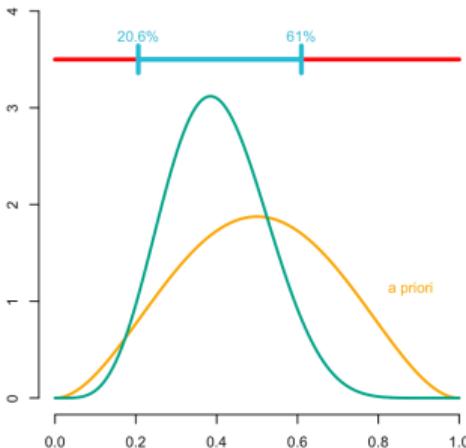
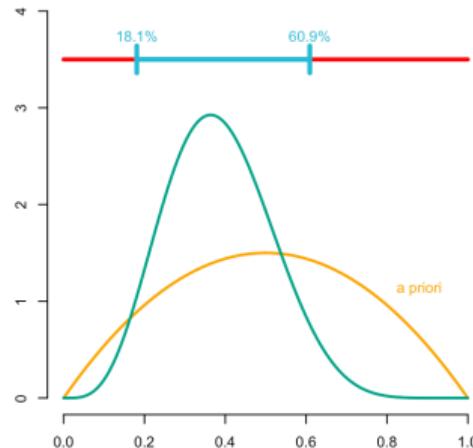
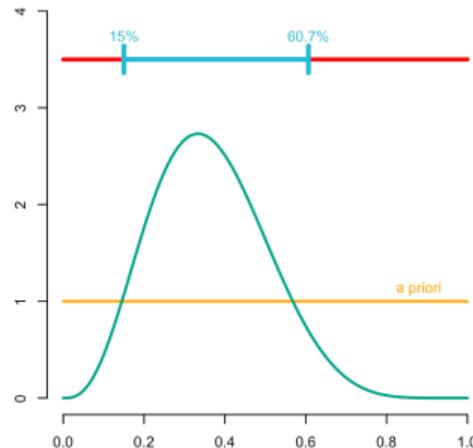


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 1, 0, \textcolor{red}{0}\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

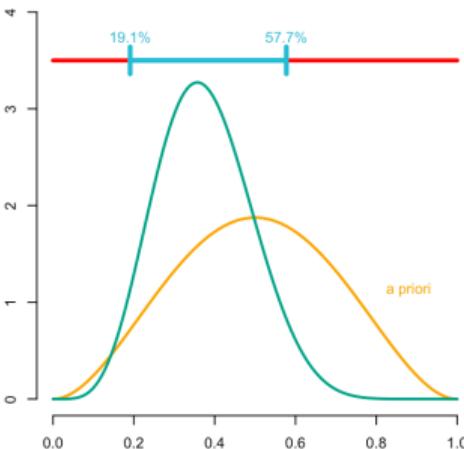
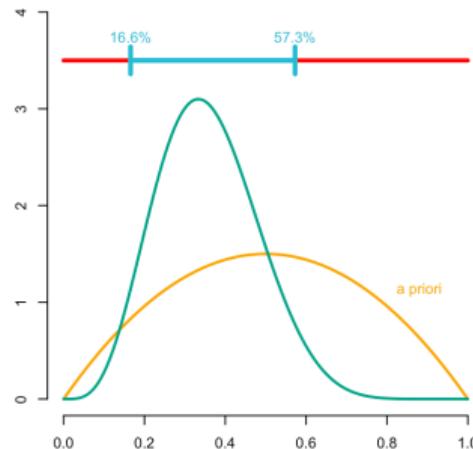
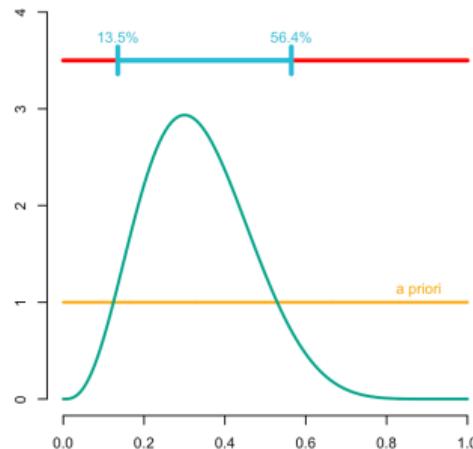


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\hat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

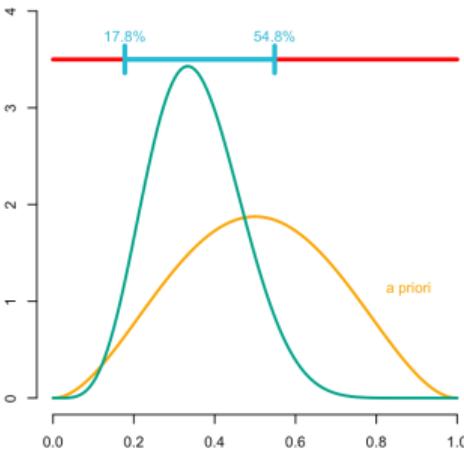
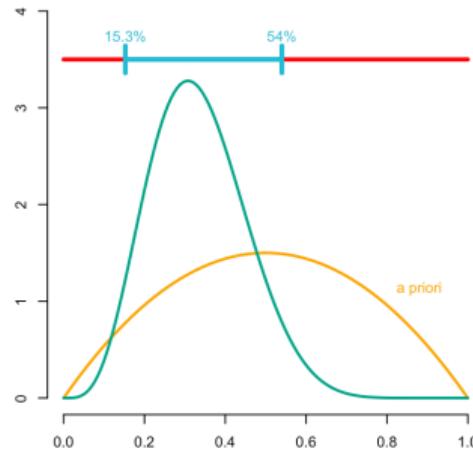
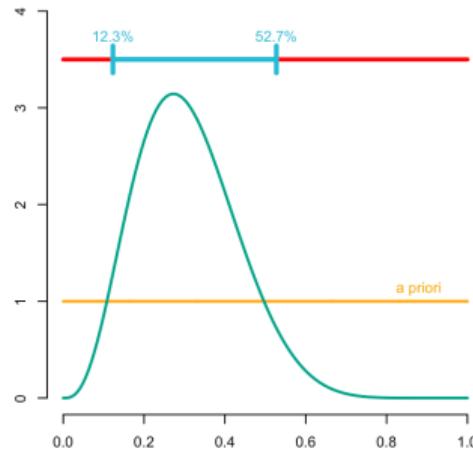


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\hat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

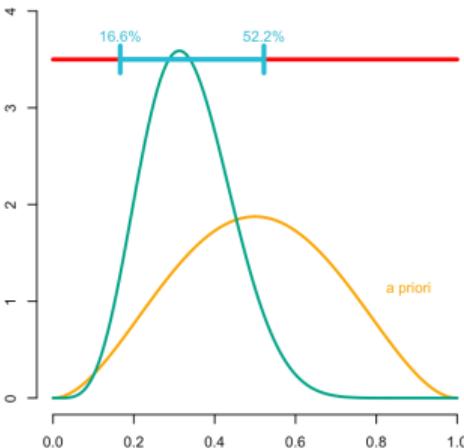
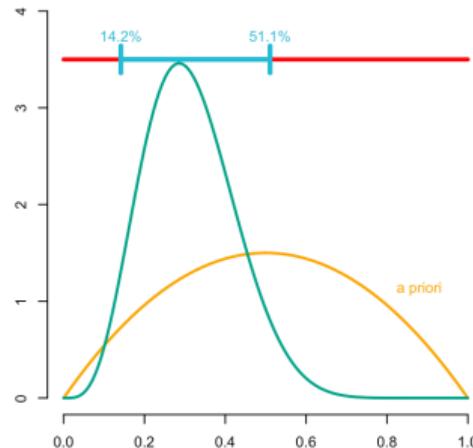
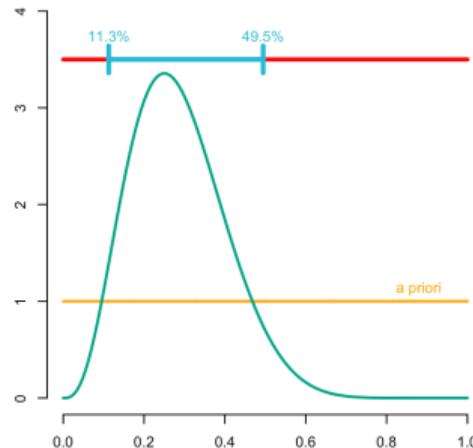


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\hat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

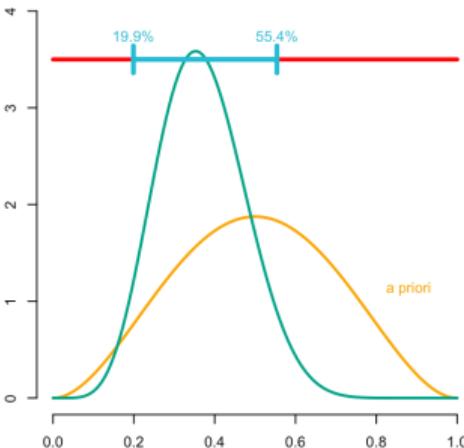
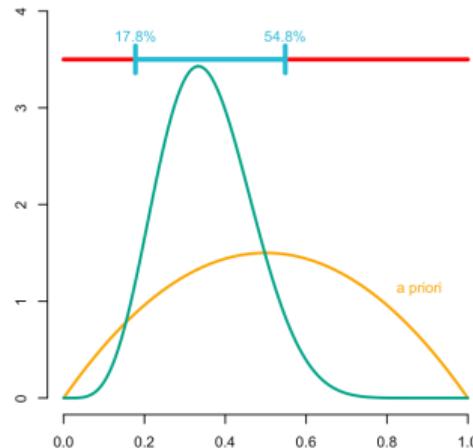
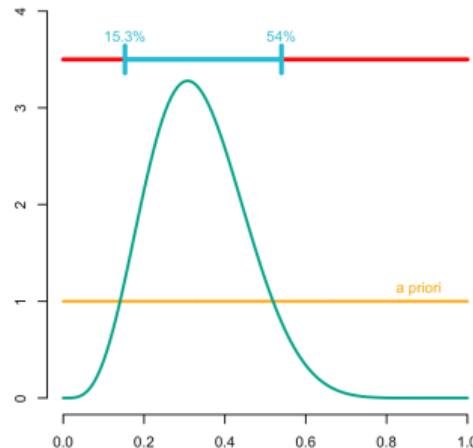


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

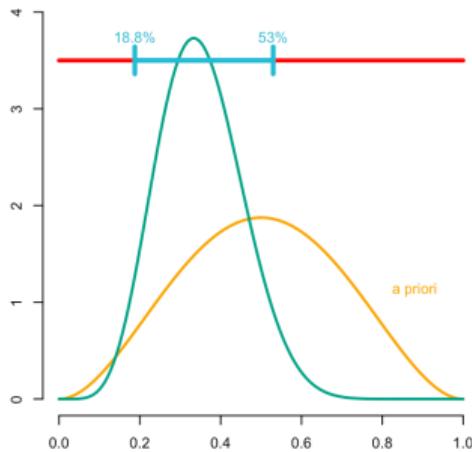
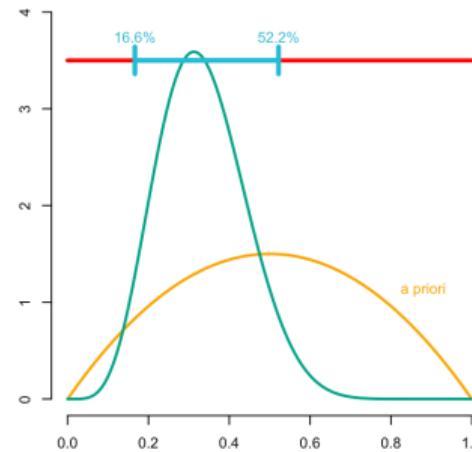
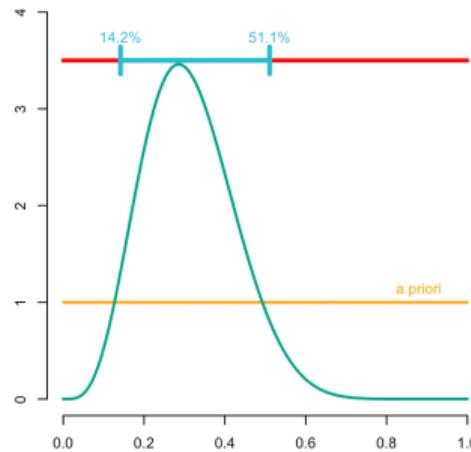


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

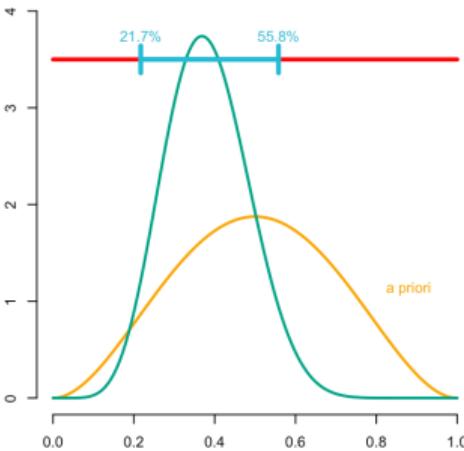
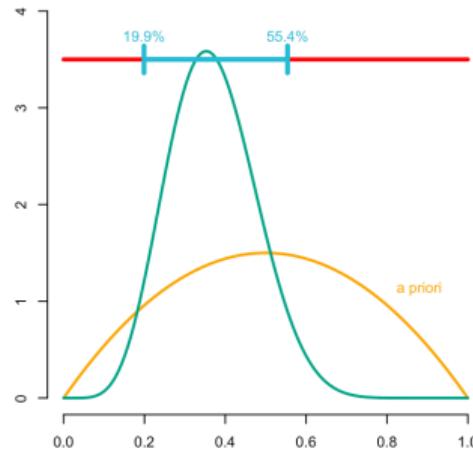
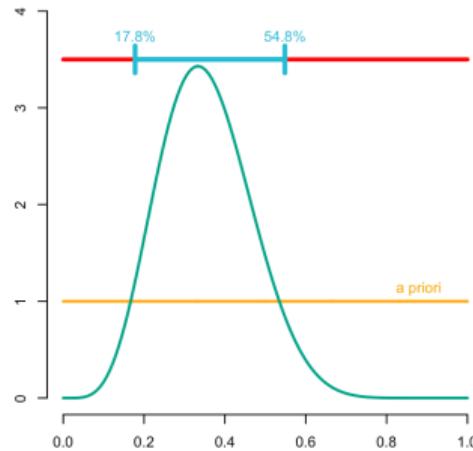


Bayesianism, statistics and calculus

- Posterior distribution

Suppose $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\widehat{\theta} | \mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

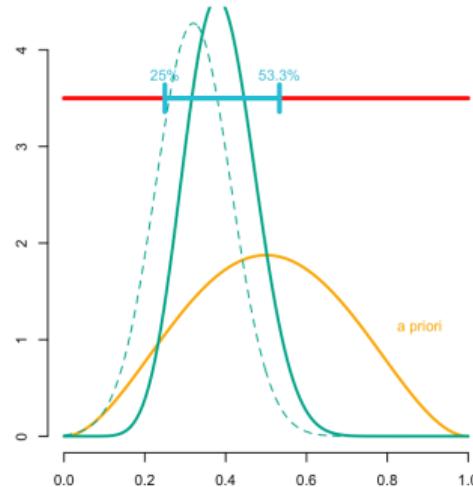
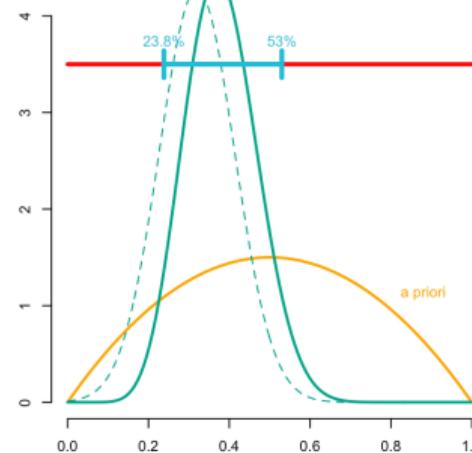
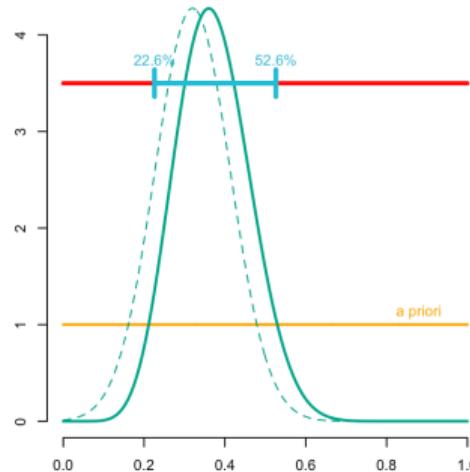


Bayesianism, statistics and calculus

- Posterior distribution

and finally $\mathbf{x} = \{0, 0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0\}$, $\mathcal{B}(\theta)$

Bayesian approach, $\hat{\theta}|\mathbf{x} \sim \text{Beta}(\alpha_0 + s, \beta_0 + n - s)$, $s = \sum_{i=1}^n x_i$

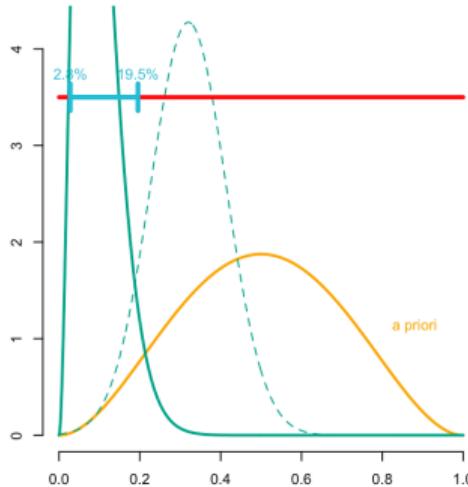
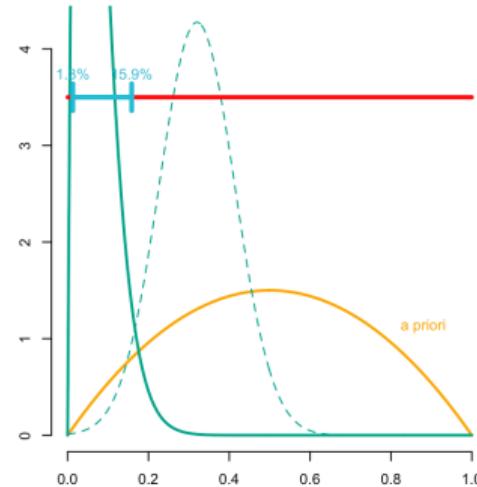
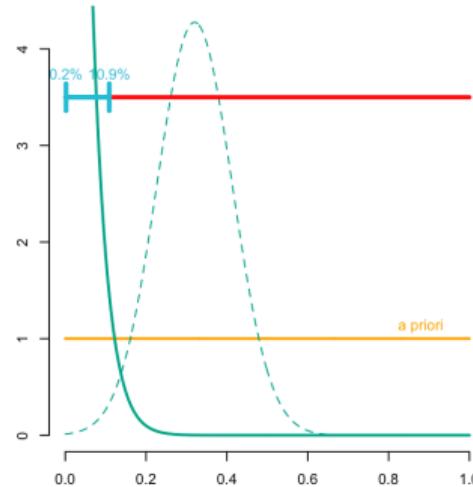


Bayesianism, statistics and calculus

- Posterior distribution

What if $x = \{0, 0\}$, $\mathcal{B}(\theta)$?

Bayesian approach, $\hat{\theta}|x \sim \text{Beta}(\alpha_0, \beta_0 + n)$, since $\sum_{i=1}^n x_i = 0$



Bayesianism, statistics and calculus

- From the distribution to the estimator

$$\begin{cases} \text{posterior average} & \hat{\theta} = \mathbb{E}[\theta|\mathcal{D}] \\ \text{maximum a posteriori (MAP)} & \hat{\theta} = \max \{\pi(\theta|\mathcal{D})\} \text{ i.e. the mode} \end{cases}$$

The average posterior is also the solution of the problem

$$\hat{\theta} = \operatorname{argmin}_{\tau} \left\{ \mathbb{E}[(\theta - \tau)^2 | \mathcal{D}] \right\} = \operatorname{argmin}_{\tau} \left\{ \int (\theta - \tau)^2 \pi(\theta | \mathcal{D}) d\theta \right\}$$

- "confidence interval" or "credibility interval"

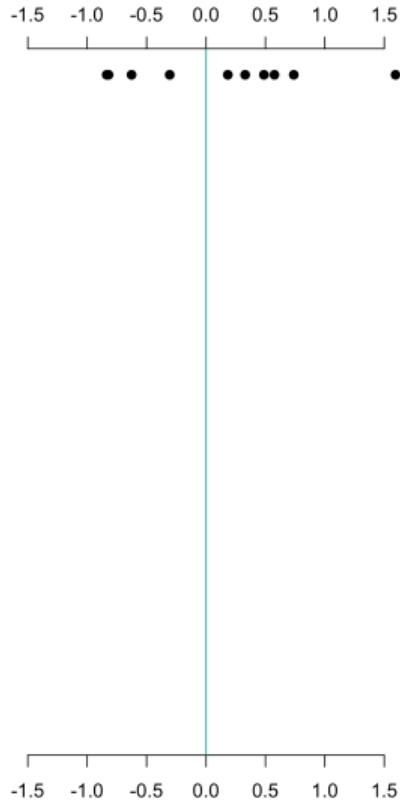
For the confidence interval, we look for $[\hat{a}_{\mathcal{D}}, \hat{b}_{\mathcal{D}}]$ such that $\mathbb{P}[\theta \in [\hat{a}_{\mathcal{D}}, \hat{b}_{\mathcal{D}}]] \geq 95\%$.

For the credibility interval, we look for $[a, b]$ such that $\mathbb{P}[\theta \in [a, b] | \mathcal{D}] \geq 95\%$.

Bayesianism, statistics and calculus

- "confidence interval"

Suppose $\mathcal{D} = \{x_1, \dots, x_n\}$, $X_i \sim \mathcal{N}(\theta, \sigma^2)$
(here $\theta = 0$)

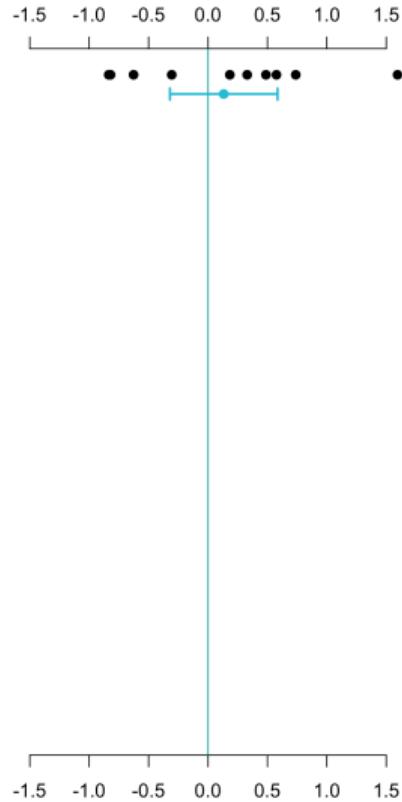


Bayesianism, statistics and calculus

- "confidence interval"

Suppose $\mathcal{D} = \{x_1, \dots, x_n\}$, $X_i \sim \mathcal{N}(\theta, \sigma^2)$
(here $\theta = 0$)

Consider $[a, b] = \left[\bar{x} \pm q_\alpha \frac{\hat{\sigma}}{\sqrt{n}} \right]$



Bayesianism, statistics and calculus

- "confidence interval"

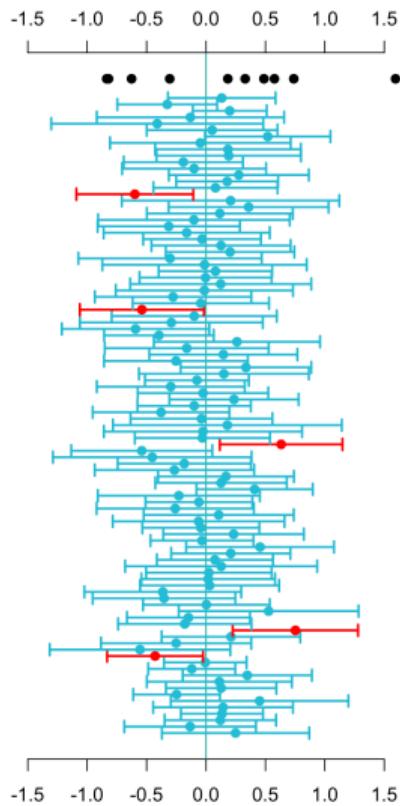
Suppose $\mathcal{D} = \{x_1, \dots, x_n\}$, $X_i \sim \mathcal{N}(\theta, \sigma^2)$
(here $\theta = 0$)

Consider $[a, b] = \left[\bar{x} \pm q_\alpha \frac{\hat{\sigma}}{\sqrt{n}} \right]$

Generate $\mathcal{D}' = \{x'_1, \dots, x'_n\}$ from $\mathcal{N}(\theta, \sigma^2)$, we want

$$\mathbb{P} \left[\theta \notin \left[\bar{x}' \pm q_\alpha \frac{\hat{\sigma}}{\sqrt{n}} \right] \right] \approx \alpha$$

interpreted as a frequency, and repeating the experience.
Here, $\alpha = 5\%$: in 5% of the simulations, 0 is not in $[a, b]$.

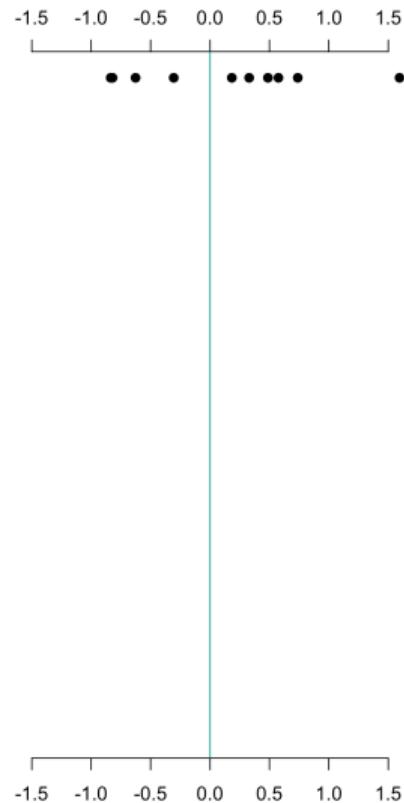


Bayesianism, statistics and calculus

- "credibility interval"

Suppose $\mathcal{D} = \{x_1, \dots, x_n\}$, $X_i \sim \mathcal{N}(\theta, \sigma^2)$

Consider some prior distribution $\pi(\cdot)$ for θ

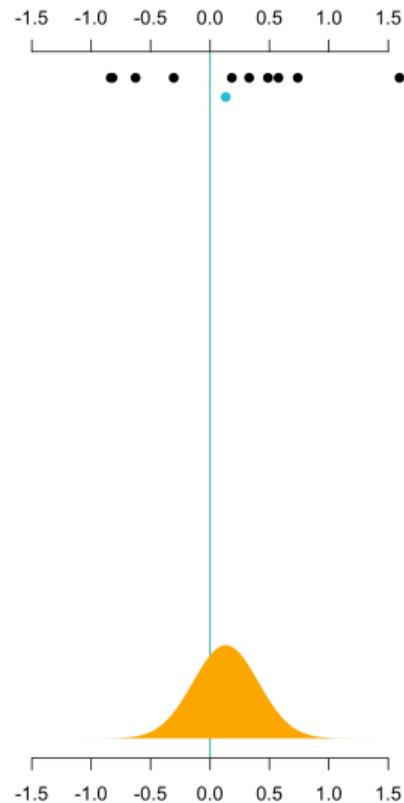


Bayesianism, statistics and calculus

- "credibility interval"

Suppose $\mathcal{D} = \{x_1, \dots, x_n\}$, $X_i \sim \mathcal{N}(\theta, \sigma^2)$

Consider some prior distribution $\pi(\cdot)$ for θ
and $\pi(\cdot | \mathcal{D})$ is the posterior distribution
(potentially complicated)



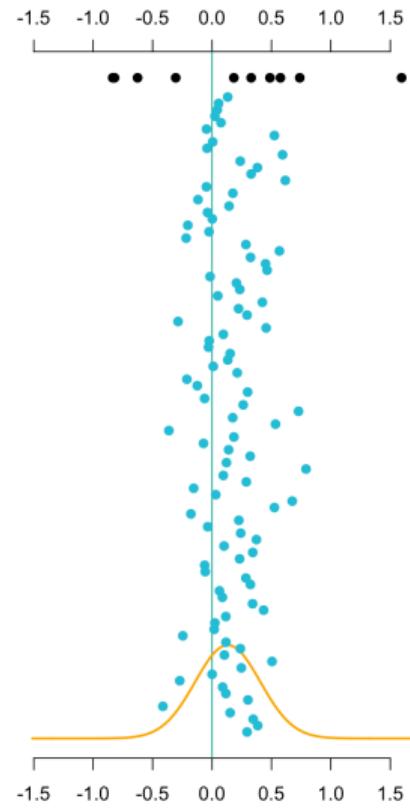
Bayesianism, statistics and calculus

- "credibility interval"

Suppose $\mathcal{D} = \{x_1, \dots, x_n\}$, $X_i \sim \mathcal{N}(\theta, \sigma^2)$

Consider some prior distribution $\pi(\cdot)$ for θ
and $\pi(\cdot|\mathcal{D})$ is the posterior distribution
(potentially complicated)

Suppose we generate $\tilde{\theta}_1, \dots, \tilde{\theta}_k$ given $\pi(\cdot|\mathcal{D})$.



Bayesianism, statistics and calculus

- "credibility interval"

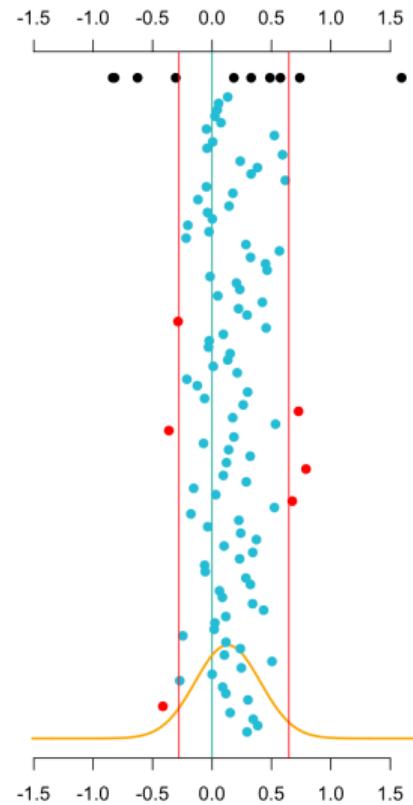
Suppose $\mathcal{D} = \{x_1, \dots, x_n\}$, $X_i \sim \mathcal{N}(\theta, \sigma^2)$

Consider some prior distribution $\pi(\cdot)$ for θ
and $\pi(\cdot|\mathcal{D})$ is the posterior distribution
(potentially complicated)

Suppose we generate $\tilde{\theta}_1, \dots, \tilde{\theta}_k$ given $\pi(\cdot|\mathcal{D})$.

Consider

$$\begin{cases} a = \hat{\Pi}^{-1}(\alpha/2|\mathcal{D}) \text{ quantile with level } \alpha/2 \\ b = \hat{\Pi}^{-1}(1 - \alpha/2|\mathcal{D}) \text{ quantile with level } 1 - \alpha/2 \end{cases}$$



Bayesianism, statistics and calculus

- "credibility interval"

Suppose $\mathcal{D} = \{x_1, \dots, x_n\}$, $X_i \sim \mathcal{N}(\theta, \sigma^2)$

Consider some prior distribution $\pi(\cdot)$ for θ
and $\pi(\cdot | \mathcal{D})$ is the posterior distribution
(potentially complicated)

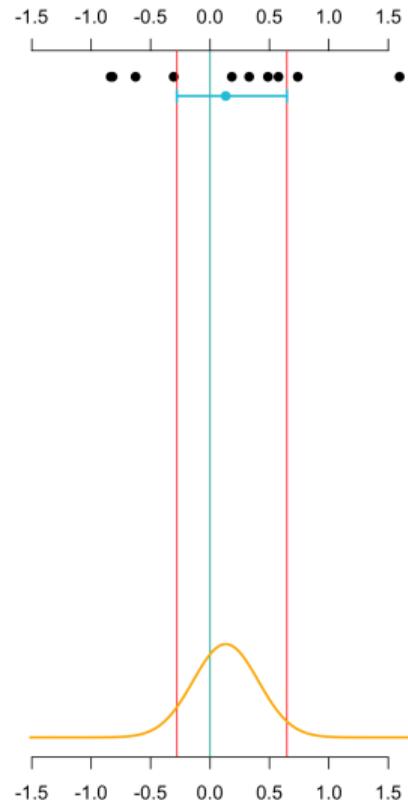
Suppose we generate $\tilde{\theta}_1, \dots, \tilde{\theta}_k$ given $\pi(\cdot | \mathcal{D})$.

Consider

$$\begin{cases} a = \hat{\Pi}^{-1}(\alpha/2 | \mathcal{D}) \text{ quantile with level } \alpha/2 \\ b = \hat{\Pi}^{-1}(1 - \alpha/2 | \mathcal{D}) \text{ quantile with level } 1 - \alpha/2 \end{cases}$$

then

$$\mathbb{P} [\theta \notin [\hat{\Pi}^{-1}(\alpha/2 | \mathcal{D}); \hat{\Pi}^{-1}(1 - \alpha/2 | \mathcal{D})]] \approx \alpha$$



Bayesianism, statistics and calculus

We can also evoke the nonparametric Bayesian modeling, Ferguson (1973). Instead of assuming $X_i \sim f \in \mathcal{F}_\Theta$ where $\mathcal{F}_\Theta = \{f_\theta : \theta \in \Theta\}$, we consider a more general family,

$$X_i \sim f \in \mathcal{F} = \left\{ f : \int_{\mathbb{R}} [f'(y)]^2 dy < \infty \right\}$$

We can always compute a posterior law,

$$\pi(f \in A | \mathcal{D}) = \mathbb{P}(X \in A | \mathcal{D}) = \frac{\int_A \mathcal{L}_n(f) d\pi(f)}{\int_{\mathcal{F}} \mathcal{L}_n(f) d\pi(f)}, \text{ where } \mathcal{L}_n(f) = \prod_{i=1}^n f(x_i)$$

where π is an a prior distribution on \mathcal{F} . Very close to the Pólya urn problems (infinite), to the Chinese restaurant process and to the Dirichlet processes, Blackwell and MacQueen (1973), Ghosh and Ramamoorthi (2003), Orbanz and Teh (2010).

Bayesianism, statistics and calculus

For example, if X_1, \dots, X_n i.i.d. of distribution F . The a priori law π is a Dirichlet process, $D(\alpha, F_0)$, where $F_0 \in \mathcal{F}$ is a prior distribution for X , while α indicates the dispersion around F_0 .

To draw according to $D(\alpha, F_0)$,

- we draw z_1, z_2, \dots according to F_0 ,
- we draw v_1, v_2, \dots according to a Beta law $\mathcal{B}(1, \alpha)$,
- we define iteratively weights, $\omega_1 = v_1$ and $\omega_j = v_j(1 - v_{j-1}) \cdots (1 - v_1)$
- $F(x) = \sum_{j \geq 1} \omega_j \mathbf{1}(x \leq z_j)$

If prior $\pi \sim D(\alpha, F_0)$, then the posterior is, $\pi | \mathcal{D} \sim D(\alpha + n, F_n)$ where

$$F_n = \frac{n}{n + \alpha} \hat{F}_n + \frac{\alpha}{n + \alpha} F_0, \text{ where } \hat{F}_n(x) = \frac{1}{n} \sum_{j=1}^n \mathbf{1}(x \leq x_j)$$

Statistics with Weights

Suppose that instead of $\{x_1, x_2, \dots, x_n\}$ we have $\{(x_1, \omega_1), (x_2, \omega_2), \dots, (x_n, \omega_n)\}$

Suppose that weights are normalized, i.e., $\omega_1 + \omega_2 + \dots + \omega_n = 1$.

