

# An introduction to Bayesian (thinking and) modeling

**Arthur Charpentier<sup>1</sup>**

<sup>1</sup> Université du Québec à Montréal

November 2022

# Agenda

Uncertainty, insurance and economics

Probabilities and random variables

Motivation with an historical perspective

Beliefs, subjective probabilities and predictive markets

Bayesianism, statistics and calculus (1)

Bayesianism, statistics and calculus (2)

Bayes and Markov property

Bayesianism and statistical learning

Bayesianism, learning and neuroscience

# Preliminaries

Keynote in 2014 at the Cass Business School (now Bayes Business School)...

## Getting into Bayesian Wizardry... (with the eyes of a muggle actuary)

Arthur Charpentier

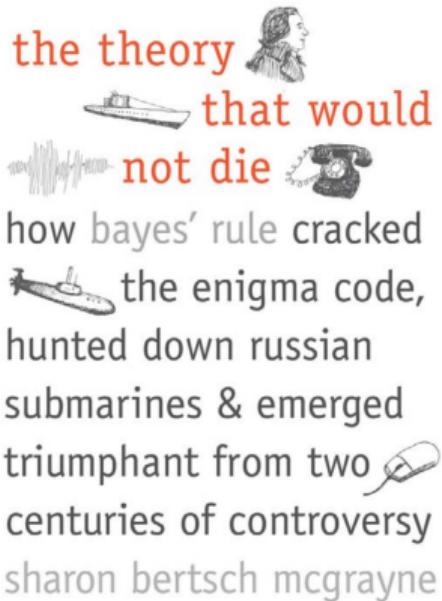
charpentier.arthur@uqam.ca

<http://freakonometrics.hypotheses.org/>

R in Insurance, London, July 2014



# A little bit of history



## contents

Preface and Note to Readers ix  
Acknowledgments xii

### Part I. Enlightenment and the Anti-Bayesian Reaction 1

1. Causes in the Air 3
2. The Man Who Did Everything 13
3. Many Doubts, Few Defenders 34

### Part II. Second World War Era 59

4. Bayes Goes to War 61
5. Dead and Buried Again 87

### Part III. The Glorious Revival 89

6. Arthur Bailey 91
7. From Tool to Theology 97
8. Jerome Cornfield, Lung Cancer, and Heart Attacks 108
9. There's Always a First Time 119
10. 46,658 Varieties 129

### Part IV. To Prove Its Worth 137

11. Business Decisions 139
12. Who Wrote *The Federalist*? 154
13. The Cold Warrior 163
14. Three Mile Island 176
15. The Navy Searches 182

## arthur bailey

### 6.

92 The Glorious Revival

with his four children, and annotating a copy of Grey's *Botany* with the locations of his favorite wild orchids. His motto was, "Some people live in the past, some people live in the future, but the wisest ones live in the present."

Settling into his new job, Bailey was horrified to see "hard-shelled underwriters" using the semi-empirical, "sledge hammer" Bayesian technique developed in 1918 for workers' compensation insurance.<sup>1</sup> University statisticians had long since virtually outlawed those methods, but as practical business people, actuaries refused to discard their prior knowledge and continued to modify their old data with new. Thus they based next year's premiums on this year's rates as refined and modified with new claims information. They did not ask what the new rates should be. Instead, they asked, "How much should the present rates be changed?" A Bayesian estimating how much ice cream someone would eat in the coming year, for example, would combine data about the individual's recent ice cream consumption with other information, such as national dessert trends.

As a modern statistical sophisticate, Bailey was scandalized. His professors, influenced by Ronald Fisher and Jerry Neyman, had taught him that Bayesian priors were "more horrid than 'split,'" in the words of a particularly polite actuary.<sup>2</sup> Statisticians should have no prior opinions about their next experiments or observations and should employ only directly relevant observations while rejecting peripheral, nonstatistical information. No standard methods even existed for evaluating the credibility of prior knowledge (about previous rates, for example) or for correlating it with additional statistical information.

Bailey spent his first year in New York trying to prove to himself that "all of the fancy actuarial [Bayesian] procedures of the casualty business were mathematically unsound."<sup>3</sup> After a year of intense mental struggle, however, he realized to his consternation that actuarial sledgehammering worked. He even preferred it to the elegance of frequentism. He passionately liked formulas that described "actual data. . . . I realized that the hard-shelled underwriters were recognizing certain facts of life neglected by the statistical theorists."<sup>4</sup> He wanted to give more weight to a large volume of data than to the frequentists' small sample, doing so fish surprisingly "logical and reasonable."<sup>5</sup> He concluded that only a "suicidal" actuary would use Fisher's method of maximum likelihood, which assigned a zero probability to nonevents.<sup>6</sup> Since many businesses file no insurance claims at all, Fisher's method would produce premiums too low to cover future losses.

Abandoning his initial suspicions of Bayes' rule, Bailey spent the Second

McGrayne (2011), that mentioned Bailey (1950) (but not Whitney (1918))

# Uncertainty, insurance and economics I



## Uncertainty, insurance and economics II



for the policyholder,  $\pi \preceq X$  (reservation price  $\geq \pi$ )

for the insurer,  $X + \sum_{i=1}^n X_i \leq \pi + \sum_{i=1}^n \pi_i$

# Uncertainty, insurance and economics III



for the policyholder,  $\pi \preceq X$  (reservation price  $\geq \pi$ )

formally,  $\preceq$  is characterized by some utility function  $u$  and beliefs  $\mathbb{Q}_p$

for the insurer,  $X + \sum_{i=1}^n X_i \leq \pi + \sum_{i=1}^n \pi_i$

formally, that inequality holds on average, or on probability

based on some beliefs  $\mathbb{Q}_i$ , e.g.  $\mathbb{Q}_i \left( X + \sum_{i=1}^n X_i \leq \pi + \sum_{i=1}^n \pi_i \right) = 90\%$

# Probabilities and random variables I

*“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)*

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

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

e.g.  $\mathbb{P}(\Omega) = 1$ ; for disjoint events, an additiviy property:  $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B)$ ; aa 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 II

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? cf "*Is the probability inherent to the event, or to our judgment?*" 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 III

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*"

Too much 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.

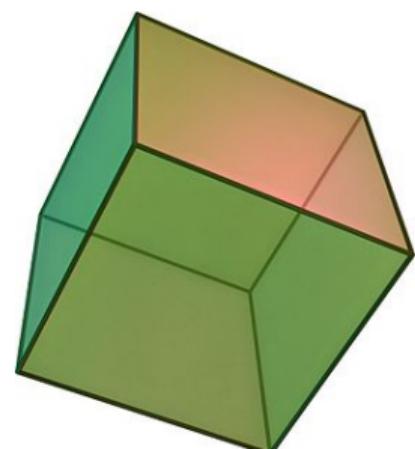
## Probabilities and random variables IV

*"The theory of probabilities is basically only common sense reduced to calculation: it makes appreciate with exactitude, what the just minds feel by a kind of instinct, without them often being able to realize it"*, Laplace (1774)

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 receive 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 V

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 VI

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."*

## Probabilities and random variables VII

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. "*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).

## Probabilities and random variables VIII

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).

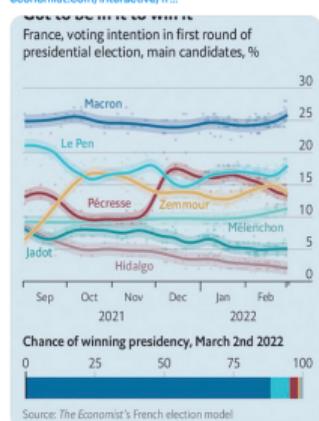
# Probability ? Probability to win an election ?

@PedderSophie (The Economist), vs @HuffPost or @tsrandall (Bloomberg)



Sophie Pedder @PedderSophie · 5 mars

With the usual caveat that one poll is only one poll, this nonetheless fits what @TheEconomist electoral forecast model has been saying for a while. It now gives Macron a 91% chance of winning the # presidency economist.com/interactive/r...



H Huffington Post @huffingtonpost

Folgen

Our @pollsterpolis model gives @HillaryClinton a 98.1% chance of winning the presidency elections.huffingtonpost.com/2016/forecast/ ...

Obersetzung anzeigen



RETWEETS 2.655

FAVORITES 2.120

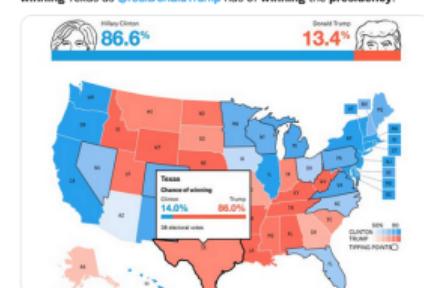
17:26 - 7. Nov. 2016

4h 23 2,7 Tsd. 2,1 Tsd. \*\*\*



Tom Randall @tsrandall

· 16 oct. 2016  
In @FiveThirtyEight's model, @HillaryClinton now has as good a chance of winning Texas as @realDonaldTrump has of winning the presidency.



Adam Singer @AdamSinger

En réponse à @BagholderQuotes

no % margin of error eh?

Traduire le Tweet

3:01 PM · 9 nov. 2016 depuis Milan, Lombardie · Twitter for Android

How to interpret this "probability of winning" ?