- Weighted mean (or barycenter)

$$\bar{x} = \operatorname{argmin} \left\{ \sum_{i=1}^n (x_i - x)^2 \right\} \text{ while } \bar{x}_\omega = \operatorname{argmin} \left\{ \sum_{i=1}^n \omega_i (x_i - x)^2 \right\}$$

- Weighted quantile

Given a dataset $\{x_1, x_2, \dots, x_n\}$ with corresponding weights $\{\omega_1, \omega_2, \dots, \omega_n\}$, the weighted quantile $Q_{\omega,p}$ for a probability p is determined as follows:

- Sort the observations in ascending order, $x_{(1)}, x_{(2)}, \dots, x_{(n)}$ with corresponding reordered weights.
- Compute the cumulative sum of weights: $W_j = \sum_{i=1}^j \omega_{(i)}$

Statistics with Weights

- **Find the quantile position:** The weighted quantile Q_p corresponds to the value where the cumulative weight first reaches or exceeds pW_n , i.e., $W_j \geq pW_n$. The value at this index is the weighted quantile: $Q_p = x_{(j)}$.
- Linear Interpolation (if needed): If pW_n falls between two cumulative weights, we interpolate:

$$Q_{\omega,p} = x_{(j-1)} + \frac{pW_n - W_{j-1}}{W_j - W_{j-1}}(x_{(j)} - x_{(j-1)})$$

```
1 > library(cNORM)
2 > set.seed(1234)
3 x = runif(10)
4 > x
5 [1] 0.114 0.622 0.609 0.623 0.861 0.640 0.009 0.233 0.666 0.514
6 > w = dbeta(x,2,2)
7 > w
8 [1] 0.605 1.410 1.428 1.409 0.718 1.382 0.056 1.071 1.334 1.499
9 weighted.quantile(x, probs = .9, weights = w)
```

Statistics with Weights

10 [1] 0.7555

Independence

Independence

Independence is a fundamental notion in probability theory, as in statistics and the theory of stochastic processes. Two events are independent if, informally speaking, the occurrence of one does not affect the probability of occurrence of the other or, equivalently, does not affect the odds. \square

Definition 1.15: Independence (dimension 2)

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

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

Definition 1.16: Linear Independence (dimension 2)

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

freakonometrics

freakonometrics.hypotheses.org

Independence

Correlation

in the broadest sense, "correlation" may indicate any type of association, in statistics it usually refers to the degree to which a pair of variables are linearly related. \mathbb{W}

Definition 1.17: Correlation (dimension 2), Pearson (1895)

X and Y are two random variables

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

where $\text{Cov}[X, Y] = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$.

Independence

From Cauchy-Schwarz theorem, $-1 \leq \text{Corr}[X, Y] \leq +1$ but those bounds are rarely sharp,

Proposition 1.5: Correlation bounds (dimension 2)

For any random variables X and Y (with finite variances),

$r_{\min} \leq \text{Corr}[X, Y] \leq r_{\max}$, where

$$r_{\min} = \frac{\text{Cov}[F_x^{-1}(U), F_y^{-1}(1 - U)]}{\sqrt{\text{Var}[X] \cdot \text{Var}[Y]}} \quad \text{and} \quad r_{\max} = \frac{\text{Cov}[F_x^{-1}(U), F_y^{-1}(U)]}{\sqrt{\text{Var}[X] \cdot \text{Var}[Y]}}$$

Maximal correlation is obtained when X and Y are comonotonic (minimal correlation when X and $-Y$ are comonotonic).

Related to optimal transport, see also [Knott and Smith \(1984\)](#).

Proposition 1.6

Consider two random variables X and Y . $X \perp\!\!\!\perp Y$ if and only if for any functions $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ and $\psi : \mathbb{R} \rightarrow \mathbb{R}$ (such that the expected values below exist and are well-defined) $\text{Cov}[\varphi(X), \psi(Y)] = 0$, i.e.,

$$\mathbb{E}[\varphi(X) \cdot \psi(Y)] = \mathbb{E}[\varphi(X)] \cdot \mathbb{E}[\psi(Y)].$$

Definition 1.18: Maximal Correlation, HGR

Consider two random variables X and Y ,

$$r^*(X, Y) = \max_{\varphi, \psi} \{ \text{Corr}[\varphi(X), \psi(Y)] \}.$$

Independence

HGR because of Hirschfeld (1935), Gebelein (1941) and Rényi (1959) (also Sarmanov (1958a,b)).

$$r^*(X, Y) = \max_{\varphi \in \mathcal{F}_x, \psi \in \mathcal{G}_y} \mathbb{E}[\varphi(X)\psi(Y)],$$

where

$$\begin{cases} \mathcal{F}_x = \{\varphi : \mathcal{X} \rightarrow \mathbb{R} : \mathbb{E}[\varphi(X)] = 0 \text{ and } \mathbb{E}[\varphi^2(X)] = 1\} \\ \mathcal{G}_y = \{\psi : \mathcal{Y} \rightarrow \mathbb{R} : \mathbb{E}[\psi(Y)] = 0 \text{ and } \mathbb{E}[\psi^2(Y)] = 1\} \end{cases}$$

See either `ccaPP` or `acepack` package,

```
1 > ccaPP::maxCorProj(x = x, y = y, method = "pearson")
2 > corstar = acepack::ace(x = x, y = y)
3 > cor(corstar$tx, corstar$ty)
```

Independence

Proposition 1.7

Consider two random variables X and Y . $X \perp\!\!\!\perp Y$ if and only if $r^*(X, Y) = 0$.

Proof: Given a random variable X , its characteristic function is $\phi_X(t) = \mathbb{E}[e^{itX}]$. Recall that

$$\begin{cases} \phi_X(t) = \phi_Y(t), \forall t \in \mathbb{R} \text{ if and only if } X \stackrel{\mathcal{L}}{=} Y \\ \phi_{X,Y}(s, t) = \mathbb{E}[e^{i(sX+tY)}] = \phi_X(s) \cdot \phi_Y(t), \forall s, t \in \mathbb{R} \text{ if and only if } X \perp\!\!\!\perp Y \end{cases}$$

If $r^*(X, Y) = 0$, let $s, t \in \mathbb{R}$ and consider $\varphi(x) = \phi_X(x) = \mathbb{E}[e^{ixX}]$ and $\psi(y) = \phi_Y(y) = \mathbb{E}[e^{iyY}]$, then $\text{Cov}[e^{isX}, e^{itY}] = \text{Cov}[X'_s, Y'_t] = 0$, i.e. $\mathbb{E}[X'_s Y'_t] = \mathbb{E}[X'_s] \mathbb{E}[Y'_t]$,

$$\underbrace{\mathbb{E}[e^{i(sX+tY)}]}_{\phi_{XY}(s,t)} = \underbrace{\mathbb{E}[e^{isX}] \cdot \mathbb{E}[e^{itY}]}_{\phi_X(s) \cdot \phi_Y(t)}, \forall s, t \in \mathbb{R} \text{ i.e. } X \perp\!\!\!\perp Y.$$

Proposition 1.8

Consider two random variables X and Y such that (X, Y) is a Gaussian vector. Then $r^*(X, Y) = |\text{Corr}[X, Y]|$.

See [Lancaster \(1957, 1958\)](#), and Gauss-Hermite decomposition

$$f(x, y) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(-\frac{x^2 - 2\rho xy + y^2}{2[1-\rho^2]}\right) = \phi(x)\phi(y) \cdot \sum_{i=0}^{\infty} r^i H_i(x)H_i(y)$$

where H_i 's are Hermite polynomial.

Independence

Instead of

$$r^*(X, Y) = \max_{\varphi \in \mathcal{F}_x, \psi \in \mathcal{G}_y} \mathbb{E}[\varphi(X)\psi(Y)],$$

where

$$\begin{cases} \mathcal{F}_x = \{\varphi : \mathcal{X} \rightarrow \mathbb{R} : \mathbb{E}[\varphi(X)] = 0 \text{ and } \mathbb{E}[\varphi^2(X)] = 1\} \\ \mathcal{G}_y = \{\psi : \mathcal{Y} \rightarrow \mathbb{R} : \mathbb{E}[\psi(Y)] = 0 \text{ and } \mathbb{E}[\psi^2(Y)] = 1\} \end{cases}$$

Definition 1.19: Constrained Maximal Correlation, Bach and Jordan (2002), Gretton et al. (2005)

Consider two random variables X and Y , as well as some Hilbert spaces $\bar{\mathcal{F}}_x \subset \mathcal{F}_x$ and $\bar{\mathcal{G}}_y \subset \mathcal{G}_y$,

$$\bar{r}^*(X, Y) = \max_{\varphi \in \bar{\mathcal{F}}_x, \psi \in \bar{\mathcal{G}}_y} \{\text{Corr}[\varphi(X), \psi(Y)]\}.$$

Independence

Kimeldorf and Sampson (1978) and Kimeldorf et al. (1982) suggested to consider for $\bar{\mathcal{F}}_x$ and $\bar{\mathcal{G}}_y$ as subsets of monotone functions.

$$\begin{cases} \bar{\mathcal{F}}_x = \{\varphi \in \mathcal{F}_x : \varphi \text{ monotone}\} \\ \bar{\mathcal{G}}_y = \{\psi \in \mathcal{G}_y : \psi \text{ monotone}\} \end{cases}$$

See Mourier (1953), Hannan (1961), Jensen and Mayer (1977) and Lin (1987).

Independence

Definition 1.20: Linear Independence

In a general context, consider two random vectors \mathbf{X} and \mathbf{Y} , in \mathbb{R}^{d_x} and \mathbb{R}^{d_y} , respectively. $\mathbf{X} \perp \mathbf{Y}$ if and only if for any $\mathbf{a} \in \mathbb{R}^{d_x}$ and $\mathbf{b} \in \mathbb{R}^{d_y}$

$$\text{Cov}[\mathbf{a}^\top \mathbf{X}, \mathbf{b}^\top \mathbf{Y}] = 0.$$

Definition 1.21: Independence

In a general context, consider two random vectors \mathbf{X} and \mathbf{Y} . $\mathbf{X} \perp\!\!\!\perp \mathbf{Y}$ if and only if for any $\mathcal{A} \subset \mathbb{R}^{d_x}$ and $\mathcal{B} \subset \mathbb{R}^{d_y}$,

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

Proposition 1.9: Independence

Consider two random vectors \mathbf{X} and \mathbf{Y} . $\mathbf{X} \perp\!\!\!\perp \mathbf{Y}$ if and only if for any functions $\varphi : \mathbb{R}^{d_x} \rightarrow \mathbb{R}$ and $\psi : \mathbb{R}^{d_y} \rightarrow \mathbb{R}$ (such that the expected values below exist and are well-defined)

$$\mathbb{E}[\varphi(\mathbf{X})\psi(\mathbf{Y})] = \mathbb{E}[\varphi(\mathbf{X})] \cdot \mathbb{E}[\psi(\mathbf{Y})],$$

or equivalently

$$\text{Cov}[\varphi(\mathbf{X}), \psi(\mathbf{Y})] = 0.$$

Independence

Definition 1.22: Mutual Independence

Let $\mathbf{Y} = (Y_1, \dots, Y_k)$ denote some random vector. All components of \mathbf{Y} are (mutually) independent if for any $\mathcal{A}_1, \dots, \mathcal{A}_k \subset \mathbb{R}$

$$\mathbb{P}\left[\{(Y_1, \dots, Y_k) \in \bigcap_{i=1}^k \mathcal{A}_i\}\right] = \prod_{i=1}^k \mathbb{P}[\{Y_i \in \mathcal{A}_i\}].$$

Definition 1.23: Conditional Independence (dimension 2)

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

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

Definition 1.24: Conditional Independence

In a general context, consider three random vectors \mathbf{X} , \mathbf{Y} and \mathbf{Z} . $(\mathbf{X} \perp\!\!\!\perp \mathbf{Y})|\mathbf{Z}$ if and only if for any $\mathcal{A} \subset \mathbb{R}^{d_x}$, $\mathcal{B} \subset \mathbb{R}^{d_y}$ and $\mathcal{C} \subset \mathbb{R}^{d_z}$,

$$\mathbb{P}[\{\mathbf{X} \in \mathcal{A}\} \cap \{\mathbf{Y} \in \mathcal{B}\} | \mathbf{Z} \in \mathcal{C}] = \mathbb{P}[\{\mathbf{X} \in \mathcal{A}\} | \mathbf{Z} \in \mathcal{C}] \cdot \mathbb{P}[\{\mathbf{Y} \in \mathcal{B}\} | \mathbf{Z} \in \mathcal{C}].$$

Proposition 1.10

Consider three random variables X , Y , and Z . If $X \perp Z$ and $Y \perp Z$, then $aX + bY \perp Z$, for any $a, b \in \mathbb{R}$.

Independence

Proposition 1.11: $X \perp Z, Y \perp Z \not\Rightarrow \psi(X, Y) \perp Z$

Consider three random variables X , Y , and Z . If $X \perp Z$ and $Y \perp Z$, it does not imply that $\psi(X, Y) \perp Z$, for any $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$.

$$(X, Y, Z) = \begin{cases} (0, 0, 0) & \text{with probability } 1/4, \\ (0, 1, 1) & \text{with probability } 1/4, \\ (1, 0, 1) & \text{with probability } 1/4, \\ (1, 1, 0) & \text{with probability } 1/4. \end{cases}$$

Proposition 1.12

Consider a random vector \mathbf{X} in \mathbb{R}^k , and a random variable Z .
 $\mathbf{X} \perp Z$ does not imply that $\psi(\mathbf{X}) \perp Z$, for any $\psi : \mathbb{R}^k \rightarrow \mathbb{R}$.

Proposition 1.13

Consider three random variables X , Y , and Z . Even if $X \perp\!\!\!\perp Z$ and $Y \perp\!\!\!\perp Z$, it does not imply either that $\psi(X, Y) \perp Z$ or that $\psi(X, Y) \perp\!\!\!\perp Z$, for any $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$.

Proposition 1.14

Consider a random vector \mathbf{X} in \mathbb{R}^k , and a random variable Z .

$\mathbf{X} \perp\!\!\!\perp Z$ does not imply either that $\psi(\mathbf{X}) \perp Z$ or $\psi(\mathbf{X}) \perp\!\!\!\perp Z$, for any $\psi : \mathbb{R}^k \rightarrow \mathbb{R}$.

Random Numbers

- In Machine Learning, simulations help estimate quantities that are hard to compute directly.
- **Monte Carlo Simulation** is a popular method for approximating complex integrals, probabilities, and expectations.
- Simulations are essential in probabilistic modeling, where we rely on random sampling to estimate model parameters and distributions.
- Example: Simulating the expected value of a random variable when analytical solutions are not feasible.
- **Monte Carlo Simulation:** A statistical technique using random sampling to approximate solutions to problems that may be deterministic in nature.
- Typical application in ML:
 - Estimating integrals, expectations, and probabilities.
 - Sampling from a probability distribution.
 - Bayesian Inference and Markov Chain Monte Carlo (MCMC) methods.

Random Numbers

- Monte Carlo methods are useful when an exact solution is computationally expensive or analytically impossible.
- **Importance Sampling:** A technique to improve the efficiency of Monte Carlo simulations.
- Instead of directly sampling from the distribution p , we sample from an easier-to-sample distribution q and re-weight the samples:

$$\mathbb{E}_p[h(x)] = \mathbb{E}_q \left[\frac{p(x)}{q(x)} h(x) \right]$$

- This is useful when the target distribution is difficult to sample from, but we can easily sample from a different distribution.
- Example: Importance sampling is used in rare-event simulation, such as estimating the tail probabilities of heavy-tailed distributions.

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities)

BY-NC 4.0 145 / 335

Random Numbers

- In Bayesian Inference, we often need to calculate posterior distributions:

$$p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta) \cdot p(\theta)}{p(\mathcal{D})}$$

- Exact computation of the posterior may not be possible, especially when the posterior is high-dimensional.
- Monte Carlo methods (e.g., [Markov Chain Monte Carlo \(MCMC\)](#)) can be used to approximate the posterior by generating samples.
- Popular MCMC algorithms:
 - Metropolis-Hastings
 - Gibbs sampling
- MCMC helps to efficiently sample from complex distributions, critical in probabilistic modeling.

Random Numbers

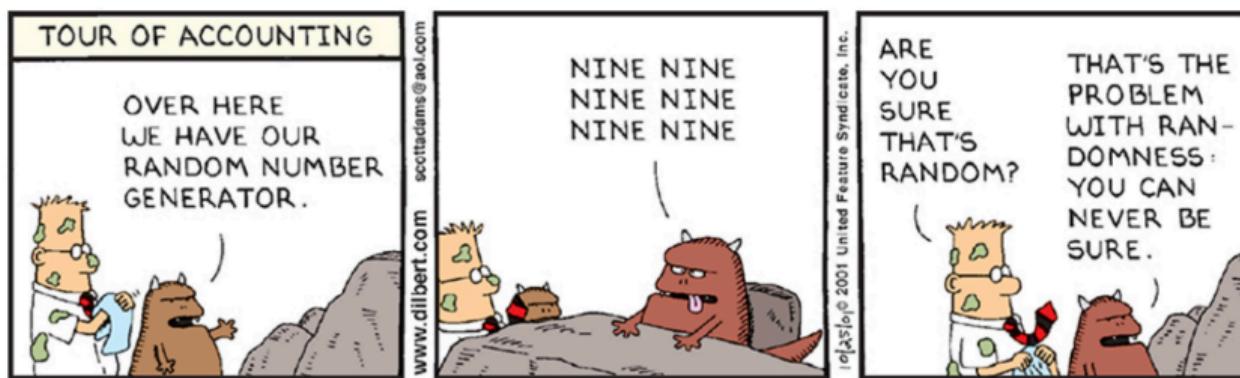
- **Bayesian Machine Learning:** Sampling from posterior distributions to make predictions or optimize hyperparameters.
- **Reinforcement Learning (RL):** Simulating environments for agent training.
- **Uncertainty Quantification:** Simulating and modeling uncertainties in predictions, important for robustness in ML models.
- **Model Evaluation:** Using simulations to estimate model performance under various conditions (cross-validation, bootstrap).
- Example: In RL, Monte Carlo methods help simulate environments and estimate the expected reward for different actions.
- **Computational Cost:** Monte Carlo methods can be computationally expensive, especially for high-dimensional problems.
- **Convergence Rate:** In some cases, Monte Carlo estimates converge slowly, requiring a large number of samples.

Random Numbers

- **Choice of Proposal Distribution (for Importance Sampling):** Poor choice of proposal distribution leads to inefficiency and high variance in estimates.
- **Curse of Dimensionality:** In high-dimensional spaces, the number of samples required grows exponentially.
- Despite these challenges, simulations are indispensable in modern ML for probabilistic modeling, optimization, and decision-making under uncertainty.

Random Numbers

Random number generation is a process by which, often by means of a random number generator (RNG), a sequence of numbers or symbols is generated that cannot be reasonably predicted better than by random chance. In common understanding, "1 2 3 4 5" is not as random as "3 5 2 1 4" and certainly not as random as "47 88 1 32 41" but "we can't say authoritatively that the first sequence is not random ... it could have been generated by chance." W



Monte Carlo

Monte Carlo methods, or Monte Carlo experiments, are a broad class of computational algorithms that rely on repeated random sampling to obtain numerical results. The underlying concept is to use randomness to solve problems that might be deterministic in principle. W

The underlying idea is that if $\mathbf{X}_1, \dots, \mathbf{X}_n, \dots$ are i.i.d. with density p , from the law of large numbers,

$$\frac{1}{n} \sum_{i=1}^n g(\mathbf{X}_i) \xrightarrow{a.s.} \mathbb{E}_p[g(\mathbf{X})] = \int p(\mathbf{x})g(\mathbf{x})d\mathbf{x}.$$

Non-unique decomposition

$$\int a(\mathbf{x})b(\mathbf{x})c(\mathbf{x})d\mathbf{x} = \mathbb{E}_a[bc(\mathbf{X})] \text{ or } \mathbb{E}_{ab}[c(\mathbf{X})] \text{ or } \dots$$

Importance sampling

if $\mathbf{X}_1, \dots, \mathbf{X}_n, \dots$ are i.i.d. with density p , from the law of large numbers,

$$\frac{1}{n} \sum_{i=1}^n g(\mathbf{X}_i) \xrightarrow{a.s.} \mathbb{E}_p[g(\mathbf{X})] = \int p(x)g(x)dx.$$

...

$$p(x) g(x) dx = q(x) \frac{p(x)}{q(x)} g(x) dx$$

Might be easier to sample from q than from p

If $\mathbf{X}_1, \dots, \mathbf{X}_n, \dots$ are i.i.d. with density q , from the law of large numbers,

$$\frac{1}{n} \sum_{i=1}^n \frac{p(\mathbf{X}_i)}{q(\mathbf{X}_i)} g(\mathbf{X}_i) \xrightarrow{a.s.} \mathbb{E}_q \left[\frac{p(\mathbf{X}_i)}{q(\mathbf{X}_i)} g(\mathbf{X}) \right] = \mathbb{E}_p[g(\mathbf{X})] = \int p(x)g(x)dx.$$

Monte Carlo Markov Chain

From the law of large numbers, if $\mathbf{X}_1, \dots, \mathbf{X}_n, \dots$ are i.i.d. with density p ,

$$\frac{1}{n} \sum_{i=1}^n g(\mathbf{X}_i) \rightarrow \mathbb{E}_p[g(\mathbf{X})] = \int p(x)g(x)dx, \quad \text{almost surely as } n \rightarrow \infty.$$

but similar results can also be obtained if $\mathbf{X}_1, \dots, \mathbf{X}_n, \dots$ are not i.i.d..

Monte Carlo Markov Chain

Definition 1.25: Markov Chains

A sequence of random variables (X_n) forms a Markov chain if

$$\mathbb{P}(X_{n+1} \in A | X_n, X_{n-1}, \dots, X_0) = \mathbb{P}(X_{n+1} \in A | X_n).$$

If the state space is not finite, the transition probabilities are given by a Markov kernel $P(x, A)$, which satisfies:

- $P(x, A)$ is a probability measure on the state space for each fixed x .
- $P(x, A)$ is measurable in x for each measurable set A .

The stationary distribution (if it exists) π satisfies $\pi(A) = \int P(x, A)\pi(x)dx$ for all measurable sets A .

Interestingly, we have an “ergodic theorem,”

Monte Carlo Markov Chain

A central concern of ergodic theory is the behavior of a dynamical system when it is allowed to run for a long time. \mathbb{W}

Proposition 1.15: Ergodic Theorem

If (X_n) is an irreducible, aperiodic, and positive recurrent Markov chain with stationary distribution π , then for any function g with $\mathbb{E}_\pi[|g(X)|] < \infty$,

$$\frac{1}{n} \sum_{i=1}^n g(X_i) \rightarrow \mathbb{E}_\pi[g(X)] = \int \pi(x)g(x)dx, \quad \text{almost surely as } n \rightarrow \infty.$$

Given a distribution π , there are algorithms that can be used to create a Markov kernel $P(x, A)$ with stationary distribution π , such as Gibbs sampler and Hastings-Metropolis.

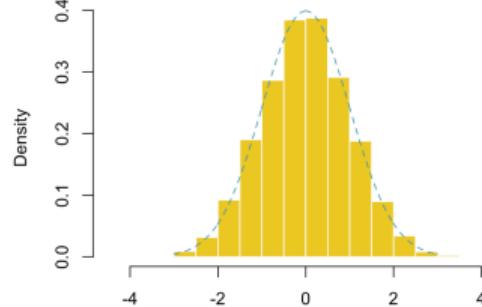
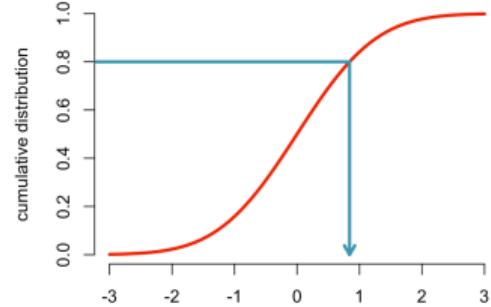
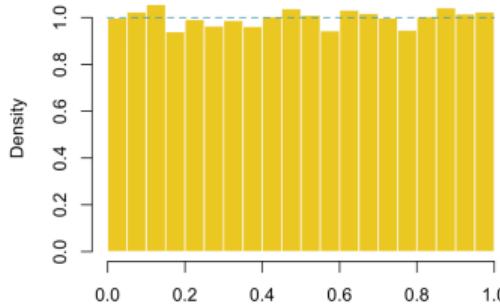
Bootstrap and Resampling

Proposition 1.16: Probability Integral Transform

If F is a cumulative distribution function $\mathbb{R} \rightarrow [0, 1]$, and if $U \sim \mathcal{U}([0, 1])$, then $X := F^{-1}(U)$ has cumulative distribution function F , i.e., $X \sim F$.

Proof Let $x \in \mathbb{R}$, then $\mathbb{P}[X \leq x]$ can be written, with $F^{-1}(u) = \inf \{x \mid F(x) \geq u\}$

$$\mathbb{P}[X \leq x] = \mathbb{P}[F^{-1}(U) \leq x] = \mathbb{P}[F(F^{-1}(U)) \leq F(x)] = \mathbb{P}[U \leq F(x)] = F(x),$$



Bootstrap and Resampling

```
1 > U <- runif(100)
2   [1] 0.29 0.79 0.41 0.88 0.94 0.05 0.53 0.89 0.55 0.46 0.96 0.45 0.68
3   [14] 0.57 0.10 0.90 0.25 0.04 0.33 0.95 0.89 0.69 0.64 0.99 0.66 0.71
4   [27] 0.54 0.59 0.29 0.15 0.96 0.90 0.69 0.80 0.02 0.48 0.76 0.22 0.32
5   [40] 0.23 0.14 0.41 0.41 0.37 0.15 0.14 0.23 0.47 0.27 0.86 0.05 0.44
6   [53] 0.80 0.12 0.56 0.21 0.13 0.75 0.90 0.37 0.67 0.09 0.38 0.27 0.81
7   [66] 0.45 0.81 0.81 0.79 0.44 0.75 0.63 0.71 0.00 0.48 0.22 0.38 0.61
8   [79] 0.35 0.11 0.24 0.67 0.42 0.79 0.10 0.43 0.98 0.89 0.89 0.18 0.13
9   [92] 0.65 0.34 0.66 0.32 0.19 0.78 0.09 0.47 0.51
```

```
1 > qnorm(U)
2   [1] -0.56 0.80 -0.23 1.19 1.56 -1.69 0.07 1.24 0.13 -0.11 1.72
3   [12] -0.12 0.46 0.18 -1.27 1.28 -0.69 -1.73 -0.45 1.69 1.22 0.50
4   [23] 0.36 2.53 0.40 0.55 0.11 0.24 -0.56 -1.05 1.79 1.29 0.50
5   [34] 0.83 -1.97 -0.06 0.70 -0.78 -0.47 -0.73 -1.07 -0.22 -0.22 -0.33
6   [45] -1.03 -1.09 -0.73 -0.09 -0.63 1.07 -1.69 -0.15 0.84 -1.17 0.15
7   [56] -0.82 -1.14 0.68 1.25 -0.32 0.43 -1.31 -0.30 -0.60 0.90 -0.13
8   [67] 0.88 0.89 0.82 -0.15 0.69 0.33 0.55 -3.23 -0.06 -0.77 -0.31
```

Bootstrap and Resampling

Let's consider a sample $\{x_1, \dots, x_n\}$ i.i.d.

of (theoretical) law $F(x) = \mathbb{P}[X_i \leq x]$,

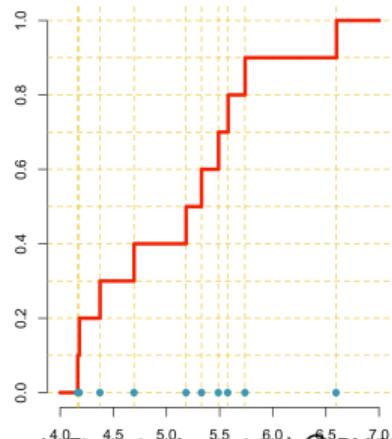
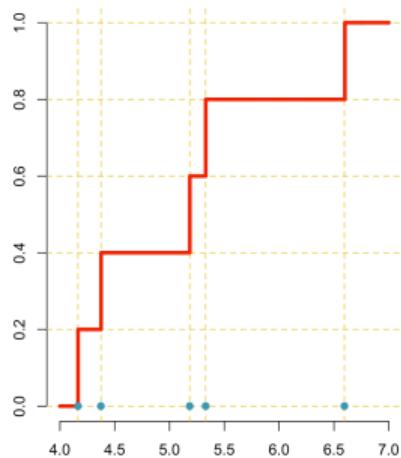
The empirical distribution function is

$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(x_i \leq x), \quad x \in \mathbb{R}$$

Glivenko-Cantelli theorem: $\widehat{F}_n \rightarrow F$ when $n \rightarrow \infty$,
or more precisely

$$\|\hat{F}_n - F\|_\infty = \sup_{x \in \mathbb{R}} |\hat{F}_n(x) - F(x)| \longrightarrow 0$$

when $n \rightarrow \infty$.



Bootstrap and Resampling

Drawing according to \hat{F}_n^{-1} means resampling in $\{x_1, \dots, x_n\}$ with probability $1/n$ (or with replacement).

```
1 > x
2 [1] 4.164 4.374 5.184 5.330 6.595
3 > Qemp = Vectorize(function(u) sort(U)[ceiling(length(x)*u)])
4 > Qemp(U)
5 [1] 4.37 5.33 5.18 6.60 6.60 4.16 5.18 6.60 5.18 5.18 6.60 5.18 5.33
6 [14] 5.18 4.16 6.60 4.37 4.16 4.37 6.60 6.60 5.33 5.33 6.60 5.33 5.33
7 [27] 5.18 5.18 4.37 4.16 6.60 6.60 5.33 5.33 4.16 5.18 5.33 4.37 4.37
8 [40] 4.37 4.16 5.18 5.18 4.37 4.16 4.16 4.37 5.18 4.37 6.60 4.16 5.18
9 [53] 5.33 4.16 5.18 4.37 4.16 5.33 6.60 4.37 5.33 4.16 4.37 4.37 6.60
10 [66] 5.18 6.60 6.60 5.33 5.18 5.33 5.33 5.33 4.16 5.18 4.37 4.37 5.33
11 [79] 4.37 4.16 4.37 5.33 5.18 5.33 4.16 5.18 6.60 6.60 6.60 4.16 4.16
12 [92] 5.33 4.37 5.33 4.37 4.16 5.33 4.16 5.18 5.18
```

see also Davison and Hinkley (1997)

```
1 > sample(x, size = 100, replace = TRUE)
```

(Linear) Model & Bootstrap (1)

Dataset $\mathcal{D}_n = \{\mathbf{z}_i = (y_i, \mathbf{x}_i)\}, i = 1, \dots, n$.

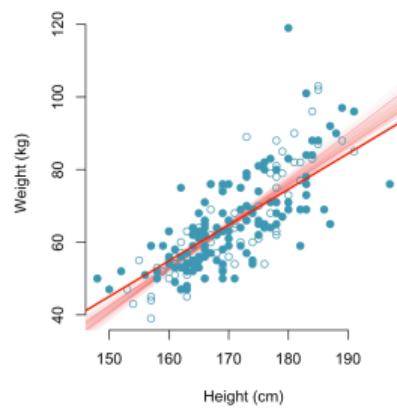
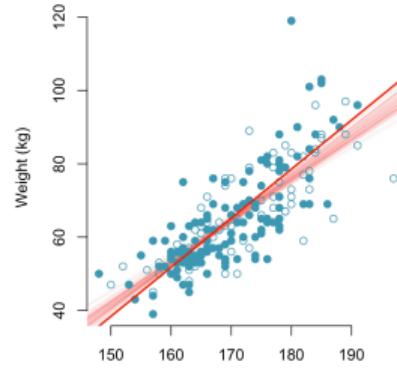
1) approach by **paired sampling**, i.e. $\{\mathbf{z}_1^*, \dots, \mathbf{z}_n^*\}$.

- draw $\{i_1^{(b)}, \dots, i_n^{(b)}\}$ randomly (with replacement) in $\{1, 2, \dots, n\}$
- consider database $(\mathbf{x}_i^{(b)}, y_i^{(b)}) = (\mathbf{x}_{i^{(b)}}, y_{i^{(b)}})$'s, and estimate a parametric model
- let $\hat{\beta}^{(b)}$ denote the estimator of β and $\hat{y}_{n+1}^{(b)}$ some prediction (associated with \mathbf{x}_{n+1})

(loop over $b = 1, \dots, B$)

(Linear) Model & Bootstrap (1)

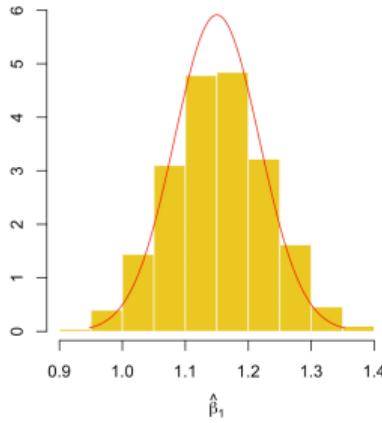
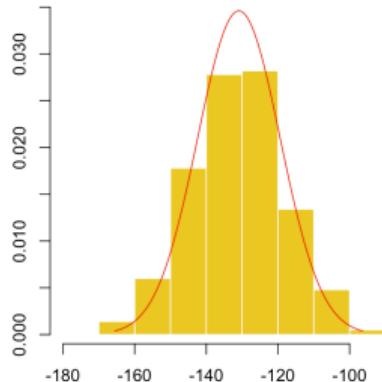
```
1 > Davis = read.table("http://freakonometrics.free.fr/  
2   /Davis.txt")  
3 > BETA = matrix(NA,1000,2)  
4 > for(s in 1:1000){  
5 +   idx = sample(1:nrow(Davis),nrow(Davis),  
6 +                 replace=TRUE)  
7 +   reg_sim = lm(weight~height, data=Davis[idx,])  
8 +   BETA[s,] = reg_sim$coefficients  
9 }  
10 > hist(BETA[,1])  
11 > hist(BETA[,2])
```



(Linear) Model & Bootstrap (1)

or use the `boot` package

```
1 > library(boot)
2 > coef = function(formula, data, indices) {
3 +   d = data[indices,]
4 +   fit = lm(formula, data=d)
5 +   return(coef(fit))
6 + }
7 > results = boot(data=Davis, statistic=coef, R=1000,
8 +                   formula=weight~height)
9 > plot(results, index=1)
10 > plot(results, index=2)
```



(Linear) Model & Bootstrap (2)

2) use **model-based resampling** (or semiparametric residual bootstrap)

- fit a (linear) model, get predictions \hat{y}_i and residuals $\hat{\varepsilon}_i$
- draw $\hat{\varepsilon}_1^{(b)}, \dots, \hat{\varepsilon}_n^{(b)}$ from the original sample $\{\hat{\varepsilon}_1, \hat{\varepsilon}_2, \dots, \hat{\varepsilon}_n\}$
- set $y_i^{(b)} = \hat{y}_i + \hat{\varepsilon}_i^{(b)}$ (from the fitted model)
- from database $(x, y^{(b)}) = (x_i, y_i^{(b)})$, fit a (linear) parametric model
- set $\hat{\beta}^{(b)}$ the estimator of β , and $\hat{y}_{n+1}^{(b)}$ the prediction (associated to x_{n+1})

Note for a simple linear regression, $\hat{\beta}_1^{(b)} = \frac{\sum[x_i - \bar{x}] \cdot y_i^{(b)}}{\sum[x_i - \bar{x}]^2} = \hat{\beta}_1 + \frac{\sum[x_i - \bar{x}] \cdot \hat{\varepsilon}_i^{(b)}}{\sum[x_i - \bar{x}]^2}$ i.e.,
 $E[\hat{\beta}_1^{(b)}] = \hat{\beta}_1$, while $Var[\hat{\beta}_1^{(b)}] = \frac{\sum[x_i - \bar{x}]^2 \cdot Var[\hat{\varepsilon}_i^{(b)}]}{(\sum[x_i - \bar{x}]^2)^2} \sim \frac{\sigma^2}{\sum[x_i - \bar{x}]^2}$.