How to interpret a "confidence interval" on that probability ? (@AdamSinger)



@freakonometrics



freakonometrics



freakonometrics.hypotheses.org

# Probability ? Probability of precipitation ? |

How to interpret the 'P.o.P.' ("Probability of Precipitation") on weather websites ?

	jeu. 14/07	ven. 15/07	sam. 16/07	dim. 17/07	lun. 18/07	mar. 19/07	mer. 20/07	Risque d'averses
Ensoleillé avec passages nuageux	Ensoleillé	Ensoleillé	Ensoleillé avec passages nuageux	Ensoleillé	Ensoleillé avec passages nuageux	Ensoleillé	Risque d'averses	
31°	27°	28°	30°	35°	38°	28°		
T. ressentie	30	26	27	28	32	35	28	
Nuit	16°	15°	15°	18°	22°	21°	19°	
P.D.P.	20 %	0 %	0 %	20 %	0 %	10 %	40 %	
Vents (km/h)	19 N-E	13 N-E	15 N-E	17 N-E	15 E	20 E	19 S-E	
Rafales (km/h)	28	19	22	25	23	30	29	
Envol. (h)	11 h	15 h	14 h	12 h	15 h	12 h	12 h	
Pluie 24 h	-	-	-	-	-	~1 mm	~1 mm	

	jeu. 14/07	ven. 15/07	sam. 16/07	dim. 17/07	lun. 18/07	mar. 19/07	mer. 20/07	Nuageux avec éclairs
Nuageux avec orages dispersés	Généralement ensoleillé	Ciel variable	Possibilité d'orages	Risque d'averses	Risque d'averses	Nuageux avec éclairs		
24°	26°	27°	28°	28°	28°	29°		
T. ressentie	27	29	31	33	35	34	35	
Nuit	15°	16°	19°	20°	20°	22°	21°	
P.D.P.	40 %	10 %	20 %	40 %	40 %	40 %	30 %	
Vents (km/h)	15 N	15 E	19 S-E	20 S-E	6 E	28 S-E	26 S-E	
Rafales (km/h)	23	23	29	30	9	42	39	
Envol. (h)	3 h	13 h	9 h	4 h	4 h	6 h	3 h	
Pluie 24 h	<1 mm	-	-	~1 mm	<1 mm	~5 mm	~5 mm	

	jeu. 14/07	ven. 15/07	sam. 16/07	dim. 17/07	lun. 18/07	mar. 19/07	mer. 20/07	Ensoleillé avec passages nuageux
Pièce	Faible pluie	Ciel variable	Nuageux	Ensoleillé avec passages nuageux	Ensoleillé	Ensoleillé	Ensoleillé	
8°	6°	9°	9°	11°	11°	17°	17°	
T. ressentie	8	6	9	9	11	17	17	
Nuit	3°	1°	4°	4°	6°	10°	7°	
P.D.P.	100 %	90 %	20 %	30 %	10 %	0 %	0 %	
Vents (km/h)	11 N	6 N-E	5 E	5 S-E	3 S	5 S-E	6 E	
Rafales (km/h)	17	8	7	8	4	7	9	
Envol. (h)	1 h	0 h	5 h	0 h	6 h	10 h	10 h	
Pluie 24 h	25 - 35 mm	5-10 mm	-	~15 mm	-	-	-	

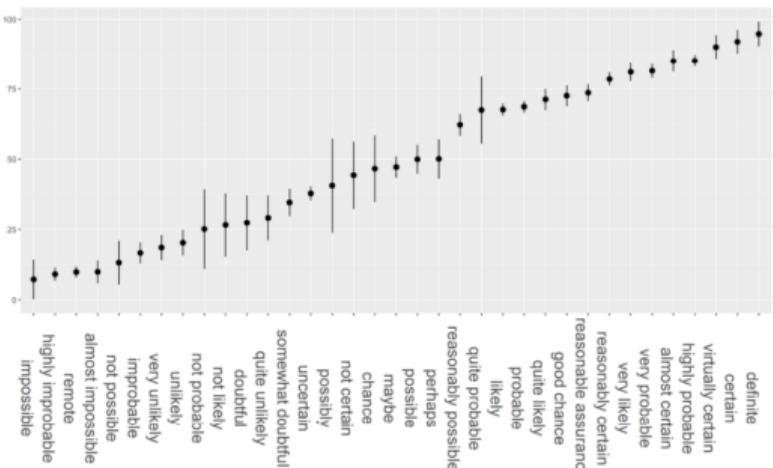
*"Out of all the times you said there was a 40 percent chance of rain, how often did rain actually occur? If, over the long run, it really did rain about 40 percent of the time, that means your forecasts were well calibrated,* Silver (2012)

Murphy and Epstein (1967), Roberts (1968)

Gneiting and Raftery (2005) on ensemble methods for weather forecasting.

# Probability ? Probability of precipitation ? II

More generally, we can think of the "probabilities" mentioned by the IPCC, Mastrandrea et al. (2010) discussed in Stoerk et al. (2020) or Kause et al. (2022)

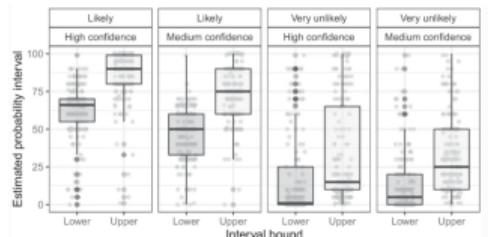


(source Vogel et al. (2022))

High agreement Limited evidence	High agreement Medium evidence	High agreement Robust evidence
Medium agreement Limited evidence	Medium agreement Medium evidence	Medium agreement Robust evidence
Low agreement Limited evidence	Low agreement Medium evidence	Low agreement Robust evidence

↑  
Agreement  
→  
Evidence (type, amount, quality, consistency) →  
Confidence Scale

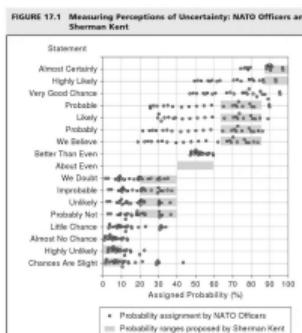
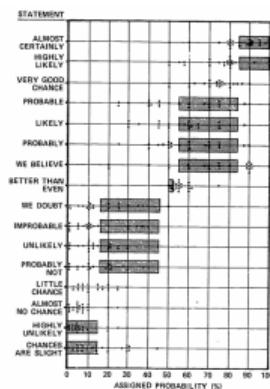
Table 1. Likelihood Scale	
Term*	Likelihood of the Outcome
Virtually certain	99-100% probability
Very likely	90-100% probability
Likely	66-100% probability
About as likely as not	33 to 66% probability
Unlikely	0-33% probability
Very unlikely	0-10% probability
Exceptionally unlikely	0-1% probability



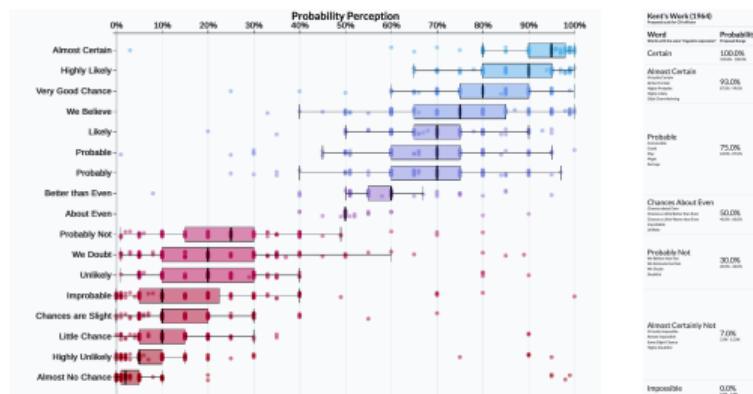
# Probability ? Probability of precipitation ? III

**Note** : “Cromwell’s rule”: one should not give a probability of 1 to an event that cannot logically be shown to be true, and one should never give a probability of 0 to an event unless it can logically be shown to be false,

Lindley (2013), Barclay et al. (1977) et Pherson and Pherson (2012).

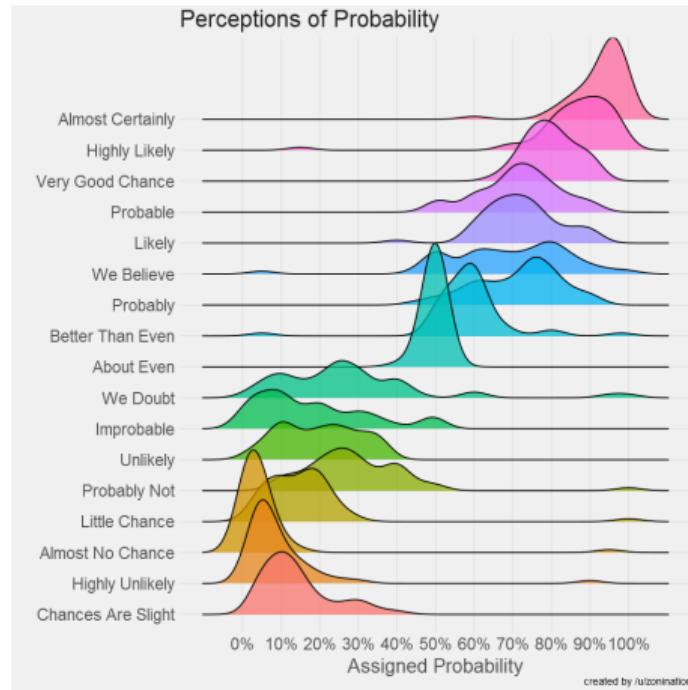
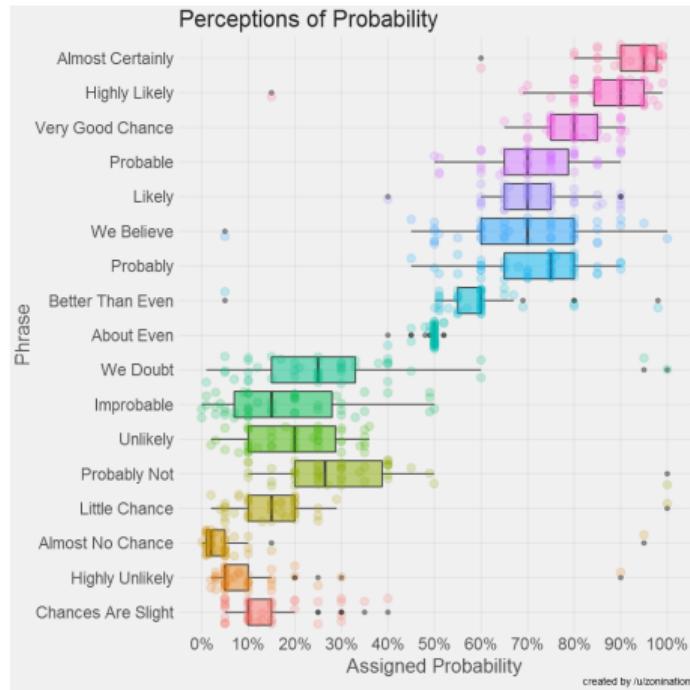