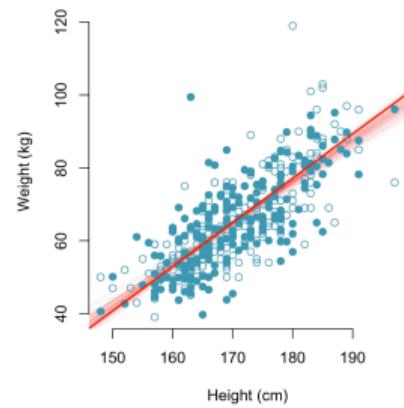
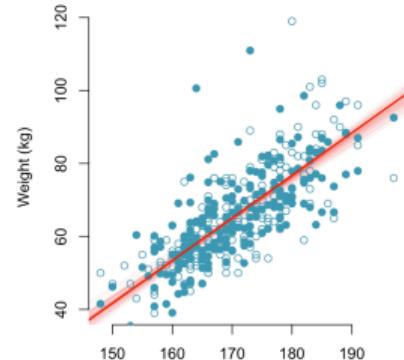
freakonometrics

freakonometrics.hypotheses.org

BY-NC 4.0 162 / 335

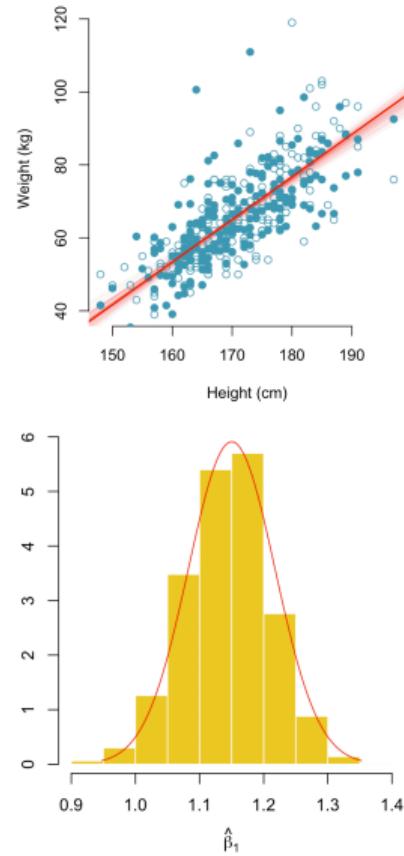
(Linear) Model & Bootstrap (2)

```
1 > BETA = matrix(NA,1000,2)
2 > reg = lm(weight~height, data=Davis)
3 > epsilon = residuals(reg)
4 > for(s in 1:1000){
5 +   eps = sample(epsilon,nrow(Davis),replace=TRUE)
6 +   Davis_s = data.frame(height = Davis$height,
7 +                         weight = predict(reg)+eps)
8 +   reg_sim = lm(weight~height, data=Davis_s)
9 +   BETA[s,] = reg_sim$coefficients
10 + }
11 > hist(BETA[,1])
12 > hist(BETA[,2])
```



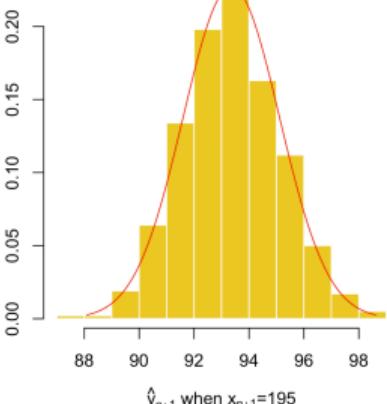
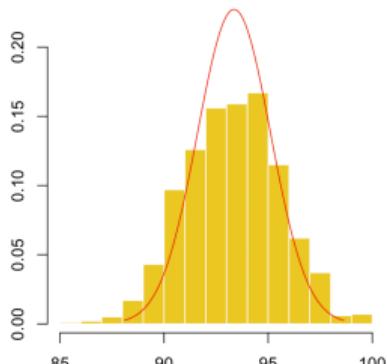
(Linear) Model & Bootstrap (2)

```
1 > BETA = matrix(NA,1000,2)
2 > reg = lm(weight~height, data=Davis)
3 > epsilon = residuals(reg)
4 > for(s in 1:1000){
5 +   eps = sample(epsilon,nrow(Davis),replace=TRUE)
6 +   Davis_s = data.frame(height = Davis$height,
7 +                         weight = predict(reg)+eps)
8 +   reg_sim = lm(weight~height, data=Davis_s)
9 +   BETA[s,] = reg_sim$coefficients
10 + }
11 > hist(BETA[,1])
12 > hist(BETA[,2])
```



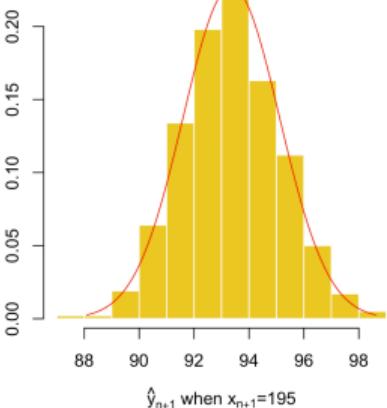
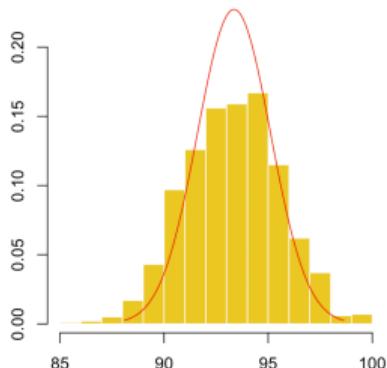
(Linear) Model & Bootstrap (1/2)

```
1 > PRED = matrix(NA,1000,2)
2 > nwDavis = data.frame(height = 195)
3 > for(s in 1:1000){
4 +   idx = sample(1:n,n,replace=TRUE)
5 +   reg_sim = lm(weight~height, data=Davis[idx,])
6 +   PRED[s,1] = predict(reg_sim, newdata=nwDavis)
7 +   eps = sample(epsilon,nrow(Davis),replace=TRUE)
8 +   Davis_s = data.frame(height = Davis$height,
9 +                         weight = predict(reg)+eps)
10 +  reg_sim = lm(weight~height, data=Davis_s)
11 +  PRED[s,2] = predict(reg_sim, newdata=nwDavis)
12 + }
```



(Linear) Model & Bootstrap (2/2)

```
1 > apply(PRED,2,function(x) quantile(x,.025))
2 [1] 89.04203 89.63030
3 > apply(PRED,2,function(x) quantile(x,.975))
4 [1] 97.60345 97.02423
5 > predict(lm(weight~height, data=Davis),
6 +           newdata=nwDavis,interval="confidence",se.
7   fit = TRUE)
8 $fit
9       fit      lwr      upr
10  1 93.35749 89.89571 96.81927
11 $se.fit
12 [1] 1.755451
```



Introduction to Metrics and Similarities

- **Metrics** and **distances** are used to quantify the similarity or dissimilarity between data points.
- They are critical in many machine learning algorithms such as:
 - Clustering (K-means)
 - Classification (K-nearest neighbors, Support Vector Machines)
 - Dimensionality reduction (PCA)
- Commonly used metrics:
 - Euclidean distance
 - Cosine similarity
 - Manhattan distance
 - Mahalanobis distance
- Understanding the properties of different metrics helps choose the right one for a given problem.

Introduction to Metrics and Similarities

- In machine learning, **loss functions** are used to quantify how well a model's predictions match the actual values.
- Loss functions are often based on some notion of "distance" between the predicted and true values, especially in regression and classification problems.
- Key idea: The loss function can be seen as a distance metric that we minimize during model optimization.
- The relationship between distances and losses is crucial in tasks like:
 - Regression: Measuring the difference between predicted and actual values.
 - Classification: Measuring the difference between predicted class probabilities and true labels.
 - Generative Models: Comparing generated data distributions to real data distributions.

Introduction to Metrics and Similarities

- **Squared Euclidean Loss (L2 Loss)**: Measures the squared Euclidean distance between predicted and true values.

$$L_2(\hat{y}, y) = \sum_{i=1}^n \ell_2(\hat{y}_i, y_i) \text{ where } \ell_2(\hat{y}_i, y_i) = (\hat{y}_i - y_i)^2$$

- **Mean Absolute Error (MAE, L1 Loss)**: Measures the absolute distance between predicted and actual values.

$$L_1(\hat{y}, y) = \sum_{i=1}^n \ell_1(\hat{y}_i, y_i) \text{ where } \ell_1(\hat{y}_i, y_i) = |\hat{y}_i - y_i|$$

- **Cross-Entropy Loss**: A loss function based on Kullback-Leibler divergence for classification tasks.

$$L_{\text{CE}}(p, q) = - \sum_i p_i \log(q_i)$$

where p is the true distribution and q is the predicted distribution.

Introduction to Metrics and Similarities

- **Hinge Loss:** Used for Support Vector Machines (SVM), measures how far the predictions are from the decision boundary.

$$\ell_{\text{hinge}}(y, \hat{y}) = \max(0, 1 - y \cdot \hat{y})$$

where y is the true label and \hat{y} is the predicted value.

- Many common loss functions can be interpreted as a distance metric:
 - **L2 loss** (Squared Euclidean loss) is a direct application of Euclidean distance between predicted and true values.
 - **L1 loss** (Mean Absolute Error) corresponds to Manhattan distance, measuring the absolute differences between predictions and true values.
 - **Cross-Entropy loss** can be seen as the Kullback-Leibler (KL) divergence between the true and predicted distributions.
- **Optimization Perspective:**

Introduction to Metrics and Similarities

- Minimizing a loss function is equivalent to minimizing a certain distance between the predicted and true values or distributions.
 - In most cases, we aim to minimize the distance, i.e., minimize the loss, to improve model performance.
- **Connection with Regularization:**
 - Regularization terms, such as L2 regularization (Ridge), are often added to the loss function to prevent overfitting, making the model less sensitive to small fluctuations in the data.
- But one can also consider “distance between distributions”

- In many machine learning and statistical problems, we are interested in comparing probability distributions rather than just individual data points.
- This is crucial in tasks such as:
 - Probabilistic forecasting and predictive modeling.

Introduction to Metrics and Similarities

- Generative models (e.g., GANs, VAEs) where we compare generated vs true distributions.
- Uncertainty quantification and model evaluation.
- Examples:
 - Comparing the predicted distribution of a model with the true distribution of the data.
 - Measuring the discrepancy between two probability distributions, e.g., in model calibration.
- Some key metrics used to measure the "distance" between distributions are:
 - Kullback-Leibler (KL) Divergence
 - Total Variation Distance
 - Wasserstein Distance (Earth Mover's Distance)

Introduction to Metrics and Similarities

- **KL Divergence:** Measures the "distance" between two probability distributions P and Q .

$$D_{KL}(P \parallel Q) = \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

- The KL divergence is not a true distance because it is asymmetric (i.e., $D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$).
 - It measures how much information is lost when Q is used to approximate P .
- Commonly used in:
 - Variational inference.
 - Machine learning models where we need to measure how much a learned distribution differs from the true one.
- Intuition:
 - If $P = Q$, then $D_{KL}(P \parallel Q) = 0$.
 - If P and Q are very different, the divergence becomes large.

Introduction to Metrics and Similarities

- **Wasserstein Distance** (also known as Earth Mover's Distance) measures the minimum amount of work required to transform one distribution into another.
- [•] The first Wasserstein distance between two probability distributions P and Q is defined as:

$$W_1(P, Q) = \inf_{\gamma \in \Gamma(P, Q)} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$

where $\Gamma(P, Q)$ is the set of all couplings between P and Q , and $\|x - y\|$ is the distance between points x and y .

- Unlike KL Divergence, the Wasserstein distance is symmetric and satisfies the triangle inequality, making it a true distance metric.
- Widely used in:
 - Generative Adversarial Networks (GANs), especially Wasserstein GANs (WGANs).

Introduction to Metrics and Similarities

- Comparing distributions with complex structures or in high-dimensional spaces.
- Intuition:
 - It measures the "cost" of transforming one distribution into another, which is especially useful when comparing distributions over different domains or supports.

Vectors

$$\vec{u} = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} = \mathbf{u} \in \mathbb{R}^n$$
$$\mathbf{u} = (u_1 \cdots u_n)^\top$$

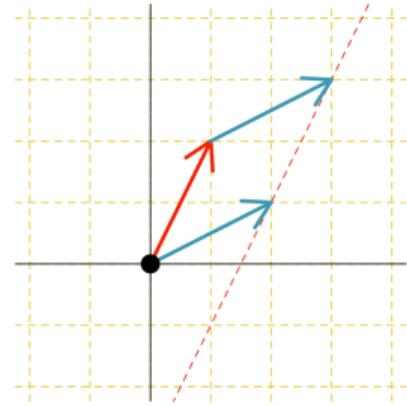
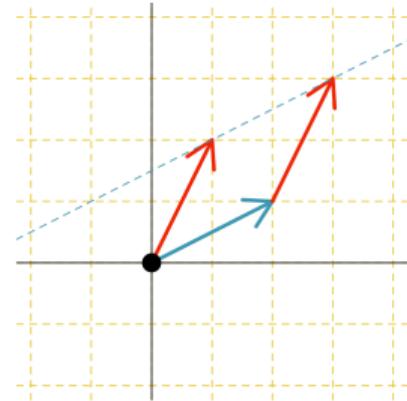
Example :

$\vec{u}_j = (178, 185, 162, 170, \dots, 169)^\top \in \mathbb{R}^{200}$
(vector of **variables**, $n = 200 = \text{sample size}$)

$\vec{v}_i = (178, 74, 1, 0, 14321, 1)^\top \in \mathbb{R}^p$
(vector of **individuals**, $p = 6 = \text{features}$)

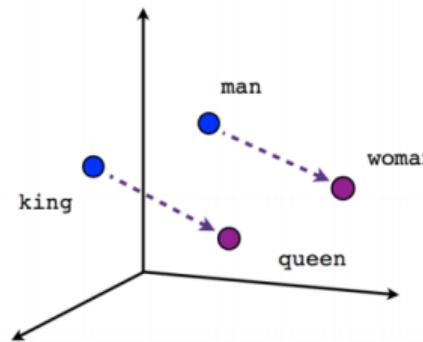
```
1 > u = [178, 185, 162, 170]
```

```
1 > u = c(178, 185, 162, 170)
```

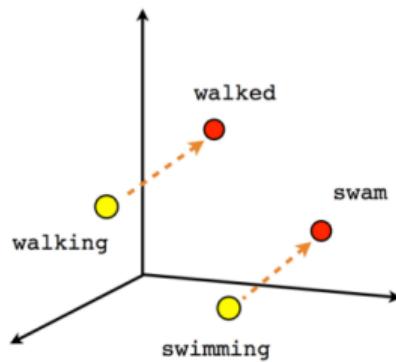


Word2Vec

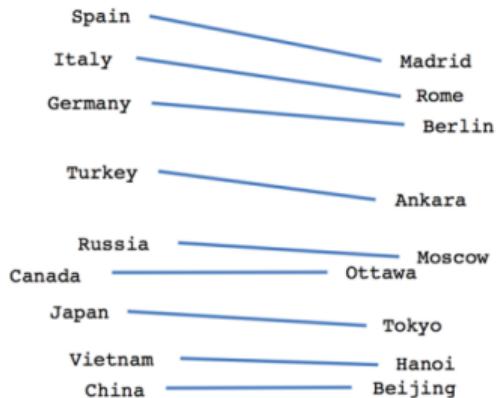
Word2vec has been one of the most popular models used to create word embeddings (words are represented in a large dimensional space... that can be projected into 2 or 3)



Male-Female



Verb tense



Country-Capital

(source <https://www.tensorflow.org>)

Scalar Product & Euclidean Distance

$$\vec{u} + \vec{v} = (u_1 + v_1, \dots, u_n + v_n)^\top \in \mathbb{R}^n$$

Note that $\vec{u} + \vec{v} = \vec{v} + \vec{u}$

Given \vec{u} and \vec{v} ,

$$\langle \vec{u}, \vec{v} \rangle = \mathbf{u} \cdot \mathbf{v} = \mathbf{u}^\top \mathbf{v} = \sum_{i=1}^n u_i v_i = u_1 v_1 + \dots + u_n v_n$$

Example Let $\mathbf{1} = (1, 1, \dots, 1)^\top$, then $\bar{u} = \frac{1}{n} \langle \mathbf{1}, \vec{u} \rangle = \frac{1}{n} \mathbf{1}^\top \mathbf{u}$

Example \mathbf{p} a probability vector, \mathbf{x} a vector of outcome, $\mathbb{E}(X) = \langle \vec{p}, \vec{x} \rangle$

$\vec{u} \perp \vec{v}$ if and only if $\langle \vec{u}, \vec{v} \rangle = 0$

$$\|\vec{u}\|^2 = \langle \vec{u}, \vec{u} \rangle = \mathbf{u} \cdot \mathbf{u} = \sum_{i=1}^n u_i^2 = u_1^2 + \dots + u_n^2$$

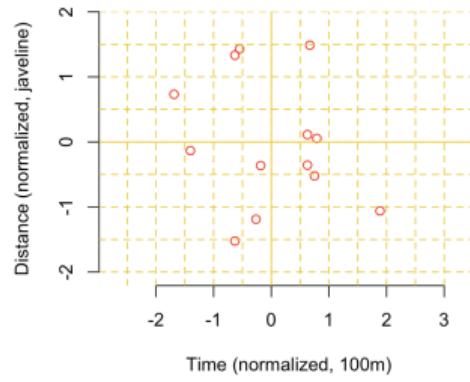
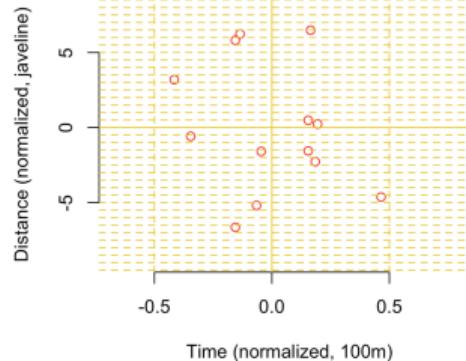
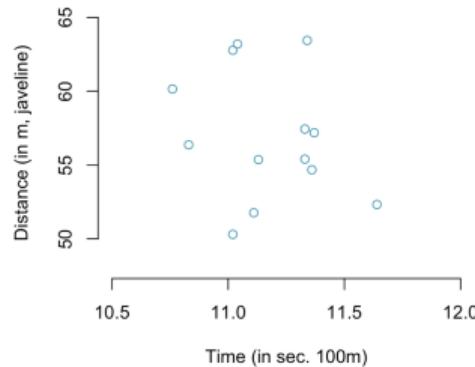
Mahalanobis Distance

The standard Euclidean distance is $\|\vec{u}\|_E^2 = \langle \vec{u}, \vec{u} \rangle = \mathbf{u}^\top \mathbf{u}$

Given some diagonal positive matrix D , and $\mu \in \mathbb{R}^n$,

$$\|\vec{u}\|_D^2 = (\mathbf{u} - \mu)^\top D^{-1}(\mathbf{u} - \mu) = \|\tilde{\mathbf{u}}\|_E^2 \text{ where } \tilde{\mathbf{u}}_i = \frac{\mathbf{u}_i - \mu_i}{\sqrt{D_{i,i}}}$$

(popular in statistics, when \mathbf{u}_i is the observation of individual i)



Euclidean Distance

Define $d(\mathbf{u}, \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|_2 = \sqrt{(\mathbf{u} - \mathbf{v})^\top (\mathbf{u} - \mathbf{v})}$

Set of points \mathbf{u} such that $d(\mathbf{0}, \mathbf{u}) = \|\mathbf{u}\|_2 = 1$

- circle (or sphere in higher dimension)

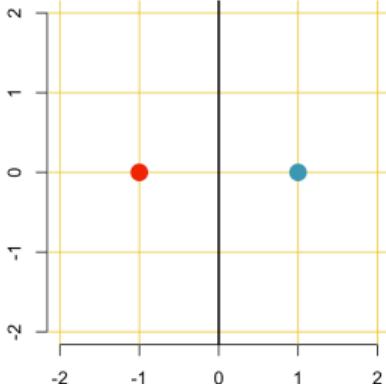
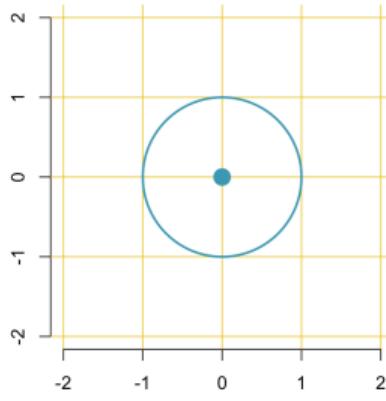
Set of points \mathbf{x} such that $d(\mathbf{x}, \mathbf{u}) = d(\mathbf{x}, \mathbf{v})$?

- straight line (or plane in higher dimension)

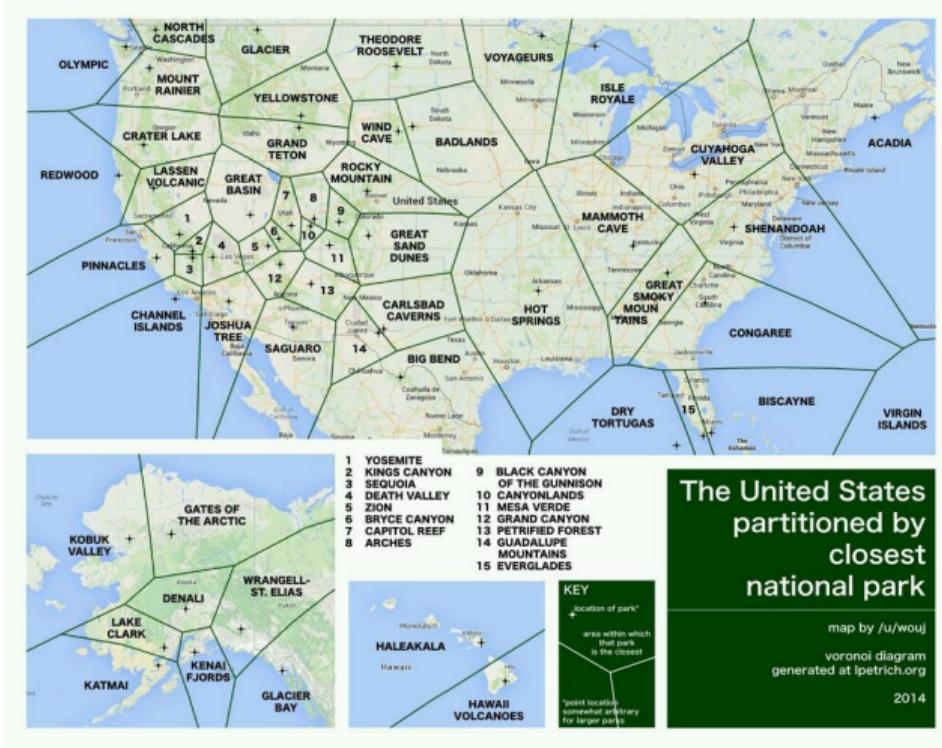
(orthogonal to $\mathbf{u} - \mathbf{v}$)

Let \mathbf{X} be a $n \times k$ matrix, n individuals, k columns

- $\mathbf{X}^\top \mathbf{X}$ is $k \times k$, it is a “**covariance**” (related) matrix
also called “Gram matrix”
- $\mathbf{X} \mathbf{X}^\top$ is $n \times n$, it is a “**similarity**” matrix



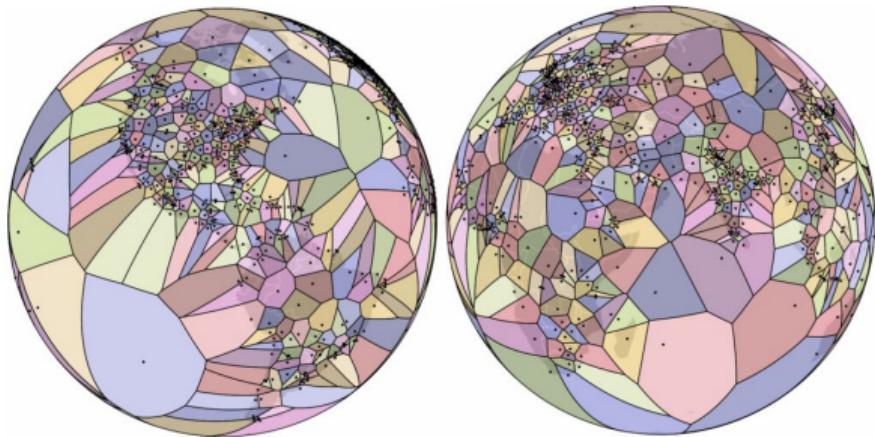
Euclidean Distance



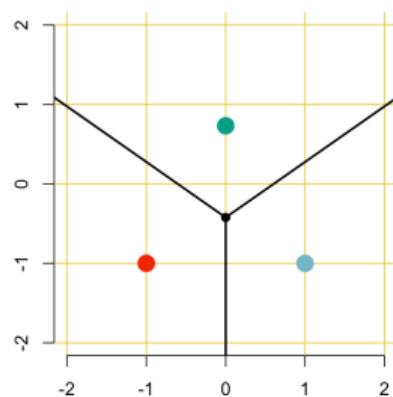
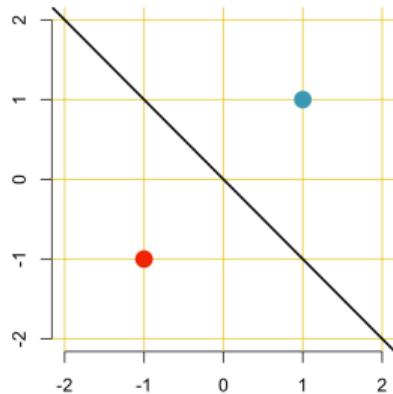
Euclidean Distance

Given a set of points $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$

What is the set of points x the closest to \mathbf{u}_i ? (see also Delaunay triangulation)



(intersections are straight lines or planes)

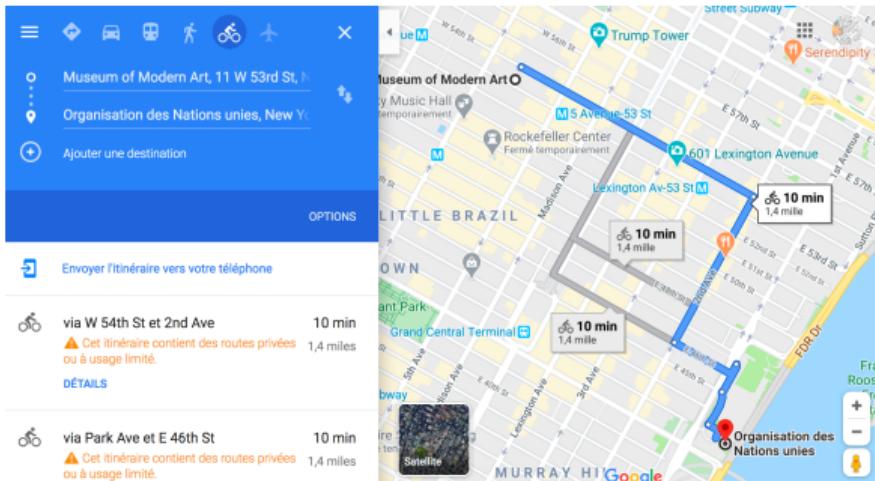


Manhattan (ℓ_1) Distance

$\|\mathbf{u}\|_1 = |u_1| + |u_2|$ is also a norm (Manhattan)

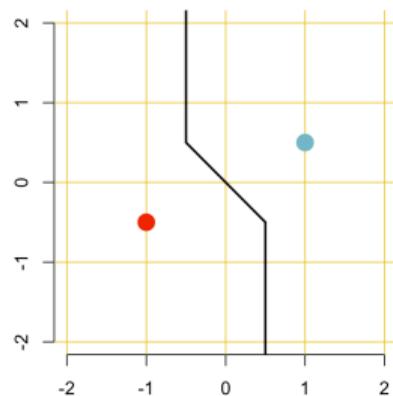
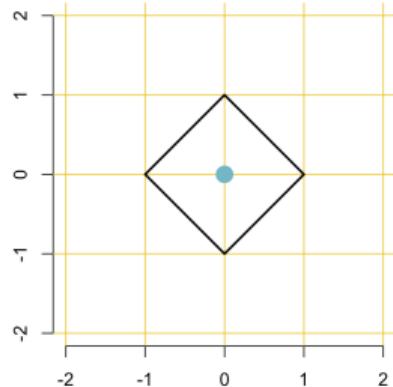
Set of points \mathbf{u} such that $d(\mathbf{0}, \mathbf{v}) = \|\mathbf{u}\| = 1$

- square (or cube in higher dimension)



Set of points \mathbf{x} such that $d(\mathbf{x}, \mathbf{u}) = d(\mathbf{x}, \mathbf{v})$?

- portions of straight lines (or planes)



Distance

...

A distance measure is an objective score that summarizes the relative difference between two objects in a problem domain.

Definition 1.26: Distance

A **distance** or **metric** over $\mathcal{X} \in \mathbb{R}^p$ is a map $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ that satisfies the following properties

- coincidence: $d(\mathbf{x}, \mathbf{y}) = 0 \iff \mathbf{x} = \mathbf{y}$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,
- symmetry $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,
- triangle inequality $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,
- non-negativity: $d(\mathbf{x}, \mathbf{y}) \geq 0$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,
- generalized triangle inequality $|d(\mathbf{x}, \mathbf{y}) - d(\mathbf{y}, \mathbf{z})| \leq d(\mathbf{x}, \mathbf{y})$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,

Distance

Definition 1.27: Manhattan-Taxicab distance (and ℓ_1 -norm)

Given $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$, $d_1(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^p |x_i - y_i|$.

Definition 1.28: Euclidean distance (and ℓ_2 -norm)

Given $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$, $d_2(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2 = \sqrt{(\mathbf{x} - \mathbf{y})^\top (\mathbf{x} - \mathbf{y})} = \sqrt{\sum_{i=1}^p (x_i - y_i)^2}$.

Definition 1.29: Mahalanobis distance

Given $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ and M a positive semidefinite matrix of dimension $p \times p$,

$$d_M(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^\top M (\mathbf{x} - \mathbf{y})}$$

Mahalanobis distances also satisfy

- homogeneousness: $d(a\mathbf{x}, a\mathbf{y}) = |a| \cdot d(\mathbf{x}, \mathbf{y})$, for every $a \in \mathbb{R}$ and $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,
- translation invariance $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y} + \mathbf{z}, \mathbf{y} + \mathbf{z})$ for every $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{X}$,

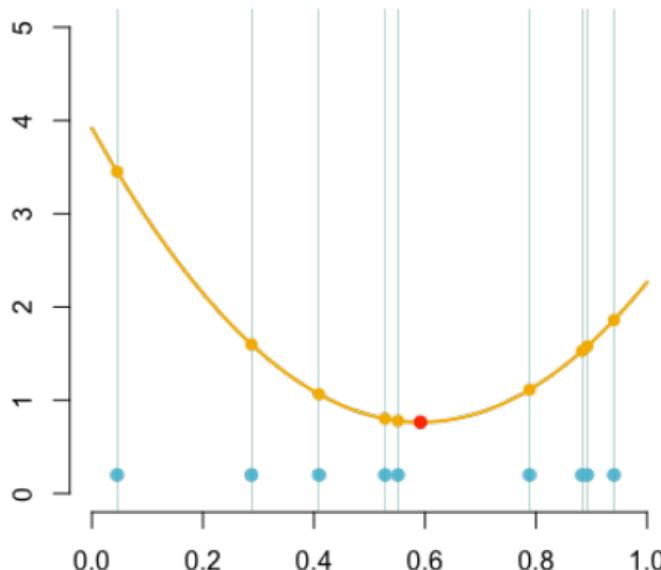
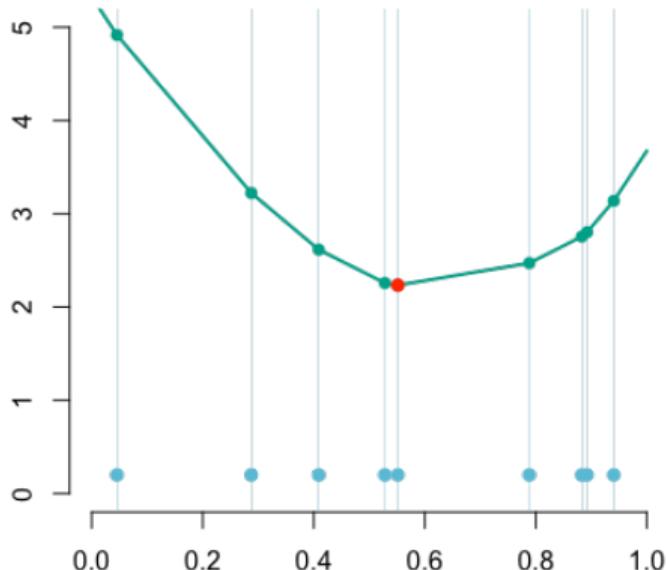
If $M = L^\top L$,

$$d_M(\mathbf{x}, \mathbf{y})^2 = (\mathbf{x} - \mathbf{y})^\top M (\mathbf{x} - \mathbf{y}) = [L(\mathbf{x} - \mathbf{y})]^\top [L(\mathbf{x} - \mathbf{y})] = \|L(\mathbf{x} - \mathbf{y})\|_2^2$$

where $\|\cdot\|_2$ is the standard Euclidean norm.

Loss

Given $\mathbf{y} = (y_1, \dots, y_n)$, find $x^* \in \operatorname{argmin}_{x \in \mathbb{R}} \left\{ \sum_{i=1}^n |x - y_i| \right\}$ or $x^* \in \operatorname{argmin}_{x \in \mathbb{R}} \left\{ \sum_{i=1}^n (x - y_i)^2 \right\}$



Loss

The τ th quantile ($\tau \in (0, 1)$) of Y is given by $q_\tau = F_Y^{-1}(\tau) = \inf \{y : F_Y(y) \geq \tau\}$,

$$q_Y(\tau) = \operatorname{argmin}_{q \in \mathbb{R}} \left\{ (\tau - 1) \int_{-\infty}^q (q - y) dF_Y(y) + \tau \int_q^\infty (y - q) dF_Y(y) \right\}.$$

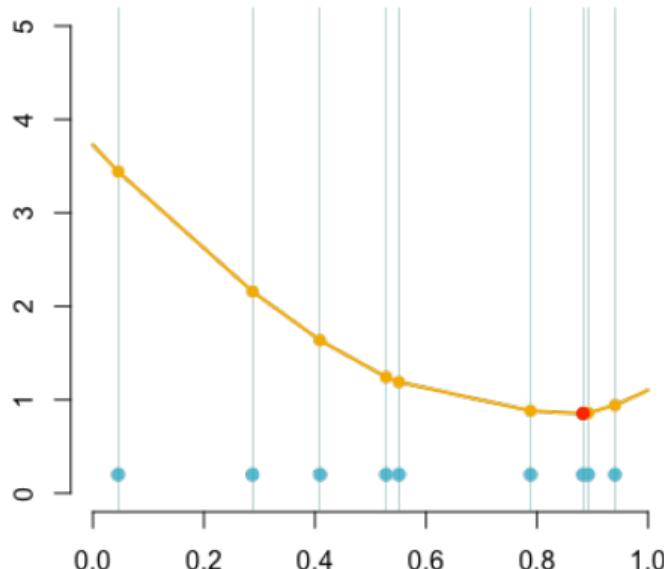
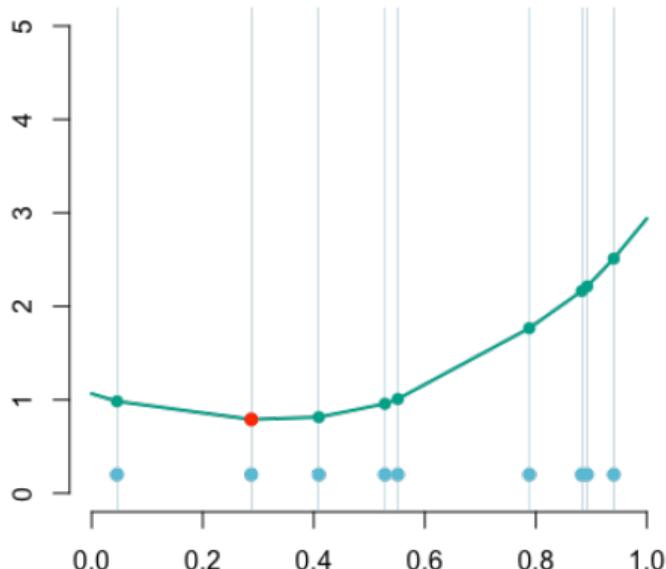
$$\hat{q}_\tau = \operatorname{argmin}_{q \in \mathbb{R}} \left\{ (\tau - 1) \sum_{y_i < q} (q - y_i) + \tau \sum_{y_i \geq q} (y_i - q) \right\},$$

also coined **pinball loss**, which is asymmetric,

$$\ell_\tau(y, \hat{y}) = \begin{cases} \tau \cdot (y - \hat{y}) & \text{if } y \geq \hat{y} \\ (1 - \tau) \cdot (\hat{y} - y) & \text{if } y < \hat{y} \end{cases}$$

Loss

Given $\mathbf{y} = (y_1, \dots, y_n)$, find x_τ^* such that $x_\tau^* \in \operatorname{argmin}_{x \in \mathbb{R}} \left\{ \sum_{i=1}^n \ell_\tau(y_i, x) \right\}$



with $\tau = 0.2$ and $\tau = 0.7$ respectively.