Source: Scott Barclay et al., "Handbook for Sherman Kent," "Ways of Estimated Decision Analysis," Institute, Vienna, Vienna, and Douglas, Inc., 1971; and Donald P. Horan et al., "Probability," in Sherman Kent and the Board of National Estimates, *Collected Essays*, Washington DC: Center for Study of Intelligence, USA, 1996.



# Probability ? Probability of precipitation ? IV

See also [@zonination](#) on "probability perceptions"



# Bayesian statistics ?

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

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

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A) \cdot \mathbb{P}(A)}{\mathbb{P}(B)}$$

“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  $A$  and  $B$  such that  $\mathbb{P}(B) \neq 0$ ,

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A) \cdot \mathbb{P}(A)}{\mathbb{P}(B)}$$

- ▶ 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  $A$  and  $B$  such that  $\mathbb{P}(B) \neq 0$ ,

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A) \cdot \mathbb{P}(A)}{\mathbb{P}(B)}$$

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

## Bayesian statistics ?

A person coughs (event  $B$ ). Which hypothesis is the most credible?  
(from Dehaene (2012))

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

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*)



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

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

# The practice of conditional probabilities

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



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

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

# 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



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

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

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

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

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

# Practice of Bayesian Statistics

*“Do doctors understand test results? ”, 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? ”, 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”.

# Axiomatic of beliefs I

Axioms of Bayesian approach, [Titelbaum \(2022a\)](#), [\(2022b\)](#), are

- ▶ step 1 : [beliefs](#)

Beliefs are quantified on a scale from 0 to 1

The "rationality of beliefs" means that beliefs are measures of probabilities (and verify the associated axioms), [Buehler \(1976\)](#).

**Note:** a weaker version of coherence can be defined using capacities (in the sense of [Choquet \(1954\)](#)), based on the axiom : if  $A \subset B$ , then  $\mathbb{Q}[A] \leq \mathbb{Q}[B]$  (and no longer the additivity of disjoint events)

## Axiomatic of beliefs II

### ► step 2 : updating beliefs

For Popper (1955), an agent who believes  $A$  to the degree  $Q[A]$ , if he learns  $B$ , he then believes  $A$  to the degree  $Q[A|B]$

$$Q[A] \mapsto Q[A|B] \cdot \underbrace{Q[B]}_{=1} + Q[A|\neg B] \cdot \underbrace{Q[\neg B]}_{=0} = Q[A|B] = Q_B[A]$$

Jeffrey (1965) proposed a generalization if  $B$  is associated with a belief  $Q'[B]$ ,

$$Q[A] \mapsto Q'[A] = Q[A|B] \cdot Q'[B] + Q[A|\neg B] \cdot Q'[\neg B]$$

In other words, "reasoning consists of graduating one's beliefs and revising one's degrees of belief by Bayesian conditionalization as new information becomes available", Drouet (2016).

## Axiomatic of beliefs III

*"La differenza essenziale da rilevare è nell'attribuzione del 'perchè': non cerco perchè IL FATTO che io prevedo accadrà, ma perchè IO prevedo che il fatto accadrà. Non sono più i fatti che hanno bisogno di una causa per prodursi : è il nostro pensiero che trova comodo di immaginare dei rapporti di causalità per spiegarli, coordinarli, e renderne possibile la previsione"*, De Finetti (1931)



"I do not seek to know why the fact that I foresee will come true, but why I foresee that the fact will come true. It is no longer the facts that need a cause to happen: it is our mind that finds it convenient to imagine causal relationships in order to explain them, to coordinate them and to make the prediction possible"

# The Dutch book I

Ramsey (1926) and De Finetti (1937) suggested to understand the rationality of beliefs with the help of bets (formalized by Lehman (1955) Kemeny (1955), Teller (1973), Lindley et al. (1979) and Skyrms (1987)) and "arbitrage" (we speak of Subjective Bayesianism).

We assign the belief  $q$  to a bet (lottery) associated to  $A$ , yielding  $a$  if  $A$  occurs and 0 otherwise if and only if the value of the lottery is  $qa$ , Hájek (2009)

The dutch book argument is that if an individual has beliefs that violate the probabilities and if he bets based on those beliefs, then he is willing to accept a set of bets that he is certain to lose, Pettigrew (2020).

**Note:** Lehman (1955) used the term "dutch book", but it corresponds to the notion of "arbitrage" in financial mathematics.

## The Dutch book II

Lehman (1955) “if a set of betting prices violate the probability calculus, then there is a Dutch Book consisting of bets at those prices.”

Kemeny (1955), “if a set of betting prices obey the probability calculus, then there does not exist a Dutch Book consisting of bets at those prices”

This characterization is also called Cox-Jaynes theorem, Cox (1946) taken up by Jaynes (1988) and Jaynes (2003) : probabilities (characterized by Kolmogorov axioms) are the only normative mechanism for plausibility induction

See also Good (1966)

or Eisenberg and Gale (1959) and Baron and Lange (2006), Chen and Pennock (2010) on parimutuel, and predictive markets

Suppose that  $I$  players bet on  $J$  horses. Each player bets  $b_i$ , and normalize  $(b_1 + \dots + b_I = 1)$ .

Player  $i$  bets  $\beta_{i,j}$  on horse  $j$  ( $b_i = \beta_{i,1} + \dots + \beta_{i,J}$ ).

## The Dutch book III

We note  $\pi_j$  the amount bet on the horse  $j$  ( $\pi_j = \beta_{1,j} + \dots + \beta_{I,j}$ ).

Since  $\pi_j \in (0, 1)$  and  $\pi_1 + \dots + \pi_J = 1$  is interpreted as a probability, describing a "collective belief".

We can also add empirical constraints, and associate the beliefs to known frequencies  
(this is called [Empirical Bayesianism](#))

[Williamson \(2004\)](#) introduced an objective Bayesianism, inspired by [Jaynes \(1957\)](#),  
based on entropy maximization (maxmin approach), associated with a precautionary principle.

# Non-boolean logic I

**Note** We can also find links with logic.

Classically, if we have the proposition "[If A is true, then B is true](#)"

- $$\begin{cases} \text{If I observe that } A \text{ is true, I conclude that } B \text{ is true} \\ \text{If I observe that } B \text{ is false, I conclude that } A \text{ is false.} \end{cases}$$

With [boolean logic](#), these are the only equivalent assertions

$$(A \Rightarrow B \text{ and } \neg B \Rightarrow \neg A)$$

But there may be some [plausible reasoning](#), Pólya (1958)

- $$\begin{cases} \text{If I observe that } A \text{ is false, it seems to me that } B \text{ becomes less plausible} \\ \text{If I observe that } B \text{ is true, it seems to me that } A \text{ becomes more plausible.} \end{cases}$$

What means "plausible" here ?

# Bayesianism, statistics and calculus I

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

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

## ► Conjugate distributions: **Binomial - Beta**

The likelihood for binomial (Bernoulli) variables

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

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

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

## Bayesianism, statistics and calculus II

### ► Conjugate distributions : **Poisson - Gamma**

The likelihood for Poisson variables is

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

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 III

## ► Conjugate distributions : **Normal - Normal**

If variance  $\Sigma$  is known

$$\begin{cases} \mathbf{x}_i | \boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma) \\ \boldsymbol{\mu} \sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0) \end{cases} \quad \text{then } \boldsymbol{\mu} | \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_x, \Sigma_x)$$

where 
$$\begin{cases} \boldsymbol{\mu}_x = (\Sigma_0^{-1} + n\Sigma^{-1})^{-1} (\Sigma_0^{-1}\boldsymbol{\mu}_0 + n\Sigma^{-1}\bar{\mathbf{x}}) \\ \Sigma_x = (\Sigma_0^{-1} + n\Sigma^{-1})^{-1} \end{cases}$$

used classically in Bayesian econometrics.

# Bayesianism, statistics and calculus IV

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

If mean  $\mu$  is known

$$\begin{cases} \mathbf{x}_i | \Sigma \sim \mathcal{N}(\mu, \Sigma) \\ \Sigma \sim IW(\nu_0, \Psi_0) \end{cases} \quad \text{then } \Sigma | \mathbf{x} \sim IW(\nu_x, \Psi_x)$$

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

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 V

Bayesian methods can be very powerful for estimating panel, hierarchical, or multilevel models, [Gelman and Hill \(2006\)](#).

## ► Hierarchical model

When the individual  $i$  belongs to the group  $j$ ,

$$y_{i,j} = \alpha_j + \mathbf{x}_i^\top \boldsymbol{\beta}_j + \varepsilon_{i,j}, \text{ where } \begin{cases} \alpha_j = a_0 + \mathbf{z}_j^\top \boldsymbol{\beta}_1 + u_j \\ \boldsymbol{\beta}_j = \mathbf{b}_0 + \mathbf{Z}_j^\top \mathbf{B}_1 + \mathbf{u}_j \end{cases}$$

with constants and slopes depending on the groups.

(usually in a GLM model).

## Bayesianism, statistics and calculus VI

Otherwise, either simulations are used (see MCMC) or simplifying assumptions are made.