Loss

```
1 > library(cNORM)
2 > set.seed(1234)
3 > n = 10
4 > x = runif(n)
5 > x
6 [1] 0.114 0.622 0.609 0.623 0.861 0.640 0.009 0.233 0.666 0.514
7 > tau = .7
8 > quantile(x,tau)
9 [1] 0.6284588
```

The linear program is

$$\min_{q^+, q^-, \mathbf{a}, \mathbf{b}} \left\{ \sum_{i=1}^n \tau a_i + (1 - \tau) b_i \right\}, \text{ where } y_i = q^+ - q^- + a_i - b_i, \forall i \in \{1, \dots, n\}$$

and where $a_i, b_i \geq 0$ as well as $q^+, q^- \geq 0$.

Loss

Set $\mathbf{z} = (q^+, q^-, \mathbf{a}, \mathbf{b})^\top \in \mathbb{R}_+^{2n+2}$,

$$\min_{\mathbf{z}} \left\{ \mathbf{c}^\top \mathbf{z} \right\}, \text{ where } \mathbf{A}\mathbf{z} = \mathbf{b} \text{ and } \mathbf{z} \geq 0,$$

where \mathbf{A} is $n \times (2n + 2)$,

$$\mathbf{A} = [\mathbf{1}_n \ -\mathbf{1}_n \ \mathbb{I}_n \ -\mathbb{I}_n] = \begin{bmatrix} 1 & -1 & 1 & 0 & \cdots & 0 & -1 & 0 & \cdots & 0 \\ 1 & -1 & 0 & 1 & \cdots & 0 & 0 & -1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & -1 & 0 & 0 & \cdots & 1 & 0 & 0 & \cdots & -1 \end{bmatrix}$$

$$\mathbf{b} = \mathbf{y} \in \mathbb{R}^n \text{ and } \mathbf{c} = [0 \ 0 \ \tau \mathbf{1}_n \ (1-\tau) \mathbf{1}_n]^\top \in \mathbb{R}^{2n+2}$$

Loss

$$\min_{\mathbf{z}} \left\{ \mathbf{c}^\top \mathbf{z} \right\}, \text{ where } \mathbf{A}\mathbf{z} = \mathbf{b} \text{ and } \mathbf{z} = 0,$$

```
1 > library(lpSolve)
2 > ones = rep(1,n)
3 > A = cbind(ones, -ones, diag(n), -diag(n))
4 > b = x
5 > c = c(rep(0,2), tau*rep(1,n),(1-tau)*rep(1,n))
6 > equal_type = rep("=", n)
7 > lp("min", c,A,equal_type,b)$solution[1]
8 [1] 0.6403106
```

Deviance loss

These different variance functions can be translated to deviance loss functions by selecting the corresponding distribution within the exponential dispersion family.

- the [Gaussian case](#) translates to the square loss $\ell(y, \hat{y}) = (y - \hat{y})^2$
- the [Poisson case](#) translates to the Poisson deviance loss

$$\ell(y, \hat{y}) = 2 \left(\hat{y} - y - y \log \frac{\hat{y}}{y} \right)$$

- the [Gamma case](#) translates to the Gamma deviance loss

$$\ell(y, \hat{y}) = 2 \left(\frac{y - \hat{y}}{\hat{y}} + y \log \frac{\hat{y}}{y} \right)$$

- the [inverse Gaussian case](#) translates to the inverse Gaussian deviance loss

$$\ell(y, \hat{y}) = \frac{(y - \hat{y})^2}{\hat{y}^2 y}$$

Deviance loss

If $Y \sim \mathcal{P}(\lambda)$, consider

$$\hat{y} \mapsto \mathbb{E}[\ell(Y, \hat{y})] = 2 \sum_{y=0}^{\infty} \left(\hat{y} - y - y \log \frac{\hat{y}}{y} \right) \frac{e^{-\lambda} \lambda^y}{y!} = 2\hat{y} - 2\lambda - 2\lambda \log(\hat{y}) + 2\mathbb{E}[Y \log Y]$$

thus, $\min \{\mathbb{E}[\ell(Y, \hat{y})]\}$ is attained when $\hat{y} = \lambda = \mathbb{E}[Y]$.

A divergence is a magnitude to measure the closeness between certain objects in a set.

Definition 1.30: Divergence

D is a divergence if it satisfies

- non negativity $D(\mathbf{x} \parallel \mathbf{y}) \geq 0$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,
- coincidence $D(\mathbf{x} \parallel \mathbf{y}) = 0 \iff \mathbf{x} = \mathbf{y}$.

Deviance loss

Definition 1.31: Divergence on \mathbb{R}^n

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

- D is separable, $\forall (\mathbf{x}, \mathbf{y}) \in E^2$, $D(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} = \mathbf{y}$,
- D admits development $\forall (\mathbf{x}, \mathbf{x} + \boldsymbol{\epsilon}) \in E^2$, $D(\mathbf{x}, \mathbf{x} + \boldsymbol{\epsilon}) = \frac{1}{2} \sum A_{i,j}(\boldsymbol{\epsilon}) \epsilon_i \epsilon_j + O(|\boldsymbol{\epsilon}|^3)$, where $A(\boldsymbol{\epsilon})$ is definite positive.

A loss ℓ is consistent for mean estimation w.r.t. \mathcal{M} if

$$\mathbb{E}[\ell(Y, \mathbb{E}[Y|\mathbf{X}])] \leq \mathbb{E}[\ell(Y, m(\mathbf{X}))] \text{ for all } m \in \mathcal{M}$$

Strictly consistent loss functions for mean estimation are the Bregman divergences, Savage (1971).

Deviance loss

Definition 1.32: Bregman Divergence, Bregman (1967)

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

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

All of these deviance loss functions are strictly consistent for mean estimation, but the one with the correct conditional variance behavior for the response Y , given \mathbf{X} , has the best finite sample properties (on average).

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

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

Proposition 1.17: Bregman Divergence

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

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

Deviance loss

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

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

(see Mahalanobis distance).

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

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

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

Deviance loss

Kullback–Leibler (KL) divergences and deviance loss functions are Bregman divergences.

Definition 1.33: Projection, Bregman (1967), Bauschke et al. (1997)

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

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

Deviance loss

Definition 1.34: Entropy or Kullback-Liebler Divergence

Kullback-Liebler divergence, between distribution f and g is

$$\text{KL}(f\|g) = \mathbb{E}_f \left[\log \frac{f(X)}{g(X)} \right] = \int_{-\infty}^{\infty} \log \left(\frac{f(x)}{g(x)} \right) f(x) dx,$$

when $f \ll g$, i.e., $g(x) = 0$ implies $f(x) = 0$.

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

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

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

Deviance loss

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

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

Proposition 1.18: Divergence for Gaussian vectors

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

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

Deviance loss

Proposition 1.19: Gibbs' inequality

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

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

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


Deviance loss

Proposition 1.20: Additivity for independence distributions

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

Proof By definition

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

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

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

Definition 1.35: Symmetric relative entropy or Jeffrey divergence

Jeffrey divergence is some sort of “symmetric Kullback-Liebler” divergence

$$JF(f, g) = KL(f\|g) + KL(g\|f).$$

Deviance loss

Definition 1.36: Divergence based inference

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

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} \{D(f, f_\theta)\}$$

unknown f $f_\theta \in \mathcal{Q}$

or its empirical version

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \{D(\hat{f}_n, f_\theta)\}$$

estimated \hat{f}_n

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities)

BY-NC 4.0 205 / 335

Deviance loss

Definition 1.37: Hellinger distance, Hellinger (1909)

For two discrete distributions f and g , Hellinger distance is

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

and for absolutely continuous distributions, if f and g are densities,

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

See Pardo (2018).

Deviance loss

Proposition 1.21: Distance between Beta variables

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

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

Proof

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

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

Deviance loss

Proposition 1.22: Distance between Gaussian vectors

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

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

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

Deviance loss

Definition 1.38: Cramér, Cramér (1928a,b) and Székely (2003)

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

$$C_k(f, g) = \left(\int_{-\infty}^{\infty} |F(x) - G(x)|^k dx \right)^{1/k}, \text{ for } k \geq 1$$

Definition 1.39: Wasserstein, Wasserstein (1969)

Consider two measures on f and g on \mathbb{R} . Then define **Wasserstein distance**

$$W_k(f, g) = \left(\int_0^1 |F^{-1}(u) - G^{-1}(u)|^k du \right)^{1/k}, \text{ for } k \geq 1$$

Deviance loss

Proposition 1.23: C_1 and W_1

Consider two measures on f and g on \mathbb{R} .

$$W_1(f, g) = \int_0^1 |F^{-1}(u) - G^{-1}(u)| du = \int_{-\infty}^{\infty} |F(x) - G(x)| dx = C_1(f, g).$$

Proof See Prokhorov (1956), Dall'Aglio (1956)

Proposition 1.24: W_2 for Gaussian / Bernoulli distributions

Consider two Gaussian distributions, then

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

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

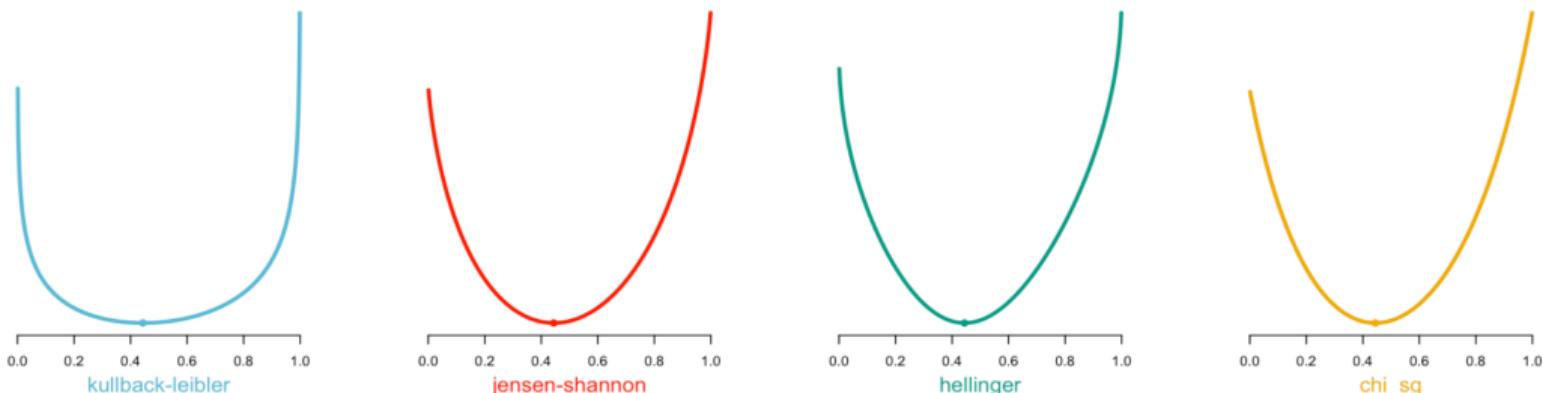
$$W_2(\mathcal{B}(p_1), \mathcal{B}(p_2)) = \sqrt{p_2 - p_1}.$$

Distribution Distance Based Inference

μ : multinomial distribution on $\{0, 1, 10\}$, with $p = (.5, .1, .4)$

ν_θ : binomial type distribution on $\{0, 10\}$, with $q_\theta = (1 - \theta, \theta)$

Let $\theta^* = \operatorname{argmin}\{d(p, q_\theta)\}$ or $\theta^* = \operatorname{argmin}\{d(p \| q_\theta)\}$



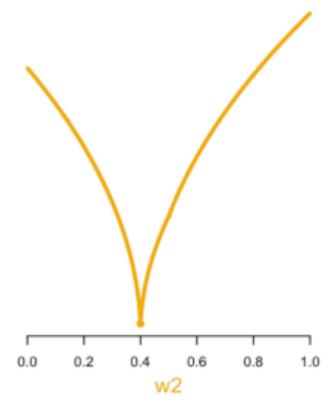
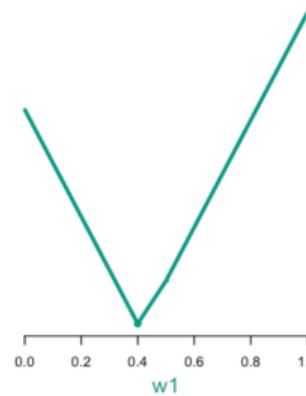
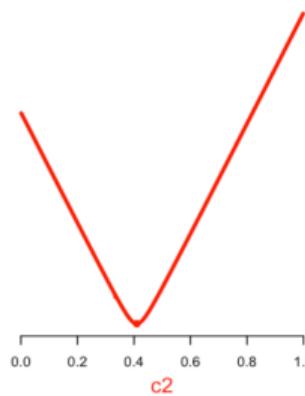
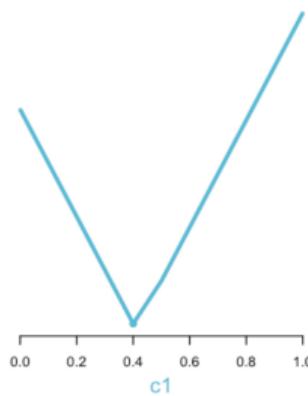
with $d_{KL}(p \| q_\theta)$, $d_{JS}(p, q_\theta)$, $d_H(p, q_\theta)$ and $d_{H_{\chi^2}}(p \| q_\theta)$.

Distribution Distance Based Inference

μ : multinomial distribution on $\{0, 1, 10\}$, with $p = (.5, .1, .4)$

ν_θ : binomial type distribution on $\{0, 10\}$, with $q_\theta = (1 - \theta, \theta)$

Let $\theta^* = \operatorname{argmin}\{d(p, q_\theta)\}$



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

Optimal Transport, the Setting ★★★

Consider measures ν_0 and ν_1 on \mathcal{Z}_0 , \mathcal{Z}_1 , compact subsets of \mathbb{R}^d .

There exists T such that $\nu_1 = T_\sharp \nu_0$, where ν_0 is atomless ($T_\sharp \nu_0(B) = \nu_0(T^{-1}(B))$).

As shown in [Villani \(2003\)](#) and [Santambrogio \(2015\)](#), we can be interested in “optimal” mappings, satisfying Monge problem, from [Monge \(1781\)](#), i.e., solutions of

$$\inf_{T_\sharp \nu_0 = \nu_1} \int_{\mathcal{Z}_0} c(\mathbf{z}_0, T(\mathbf{z}_0)) \nu_0(d\mathbf{z}_0), \quad (1)$$

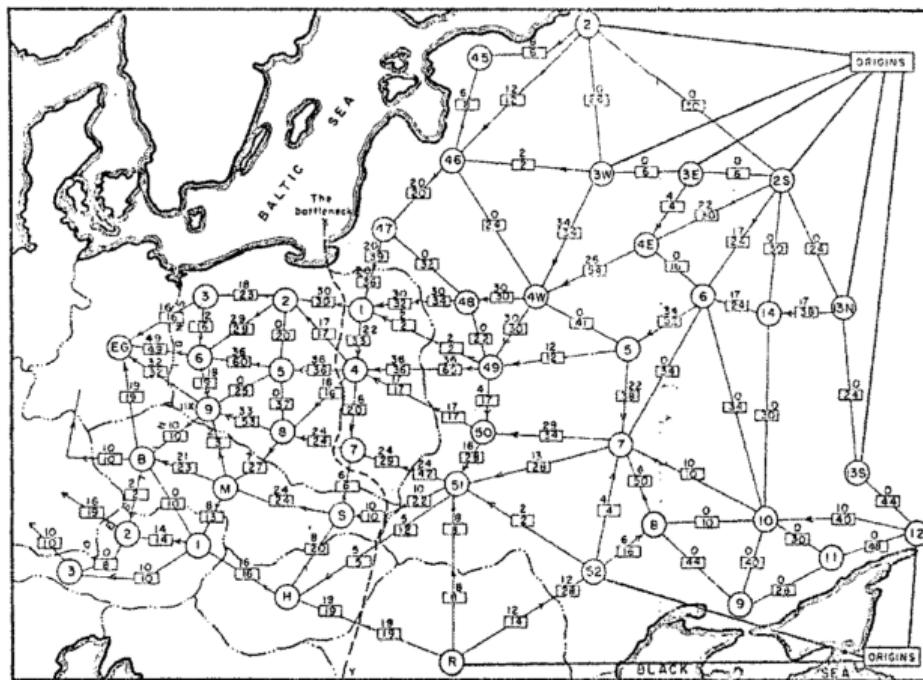
for some positive ground cost function $c : \mathcal{Z}_0 \times \mathcal{Z}_1 \rightarrow \mathbb{R}_+$. If ν_0 and ν_1 are not absolutely continuous (with respect to Lebesgue measure), there might not be such deterministic mapping T . This limitation motivates a relaxation of Monge’s problem, as considered in [Kantorovich \(1942\)](#),

$$\inf_{\pi \in \Pi(\nu_0, \nu_1)} \int_{\mathcal{Z}_0 \times \mathcal{Z}_1} c(\mathbf{z}_0, \mathbf{z}_1) \pi(d\mathbf{z}_0, d\mathbf{z}_1), \quad (2)$$

freakonometrics

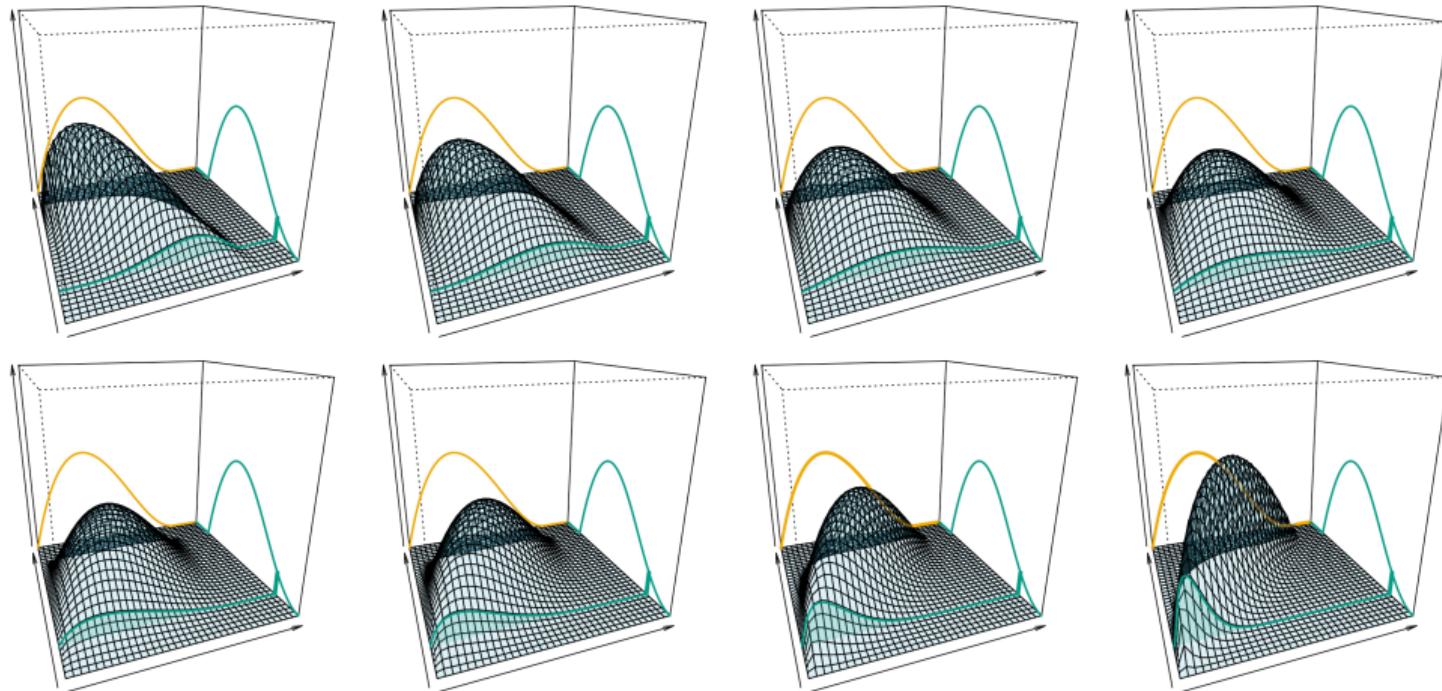
freakonometrics.hypotheses.org

Optimal Transport, the Setting ★★★



(flow on networks, via [Harris and Ross \(1955\)](#))

Optimal Transport, the Setting ★★★



Optimal Transport, the Gaussian case ★★★

With a quadratic cost (and Euclidean distance in \mathbb{R}^d) and absolutely continuous measures, the optimal Monge map T^* is unique, and it is the gradient of a convex function, $T^* = \nabla\varphi$, see [Brenier \(1991\)](#).

If $\nu_j \sim \mathcal{N}(\mu_j, \Sigma_j)$,

$$z_1 = T^*(z_0) = \mu_1 + M(z_0 - \mu_0),$$

where M is a symmetric positive matrix that satisfies $M\Sigma_0M = \Sigma_1$, which has a unique solution given by $M = \Sigma_0^{-1/2}(\Sigma_0^{1/2}\Sigma_1\Sigma_0^{1/2})^{1/2}\Sigma_0^{-1/2}$, where $M^{1/2}$ denotes the square root of the square (symmetric) positive matrix M based on the Schur decomposition ($M^{1/2}$ is a positive symmetric matrix), as described in [Higham \(2008\)](#).

In the univariate case

$$z_1 = T^*(z_0) = \mu_1 + \sigma_1\sigma_0^{-1}(z_0 - \mu_0),$$

that is a linear non-decreasing mapping $\mathbb{R} \rightarrow \mathbb{R}$.

Optimal Transport, the discrete case ★★★

Consider two samples in the \mathbb{R}^d , $\{\mathbf{z}_{0,1}, \dots, \mathbf{x}_{z,n_0}\}$ and $\{\mathbf{z}_{1,1}, \dots, \mathbf{z}_{1,n_1}\}$. The discrete version of the Kantorovich problem (Equation 2) is

$$\min_{P \in U(n_0, n_1)} \left\{ \sum_{i=1}^{n_0} \sum_{j=1}^{n_1} P_{i,j} C_{i,j} \right\} = \min_{P \in U(n_0, n_1)} \{ \langle P, C \rangle \} \quad (3)$$

where, as in [Brualdi \(2006\)](#), $U(n_0, n_1)$ is the set of $n_0 \times n_1$ matrices corresponding to the convex transportation polytope

$$U(n_0, n_1) = \left\{ P : P\mathbf{1}_{n_1} = \mathbf{1}_{n_0} \text{ and } P^\top \mathbf{1}_{n_0} = \frac{n_0}{n_1} \mathbf{1}_{n_1} \right\},$$

and where C denotes the $n_0 \times n_1$ cost matrix, $C_{i,j} = c(\mathbf{z}_i, \mathbf{z}_j)$, associated with cost c .

freakonometrics.hypotheses.org

Optimal Transport, the discrete case ★★★

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

$$\underset{P \in U(1_n, 1_n)}{\operatorname{argmin}} \left\{ \langle P, C \rangle \right\}$$

over the set of **permutation matrices**.

minimal cost 0.2116

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

Optimal Transport, the discrete case ★★★

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

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

Optimal Transport, the discrete case ★★★

$$\underset{P \in U(\mathbf{1}_{n_0}, \mathbf{1}_{n_1})}{\operatorname{argmin}} \left\{ \langle P, C \rangle - \gamma \cdot \text{entropy}(P) \right\}$$

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

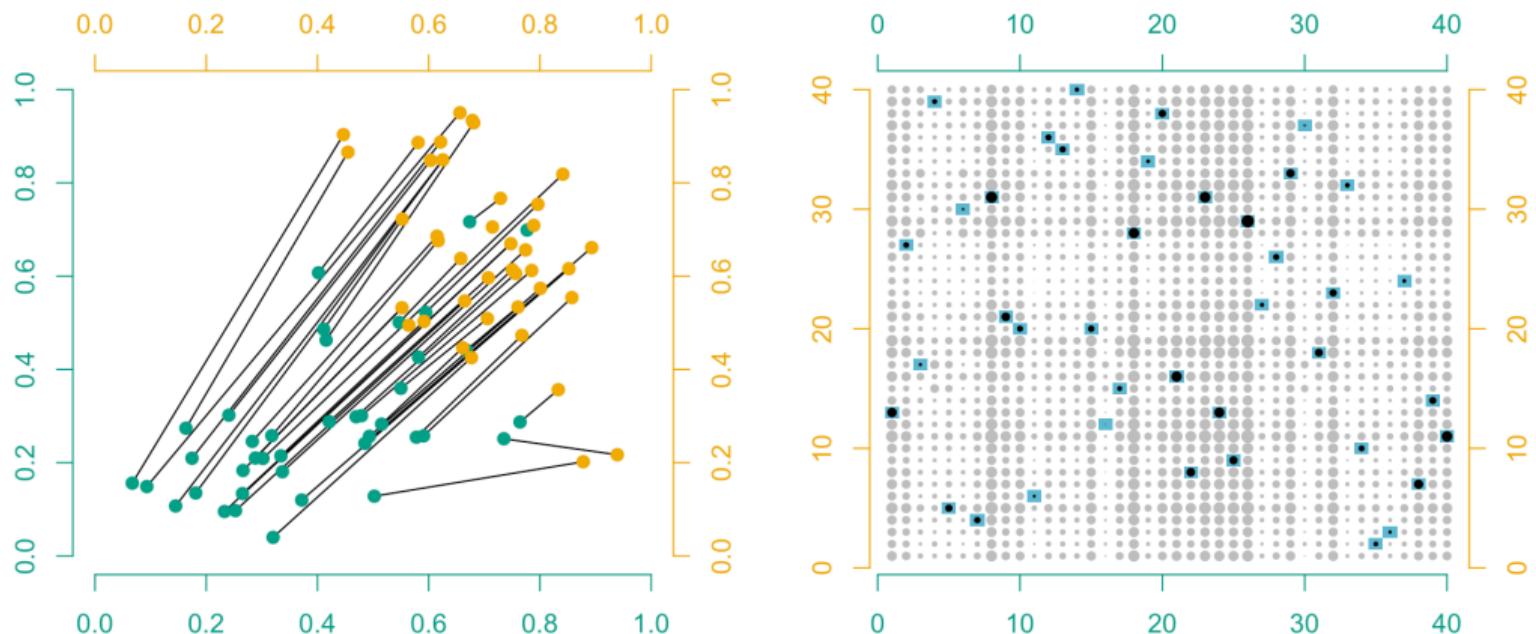
(regularization via Sinkhorn algorithm,
associated to “Matrix Scaling Problem”
in [Sinkhorn \(1962\)](#)) or

$$\underset{P \in U(\mathbf{1}_{n_0}, \mathbf{1}_{n_1})}{\operatorname{argmin}} \left\{ \langle P, C \rangle + \gamma \cdot d_{KL}(P || \mathbf{1}_{n_0} \otimes \mathbf{1}_{n_1}^\top) \right\}$$

minimal cost 0.2189 (> 0.2116)

	7	8	9	10	11	12	
1	.	.	0.35	.	0.59	0.06	$1 \leftrightarrow 11 + \dots$
2	0.12	0.82	.	0.07	.	.	$2 \leftrightarrow 8 + \dots$
3	.	.	0.55	.	0.41	0.04	$3 \leftrightarrow 9 + \dots$
4	0.85	.	.	0.14	.	.	$4 \leftrightarrow 7 + \dots$
5	0.03	0.18	.	0.79	.	.	$5 \leftrightarrow 10 + \dots$
6	.	.	0.10	.	.	0.90	$6 \leftrightarrow 12 + \dots$

Optimal Transport, the discrete case ★★★



Discrete matching, $n_0 = n_1 = 40$ individuals in two groups, $\{z_i, s_i = 0\}$ and $\{z_i, s_i = 1\}$, with matrices C ($C_{i,j} = \bullet$) and P^* (\blacksquare if $P_{i,j}^* = 1$)

Optimal Transport, the univariate case ★★★

Wasserstein distance is defined as the optimal transportation cost, when c is the ℓ_k distance, for $k \geq 1$, [Wasserstein \(1969\)](#),

$$W_k(\nu_0, \nu_1)^k := \inf_{T \sharp \nu_0 = \nu_1} \int_{\mathbb{R}} |z_0 - T(z_0)|^k \nu_0(dz_0) = \int_0^1 |Q_1(u) - Q_0(u)|^k du$$

i.e., the optimal mapping T^* (if $k > 1$) is non-decreasing, i.e.,

$$z_1 = T^*(z_0) = Q_1 \circ F_0(z_0) \text{ where } \begin{cases} F_j(z) = \mu_j((-\infty, z]) \\ Q_j(u) = F_j^{-1}(u) = \inf\{t \in \mathbb{R} : F_j(t) \geq u\} \end{cases}$$

$T^* := Q_1 \circ F_0$ is a non-decreasing mapping $\mathbb{R} \rightarrow \mathbb{R}$.

In higher dimension, some multivariate extensions of quantiles can be introduced to keep this construction, see [Hallin et al. \(2021\)](#), [Hallin and Konen \(2024\)](#)

Optimal Transport, the univariate case ★★★

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

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

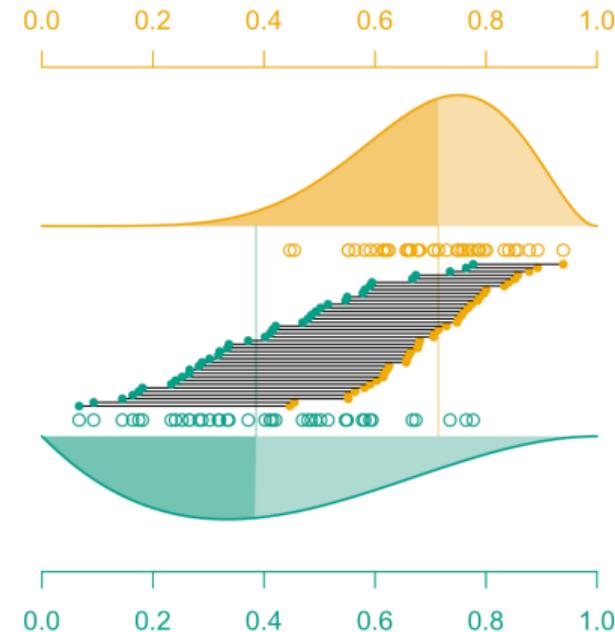
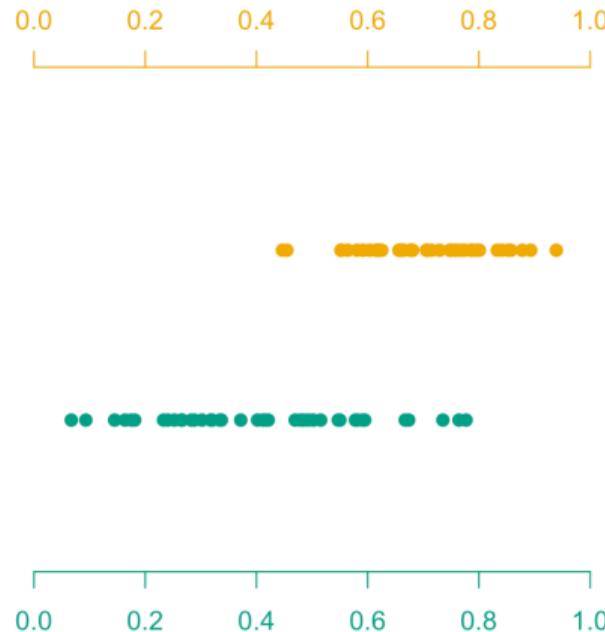
while if $\Phi : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is submodular,

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

see [Hardy et al. \(1952\)](#).

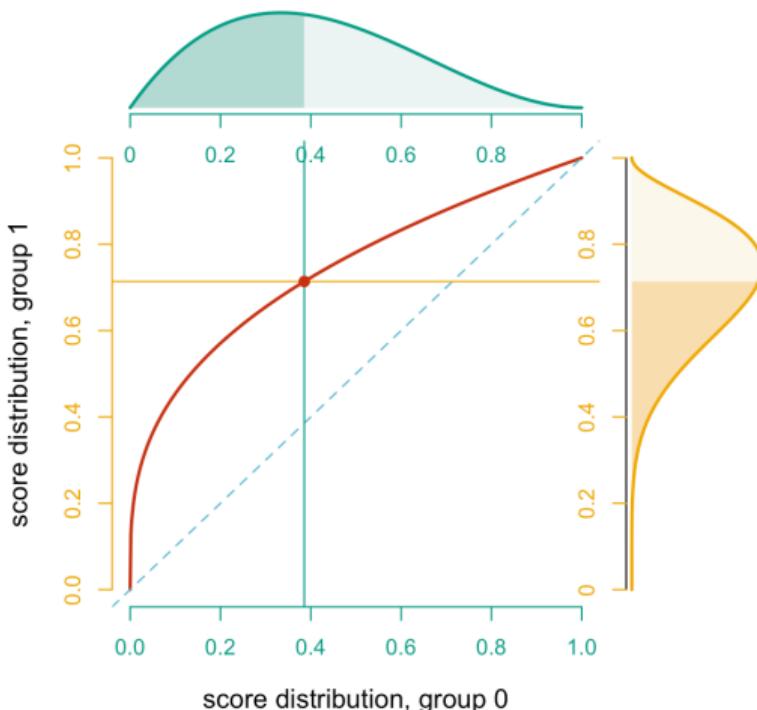
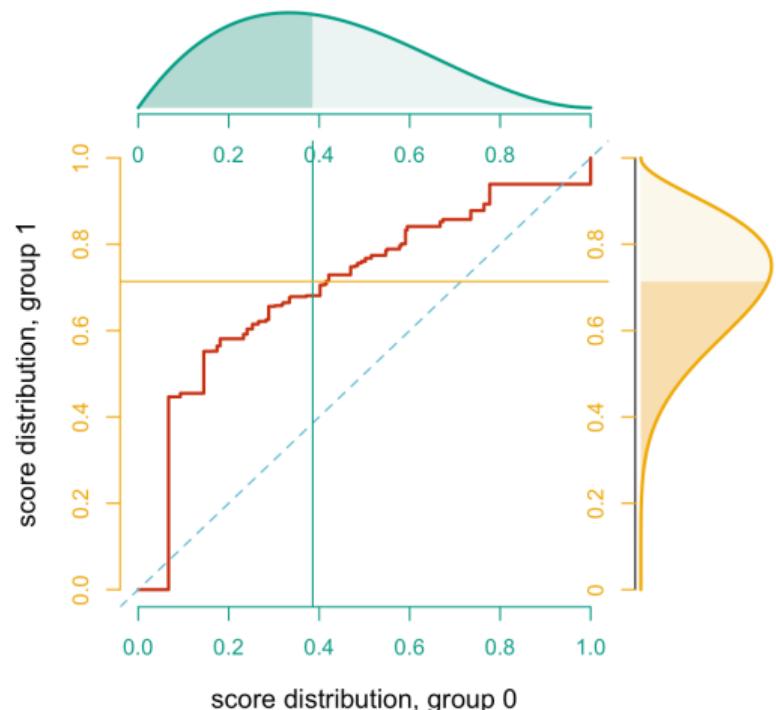
See related work on risk measures, [Denuit et al. \(2006\)](#), [Galichon \(2016\)](#), [Carlier et al. \(2016\)](#).