Consider symptoms  $s_1, \dots, s_k$  and diseases  $m_1, \dots, m_j$  (in  $\{0, 1\}$ )

$$\mathbb{P}[\mathbf{M} = \mathbf{m} | \mathbf{S} = \mathbf{s}] = \frac{\mathbb{P}[\mathbf{M} = \mathbf{m}] \cdot \mathbb{P}[\mathbf{S} = \mathbf{s} | \mathbf{M} = \mathbf{m}]}{\sum_{\mathbf{x}} \mathbb{P}[\mathbf{M} = \mathbf{x}] \cdot \mathbb{P}[\mathbf{S} = \mathbf{s} | \mathbf{M} = \mathbf{x}]}$$

“Naïve Bayes” relies on assumptions (Spiegelhalter et al. (1993))

- ▶ diseases are mutually exclusive  $\mathbb{P}[\mathbf{M} = \mathbf{m} | \mathbf{S} = \mathbf{s}] = 0$  si  $\mathbf{m}^\top \mathbf{1} > 1$ ,
- ▶ the symptoms are conditionally independent

$$\mathbb{P}[\mathbf{S} = \mathbf{s} | M_i = m_i] = \prod_{j=1}^k \mathbb{P}[S_j = s_j | M_i = m_i]$$

## Bayesianism, statistics and calculus VII

In that case

$$\mathbb{P}[M_i = m_i | \mathbf{S} = \mathbf{s}] = \frac{\mathbb{P}[M_i = m_i] \cdot \prod_{j=1}^k \mathbb{P}[S_j = s_j | M_i = m_i]}{\mathbb{P}[M_i = 0] \cdot \prod_{j=1}^k \mathbb{P}[S_j = s_j | M_i = 0] + \mathbb{P}[M_i = 1] \cdot \prod_{j=1}^k \mathbb{P}[S_j = s_j | M_i = 1]}$$

We can improve the model by using a [Bayesian network](#) (we will talk about it later).

## Bayesianism, statistics and calculus VIII

To determine  $\mathbb{P}[M_i = m_i | \mathbf{S} = \mathbf{s}]$ , we need to know

- ▶ prevalence of disease  $\mathbb{P}[M_i = 1]$
- ▶ sensitivity  $\mathbb{P}[S_j = 1 | M_i = 1]$
- ▶ specificity  $\mathbb{P}[S_j = 0 | M_i = 0]$

for all symptoms  $S_j$  and all disease  $M_i$ .

Note that  $\mathbb{P}[S_j = s_j | M_i = m_i]$  have a causal interpretation: it is the diseases that cause the symptoms.

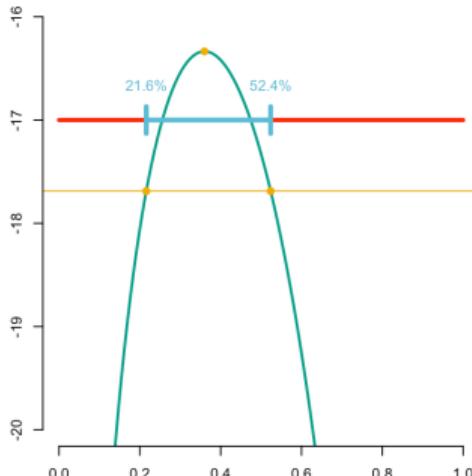
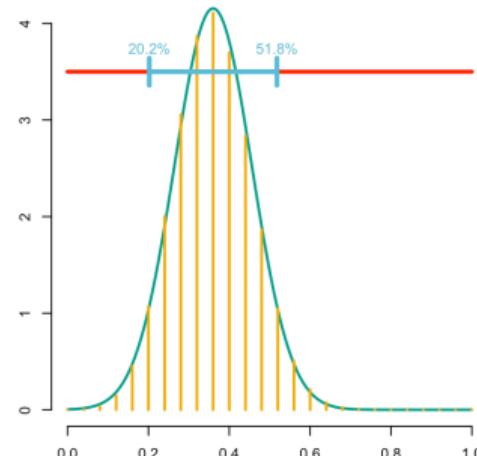
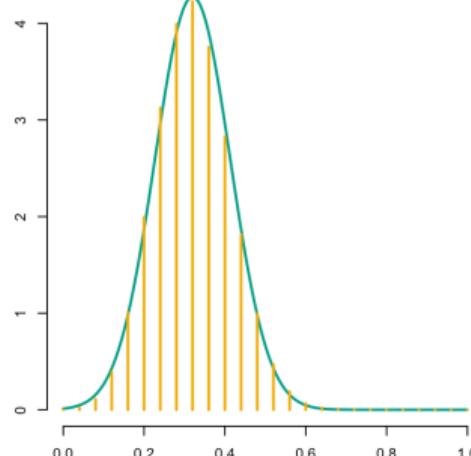
See [Sadegh-Zadeh \(1980\)](#) on Bayesian diagnostics, or [Donnat et al. \(2020\)](#).