Optimal Transport with Pictures, the univariate case ★★★



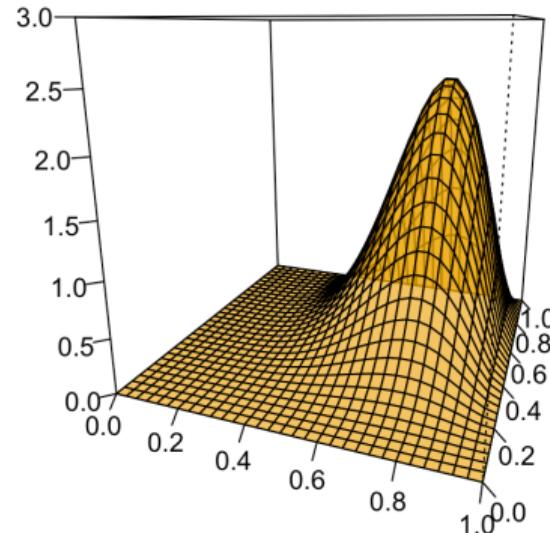
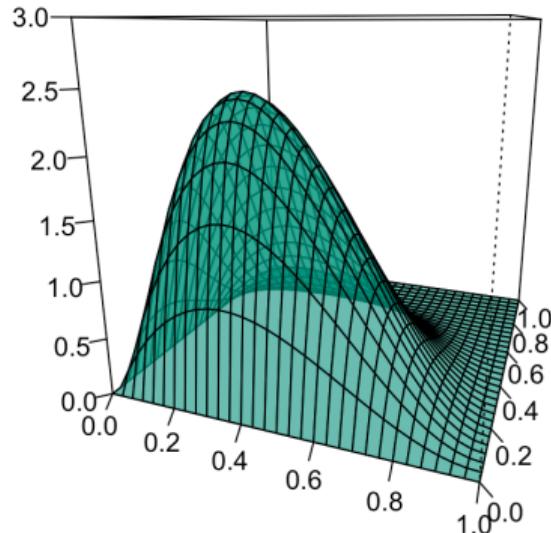
Univariate predictions for n individuals in two groups, $\{z_i, s_i = 0\}$ and $\{z_i, s_i = 1\}$, or measures ν_0 and ν_1 . $T^*(\cdot) = Q_1 \circ F_0(\cdot)$.

Optimal Transport with Pictures, the univariate case ★★★



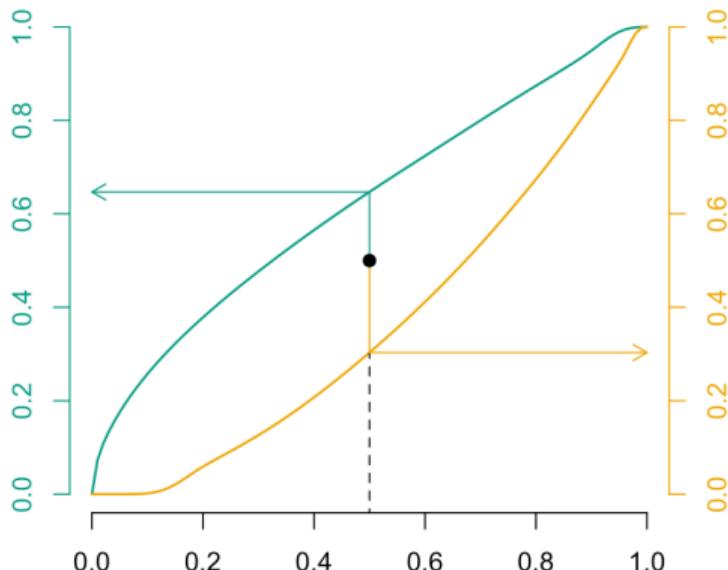
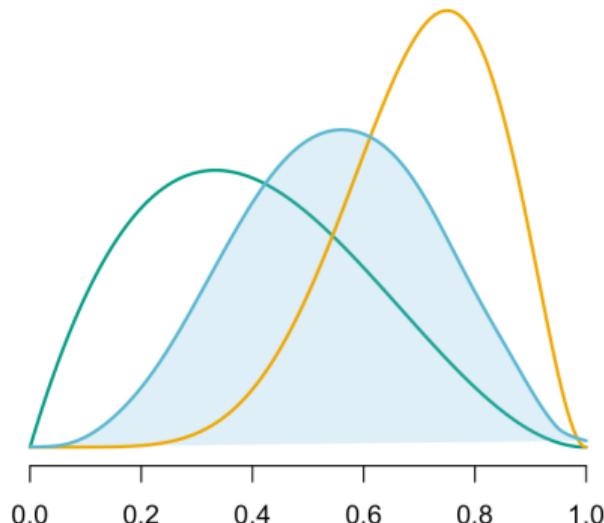
Univariate predictions for n individuals in two groups, $T^*(\cdot) = Q_1 \circ F_0(\cdot)$.

Optimal Transport with Pictures, the univariate case ★★★



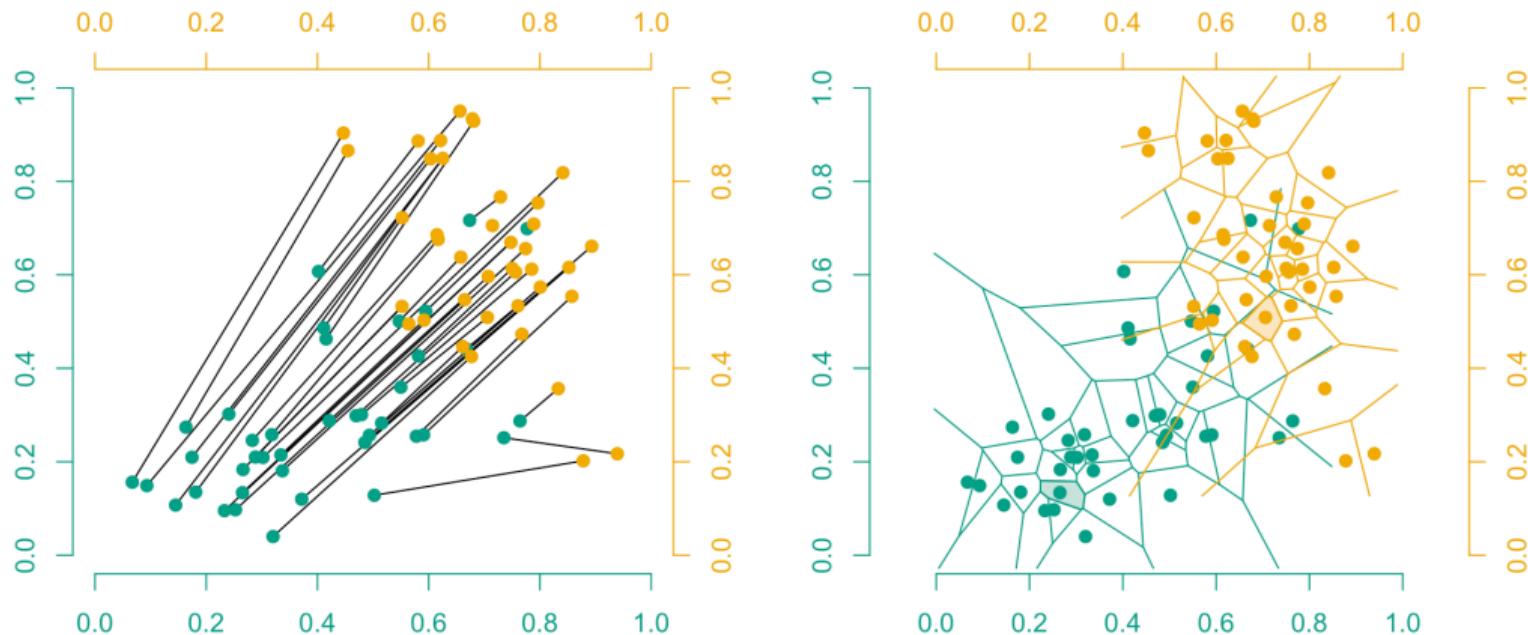
Multivariate observations for n individuals in two groups, ν_0 and ν_1 (densities).

Optimal Transport with Pictures, the univariate case ★★★



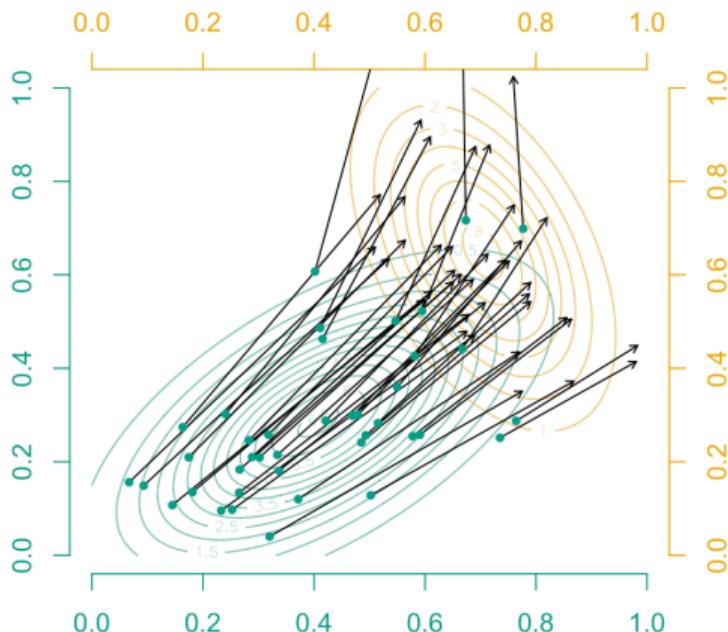
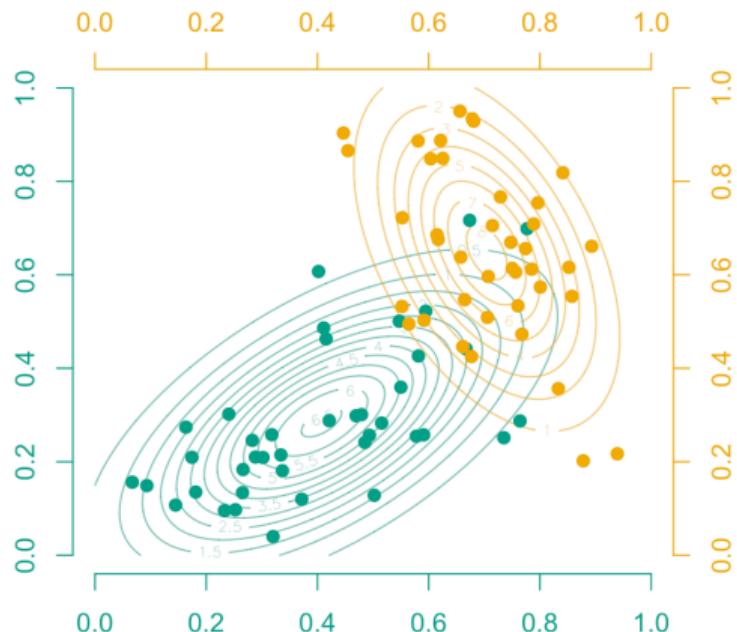
Bivariate observations for n individuals in two groups, $\{z_i, s_i = 0\}$ and $\{z_i, s_i = 1\}$.

Optimal Transport with Pictures, the univariate case ★★★



Discrete matching, n individuals in two groups, $\{z_i, s_i = 0\}$ and $\{z_i, s_i = 1\}$.

Optimal Transport with Pictures, the univariate case ★★★



Optimal transport with Gaussian joint distributions, $\mathcal{N}(\mu_0, \Sigma_0)$ and $\mathcal{N}(\mu_1, \Sigma_1)$.

Definition 1.40: Distance

A **distance** or **metric** over $\mathcal{X} \in \mathbb{R}^p$ is a map $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ that satisfies the following properties

- coincidence: $d(\mathbf{x}, \mathbf{y}) = 0 \iff \mathbf{x} = \mathbf{y}$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,
- symmetry $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,
- triangle inequality $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$ for every $\mathbf{x}, \mathbf{y} \in \mathcal{X}$,

Definition 1.41: Loss ℓ

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

Standard loss functions are usually related to divergences, distances, or powers of distances, e.g.

$$\ell_1(y, y') = d_1(y, y') = |y - y'| \text{ or } \ell_2(y, y') = d_1(y, y')^2 = (y - y')^2.$$

Definition 1.42: Risk \mathcal{R}

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

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

Definition 1.43: Empirical risk $\hat{\mathcal{R}}_n$

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

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

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

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

Proposition 1.25: Consistency of the empirical risk

The empirical risk, associated with ℓ is a consistent estimator of the ℓ -risk, i.e.

$$\lim_{n \rightarrow \infty} \{\widehat{\mathcal{R}}_n(m)\} = \mathcal{R}(m), \text{ or } \lim_{n \rightarrow \infty} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(m(\mathbf{X}_i), Y_i) \right\} = \mathbb{E}[\ell(m(\mathbf{X}), Y)], \text{ in probability.}$$

Proposition 1.26: Optimal Decision, "*Bayes decision rule*"

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

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

Risk

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

Definition 1.44: Misclassification loss, $\ell_{0/1}$

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

In the case of a binary classifier, observe that

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

Risk

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

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

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

Definition 1.45: Excess of risk of \hat{m}

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

For a classifier

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

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

Definition 1.46: Plug-in Estimator

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

Proposition 1.27

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

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

Proof We have seen that

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

Risk

But

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

(from the definition of m^*)

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

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

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

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

Risk

i.e.

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

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

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

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

Definition 1.47: Quadratic loss, ℓ_2

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

Risk

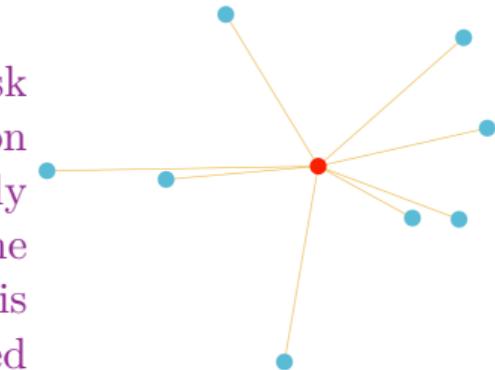
Observe that

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

The expected value is “**elicitable**” (for the ℓ_2 loss).

“The elicitability of a risk measure means that the risk measure can be obtained by minimizing the expectation of a forecasting objective function. Elicitability is closely related to backtesting, whose objective is to evaluate the performance of a risk forecasting model. If a risk measure is elicitable, then the sample average forecasting error based on the objective function can be used for backtesting the risk measure,” He et al. (2022).

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



Proposition 1.28: Optimal Decision, "Bayes decision rule", ℓ_2

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

Proposition 1.29: Optimal Decision, "Bayes decision rule", ℓ_{0-1}

For the binary loss ℓ_{0-1} , Bayes decision rule is the majority class
 $m_x^* = \operatorname{argmax}_{y \in \{0,1\}} \{\mathbb{P}[Y = y | \mathbf{X} = \mathbf{x}]\}$.

Risk

For the quadratic loss ℓ_2 , Bayes decision rule is the (conditional) expected value,

$$m_{\mathbf{x}}^* = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}] = \mu(\mathbf{x}) = \operatorname{argmin}_{m \in \mathcal{M}} \{\mathbb{E}[\ell_2(Y, m(\mathbf{X}))]\},$$

For the pinball loss $\ell_{q,\tau}$, Bayes decision rule is the quantile regression,

$$q_{\tau,\mathbf{x}}^* = Q_\tau[Y|\mathbf{X} = \mathbf{x}] = \operatorname{argmin}_{m \in \mathcal{M}} \{\mathbb{E}[\ell_{q,\tau}(Y, m(\mathbf{X}))]\},$$

where $\ell_{q,\tau}(y, y') = \tau|y - y'|\mathbf{1}_{y > y'} + (1 - \tau)|y - y'|\mathbf{1}_{y < y'}$.

For any $\alpha \in (0, 1)$,

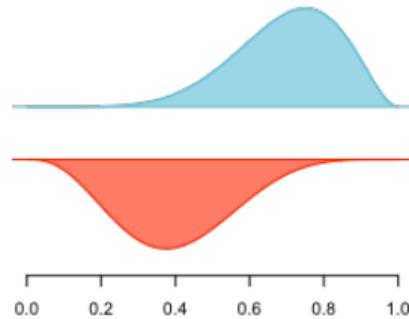
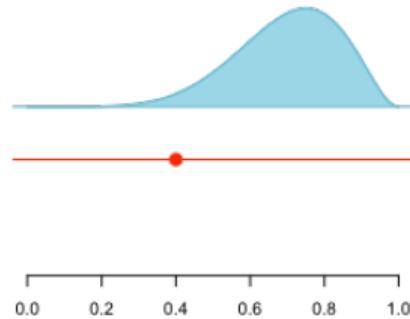
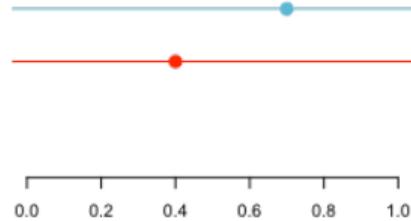
$$\mathbb{P}\left(Y \in [Q_{\alpha/2}(Y|\mathbf{X}); Q_{1-\alpha/2}(Y|\mathbf{X})]\right) = 1 - \alpha.$$

however, in general, with a finite number of observations

$$\mathbb{P}\left(Y \in [\hat{Q}_{\alpha/2}(Y|\mathbf{X}); \hat{Q}_{1-\alpha/2}(Y|\mathbf{X})]\right) \neq 1 - \alpha.$$

(see conformal prediction).

Loss



Loss

A statistic θ is elicitable if it can be expressed as a risk minimization for a given loss.

Definition 1.48: Elicitation, Brier et al. (1950), Good (1952)

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

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

A statistics $\hat{\theta} = \hat{\mathcal{I}}(\mathbf{y})$ (where $\mathbf{y} = \{y_1, \dots, y_n\}$) is said to be elicitable if it minimizes expected loss for some loss function ℓ , in the sense that

$$\hat{\mathcal{I}}(\mathbf{y}) = \operatorname{argmin}_{y \in \mathbb{R}} \left\{ \sum_{i=1}^n \ell(y_i, y) \right\}$$

Loss

E.g.

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \operatorname{argmin}_{m \in \mathbb{R}} \{\hat{\mathcal{R}}_2(m)\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \ell_2(y_i, m) \right\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n (y_i - m)^2 \right\}.$$

for any distribution of Y_i 's, or, e.g.,

$$\text{median}(y) \in \operatorname{argmin}_{m \in \mathbb{R}} \{\hat{\mathcal{R}}_1(m)\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \ell_1(y_i, m) \right\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n |y_i - m| \right\}.$$

again, for any distribution of Y_i 's.

But if Y_i 's are i.i.d. Bernoulli $\mathcal{B}(p)$ variables,

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \operatorname{argmin}_{m \in \mathbb{R}} \{\hat{\mathcal{R}}(m)\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \ell(y_i, m) \right\}.$$

Loss

where

$$\ell(y_i, \hat{y}_i) = \begin{cases} -2 \log(1 - \hat{y}_i) & y_i = 0 \\ -2 \log(\hat{y}_i) & y_i = 1. \end{cases}$$

if Y_i 's are i.i.d. Poisson $\mathcal{P}(\lambda)$ variables,

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \underset{m \in \mathbb{R}}{\operatorname{argmin}} \{ \widehat{\mathcal{R}}(m) \} = \underset{m \in \mathbb{R}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \ell(y_i, m) \right\}.$$

where

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

freakonometrics

freakonometrics.hypotheses.org

Definition 1.49: Scoring Rule

In decision theory, a scoring rule provides evaluation metrics for probabilistic predictions or forecasts. $\ell : (y, F_y)$ for a predictive distribution F_y

Popular in Bayesian regression

Quick Overview on Optimization

- Most ML models involve learning parameters to minimize a cost, or maximize a reward. For examples:
 - Linear regression: minimize mean squared error.
 - Logistic regression: minimize cross-entropy loss.
 - Neural networks: minimize empirical loss via backpropagation.
 - Reinforcement learning: maximize expected future rewards.
- Optimization is the engine that drives model training.
- Understanding optimization is crucial for developing efficient, robust learning algorithms.
- General form: $\min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta)$, where $\mathcal{L}(\theta)$ is the loss function and θ are the parameters.
- Often involves:
 - High-dimensional parameter spaces
 - Non-convex loss functions (e.g., deep nets)

Quick Overview on Optimization

- Large-scale datasets
- Goal: find parameters θ^* that minimize training loss and ideally generalize well.
- **Unconstrained vs. Constrained**

$$\min_{\theta} \mathcal{L}(\theta) \quad \text{vs.} \quad \min_{\theta \in C} \mathcal{L}(\theta)$$

- **Convex vs. Non-convex**
 - Convex: single global minimum, easier to optimize.
 - Non-convex: multiple minima, may require heuristics.
- **Deterministic vs. Stochastic**
 - Full dataset vs. mini-batches (SGD)
- **Smooth vs. Non-smooth**
 - Smooth: gradients exist everywhere.
 - Non-smooth: may need subgradient or proximal methods.

Quick Overview on Optimization

- **Stochastic methods** scale to large datasets:
 - SGD, Mini-batch SGD, Adam
- **Regularization** improves generalization:

$$\min_{\theta} \mathcal{L}(\theta) + \lambda R(\theta)$$

- **Constraints and projections:** Useful for interpretability, sparsity, and physical constraints.
- **Non-smooth optimization:**
 - Lasso (L1 regularization), hinge loss in SVMs
 - Meta-learning, Bayesian optimization, reinforcement learning: Advanced optimization in frontier ML tasks.
- **Convex optimization:**
 - Global minimum is also local minimum.

Quick Overview on Optimization

- Easier to analyze and solve.
- **Non-convex optimization:**
 - Many local minima and saddle points.
 - Common in deep learning.
- Visualization: think "bowl-shaped" vs "rugged terrain".
- Most ML models are trained using gradients.
- **Gradient Descent:** from x_0

$$x_{t+1} = x_t - \eta \nabla \mathcal{L}(x_t)$$

- **Variants:** Stochastic Gradient Descent
- Efficient and scalable for large datasets and high-dimensional problems.
- Not all problems are best solved with plain gradients.
- Alternatives include:

Quick Overview on Optimization

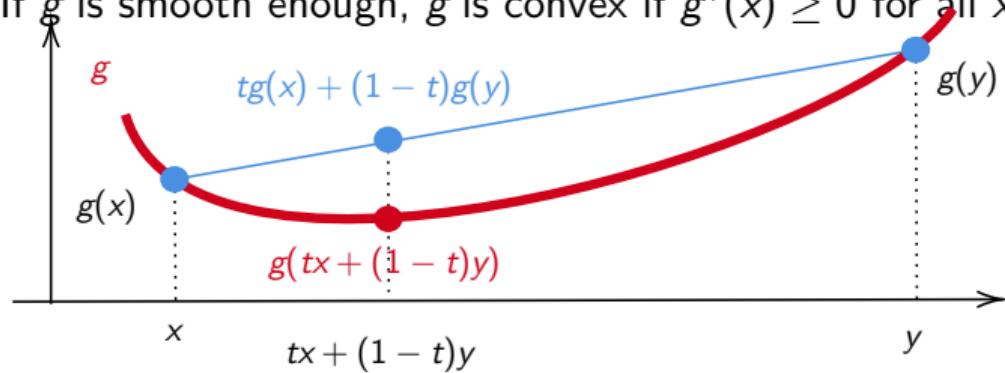
- **Newton's Method**: uses second derivatives (Hessian).
- **Coordinate Descent**: optimize one variable at a time.
- **Evolutionary algorithms**: e.g., genetic algorithms.
- **Bayesian Optimization**: for black-box functions.
- Tool choice depends on model complexity, smoothness, constraints, and noise.

Definition 1.50: Convex Function

$g : \mathbb{R} \rightarrow \mathbb{R}$ is convex if for any $t \in [0, 1]$, $g(tx + (1 - t)y) \leq tg(x) + (1 - t)g(y)$.

Example: $g(x) = \exp(x)$, $g(x) = |x|$ or $g(x) = x^2$

If g is smooth enough, g is convex if $g''(x) \geq 0$ for all x



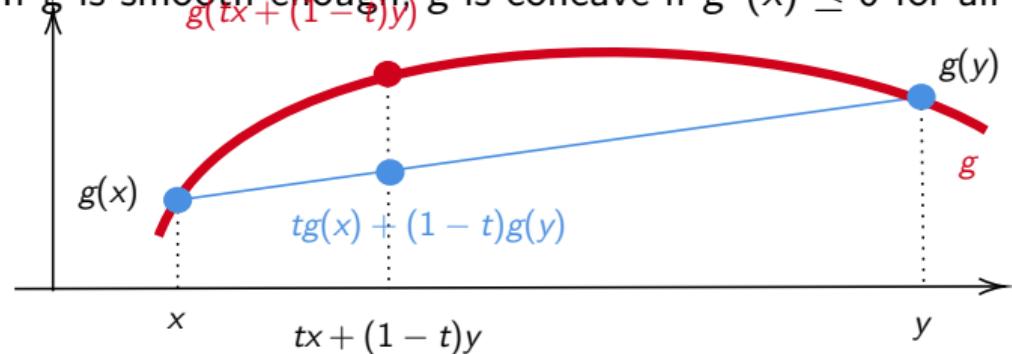
Convex and Concave

Definition 1.51: Concave Function

$g : \mathbb{R} \rightarrow \mathbb{R}$ is concave if for any $t \in [0, 1]$, $g(tx + (1 - t)y) \geq tg(x) + (1 - t)g(y)$.

Example: $g(x) = \log(x)$

If g is smooth enough, g is concave if $g''(x) \leq 0$ for all x



Approximation

Proposition 1.30: Taylor First Order Approximation

If g is continuously differentiable in the neighbourhood of a

$$g(x) \approx g(a) + g'(a) \cdot (x - a) \text{ or } g(a + h) \approx g(a) + g'(a) \cdot h$$

Sketch of the proof (1): definition of the derivative (in a)

$$g'(a) = \lim_{h \rightarrow 0} \frac{g(a + h) - g(a)}{h} \approx \frac{g(a + h) - g(a)}{h}$$

Sketch of the proof (2): fundamental theorem of analysis

$$g(x) = g(a) + \int_a^x g'(y) dy \approx g(a) + \int_a^x g'(a) dy \approx g(a) + g'(a) \cdot (x - a)$$

Approximation

since $\int_a^x g'(y)dy \approx g'(a) \cdot (x - a)$ but one could be more precise

$$\int_a^x g'(y)dy = \left[-(x - y)g'(y) \right]_a^x + \int_a^x (x - y)g''(y)dy$$

$$\int_a^x g'(y)dy = g'(a)(x - a) + \int_a^x (x - y)g''(y)dy$$

then iterate...

$$\int_a^x (x - y)g''(y)dy = \left[-\frac{(x - y)^2}{2}g''(y) \right]_a^x + \int_a^x \frac{(x - y)^2}{2}g^{(3)}(y)dy$$

$$\int_a^x (x - y)g''(y)dy \approx \frac{(x - a)^2}{2}g''(a)$$

Approximation

Proposition 1.31: Taylor Second Order Approximation

If g is twice continuously differentiable in the neighborhood of a

$$g(x) \approx g(a) + g'(a)(x - a) + \frac{g''(a)}{2}(x - a)^2$$

$$g(a + h) \approx g(a) + g'(a)h + \frac{g''(a)}{2}h^2$$

More generally, integration by parts leads to

$$\int_a^x \frac{(x - y)^{i-1}}{(i-1)!} g^{(i)}(y) dy = \left[-\frac{(x - y)^i}{(i)!} g^{(i)}(y) \right]_a^x + \int_a^x \frac{(x - y)^i}{(i)!} f^{(i+1)}(y) dy$$

Jensen Inequality

If g is twice continuously derivable in the neighborhood of a , then

$$g(x) \approx g(a) + g'(a)(x - a) + \frac{g''(a)}{2}(x - a)^2$$

Let X be a random variable, with expectation $\mathbb{E}[X] = a$, then

$$g(X) \approx g(\mathbb{E}[X]) + g'(\mathbb{E}[X])(X - \mathbb{E}[X]) + \frac{g''(\mathbb{E}[X])}{2}(X - \mathbb{E}[X])^2$$

so

$$\mathbb{E}[g(X)] \approx g(\mathbb{E}[X]) + \frac{g''(\mathbb{E}[X])}{2} \text{Var}[X]$$

Proposition 1.32: Jensen Inequality

Let g be a convex function, and X a random variable, then $\mathbb{E}[g(X)] \geq g(\mathbb{E}[X])$ (if the expected value exists).

Approximation in Higher Dimension

Proposition 1.33: Taylor Second Order Approximation, Univariate

If $g : \mathbb{R} \rightarrow \mathbb{R}$ is twice continuously differentiable in the neighborhood of a

$$g(a + h) \approx g(a) + g'(a)h + \frac{g''(a)}{2}h^2$$

Proposition 1.34: Taylor Second Order Approximation, Multivariate

If $g : \mathbb{R}^d \rightarrow \mathbb{R}$ is twice continuously differentiable in the neighbourhood of \mathbf{a}

$$g(\mathbf{a} + \mathbf{h}) \approx g(\mathbf{a}) + \nabla g(\mathbf{a})^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \mathbf{H}g(\mathbf{a})\mathbf{h}$$

where ∇g is the gradient of g and $\mathbf{H}g$ is the Hessian matrix, evaluated in \mathbf{a} .

Approximation

If $g : \mathbb{R}^2 \rightarrow \mathbb{R}$, the gradient is a vector and the Hessian a 2×2 matrix

$$\nabla g = \begin{pmatrix} \frac{\partial g}{\partial x} \\ \frac{\partial g}{\partial y} \end{pmatrix} \text{ and } \mathbf{H}g = \begin{pmatrix} \frac{\partial^2 g}{\partial x^2} & \frac{\partial^2 g}{\partial x \partial y} \\ \frac{\partial^2 g}{\partial y \partial x} & \frac{\partial^2 g}{\partial y^2} \end{pmatrix}$$

$$g(\mathbf{a} + \mathbf{h}) \approx g(\mathbf{a}) + \nabla g(\mathbf{a})^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \mathbf{H}g(\mathbf{a}) \mathbf{h}$$

Example: $g(\mathbf{x}) = \mathbf{x}^\top \mathbf{x}$, then

$$g(\mathbf{a} + \mathbf{h}) = (\mathbf{a} + \mathbf{h})^\top (\mathbf{a} + \mathbf{h}) = \mathbf{a}^\top \mathbf{a} + \mathbf{a}^\top \mathbf{h} + \mathbf{h}^\top \mathbf{a} + \mathbf{h}^\top \mathbf{h}$$

or $\mathbf{a}^\top \mathbf{h} = \mathbf{h}^\top \mathbf{a}$ so

$$g(\mathbf{a} + \mathbf{h}) = \underbrace{\mathbf{a}^\top \mathbf{a}}_{=g(\mathbf{a})} + \underbrace{2\mathbf{a}^\top}_{=\nabla g(\mathbf{a})^\top} \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \underbrace{2\mathbb{I}}_{=\mathbf{H}g(\mathbf{a})} \mathbf{h}$$

Gradients and Differentials

- The **derivative** of a $\mathbb{R} \rightarrow \mathbb{R}$ measures how a function changes with respect to its input.
- It represents the slope or rate of change of a function.
- **Notation:** $g'(x)$, $\frac{df}{dx}$ or $Dg(x)$
- **Interpretation:**
 - $g'(x) > 0$ means the function is increasing.
 - $g'(x) < 0$ means the function is decreasing.
 - $g'(x) = 0$ may indicate a maximum, minimum, or inflection point.
- **Example:** If $g(x) = x^2$, then $g'(x) = 2x$
- A **differential** represents a small change in a variable.
- $dg \approx g'(x) dx$ gives a linear approximation of change in g .
- Useful for approximating function values near a point.

Gradients and Differentials

- **Example:** If $g(x) = x^2$, then

$$dg = 2x \, dx$$

- This tells us how a small change dx affects $g(x)$.
- A **gradient** generalizes the derivative to multivariable functions $\mathbb{R}^d \rightarrow \mathbb{R}$.
- For $g(x, y)$, the gradient is a vector:

$$\nabla g = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

- Points in the direction of greatest increase of the function.
- Magnitude of ∇f shows how steep the increase is.
- **Example:** $g(x, y) = x^2 + y^2 \Rightarrow \nabla f = (2x, 2y)$
- **Example:** $g(x, y) = x + y \Rightarrow \nabla f = (1, 1)$ everywhere, $\Pi_k = \{(x, y) : g(x, y) = k\}$ it is a plane Π (linear), and $\vec{u} = \nabla f \perp \Pi_k$

Gradients and Differentials

- The **chain rule** handles compositions of functions.
- If $y = g(u)$ and $u = g(x)$, then:

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \cdot \frac{\partial u}{\partial x} \text{ or } \frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

- **Example:** If

$$y = \log(\sqrt{1 - x^2})$$

- Let $u = \sqrt{1 - x^2}$, then:

$$\frac{dy}{dx} = \frac{1}{u} \cdot \frac{du}{dx} = \frac{1}{\sqrt{1 - x^2}} \cdot \left(\frac{-x}{\sqrt{1 - x^2}} \right) = \frac{-x}{1 - x^2}$$

- So,

$$\frac{d}{dx} [\log(\sqrt{1 - x^2})] = \frac{-x}{1 - x^2}$$

- In higher dimensions: use Jacobians or total derivatives.

Optimization

Machine learning has a lot to do with optimization.

And a lot of concepts can be related to optimization.

E.g.

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \operatorname{argmin}_{m \in \mathbb{R}} \{\hat{\mathcal{R}}_2(m)\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \ell_2(y_i, m) \right\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n (y_i - m)^2 \right\}.$$

or, e.g.,

$$\text{median}(y) \in \operatorname{argmin}_{m \in \mathbb{R}} \{\hat{\mathcal{R}}_1(m)\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \ell_1(y_i, m) \right\} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n |y_i - m| \right\}.$$

Optimisation: continuous (differentiable) case

The problem is to solve $\min_{y \in \mathbb{R}} \{g(y)\}$

Note: $\min_{y \in \mathbb{R}} \{g(y)\} = \max_{y \in \mathbb{R}} \{-g(y)\}$

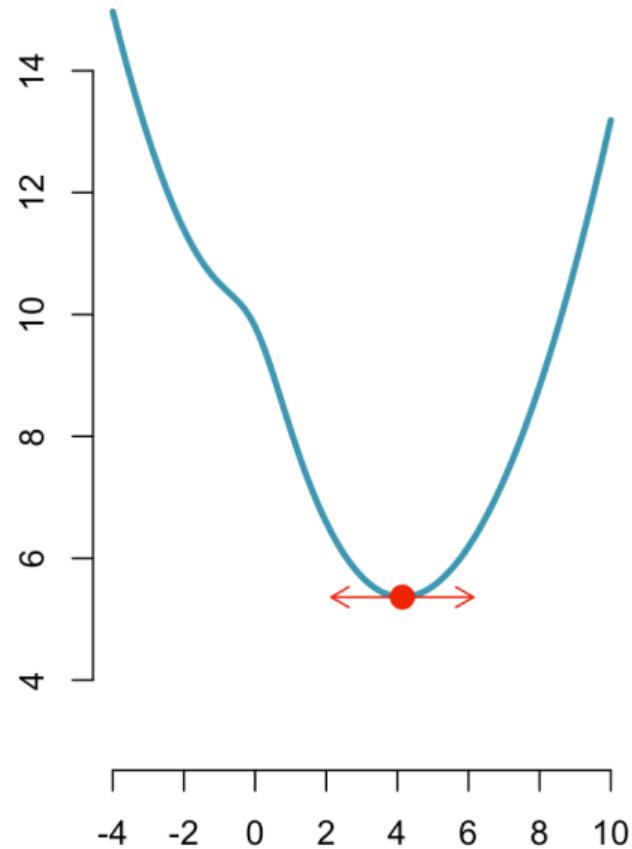
Note: $y^* \in \operatorname{argmin}_{y \in \mathbb{R}} \{g(y)\}$

and $\min_{y \in \mathbb{R}} \{g(y)\} = g(y^*)$.

First order condition

$$g'(y^*) = \frac{\partial g(y)}{\partial y} \Big|_{y=y^*} = 0$$

(necessary condition)



Optimisation: continuous (differentiable) case

First order condition

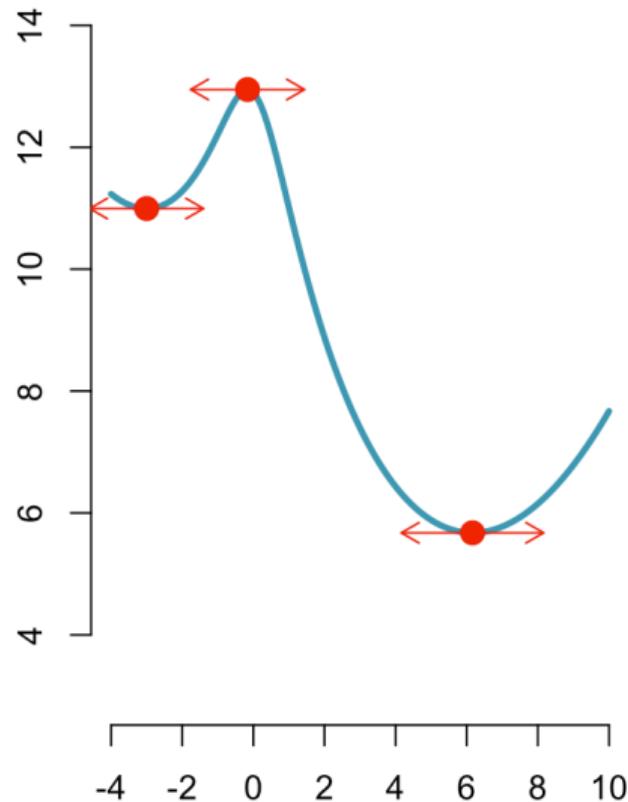
$$g'(y^*) = \frac{\partial g(y)}{\partial y} \Big|_{y=y^*} = 0$$

might be not sufficient

$$g''(y^*) = \frac{\partial^2 f}{\partial y^2} \Big|_{y=y^*} > 0 : \text{minimum}$$

$$g''(y^*) = \frac{\partial^2 f}{\partial y^2} \Big|_{y=y^*} < 0 : \text{maximum}$$

can be a local minimum...



Optimisation: continuous (differentiable) case

Example : $\{y_1, \dots, y_n\}$ in \mathbb{R} , let

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

$$\frac{\partial g(y)}{\partial y} = \frac{\partial}{\partial y} \sum_{i=1}^n (y_i - y)^2 = \sum_{i=1}^n \frac{\partial (y_i - y)^2}{\partial y} = \sum_{i=1}^n -2(y_i - y)$$

so

$$\left. \frac{\partial g(y)}{\partial y} \right|_{y=y^*} = 0 \text{ if and only if } \sum_{i=1}^n (y_i - y^*) = 0 \text{ or } \sum_{i=1}^n y_i = ny^*$$

i.e. $y^* = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$.

Optimisation: continuous (differentiable) case

Solving $g'(y^*) = 0$ numerically

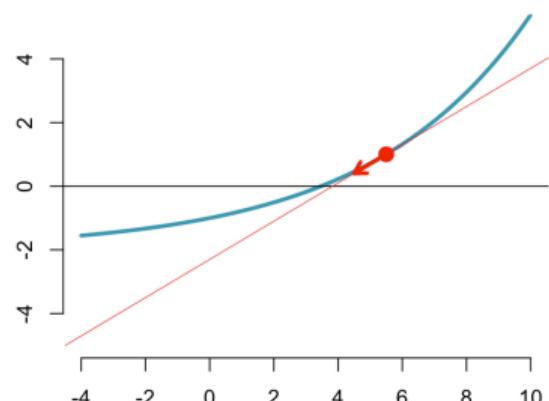
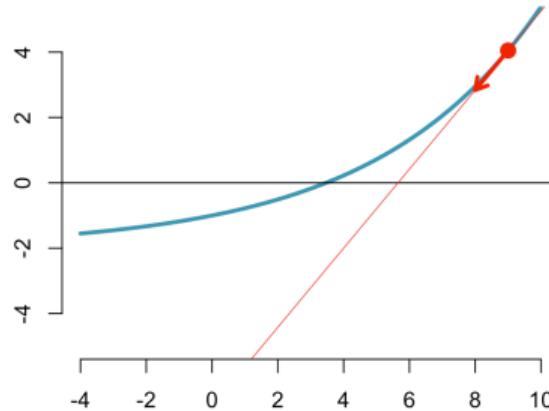
Newton's method: solve $g(y^*) = 0$

$$g(y) \simeq g(y_0) + g'(y_0)(y - y_0)$$

If $g(y) \simeq 0$, $g(y_0) + g'(y_0)(y - y_0) \simeq 0$

Start from y_0 , then

$$y_{k+1} = y_k - \frac{g(y_k)}{g'(y_k)}$$



Optimisation: continuous (differentiable) case

To solve $g'(y^*) = 0$ numerically

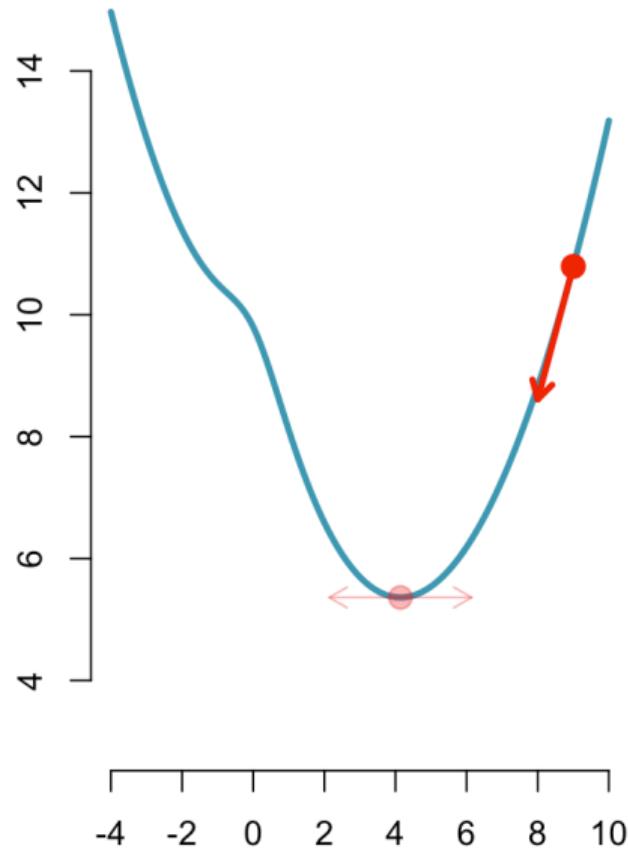
Start from y_0 , then

$$y_{k+1} = y_k - \frac{g'(y_k)}{g''(y_k)}$$

$g'(y_k)$ gives the direction

$g''(y_k)$ gives the speed of convergence

(close to the minimum $g''(y_k) > 0$)



Optimisation: continuous (differentiable) case

In python

```
1 > import statistics as stat
2 > v=[0.89367,-1.04729,1.97133,-0.38363,1.65414]
3 > stat.mean(v)
4 0.617644
5 > import numpy as np
6 > def f0(x):
7 ...     return np.sum((np.array(v)-x)**2)
8 > f = np.vectorize(f0)
9 > from scipy.optimize import fminbound
10 > fminbound(f, -1, 1)
11 0.6176439999999992
```

or

Optimisation: continuous (differentiable) case

```
1 > from scipy.optimize import minimize
2 > minimize(f, 0, method='nelder-mead')
3   final_simplex: (array([[0.617625],
4     [0.6176875]]), array([6.75753495, 6.75753495]))
5     fun: 6.757534946524999
6     message: 'Optimization terminated successfully.'
```

or

```
1 > def g(x):
2 ...   s=[0]*len(v)
3 ...   for i in range(len(v)):
4 ...     s[i]=((v[i]-x)**2)
5 ...   return sum(s)
6 > fminbound(f, -1, 1)
7 0.6176439999999992
```

Optimisation: continuous (differentiable) case

and in R

```
1 > v = c(0.89367,-1.04729,1.97133,-0.38363,1.65414)
2 > mean(v)
3 [1] 0.617644
4 > f = function(x) sum((v-x)^2)
5 > optim(0, f)
6 $par
7 [1] 0.6175781
8 $value
9 [1] 6.757535
```

Optimisation: continuous (differentiable) case

The problem is $\min_{\mathbf{y} \in \mathbb{R}^p} \{g(\mathbf{y})\}$ or $\min_{(y_1, \dots, y_p) \in \mathbb{R}^p} \{g(y_1, \dots, y_p)\}$

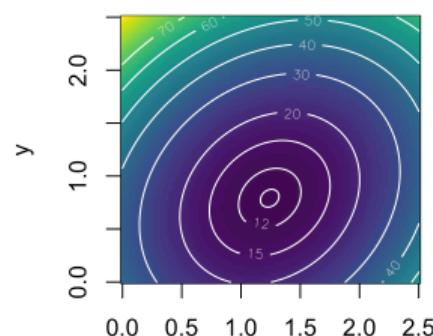
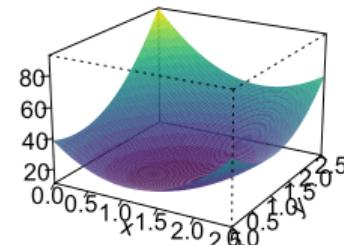
First order conditions: $\nabla g(\mathbf{y}^*) = \mathbf{0}$,

$$\frac{\partial g(y_1, y_2, \dots, y_p)}{\partial y_1} \Big|_{\mathbf{y}=\mathbf{y}^*} = 0$$

$$\frac{\partial g(y_1, y_2, \dots, y_p)}{\partial y_2} \Big|_{\mathbf{y}=\mathbf{y}^*} = 0$$

⋮

$$\frac{\partial g(y_1, y_2, \dots, y_p)}{\partial y_p} \Big|_{\mathbf{y}=\mathbf{y}^*} = 0$$



Optimisation: continuous (differentiable) case

Example : $\{(x_1, y_1), \dots, (x_n, y_n)\}$ in \mathbb{R}^2 , let

$$g(a, b) = \sum_{i=1}^n (y_i - [a + bx_i])^2$$

$$\frac{\partial g(a, b)}{\partial a} = -2 \sum_{i=1}^n (y_i - [a + bx_i]) = -2(n\bar{y} - [a + bn\bar{y}])$$

$$\frac{\partial g(a, b)}{\partial b} = -2 \sum_{i=1}^n (y_i - [a + bx_i])x_i$$

$$\left. \frac{\partial g(a, b)}{\partial a} \right|_{(a,b)=(a^*,b^*)} = 0 \text{ means that } \bar{y} = a^* + b^*\bar{x},$$

$$\left. \frac{\partial g(a, b)}{\partial b} \right|_{(a,b)=(a^*,b^*)} = 0 \text{ means that } \hat{\varepsilon} \perp \mathbf{x}, \hat{\varepsilon}_i = y_i - [a^* + b^*x_i],$$

Optimisation: continuous (differentiable) case

To solve $\nabla g(\mathbf{y}^*) = \mathbf{0}$ numerically

Start from \mathbf{y}_0 , then

$$\mathbf{y}_{k+1} = \mathbf{y}_k - \mathbf{H}_k^{-1} \nabla g(\mathbf{y}_k)$$

$\nabla g(\mathbf{y}_k)$ gives the direction

\mathbf{H}_k^{-1} gives the speed of convergence

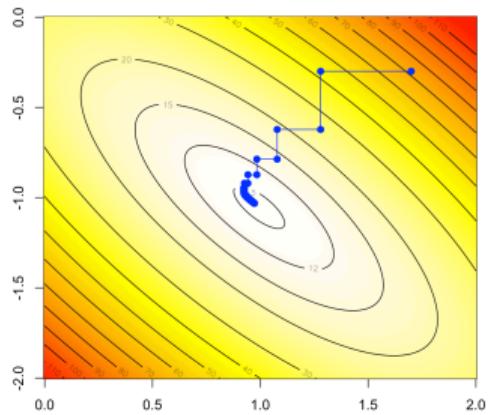
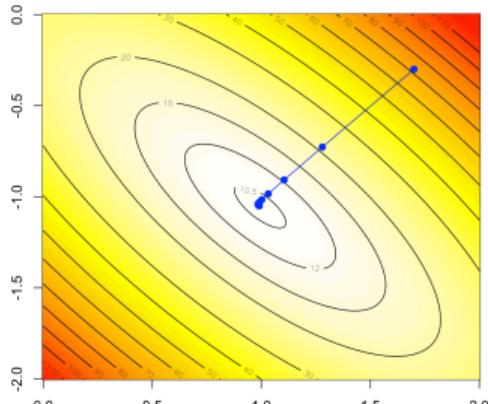
\mathbf{H}_k^{-1} is the inverse of the Hessian matrix

One can also consider some numerical tricks,

see **coordinate descent**

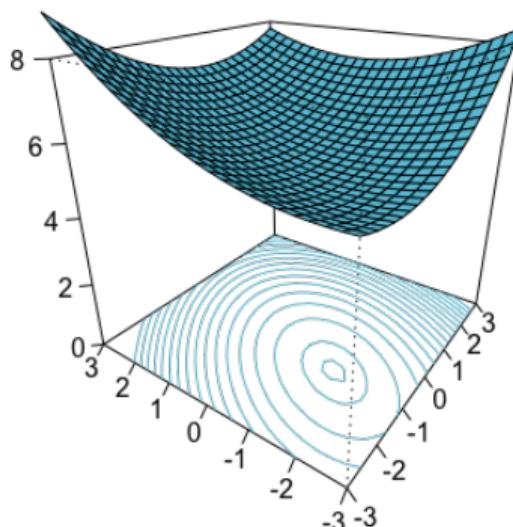
where we iterate on the dimension

(univariate optimisation problems)



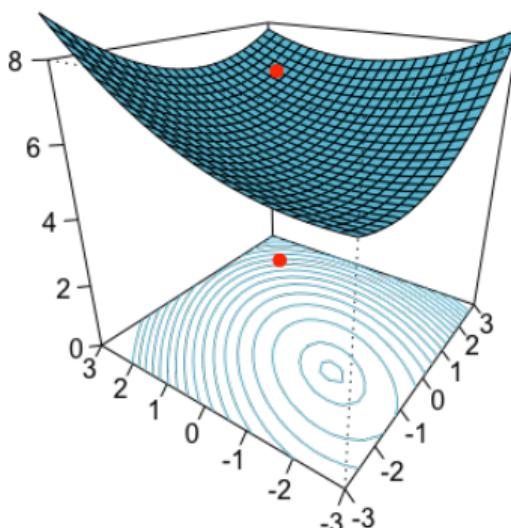
Optimisation: continuous (differentiable) case

The problem is to find $\mathbf{y}^* = \underset{\mathbf{y} \in \mathbb{R}^2}{\operatorname{argmin}} \{g(\mathbf{y})\}$,
for some convex function g
corresponding to $\nabla g(\mathbf{y}^*) = \mathbf{0}$



Optimisation: continuous (differentiable) case

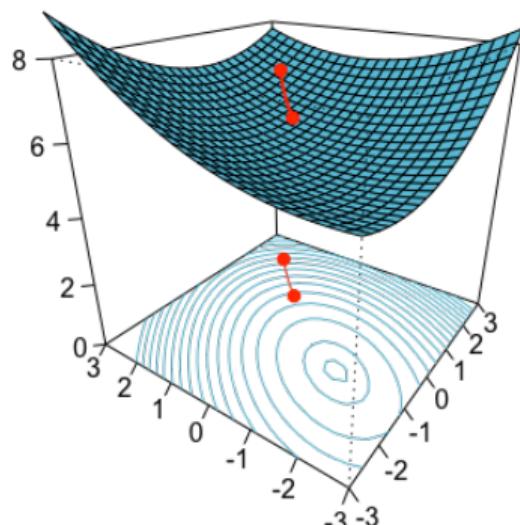
The problem is to find $\mathbf{y}^* = \underset{\mathbf{y} \in \mathbb{R}^2}{\operatorname{argmin}} \{g(\mathbf{y})\}$,



for some convex function g
corresponding to $\nabla g(\mathbf{y}^*) = \mathbf{0}$

Start from \mathbf{y}_0 ,

Optimisation: continuous (differentiable) case



The problem is to find $\mathbf{y}^* = \underset{\mathbf{y} \in \mathbb{R}^2}{\operatorname{argmin}} \{g(\mathbf{y})\}$,
for some convex function g
corresponding to $\nabla g(\mathbf{y}^*) = \mathbf{0}$

Start from \mathbf{y}_0 , then

$$\mathbf{y}_1 = \mathbf{y}_0 - \mathbf{H}_0^{-1} \nabla g(\mathbf{y}_0)$$

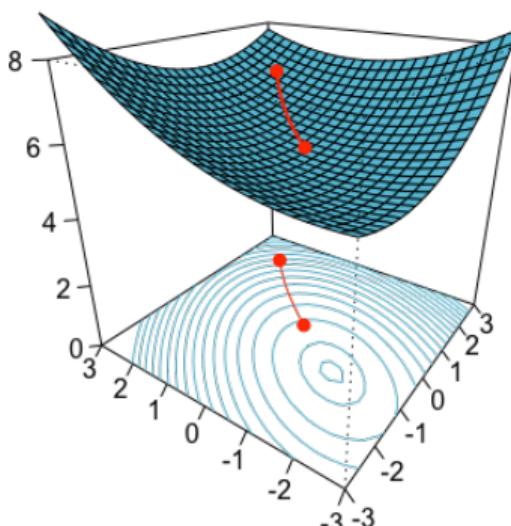
where $\nabla g(\mathbf{y}_k)$ and \mathbf{H}_k are

$$\left(\begin{array}{c} \frac{\partial g(y_1, y_2)}{\partial y_1} \\ \frac{\partial g(y_1, y_2)}{\partial y_2} \end{array} \right) \Bigg|_{\mathbf{y}=\mathbf{y}_0} \quad \left(\begin{array}{cc} \frac{\partial^2 g(y_1, y_2)}{\partial y_1^2} & \frac{\partial^2 g(y_1, y_2)}{\partial y_1 \partial y_2} \\ \frac{\partial^2 g(y_1, y_2)}{\partial y_2 \partial y_1} & \frac{\partial^2 g(y_1, y_2)}{\partial y_2^2} \end{array} \right)$$

that are either known, or approximated

Optimisation: continuous (differentiable) case

The problem is to find $\mathbf{y}^* = \underset{\mathbf{y} \in \mathbb{R}^2}{\operatorname{argmin}} \{g(\mathbf{y})\}$,



for some convex function g
corresponding to $\nabla g(\mathbf{y}^*) = \mathbf{0}$

Start from \mathbf{y}_0 , then

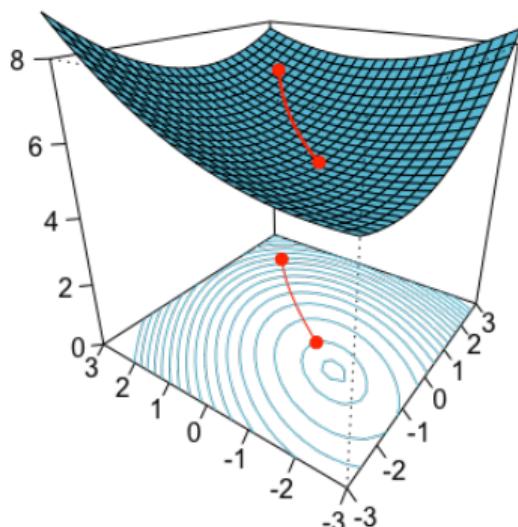
$$\mathbf{y}_2 = \mathbf{y}_1 - \mathbf{H}_1^{-1} \nabla g(\mathbf{y}_1)$$

where $\nabla g(\mathbf{y}_1)$ and \mathbf{H}_1 are

$$\left(\begin{array}{c} \frac{\partial g(y_1, y_2)}{\partial y_1} \\ \frac{\partial g(y_1, y_2)}{\partial y_2} \end{array} \right) \Bigg|_{\mathbf{y}=\mathbf{y}_1} \quad \left(\begin{array}{cc} \frac{\partial^2 g(y_1, y_2)}{\partial y_1^2} & \frac{\partial^2 g(y_1, y_2)}{\partial y_1 \partial y_2} \\ \frac{\partial^2 g(y_1, y_2)}{\partial y_2 \partial y_1} & \frac{\partial^2 g(y_1, y_2)}{\partial y_2^2} \end{array} \right)$$

that are either known, or approximated

Optimisation: continuous (differentiable) case



The problem is to find $\mathbf{y}^* = \underset{\mathbf{y} \in \mathbb{R}^2}{\operatorname{argmin}} \{g(\mathbf{y})\}$,

for some convex function g
corresponding to $\nabla g(\mathbf{y}^*) = \mathbf{0}$

Start from \mathbf{y}_0 , then iterate

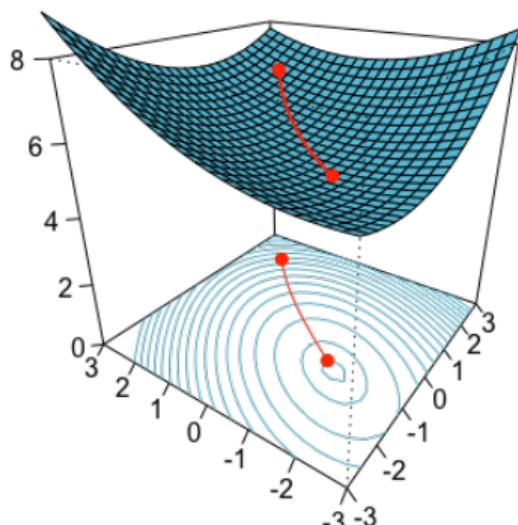
$$\mathbf{y}_{k+1} = \mathbf{y}_k - \mathbf{H}_k^{-1} \nabla g(\mathbf{y}_k)$$

where $\nabla g(\mathbf{y}_k)$ and \mathbf{H}_k are

$$\left(\begin{array}{c} \frac{\partial g(y_1, y_2)}{\partial y_1} \\ \frac{\partial g(y_1, y_2)}{\partial y_2} \end{array} \right) \Big|_{\mathbf{y}=\mathbf{y}_k} \quad \left(\begin{array}{cc} \frac{\partial^2 g(y_1, y_2)}{\partial y_1^2} & \frac{\partial^2 g(y_1, y_2)}{\partial y_1 \partial y_2} \\ \frac{\partial^2 g(y_1, y_2)}{\partial y_2 \partial y_1} & \frac{\partial^2 g(y_1, y_2)}{\partial y_2^2} \end{array} \right)$$

that are either known, or approximated

Optimisation: continuous (differentiable) case



The problem is to find $\mathbf{y}^* = \underset{\mathbf{y} \in \mathbb{R}^2}{\operatorname{argmin}} \{g(\mathbf{y})\}$,

for some convex function g

corresponding to $\nabla g(\mathbf{y}^*) = \mathbf{0}$

Start from \mathbf{y}_0 , then iterate

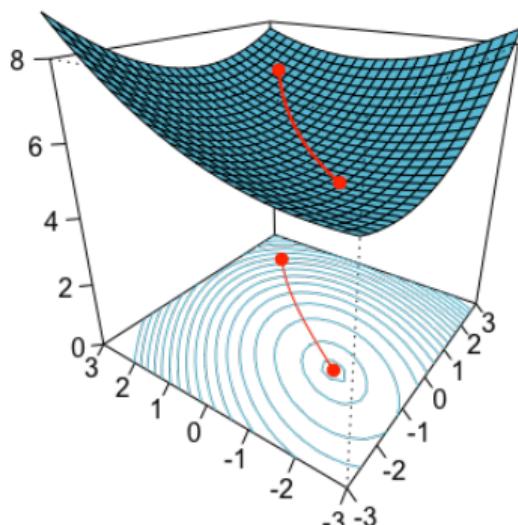
$$\mathbf{y}_{k+1} = \mathbf{y}_k - \mathbf{H}_k^{-1} \nabla g(\mathbf{y}_k)$$

where $\nabla g(\mathbf{y}_k)$ and \mathbf{H}_k are

$$\left(\begin{array}{c} \frac{\partial g(y_1, y_2)}{\partial y_1} \\ \frac{\partial g(y_1, y_2)}{\partial y_2} \end{array} \right) \Big|_{\mathbf{y}=\mathbf{y}_k} \quad \left(\begin{array}{cc} \frac{\partial^2 g(y_1, y_2)}{\partial y_1^2} & \frac{\partial^2 g(y_1, y_2)}{\partial y_1 \partial y_2} \\ \frac{\partial^2 g(y_1, y_2)}{\partial y_2 \partial y_1} & \frac{\partial^2 g(y_1, y_2)}{\partial y_2^2} \end{array} \right)$$

that are either known, or approximated

Optimisation: continuous (differentiable) case



The problem is to find $\mathbf{y}^* = \underset{\mathbf{y} \in \mathbb{R}^2}{\operatorname{argmin}} \{g(\mathbf{y})\}$,
for some convex function g
corresponding to $\nabla g(\mathbf{y}^*) = \mathbf{0}$

Start from \mathbf{y}_0 , then

$$\mathbf{y}_{k+1} = \mathbf{y}_k - \mathbf{H}_k^{-1} \nabla g(\mathbf{y}_k)$$

where $\nabla g(\mathbf{y}_k)$ and \mathbf{H}_k are

$$\left(\begin{array}{c} \frac{\partial g(y_1, y_2)}{\partial y_1} \\ \frac{\partial g(y_1, y_2)}{\partial y_2} \end{array} \right) \Big|_{\mathbf{y}=\mathbf{y}_k} \quad \left(\begin{array}{cc} \frac{\partial^2 g(y_1, y_2)}{\partial y_1^2} & \frac{\partial^2 g(y_1, y_2)}{\partial y_1 \partial y_2} \\ \frac{\partial^2 g(y_1, y_2)}{\partial y_2 \partial y_1} & \frac{\partial^2 g(y_1, y_2)}{\partial y_2^2} \end{array} \right)$$

that are either known, or approximated

Constrained Optimisation: continuous case

Constrained optimization seeks to find the best solution to an objective function subject to constraints:

$$\max_x \{f(x)\} \quad \text{subject to} \quad \begin{cases} g_i(x) \leq 0, \quad \forall i \\ h_j(x) = 0, \quad \forall j, \end{cases}$$

where:

- $f(x)$ is the objective function.
- $g_i(x)$ are inequality constraints.
- $h_j(x)$ are equality constraints.

The Lagrangian function transforms a constrained problem into an unconstrained one:

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$$

Constrained Optimisation: continuous case

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$$

where:

- λ_i are Lagrange multipliers for inequality constraints.
- μ_j are Lagrange multipliers for equality constraints.

Necessary conditions for a local optimum in constrained optimization:

- Stationarity: $\nabla f(x) + \sum_i \lambda_i \nabla g_i(x) + \sum_j \mu_j \nabla h_j(x) = 0$
- Primal feasibility: $g_i(x) \leq 0, \quad h_j(x) = 0$
- Dual feasibility: $\lambda_i \geq 0$
- Complementary slackness: $\lambda_i g_i(x) = 0$

so-called **Karush-Kuhn-Tucker** (KKT) conditions.

Constrained Optimisation: continuous case

For equality constraints only:

- Construct the Lagrangian:

$$\mathcal{L}(x, \mu) = f(x) + \sum_j \mu_j h_j(x)$$

- Solve $\nabla \mathcal{L} = 0$ to find x and μ .

For inequality constraints:

$$\lambda_i g_i(x) = 0$$

This means:

- If $g_i(x) < 0$, then $\lambda_i = 0$ (inactive constraint).
- If $g_i(x) = 0$, then $\lambda_i \geq 0$ (active constraint).

Constrained Optimisation: continuous case

The dual function:

$$q(\lambda, \mu) = \max_x \mathcal{L}(x, \lambda, \mu)$$

gives a lower bound on the optimal value:

$$q(\lambda, \mu) \leq f^*$$

- Weak Duality: $q(\lambda, \mu) \leq f^*$
- Strong Duality: $q(\lambda^*, \mu^*) = f^*$ (holds if Slater's condition is satisfied)

Slater's condition ensures strong duality for convex optimization problems. It states that if the problem:

$$\min g(x) \quad \text{subject to } g_i(x) \leq 0$$

is convex and there exists a strictly feasible point x_0 such that:

$$g_i(x_0) < 0 \quad \forall i$$

Constrained Optimisation: continuous case

then strong duality holds, meaning:

$$\max_{\lambda \geq 0} \min_x \{\mathcal{L}(x, \lambda)\} = \min_x \{f(x)\}$$

This guarantees that the primal and dual problems have the same optimal value.

Constrained Optimisation: continuous case

Quadratic Programming corresponds to minimizing

$$f(x) = \frac{1}{2}x^T Qx + c^T x$$

subject to:

$$Ax \leq b$$

Solution obtained using KKT conditions or interior-point methods.

- Constrained optimization incorporates constraints into objective minimization.
- The Lagrangian method helps transform constrained problems into unconstrained ones.
- KKT conditions provide necessary conditions for optimality.
- Duality theory helps analyze convex problems.

Convex Optimization Problem: Ridge

$$\mathcal{L}(\beta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^\top \beta - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 = \frac{1}{2} (\mathbf{X}\beta - \mathbf{y})^\top (\mathbf{X}\beta - \mathbf{y}) + \frac{1}{2} \beta^\top \beta$$

Primal problem

$$\min_{\beta} \left\{ \frac{1}{2} \beta^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p) \beta - \beta^\top \mathbf{X}^\top \mathbf{y} \right\}$$

The solution of the unconstrained quadratic program is $\hat{\beta} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$

Dual problem

$$\min_{\alpha} \left\{ \frac{1}{2} \alpha^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbb{I}_n) \alpha - \lambda \alpha^\top \mathbf{y} \right\}$$

The solution is $\hat{\alpha} = \lambda (\mathbf{X} \mathbf{X}^\top + \lambda \mathbb{I}_n)^{-1} \mathbf{y}$, i.e. $\hat{\beta} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top + \lambda \mathbb{I}_n)^{-1} \mathbf{y}$

Differentials and Gradients

Let $f: \mathbb{R}^n \rightarrow \mathbb{R}^k$, the “derivative” is $\nabla f(\mathbf{x})$, a $n \times k$ matrix, and

$$f(\mathbf{x} + \mathbf{u}) \sim f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \mathbf{u} = f(\mathbf{x}) + df(\mathbf{x}; \mathbf{u})$$

where the differential df is linear in \mathbf{u} .

Chain rule for differentials

If $h(\mathbf{x}) = g(f(\mathbf{x}))$

$$dh(\mathbf{x}; \mathbf{u}) = dg(f(\mathbf{x}); df(\mathbf{x}; \mathbf{u})) .$$

From chain rule to back-propagation

Differentials and Gradients

If $h(\mathbf{x}) = f_n(\cdots f_2(f_1(\mathbf{x}, \boldsymbol{\theta}_1), \boldsymbol{\theta}_2), \cdots, \boldsymbol{\theta}_n)$. Let

$$\mathbf{y}_0 = \mathbf{x}, \quad \mathbf{y}_i = f_i(\mathbf{y}_{i-1}, \boldsymbol{\theta}_i), \text{ and } \mathbf{y}_n = h(\mathbf{x}).$$

Step 1

$$\frac{\partial h}{\partial \mathbf{y}_{n-1}} = \frac{\partial f_n(\mathbf{y}_{n-1}, \boldsymbol{\theta}_n)}{\partial \mathbf{y}_{n-1}} \text{ and } \frac{\partial h}{\partial \boldsymbol{\theta}_n} = \frac{\partial f_n(\mathbf{y}_{n-1}, \boldsymbol{\theta}_n)}{\partial \boldsymbol{\theta}_n}$$

then loop.

Back-propagation proceeds by obtaining the derivatives with respect to $\boldsymbol{\theta}_i$ and \mathbf{y}_{i-1} from the derivative with respect to \mathbf{y}_i

$$\frac{\partial h}{\partial \boldsymbol{\theta}_i} = \frac{\partial h}{\partial \mathbf{y}_i} \cdot \frac{\partial \mathbf{y}_i}{\partial \boldsymbol{\theta}_i} = \frac{\partial h}{\partial \mathbf{y}_i} \cdot \frac{\partial f_i(\mathbf{y}_{i-1}, \boldsymbol{\theta}_i)}{\partial \boldsymbol{\theta}_i}$$

$$\frac{\partial h}{\partial \mathbf{y}_{i-1}} = \frac{\partial h}{\partial \mathbf{y}_i} \cdot \frac{\partial \mathbf{y}_i}{\partial \mathbf{y}_{i-1}} = \frac{\partial h}{\partial \mathbf{y}_i} \cdot \frac{\partial f_i(\mathbf{y}_{i-1}, \boldsymbol{\theta}_i)}{\partial \mathbf{y}_{i-1}}$$

(etc)

Differentials and Gradients

For **back-propagation**, compute

$$\frac{\partial h}{\partial \theta_i} = \left(\left(\cdots \left(\frac{\partial h}{\partial y_{n-1}} \times \frac{\partial y_{n-1}}{\partial y_{n-2}} \right) \times \cdots \right) \times \frac{\partial y_{i+1}}{\partial y_i} \right) \times \frac{\partial y_i}{\partial \theta_i},$$

For **forward-propagation**, compute

$$\frac{\partial h}{\partial \theta_i} = \frac{\partial h}{\partial y_{n-1}} \times \left(\frac{\partial y_{n-1}}{\partial y_{n-2}} \times \left(\cdots \times \left(\frac{\partial y_{i+1}}{\partial y_i} \times \frac{\partial y_i}{\partial \theta_i} \right) \cdots \right) \right)$$

Convex Optimization Algorithms

$$\min_{\mathbf{x}} \{g(\mathbf{x})\}$$

with g convex, and differentiable.

Algorithm 1: Gradient Descent

- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} - h_t \nabla g(\mathbf{x}^{(t-1)})$
-

Heuristics: Taylor expansion

$$g(\mathbf{y}) \sim g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2h} \|\mathbf{y} - \mathbf{x}\|^2$$

From Gradient Descent to Newton's Method

$$\min_{\mathbf{x}} \{g(\mathbf{x})\}$$

with g convex, twice differentiable.

Algorithm 2: Newton's Method

- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 $\mathbf{H}_t \leftarrow \nabla^2 g(\mathbf{x}^{(t-1)})$;
 - 4 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} - \mathbf{H}_t^{-1} \nabla g(\mathbf{x}^{(t-1)})$
-

Use a better quadratic approximation in Taylor expansion – $\frac{1}{h} \mathbb{I} \rightarrow H$,

$$g(\mathbf{y}) \sim g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^\top \mathbf{H} g(\mathbf{x}) (\mathbf{y} - \mathbf{x})$$

Convex Optimization Problem

$$\min_x \{g(\mathbf{x})\}$$

with g convex, but non-differentiable.

Algorithm 3: Subgradient 'Descent'

- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 $\mathbf{g}^{(t-1)} \in \partial g(\mathbf{x}^{(t-1)})$;
 - 4 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} - h_t \mathbf{g}^{(t-1)}$
-

Note that it is not necessarily a descent, so pick

$$\mathbf{x}^* = \operatorname{argmin}\{g(\mathbf{x}^{(0)}), g(\mathbf{x}^{(1)}), g(\mathbf{x}^{(2)}), \dots\}$$

Convex Optimization Problem

$$\min_{\mathbf{x}} \{g(\mathbf{x})\} \text{ or } \min_{\mathbf{x}} \{f_1(\mathbf{x}) + g_2(\mathbf{x})\}$$

with f_1 and f_2 convex, but f_2 non-differentiable.

Algorithm 4: Proximal Gradient 'Descent'

- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 $\gamma_h(x) = \frac{1}{h}(\mathbf{x} - \text{proximal}_{h,f_2}(\mathbf{x} - h\nabla f_1(\mathbf{x})))$;
 - 4 $\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} - h_t \gamma_{h_{t-1}}(\mathbf{x}^{(t-1)}) = \text{proximal}_{h,f_2}(\mathbf{x}^{(t-1)} - h\nabla f_1(\mathbf{x}^{(t-1)}))$
-

γ_h is the "generalized gradient of g ".

The trick is that $\text{proximal}_{h,f_2}(\cdot)$ usually has a closed form in most applications.

Constrained Convex Optimization

Let $C \subset \mathbb{R}^n$ denote a convex set, with g convex, and differentiable.

$$\min_{\mathbf{x} \in C} \{g(\mathbf{x})\} \iff \min_{\mathbf{x}} \{g(\mathbf{x}) + I_C(\mathbf{x})\} \text{ where } I_C(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{x} \in C \\ \infty & \text{if } \mathbf{x} \notin C \end{cases}$$

$$\text{proximal}_h(\mathbf{x}) = \operatorname{argmin}_{\mathbf{z}} \left\{ \frac{1}{2h} \|\mathbf{x} - \mathbf{z}\|_2^2 + \lambda I_C(\mathbf{z}) \right\} = \operatorname{argmin}_{\mathbf{z} \in C} \left\{ \frac{1}{2h} \|\mathbf{x} - \mathbf{z}\|_2^2 \right\}$$

i.e. $\text{proximal}_h(\mathbf{x})$ is the projection operator onto C , $\Pi_C(\mathbf{x})$

Algorithm 5: Projected Gradient Descent

- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 $\mathbf{x}^{(t)} \leftarrow \Pi_C(\mathbf{x}^{(t-1)} - h_t \nabla f(\mathbf{x}^{(t-1)}))$
-

Coordinate Descent

Let $\{\vec{e}_1, \dots, \vec{e}_n\}$ denote the standard basis in \mathbb{R}^n ,

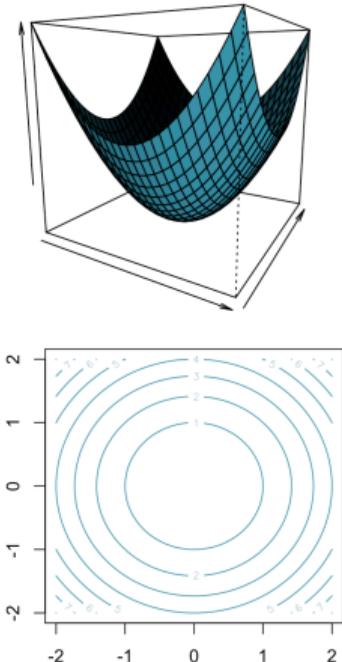
$$\vec{e}_i = (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^n$$

Proposition 1.35

If $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, differentiable,

$$g(\mathbf{x}) \leq g(\mathbf{x} + \delta \vec{e}_i), \forall i \implies g(\mathbf{x}) = \min\{g\}$$

i.e. if we are at a point \mathbf{x} such that $g(\mathbf{x})$ is minimized along each coordinate axis, then we have found a global minimizer.



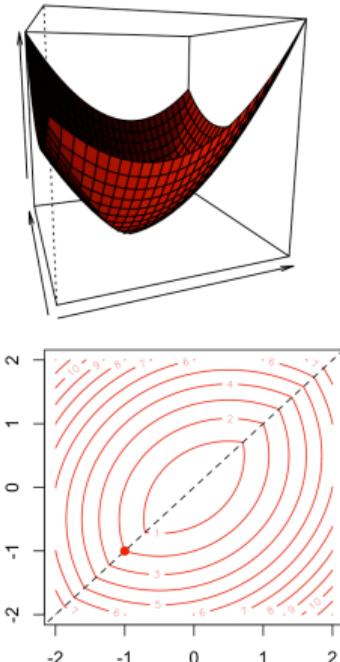
Coordinate Descent

Proposition 1.36

If $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, but **not differentiable**,

$$g(\mathbf{x}) \leq g(\mathbf{x} + \delta \vec{\mathbf{e}}_i), \quad \forall i \not\Rightarrow g(\mathbf{x}) = \min\{g\} = g^*$$

i.e. if we are at a point \mathbf{x} such that $g(\mathbf{x})$ is minimized along each coordinate axis, then we have **not** found a global minimizer.



Coordinate Descent

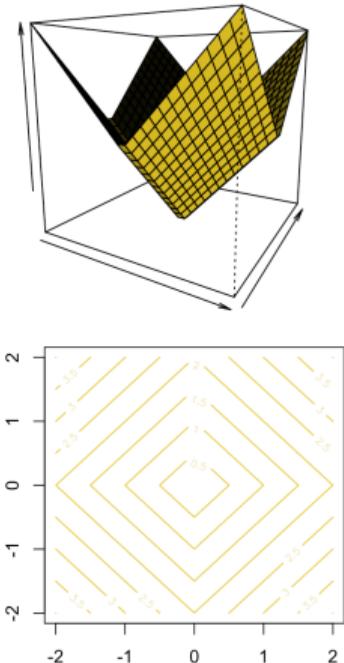
Proposition 1.37

If $g : \mathbb{R}^n \rightarrow \mathbb{R}$ can be written, for f_0, f_i convex,

$$g(\mathbf{x}) = f_0(\mathbf{x}) + \underbrace{\sum_{i=1}^n f_i(x_i)}_{\text{separable}}, \quad \text{where } \begin{cases} f_0 \text{ differentiable} \\ f_i \text{ non-differentiable} \end{cases}$$

$$g(\mathbf{x}) \leq g(\mathbf{x} + \delta \vec{e}_i), \quad \forall i \implies g(\mathbf{x}) = \min\{g\} = g^*$$

i.e. if we are at a point \mathbf{x} such that $g(\mathbf{x})$ is minimized along each coordinate axis, then we have found a global minimizer.



Coordinate Descent

If we want to solve $\min\{g(\mathbf{x})\}$ for some $g : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$g(\mathbf{x}) = f_0(\mathbf{x}) + \underbrace{\sum_{i=1}^n f_i(x_i)}_{\text{separable}}, \quad \text{where } \begin{cases} f_0 \text{ convex and differentiable} \\ f_i \text{ convex and non-differentiable} \end{cases}$$

we can use a **coordinate descent algorithm**

Algorithm 6: Coordinate Descent

- 1 initialization : $\mathbf{x}^{(0)}$;
 - 2 **for** $t=1,2,\dots$ **do**
 - 3 **for** $j=1,2,\dots,n$ **do**
 - 4 $\mathbf{x}_j^{(t)} \leftarrow \operatorname{argmin}\{g(\mathbf{x}_1^{(t)}, \dots, \mathbf{x}_{j-1}^{(t)}, x_j, \mathbf{x}_{j+1}^{(t-1)}, \dots, \mathbf{x}_n^{(t-1)})\}$
-

freakonometrics

freakonometrics.hypotheses.org

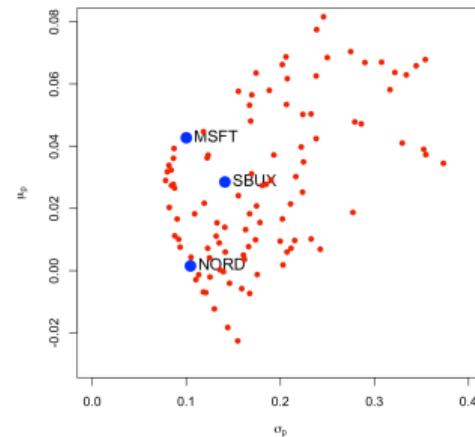
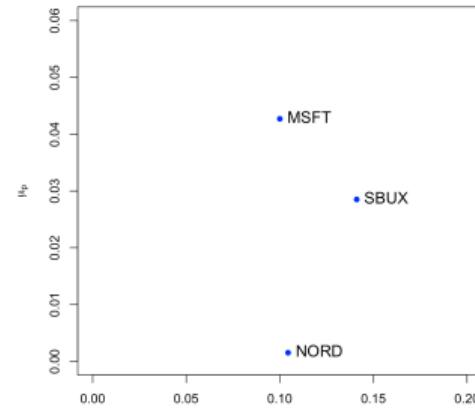
Portfolio Optimization

Consider n assets, with expected return μ and covariance matrix Σ

```
1 > asset.names = c("MSFT", "NORD", "SBUX")
2 > mu.vec = c(0.0427, 0.0015, 0.0285)
3 > names(mu.vec) = asset.names
4 > sigma.mat = matrix(c(0.0100, 0.0018,
   0.0113, 0.0018, 0.0109, 0.0026, 0.0113,
   0.0026, 0.0199), nrow=3, ncol=3)
5 > dimnames(sigma.mat) = list(asset.names,
   asset.names)
```

A **portfolio** is $\omega \in \mathbb{R}^n$ with $\omega^\top \mathbf{1} = 1$,

$$\mathbb{E}(P) = \omega^\top \mu, \text{Var}(P) = \omega^\top \Sigma \omega$$



Portfolio Optimization – Markowitz (1952)

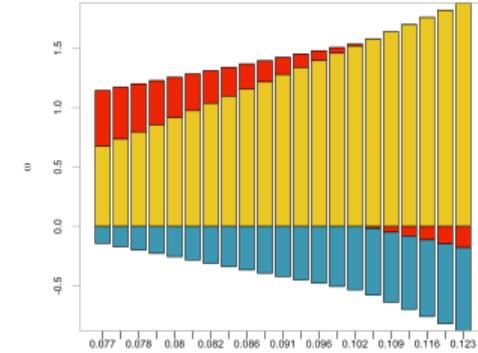
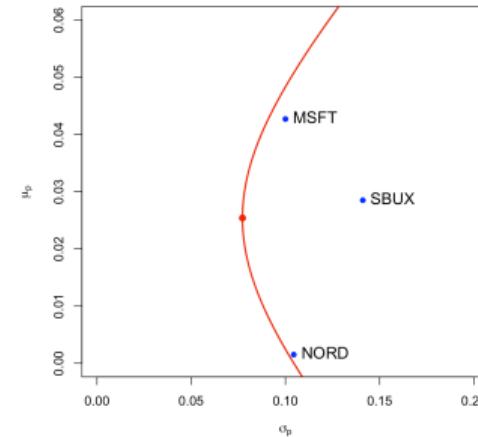
Efficient portfolio: ω_r^* ,

$$\omega_r^* = \operatorname{argmin}\{\omega^\top \Sigma \omega\} \text{ s.t. } \begin{cases} \omega^\top \mathbf{1} = 1 \\ \omega^\top \mu = r \end{cases}$$

or minimum variance portfolio

$$\omega^* = \operatorname{argmin}\{\omega^\top \Sigma \omega\} \text{ s.t. } \omega^\top \mathbf{1} = 1$$

```
1 > one.vec = rep(1, 3)
2 > sigma.inv.mat = solve(sigma.mat)
3 > top.mat = sigma.inv.mat%*%one.vec
4 > bot.val = as.numeric((t(one.vec)%*%sigma.
   inv.mat%*%one.vec))
5 > m.mat = top.mat/bot.val
6 > m.mat[, 1]
7 [1]  0.674  0.470 -0.144
```



Possible only if short sells are possible,

MSFT NORD SBUX

Minimum Variance Portfolio – Markowitz (1952)

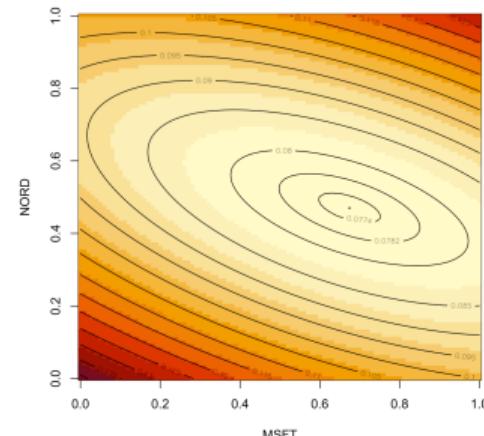
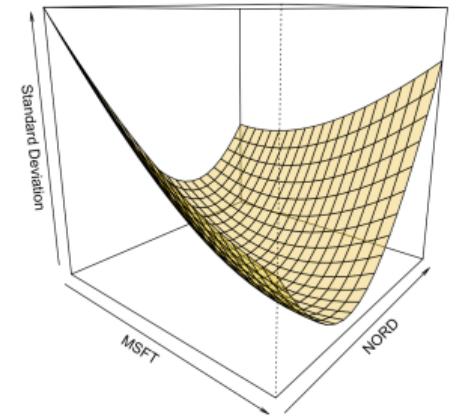
minimum variance portfolio

$$\omega^* = \operatorname{argmin}\{\omega^\top \Sigma \omega\} \text{ s.t. } \omega^\top \mathbf{1} = 1$$

$$\begin{pmatrix} \omega_M \\ \omega_N \\ 1 - \omega_M - \omega_N \end{pmatrix}^\top \Sigma \begin{pmatrix} \omega_M \\ \omega_N \\ 1 - \omega_M - \omega_N \end{pmatrix}$$

```
1 > var_poids = function(z){  
2 +   z.vec = cbind(z[1],z[2],1-z[1]-z[2])  
3 +   sqrt((z.vec) %*% sigma.mat %*% t(z.vec))  
4 + }  
5 > optim(par=c(1,1),var_poids)  
6 $par  
7 [1] 0.675 0.470
```

i.e. $\omega^* = (0.675, 0.470, -0.145)$

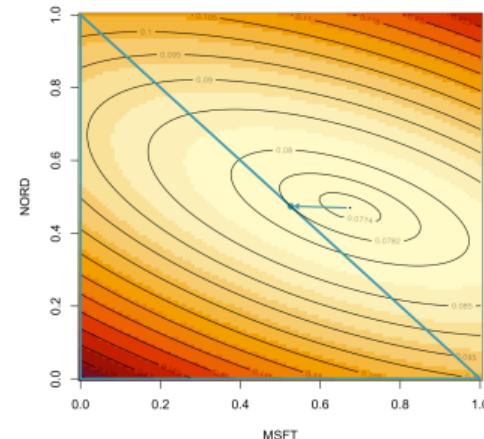


Minimum Variance Portfolio – Markowitz (1952)

with no-short sales,

$$\omega_{NSS}^* = \operatorname{argmin}_{\omega} \{\omega^\top \Sigma \omega\} \text{ s.t. } \begin{cases} \omega^\top \mathbf{1} = 1 \\ \omega \geq \mathbf{0} \\ \omega \leq \mathbf{1} \end{cases}$$

```
1 > grd_var_poids= function(z){  
2 +   h=1e-5  
3 +   c((var_poids(c(z[1]+h,z[2]))-var_poids(c  
+     (z[1]-h,z[2])))/(2*h), (var_poids(c(z  
[1],z[2]+h))-var_poids(c(z[1],z[2]-h))  
/ (2*h)) }
```

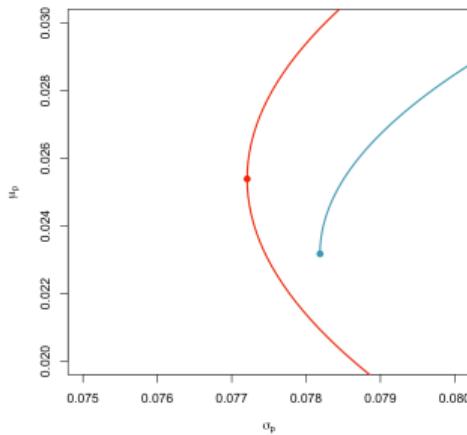


```
1 > constrOptim(c(.2,.2),var_poids,grd_var_poids, ui=matrix(c  
+   (1,0,0,1,-1,0,0,-1,-1,-1),5,2,byrow = TRUE), ci=c(0,0,-1,-1,-1))  
2 $par  
3 [1] 0.526 0.474
```

Portfolio Optimization – Markowitz (1952)

Efficient portfolio: $\omega_{r:NSS}^*$,

$$\omega_{r:NSS}^* = \operatorname{argmin}_{\omega} \{\omega^\top \Sigma \omega\} \text{ s.t. } \begin{cases} \omega^\top \mathbf{1} = 1 \\ \omega^\top \mu = r \\ \omega \geq \mathbf{0} \\ \omega \leq \mathbf{1} \end{cases}$$



Support Vector Machines

SVMs are a classification method that finds the optimal separating hyperplane:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to:

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad \forall i$$

where \mathbf{w} is the weight vector, and b is the bias.

Using Lagrange multipliers α_i , we define the Lagrangian:

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i [y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1]$$

By solving for \mathbf{w} and b , we get the dual problem:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{w}_i^\top \mathbf{w}_j$$

Support Vector Machines

subject to:

$$\sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0$$

Support vectors are the data points for which $\alpha_i > 0$. They determine the decision boundary:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

For non-linearly separable data, we use kernel functions $K(\mathbf{x}_i, \mathbf{x}_j)$ to transform the input space:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

Common kernels:

- Polynomial: $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + c)^d$
- Gaussian (RBF): $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$

Support Vector Machines

- SVMs rely on constrained optimization to find an optimal hyperplane.
- The Lagrangian formulation leads to the dual problem, reducing complexity.
- The kernel trick extends SVMs to non-linear classification problems.

Linear Algebra and Matrices

- A vector is an ordered list of numbers (or entries). It is represented as a column or as a row

$$\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \text{ or } \mathbf{v} = [v_1 \quad v_2 \quad \cdots \quad v_n]$$

- Vectors are used to represent features in machine learning models (e.g., data points in multi-dimensional space).
- Common operations on vectors:
 - **Addition:** Adding two vectors element-wise.
 - **Dot product:** Computes a scalar value. $\mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^n v_i w_i$.
 - **Scalar multiplication:** Multiplying a vector by a scalar.

Linear Algebra and Matrices

- In machine learning, vectors represent data points, and operations on vectors allow us to compute distances, projections, and transformations.
- A matrix is a two-dimensional array of numbers. It has rows and columns.

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1} & A_{m2} & \cdots & A_{mn} \end{bmatrix}$$

- In machine learning, matrices represent datasets, where each row corresponds to an observation, and each column corresponds to a feature.
- Key operations on matrices:
 - **Matrix multiplication:** \mathbf{AB} , multiplying rows of \mathbf{A} with columns of \mathbf{B} .
 - **Transpose:** The transpose of matrix \mathbf{A} , denoted \mathbf{A}^\top , swaps its rows and columns.

Linear Algebra and Matrices

- **Inverse:** If \mathbf{A} is square, its inverse \mathbf{A}^{-1} satisfies $\mathbf{AA}^{-1} = \mathbb{I}$, where \mathbb{I} is the identity matrix.
- Matrices are essential for representing data transformations and for operations in machine learning models.
- A **linear transformation** is a function that maps vectors to other vectors while preserving vector addition and scalar multiplication.
- Matrices can represent linear transformations. For example, applying matrix \mathbf{A} to a vector \mathbf{x} transforms the vector \mathbf{x} into a new vector \mathbf{y} :

$$\mathbf{y} = \mathbf{Ax}$$

- In machine learning, linear transformations are used to:
 - Map data points from one space to another (e.g., in neural networks, during the forward pass).
 - Transform features in PCA or other dimensionality reduction techniques.

Linear Algebra and Matrices

- Key properties of linear transformations:
 - Preserves vector addition: $\mathbf{T}(\mathbf{x} + \mathbf{y}) = \mathbf{T}(\mathbf{x}) + \mathbf{T}(\mathbf{y})$.
 - Preserves scalar multiplication: $\mathbf{T}(c\mathbf{x}) = c \cdot \mathbf{T}(\mathbf{x})$.
- An **eigenvector** \mathbf{v} of a matrix \mathbf{A} is a non-zero vector that satisfies the equation:

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$$

where λ is the corresponding **eigenvalue**.

- In machine learning, eigenvectors and eigenvalues are used in:
 - Principal Component Analysis (PCA) for dimensionality reduction.
 - Neural networks for understanding weight updates and activations.
- The matrix \mathbf{A} acts by scaling the eigenvector \mathbf{v} by the eigenvalue λ , without changing its direction.
- **Geometric interpretation:** Eigenvectors point in the direction of maximum variance or stretching, and eigenvalues measure the scaling factor.

Linear Algebra and Matrices

- **Matrix Decomposition** refers to breaking down a matrix into simpler components to facilitate computations and analysis.
- Common matrix decompositions include:
 - **Singular Value Decomposition (SVD)**: Decomposes a matrix into three matrices $\mathbf{A} = \mathbf{U}\Lambda\mathbf{V}^\top$, used in PCA, recommendation systems, and low-rank approximations.
 - **QR Decomposition**: Decomposes a matrix into an orthogonal matrix \mathbf{Q} and an upper triangular matrix \mathbf{R} , useful in solving linear systems and optimization.
 - **LU Decomposition**: Decomposes a matrix into a lower triangular matrix \mathbf{L} and an upper triangular matrix \mathbf{U} , used in solving systems of linear equations.
- Applications in machine learning:
 - Dimensionality reduction via PCA (using SVD).
 - Solving linear systems in regression analysis.
 - Optimization algorithms for deep learning.

Determinant of square matrices (2×2)

Let $\mathbf{M} = \begin{pmatrix} a & c \\ b & d \end{pmatrix}$, then

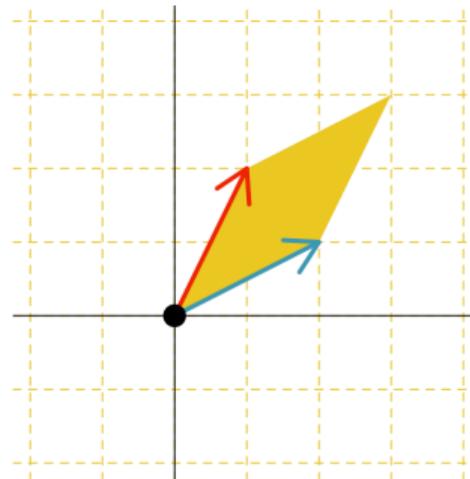
$$\det(\mathbf{M}) = |\mathbf{M}| = \begin{vmatrix} a & c \\ b & d \end{vmatrix} = ad - bc$$

Note that $|\det(\mathbf{M})|$ is the area of the parallelogram, $\vec{u} = \begin{pmatrix} a \\ b \end{pmatrix}$ and $\vec{v} = \begin{pmatrix} c \\ d \end{pmatrix}$

Thus, $\det(\mathbf{M}) = 0$ if and only if $\vec{u} \parallel \vec{v}$

Let θ denote the angle (\vec{u}, \vec{v}) ,

$$\cos(\theta) = \frac{\langle \vec{u}, \vec{v} \rangle}{\|\vec{u}\| \cdot \|\vec{v}\|} \text{ and } \sin(\theta) = \frac{\det(\mathbf{M})}{\|\vec{u}\| \cdot \|\vec{v}\|}, \text{ where } M = (\vec{u} \ \vec{v})$$



Inverse of square matrices (2×2)

If $\det(\mathbf{M}) \neq 0$, then \mathbf{M} has an **inverse**, denoted \mathbf{M}^{-1} ,

$$\mathbf{MM}^{-1} = \mathbf{M}^{-1}\mathbf{M} = \mathbb{I}$$

Example: $\mathbf{M} = \begin{pmatrix} 1 & 2 \\ -2 & 0 \end{pmatrix}$ then $\mathbf{M}^{-1} = \frac{1}{4} \begin{pmatrix} 0 & -1 \\ 2 & 1 \end{pmatrix}$

Important for linear systems: find \mathbf{x} such that $\mathbf{Mx} = \mathbf{a}$,
then the unique solution is $\mathbf{x} = \mathbf{M}^{-1}\mathbf{a}$

Example: $\boldsymbol{\Sigma} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ then $\boldsymbol{\Sigma}^{-1} = \frac{1}{1-\rho^2} \begin{pmatrix} 1 & -\rho \\ -\rho & 1 \end{pmatrix}$, if $|\rho| < 1$.
(here $\boldsymbol{\Sigma}$ is a variance matrix)

A projection matrix cannot be inverted

Eigenvalues and Eigenvectors of square matrices

Let $\mathbf{M} = \begin{pmatrix} a & c \\ b & d \end{pmatrix}$, $\vec{u} = \begin{pmatrix} x \\ y \end{pmatrix}$ and λ such that $\mathbf{M}\vec{u} = \lambda\vec{u}$,

\vec{u} is the **eigenvector** associated with **eigenvalue** λ .

$$\begin{pmatrix} a - \lambda & c \\ b & d - \lambda \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \text{ implies } \det \begin{pmatrix} a - \lambda & c \\ b & d - \lambda \end{pmatrix} = 0.$$

i.e. solve (in λ) $\det(M - \lambda\mathbb{I}) = 0$.

This can be extended in higher dimension

If $\det(\mathbf{M}) = 0$ then $\lambda = 0$ is an eigenvalue.

If \mathbf{M} is symmetric, eigenvectors are orthogonal.

\mathbf{M} is said to be **positive** if all eigenvalues are positive, and then $\mathbf{z}^\top \mathbf{M} \mathbf{z} \geq 0, \forall \mathbf{z}$.

Example: Let \mathbf{X} be a $n \times m$ matrix, then $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{X} \mathbf{X}^\top$ are symmetric and positive matrices (think of $\mathbf{X}^\top \mathbf{X}$ as a covariance matrix).

Quadratic Forms

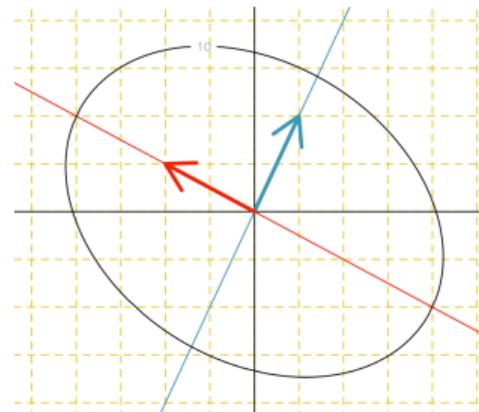
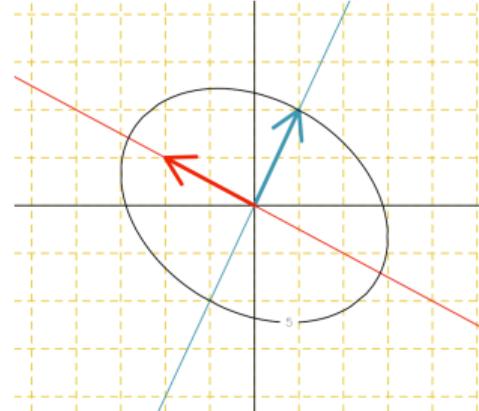
Consider $\mathbf{M} = \begin{pmatrix} a & b \\ b & c \end{pmatrix}$,
and function $\mathbf{z} \mapsto \mathbf{z}^\top \mathbf{M} \mathbf{z}$, i.e.

$$f: \begin{pmatrix} x \\ y \end{pmatrix} \mapsto \begin{pmatrix} x & y \end{pmatrix} \begin{pmatrix} a & b \\ b & c \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

or $ax^2 + 2bxy + cy^2$ is a quadratic form.

If $\det(\mathbf{M}) > 0$, points $\mathbf{z} = (x, y)$ such that $\mathbf{z}^\top \mathbf{M} \mathbf{z} = \gamma$, for some $\gamma > 0$, are on an ellipse (centered on $\mathbf{0}$)

Let $\lambda_1 \geq \lambda_2 > 0$ denote the eigenvalues of \mathbf{M} and \vec{v}_1 and \vec{v}_2 denote the eigenvectors.



Quadratic Forms

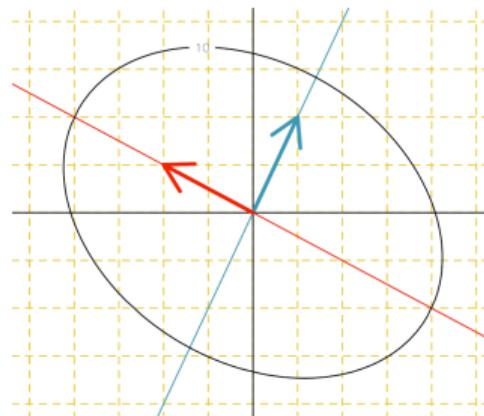
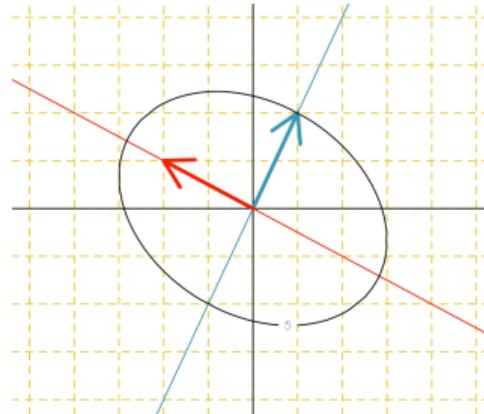
On the picture, $\mathbf{M} = \begin{pmatrix} 0.6 & 0.2 \\ 0.2 & 0.9 \end{pmatrix}$

in Python

```
1 > import numpy as np
2 > from numpy import linalg
3 > M = np.array([[.6,.2], [.2, .9]])
4 > l, v = linalg.eig(M)
5 > print(l)
6 [ 0.5  1. ]
7 > print(v)
8 [[-0.89442719 -0.4472136 ]
9 [ 0.4472136   -0.89442719]]
```

i.e. $\lambda_1 = 1/2$ and $\lambda_2 = 1$, and

$$\vec{v}_1 = \begin{pmatrix} -2\sqrt{5} \\ \sqrt{5} \end{pmatrix} = \sqrt{5} \begin{pmatrix} -2 \\ 1 \end{pmatrix}, \quad \vec{v}_2 = \sqrt{5} \begin{pmatrix} -1 \\ -2 \end{pmatrix}$$

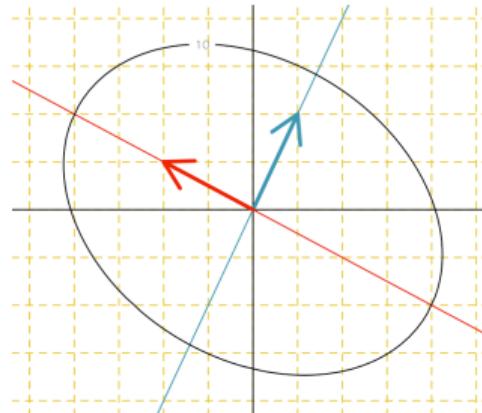
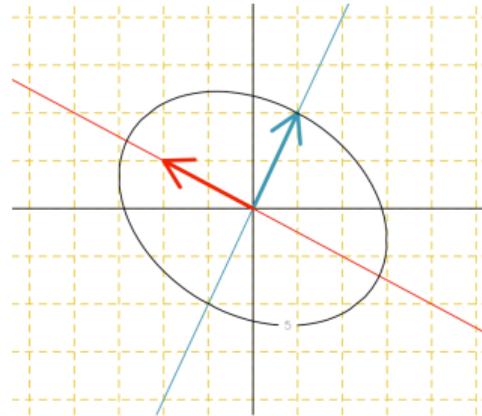


Quadratic Forms

On the picture, $\mathbf{M} = \begin{pmatrix} 0.6 & 0.2 \\ 0.2 & 0.9 \end{pmatrix}$
in \mathbb{R}^2

```
1 > M=matrix(c(.6,.2,.2,.9),2,2)
2 > eigen(M)
3 eigen() decomposition
4 $values
5 [1] 1.0 0.5
6 $vectors
7      [,1]      [,2]
8 [1,] 0.4472136 -0.8944272
9 [2,] 0.8944272  0.4472136
```

Note that $\|\vec{v}_1\| = \|\vec{v}_2\| = 1$ and $\vec{v}_1 \perp \vec{v}_2$



Spectral & Singular Value Decomposition

Let \mathbf{M} denote a symmetric $d \times d$, with eigenvalues $\lambda_1, \dots, \lambda_d$ and eigenvectors $\vec{\mathbf{u}}_1, \dots, \vec{\mathbf{u}}_d$ with $\|\vec{\mathbf{u}}_j\| = 1$,

$$\mathbf{M} = \begin{pmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_d \\ | & | & & | \end{pmatrix}_{d \times d} \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \lambda_d \end{pmatrix}_{d \times d} \begin{pmatrix} -\mathbf{u}_1^\top \\ -\mathbf{u}_2^\top \\ \vdots \\ -\mathbf{u}_d^\top \end{pmatrix}_{d \times d}$$

i.e. $\mathbf{M} = \mathbf{U}\Lambda\mathbf{U}^\top$. Let \mathbf{M} denote $p \times q$, with $p \leq q$, $\mathbf{M} = \mathbf{U}\Lambda\mathbf{V}^\top$, with $\mathbf{u}_1, \dots, \mathbf{u}_p \in \mathbb{R}^p$ and $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^q$

$$\mathbf{M} = \begin{pmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_p \\ | & | & & | \end{pmatrix}_{p \times p} \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \lambda_p \end{pmatrix}_{p \times p} \begin{pmatrix} -\mathbf{v}_1^\top \\ -\mathbf{v}_2^\top \\ \vdots \\ -\mathbf{v}_p^\top \end{pmatrix}_{p \times q}$$

Trace of square matrices

Definition 1.52: Trace of square matrix

The trace of a square matrix \mathbf{A} , denoted $\text{trace}(\mathbf{A})$ is the sum of the elements on its main diagonal,

$$\text{trace}(\mathbf{A}) = A_{11} + A_{22} + \cdots + A_{nn}.$$

The trace of a matrix is the sum of its eigenvalues.

The trace of a projection matrix is the dimension of the target space.

$$\mathbf{P}_\mathbf{X} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \implies \text{trace}(\mathbf{P}_\mathbf{X}) = \text{rank}(\mathbf{X})$$

Definition 1.53: Frobenius inner product and norm

$$\langle \mathbf{A}, \mathbf{B} \rangle = \text{trace}(\mathbf{A}^\top \mathbf{B}) \text{ the Frobenius inner product of } \mathbf{A} \text{ and } \mathbf{B}$$

Approximations

Consider $\mathbf{M} = \mathbf{U}\Lambda\mathbf{V}^\top$ et $\tilde{\mathbf{M}} = \mathbf{U}\tilde{\Lambda}\mathbf{V}^\top$ where $\tilde{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \mathbf{0})$.

$$\begin{pmatrix} 91 & 107 & 136 & 4 \\ 91 & 82 & 86 & 81 \\ 31 & 33 & -52 & 300 \\ 105 & 86 & 97 & 28 \\ 84 & 80 & 78 & 76 \\ 96 & 97 & 124 & 15 \end{pmatrix} \simeq \begin{pmatrix} 103.5 & 101.4 & 130.5 & 2.2 \\ 87.9 & 85.7 & 85.6 & 80.9 \\ 33.5 & 31.7 & -53.0 & 299.7 \\ 91.1 & 89.1 & 105.6 & 30.8 \\ 82.0 & 80.0 & 79.6 & 76.5 \\ 99.3 & 97.2 & 121.3 & 14.1 \end{pmatrix}$$

The old dataset \mathbf{M} is $n = 6$ observations, in dimension 4

The new dataset $\tilde{\mathbf{M}}$ is $n = 6$ observations, in (real) dimension 2

Réduction de dimension, cas linéaire

- If \mathbf{M} is a matrix $n \times n$, \vec{v} is a proper vector for \mathbf{M} (associated to λ) if $\mathbf{M}\vec{v} = \lambda\vec{v}$
- If \mathbf{M} is diagonalize, $\mathbf{M} = \mathbf{V}\Lambda\mathbf{V}^\top$ where $\Lambda = \text{diag}(\lambda)$, et $\mathbf{V} = [\vec{v}_1, \dots, \vec{v}_n]$.
- If \mathbf{M} is invertible, $\mathbf{M}^{-1} = \mathbf{V}\Lambda^{-1}\mathbf{V}^\top$ where $\Lambda = \text{diag}(\lambda)$, et $\mathbf{V} = [\vec{v}_1, \dots, \vec{v}_n]$.
- **Eckart-Young-Mirsky** theorem (Eckart and Young (1936), Mirsky (1960)):

$$\mathbf{M}^* \in \underset{\mathbf{H}:n \times n}{\operatorname{argmin}} \{ \|\mathbf{M} - \mathbf{H}\|_F \text{ such that } \text{rank}(\mathbf{H}) \leq r \}$$

$$\text{If } \mathbf{M} = \mathbf{V}\Lambda\mathbf{V}^\top = [\mathbf{V}_r \ \ \mathbf{V}_{n-r}] \begin{bmatrix} \Lambda_r & \mathbf{0} \\ \mathbf{0} & \Lambda_{n-r} \end{bmatrix} [\mathbf{V}_r \ \ \mathbf{V}_{n-r}]^\top, \mathbf{M}^* = \mathbf{V}_r \Lambda_r \mathbf{V}_r^\top$$

(which will be unique if $\lambda_r > \lambda_{r+1}$)

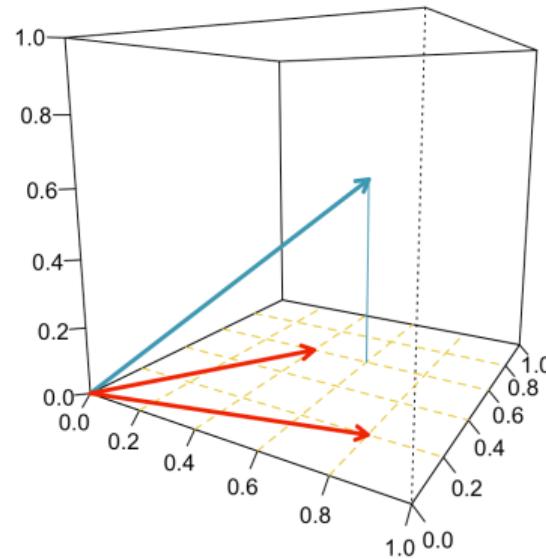
- The approximation error verifies $\|\mathbf{M} - \mathbf{M}^*\|_F = \sqrt{\lambda_{r+1}^2 + \dots + \lambda_n^2}$.

Projection

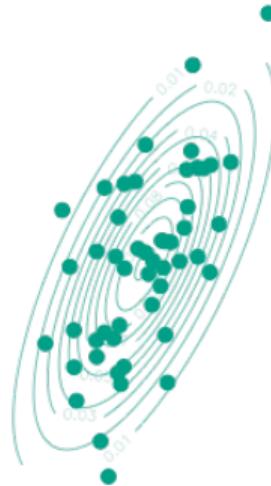
Consider the projection (in \mathbb{R}^3) on $\{\vec{x}_1, \vec{x}_2\}$. Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2]$.

$\mathbf{P} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is the (orthogonal) projection on $\{\vec{x}_1, \vec{x}_2\}$

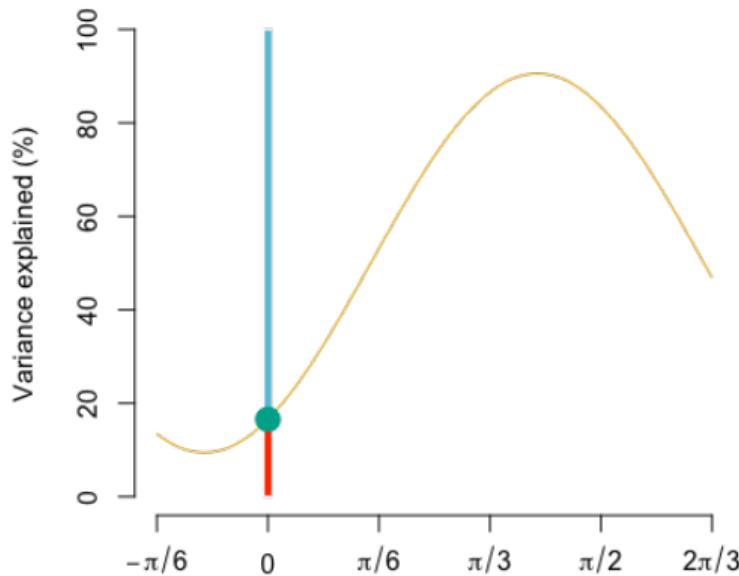
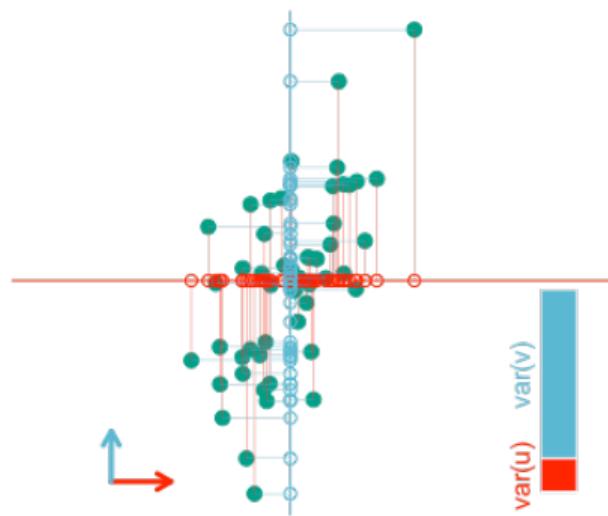
```
1 > X=cbind(c(.8,.2,0),c(.4,.6,0))
2 > P=X %*% solve(t(X)%*%X) %*% t(X)
3 > P
4      [,1] [,2] [,3]
5 [1,]    1    0    0
6 [2,]    0    1    0
7 [3,]    0    0    0
8 > P %*% c(.6,.6,0.6)
9      [,1]
10 [1,]   0.6
11 [2,]   0.6
12 [3,]   0.0
```



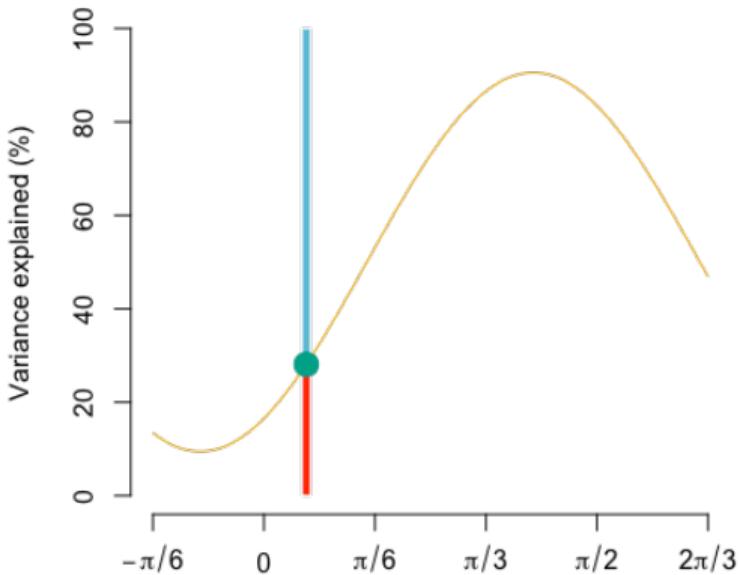
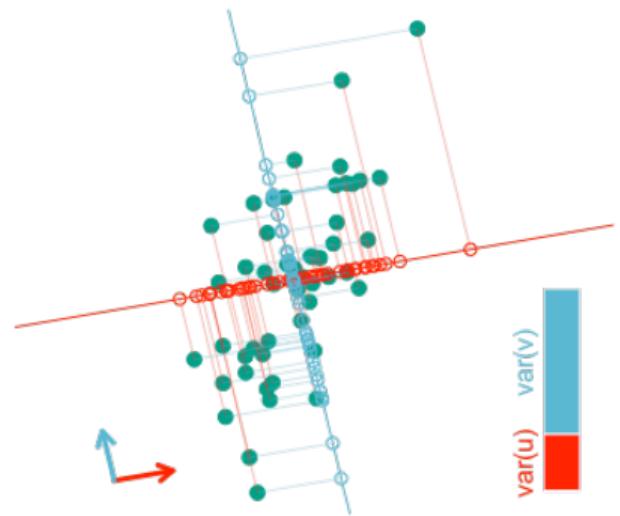
Projection and Variance



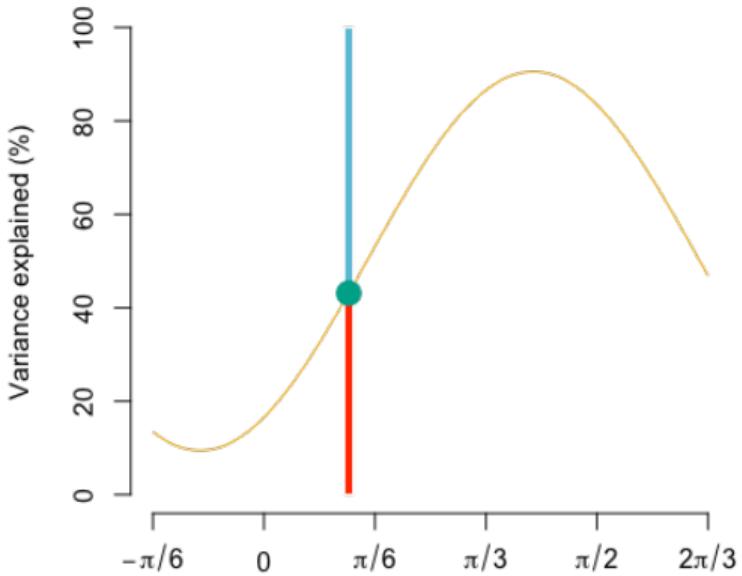
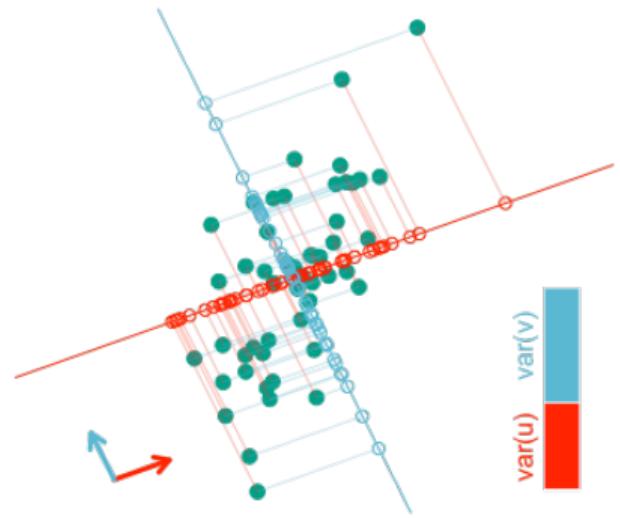
Projection and Variance



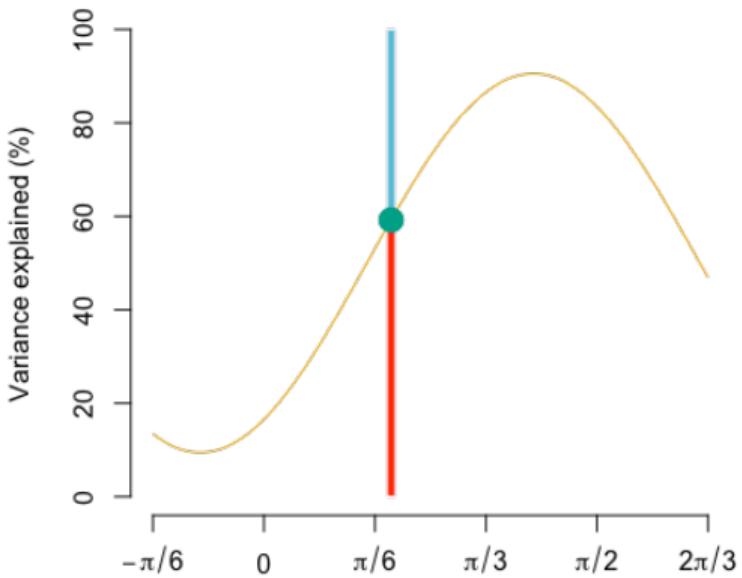
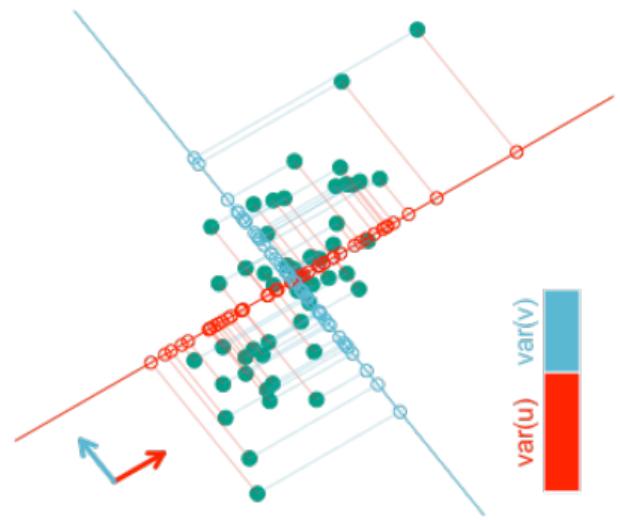
Projection and Variance



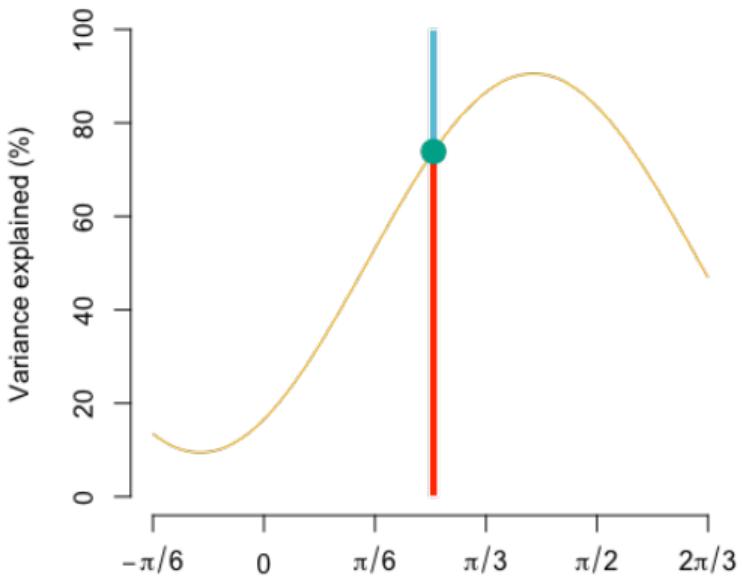
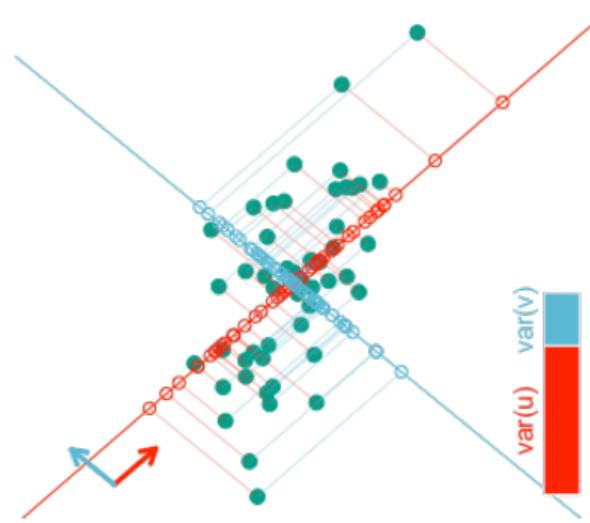
Projection and Variance



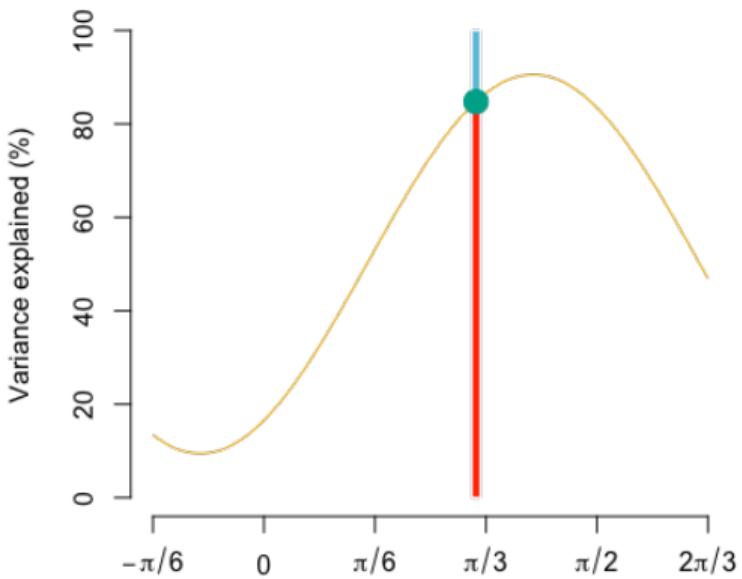
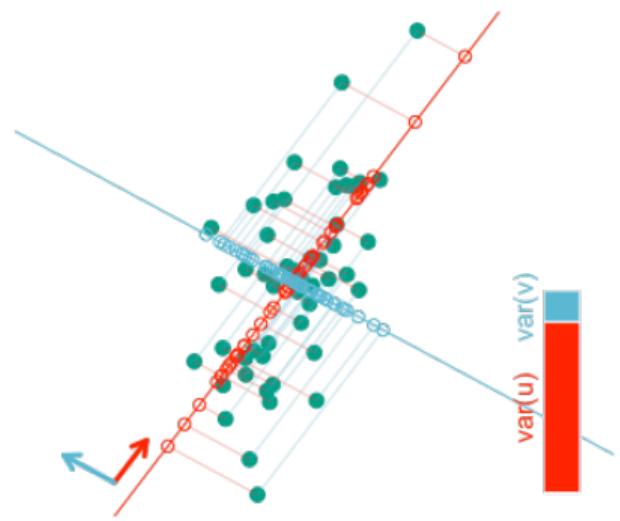
Projection and Variance



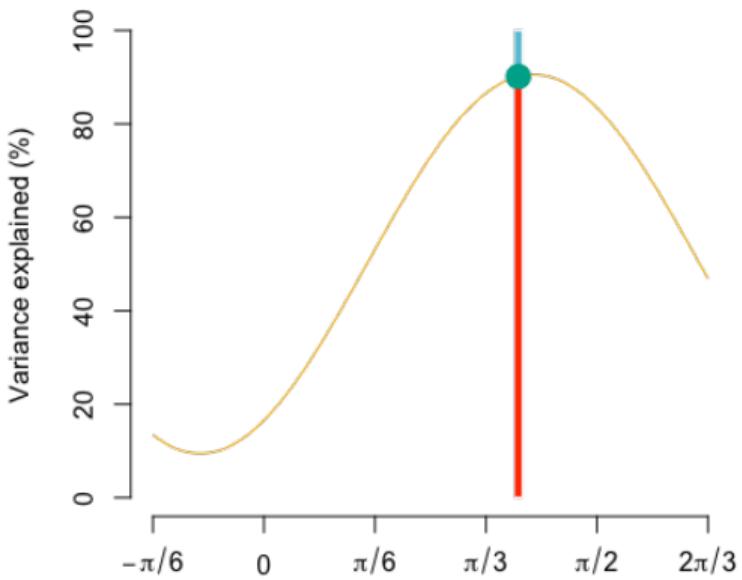
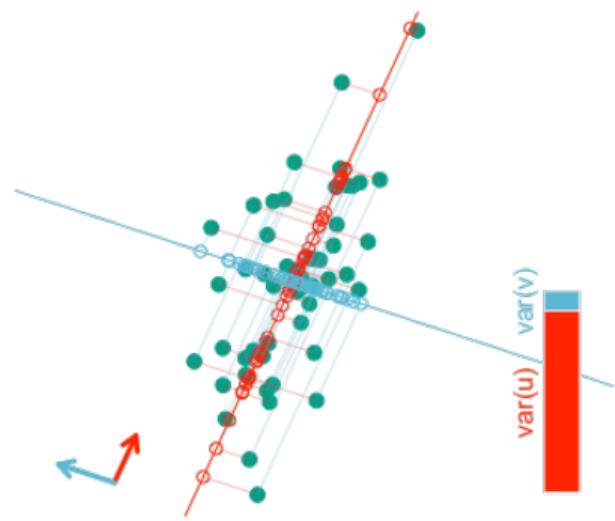
Projection and Variance



Projection and Variance



Projection and Variance



Proposition 1.38: Eigenvectors, optimization and projection

iven a symmetric $n \times n$ matrix \mathbf{M} , the eigenvector \vec{v}_1 corresponding to the largest eigenvalue λ_1 can be found as the solution to the following optimization problem

$$\max_{\vec{v}: \|\vec{v}\|=1} \{\mathbf{v}^\top \mathbf{M} \mathbf{v}\}.$$

Rotation & Orthogonal Matrices

Orthogonal matrix

In linear algebra, an orthogonal matrix, or orthonormal matrix, is a real square matrix whose columns and rows are orthonormal vectors. One way to express this is $\mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T = \mathbb{I}$. W

One can use matrices to transform vectors, e.g. $\vec{\mathbf{y}} = \mathbf{A} \vec{\mathbf{x}}$, with $\vec{\mathbf{x}}, \vec{\mathbf{y}} \in \mathbb{R}^n$, and \mathbf{A} is some $n \times n$ matrix.

Example:

$$\mathbf{A} \vec{\mathbf{x}} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x \cos \theta - y \sin \theta \\ x \sin \theta + y \cos \theta \end{bmatrix}.$$

If $\mathbf{A} = R_0(\theta)$, $\mathbf{A}^T = R_0(-\theta) = \mathbf{A}^{-1}$

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities)

BY-NC 4.0 329 / 335

Rotation & Orthogonal Matrices

Example: $\mathbf{A} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} = R_0(\theta)$ and $\mathbf{B} = \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{bmatrix} = R_0(\phi)$ then

$$\mathbf{AB} = \begin{bmatrix} \cos \theta \cos \phi - \sin \theta \sin \phi & -\cos \theta \sin \phi - \sin \theta \cos \phi \\ \sin \theta \cos \phi + \cos \theta \sin \phi & -\sin \theta \sin \phi + \cos \theta \cos \phi \end{bmatrix} \text{ i.e.}$$

$$\mathbf{AB} = \begin{bmatrix} \cos(\theta + \phi) & -\sin(\theta + \phi) \\ \sin(\theta + \phi) & \cos(\theta + \phi) \end{bmatrix} = R_0(\theta + \phi)$$

In higher dimension n , a $n \times n$ matrix \mathbf{A} is **orthogonal** if its columns and rows are orthogonal unit vectors, i.e.

$$\mathbf{A}^\top \mathbf{A} = \mathbf{AA}^\top = \mathbb{I}$$

or equivalently, $\mathbf{A}^{-1} = \mathbf{A}^\top$.

In dimension 2,

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} = \text{rotation, and } \begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{bmatrix} = \text{reflection}$$

Other Important Matrices

- A **row-stochastic matrix** is a square matrix where each row sums to 1, and all entries are non-negative.

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1} & A_{m2} & \cdots & A_{mn} \end{bmatrix} \quad \text{with} \quad \sum_{j=1}^n A_{ij} = 1 \quad \forall i$$

- Key properties:
 - Non-negative entries ($A_{ij} \geq 0$).
 - Row sums equal to 1: $\sum_j A_{ij} = 1$ for each row.
- **Application in Machine Learning:**
 - **Markov Chains:** Transition probability matrices are row-stochastic matrices. Each row represents the probability of transitioning from one state to another.

Other Important Matrices

- A **doubly stochastic matrix** is a square matrix where:
 - Each row sums to 1 (row-stochastic).
 - Each column sums to 1 (column-stochastic).

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1} & A_{m2} & \cdots & A_{mn} \end{bmatrix} \quad \text{with} \quad \sum_{j=1}^n A_{ij} = 1 \quad \forall i, \quad \sum_{i=1}^m A_{ij} = 1 \quad \forall j$$

- Key properties:
 - Non-negative entries.
 - Both row sums and column sums equal to 1.
- Application in Machine Learning:
 - **Sinkhorn Algorithm:** Used in optimal transport problems, where doubly stochastic matrices are used to balance marginal distributions.

Other Important Matrices

- **Softmax**: In the softmax function used in classification tasks, the output probabilities can be represented by a doubly stochastic matrix.
- A **permutation matrix** is a square matrix that has exactly one entry of 1 in each row and column, with all other entries being 0.

$$\boldsymbol{P} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

- Key properties:
 - The rows and columns are permutations of the standard basis vectors.
 - $\boldsymbol{P}^\top = \boldsymbol{P}^{-1}$, meaning the inverse of a permutation matrix is its transpose.
- Application in Machine Learning:
 - **Data Shuffling**: Permutation matrices are used in data preprocessing to shuffle datasets randomly while preserving the structure of the data.

Other Important Matrices

- **Attention Mechanisms:** In transformer models, permutation matrices are used to shuffle input sequences or to reorder elements based on attention scores.

- A **projection matrix** projects vectors onto a subspace.

$$\mathbf{P}^2 = \mathbf{P}, \quad \mathbf{P} = \mathbf{P}^\top \quad (\text{for an orthogonal projection matrix})$$

- Key properties:

- \mathbf{P} is idempotent: $\mathbf{P}^2 = \mathbf{P}$.
- For orthogonal projections, \mathbf{P} is symmetric: $\mathbf{P} = \mathbf{P}^\top$.

- Application in Machine Learning:

- **Linear Regression:** The projection matrix in ordinary least squares (OLS) regression projects the target vector onto the subspace spanned by the predictors.
- **PCA:** The matrix that projects data onto principal components is a projection matrix.

Other Important Matrices

- Both row-stochastic and doubly stochastic matrices are used to represent graph-based data:
 - A row-stochastic matrix can represent the transition probabilities in a Markov chain on a graph.
 - A doubly stochastic matrix can be used to represent balanced flows in networks, where both rows and columns represent marginal distributions.
- Application in Machine Learning:
 - [Graph Neural Networks \(GNNs\)](#): Markov chains and transition matrices (row-stochastic) are foundational in GNNs for message passing.
 - [Optimal Transport](#): The Sinkhorn algorithm uses doubly stochastic matrices to solve optimal transport problems, which is applied in generative models and domain adaptation.

References

- Adjemian, S. and Pelgrin, F. (2008). Un regard bayésien sur les modèles dynamiques de la macroéconomie. *Economie prevision*, (2):127–152.
- Anscombe, F. J., Aumann, R. J., et al. (1963). A definition of subjective probability. *Annals of mathematical statistics*, 34(1):199–205.
- Bach, F. R. and Jordan, M. I. (2002). Kernel independent component analysis. *Journal of machine learning research*, 3(Jul):1–48.
- Banerjee, A., Merugu, S., Dhillon, I. S., Ghosh, J., and Lafferty, J. (2005). Clustering with bregman divergences. *Journal of machine learning research*, 6(10).
- Bartlett, P. L., Jordan, M. I., and McAuliffe, J. D. (2006). Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156.
- Bauschke, H. H., Borwein, J. M., et al. (1997). Legendre functions and the method of random bregman projections. *Journal of convex analysis*, 4(1):27–67.
- Bayarri, M. J. and Berger, J. O. (2004). The interplay of bayesian and frequentist analysis. *Statistical Science*, 19(1):58–80.
- Bayes, T. (1763). An essay towards solving a problem in the doctrine of chances. *Philosophical transactions of the Royal Society of London*, (53):370–418.

freakonometrics

freakonometrics.hypotheses.org

References

- Bell, E. T. (1945). *The development of mathematics*. Courier Corporation.
- Bernoulli, J. (1713). *Ars conjectandi: opus posthumum: accedit Tractatus de seriebus infinitis; et Epistola gallice scripta de ludo pilae reticularis*. Impensis Thurnisiorum.
- Black, F. and Litterman, R. (1990). Asset allocation: combining investor views with market equilibrium. *Goldman Sachs Fixed Income Research*, 115.
- Black, F. and Litterman, R. (1992). Global portfolio optimization. *Financial analysts journal*, 48(5):28–43.
- Blackwell, D. and MacQueen, J. B. (1973). Ferguson distributions via pólya urn schemes. *The annals of statistics*, 1(2):353–355.
- Bregman, L. M. (1967). The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics*, 7(3):200–217.
- Brenier, Y. (1991). Polar factorization and monotone rearrangement of vector-valued functions. *Communications on pure and applied mathematics*, 44(4):375–417.
- Brier, G. W. et al. (1950). Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3.
- Brualdi, R. A. (2006). *Combinatorial matrix classes*, volume 13. Cambridge University Press.

References

- Bühlmann, H. (1967). Experience rating and credibility. *ASTIN Bulletin: The Journal of the IAA*, 4(3):199–207.
- Cardano, G. (1564). *Liber de ludo aleae*. Franco Angeli.
- Carlier, G., Chernozhukov, V., and Galichon, A. (2016). Vector quantile regression: an optimal transport approach. *The Annals of Statistics*, 44.
- Cournot, A. A. (1843). *Exposition de la théorie des chances et des probabilités*. Hachette.
- Cramér, H. (1928a). On the composition of elementary errors: First paper: Mathematical deductions. *Scandinavian Actuarial Journal*, 1928(1):13–74.
- Cramér, H. (1928b). On the composition of elementary errors: second paper: statistical applications. *Scandinavian Actuarial Journal*, 1928(1):141–180.
- Dall'Aglio, G. (1956). Sugli estremi dei momenti delle funzioni di ripartizione doppia. *Annali della Scuola Normale Superiore di Pisa-Classe di Scienze*, 10(1-2):35–74.
- Davison, A. C. and Hinkley, D. V. (1997). *Bootstrap methods and their application*. Number 1. Cambridge university press.
- De Finetti, B. (1937). La prévision: ses lois logiques, ses sources subjectives. In *Annales de l'institut Henri Poincaré*, volume 7, pages 1–68.

freakonometrics

freakonometrics.hypotheses.org

References

- Denuit, M., Dhaene, J., Goovaerts, M., and Kaas, R. (2006). *Actuarial theory for dependent risks: measures, orders and models*. John Wiley & Sons.
- Eckart, C. and Young, G. (1936). The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218.
- Fenton, N. and Neil, M. (2018). *Risk assessment and decision analysis with Bayesian networks*. CRC Press.
- Fenton, N., Neil, M., and Berger, D. (2016). Bayes and the law. *Annual Review of Statistics and Its Application*, 3:51.
- Ferguson, T. S. (1973). A bayesian analysis of some nonparametric problems. *The annals of statistics*, pages 209–230.
- Fisher, R. (1921). Studies in crop variation. i. an examination of the yield of dressed grain from broadbalk. *The Journal of Agricultural Science*, 11(2):107–135.
- Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical transactions of the Royal Society of London. Series A*, 222(594-604):309–368.
- Galichon, A. (2016). *Optimal transport methods in economics*. Princeton University Press.

References

- Gebelein, H. (1941). Das statistische problem der korrelation als variations- und eigenwertproblem und sein zusammenhang mit der ausgleichsrechnung. *ZAMM - Journal of Applied Mathematics and Mechanics / Zeitschrift für Angewandte Mathematik und Mechanik*, 21(6):364–379.
- Ghosh, J. K. and Ramamoorthi, R. V. (2003). *Bayesian Nonparametrics*. Springer Verlag.
- Gigerenzer, G. and Hoffrage, U. (1995). How to improve bayesian reasoning without instruction: frequency formats. *Psychological review*, 102(4):684.
- Good, I. J. (1952). Rational decisions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 14(1):107–114.
- Gretton, A., Smola, A., Bousquet, O., Herbrich, R., Belitski, A., Augath, M., Murayama, Y., Pauls, J., Schölkopf, B., and Logothetis, N. (2005). Kernel constrained covariance for dependence measurement. In *International Workshop on Artificial Intelligence and Statistics*, pages 112–119. PMLR.
- Hájek, A. (2002). Interpretations of probability. *Stanford Encyclopedia of Philosophy*.
- Hallin, M., del Barrio, E., Cuesta-Albertos, J., and Matrán, C. (2021). Distribution and quantile functions, ranks and signs in dimension d : A measure transportation approach. *The Annals of Statistics*, 49(2):1139–1165.

References

- Hallin, M. and Konen, D. (2024). Multivariate quantiles: Geometric and measure-transportation-based contours.
- Hannan, E. J. (1961). The general theory of canonical correlation and its relation to functional analysis. *Journal of the Australian Mathematical Society*, 2(2):229–242.
- Hardy, G. H., Littlewood, J. E., and Pólya, G. (1952). *Inequalities*. Cambridge university press.
- Harris, T. and Ross, F. (1955). Fundamentals of a method for evaluating rail net capacities. Technical report.
- He, X. D., Kou, S., and Peng, X. (2022). Risk measures: robustness, elicability, and backtesting. *Annual Review of Statistics and Its Application*, 9:141–166.
- Hellinger, E. (1909). Neue begründung der theorie quadratischer formen von unendlichvielen veränderlichen. *Journal für die reine und angewandte Mathematik*, 1909(136):210–271.
- Higham, N. J. (2008). *Functions of matrices: theory and computation*. SIAM.
- Hirschfeld, H. O. (1935). A connection between correlation and contingency. *Mathematical Proceedings of the Cambridge Philosophical Society*, 31(4):520–524.
- Hunt, I. and Mostyn, J. (2020). Probability reasoning in judicial fact-finding. *The international Journal of evidence & Proof*, 24(1):75–94.

References

- Jeffrey, R. (2004). *Subjective probability: The real thing*. Cambridge University Press.
- Jensen, D. and Mayer, L. (1977). Some variational results and their applications in multiple inference. *The Annals of Statistics*, pages 922–931.
- Jonakait, R. N. (1983). When blood is their argument: probabilities in criminal cases, genetic markers, and, once again, bayes' theorem. *University of Illinois Law Review*, page 369.
- Kahneman, D. and Tversky, A. (1972). Subjective probability: A judgment of representativeness. *Cognitive psychology*, 3(3):430–454.
- Kantorovich, L. V. (1942). On the translocation of masses. In *Doklady Akademii Nauk USSR*, volume 37, pages 199–201.
- Kimeldorf, G., May, J. H., and Sampson, A. R. (1982). Concordant and discordant monotone correlations and their evaluation by nonlinear optimization. *Studies in the Management Sciences*, 19:117–130.
- Kimeldorf, G. and Sampson, A. R. (1978). Monotone dependence. *The Annals of Statistics*, pages 895–903.
- Klugman, S. A. (1991). *Bayesian statistics in actuarial science: with emphasis on credibility*, volume 15. Springer Science & Business Media.

References

- Knott, M. and Smith, C. S. (1984). On the optimal mapping of distributions. *Journal of Optimization Theory and Applications*, 43:39–49.
- Kolmogorov, A. (1933). Grundbegriffe der wahrscheinlichkeitsrechnung.
- Kremer, W. (2014). Do doctors understand test results. *BBC World Service*.
- Lancaster, H. O. (1957). Some properties of the bivariate normal distribution considered in the form of a contingency table. *Biometrika*, 44(1/2):289–292.
- Lancaster, H. O. (1958). The Structure of Bivariate Distributions. *The Annals of Mathematical Statistics*, 29(3):719 – 736.
- Laplace, P. S. (1774). Mémoire sur la probabilité de causes par les événements. *Mémoire de l'académie royale des sciences*.
- Lebesgue, H. (1918). Remarques sur les théories de la mesure et de l'intégration. *Annales scientifiques de l'École Normale Supérieure*, 35:191–250.
- Lin, P.-E. (1987). Measures of association between vectors. *Communications in Statistics-Theory and Methods*, 16(2):321–338.
- Loève, M. (1977). *Probability Theory*. Springer.
- Longley-Cook, L. H. (1962). An introduction to credibility theory. Casualty Actuarial Society.

References

- Martin, T. (2009). La probabilité, un concept pluriel. *Pour la science*, (385):46–50.
- Mirsky, L. (1960). Symmetric gauge functions and unitarily invariant norms. *The quarterly journal of mathematics*, 11(1):50–59.
- Monge, G. (1781). Mémoire sur la théorie des déblais et des remblais. *Histoire de l'Académie Royale des Sciences de Paris*.
- Mourier, E. (1953). Eléments aléatoires dans un espace de banach. In *Annales de l'institut Henri Poincaré*, volume 13, pages 161–244.
- Neyman, J. (1977). Frequentist probability and frequentist statistics. *Synthese*, pages 97–131.
- Orbánz, P. and Teh, Y. W. (2010). Bayesian nonparametric models. *Encyclopedia of machine learning*, 1.
- Pardo, L. (2018). *Statistical inference based on divergence measures*. CRC press.
- Pearson, K. (1895). Note on regression and inheritance in the case of two parents. *Proceedings of the royal society of London*, 58(347-352):240–242.
- Pettigrew, R. (2016). *Accuracy and the Laws of Credence*. Oxford University Press.
- Popper, K. R. (1959). The propensity interpretation of probability. *The British journal for the philosophy of science*, 10(37):25–42.

References

- Prokhorov, Y. V. (1956). Convergence of random processes and limit theorems in probability theory. *Theory of Probability & Its Applications*, 1(2):157–214.
- Reichenbach, H. (1971). *The theory of probability*. University of California Press.
- Rényi, A. (1959). On measures of dependence. *Acta mathematica hungarica*, 10(3-4):441–451.
- Saini, A. (2011). A formula for justice. *The Guardian*, October 2nd.
- Santambrogio, F. (2015). *Optimal transport for applied mathematicians*, volume 87. Springer.
- Sarmanov, O. (1958a). Maximum correlation coefficient (non-symmetrical case). *Doklady Akademii Nauk SSSR*, 121(1):52–55.
- Sarmanov, O. V. (1958b). The maximum correlation coefficient (symmetrical case). *Doklady Akademii Nauk SSSR*, 120(4):715–718.
- Satchell, S. and Scowcroft, A. (2000). A demystification of the black–litterman model: Managing quantitative and traditional portfolio construction. *Journal of Asset Management*, 1(2):138–150.
- Savage, L. J. (1971). Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336):783–801.
- Savage, L. J. (1972). *The foundations of statistics*. Courier Corporation.

freakonometrics

freakonometrics.hypotheses.org

References

- Sinkhorn, R. (1962). On the factor spaces of the complex doubly stochastic matrices. *Notices of the American Mathematical Society*, 9:334–335.
- Székely, G. J. (2003). E-statistics: The energy of statistical samples. *Bowling Green State University, Department of Mathematics and Statistics Technical Report*, 3(05):1–18.
- Vapnik, V. (1991). Principles of risk minimization for learning theory. *Advances in neural information processing systems*, 4.
- Villani, C. (2003). *Topics in optimal transportation*, volume 58. American Mathematical Society.
- von Mises, R. (1928). *Wahrscheinlichkeit Statistik und Wahrheit*. Springer.
- von Mises, R. (1939). *Probability, statistics and truth*. Macmillan.
- Wasserstein, L. N. (1969). Markov processes over denumerable products of spaces, describing large systems of automata. *Problemy Peredachi Informatsii*, 5(3):64–72.
- Whitney, A. W. (1918). Theory of experience rating. In *Proceedings of the Casualty Actuarial society*, volume 1, pages 274–292.
- Wilks, S. S. (1938). The large-sample distribution of the likelihood ratio for testing composite hypotheses. *The annals of mathematical statistics*, 9(1):60–62.

freakonometrics

freakonometrics.hypotheses.org

– Arthur Charpentier, April 2025 (Bermuda Financial Authorities) BY-NC 4.0 335 / 335