# Bayesianism, statistics and calculus I

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

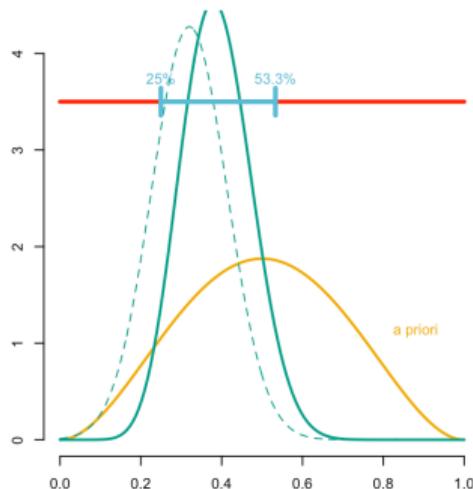
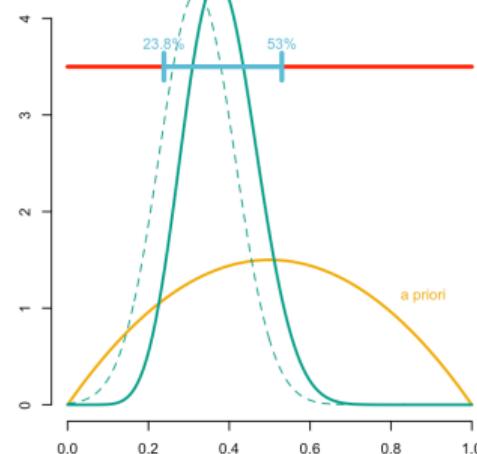
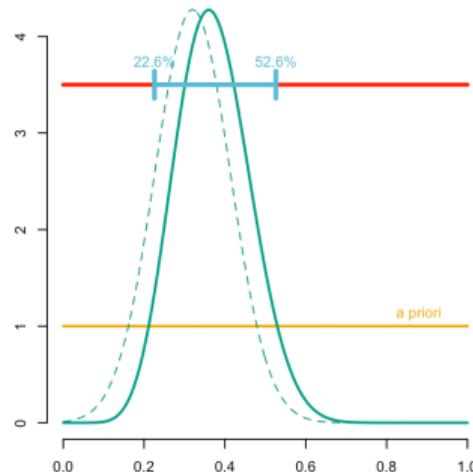


# Bayesianism, statistics and calculus II

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

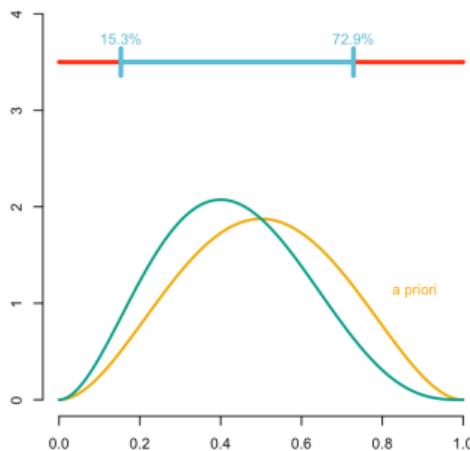
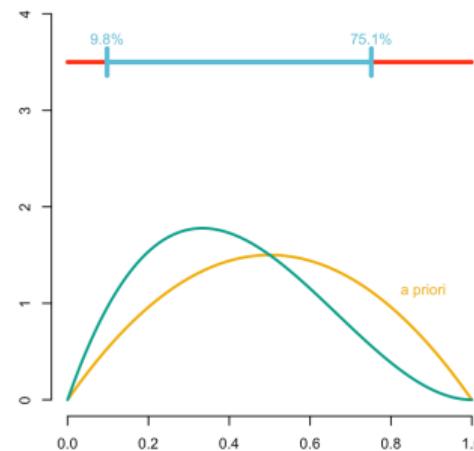
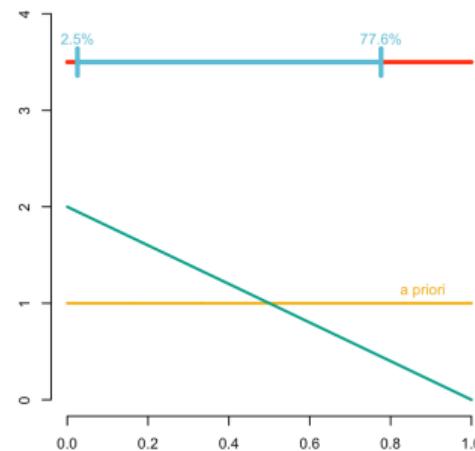


# Bayesianism, statistics and calculus III

## ► Posterior distribution

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

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

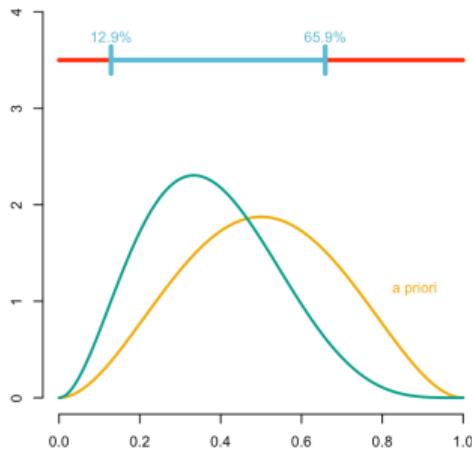
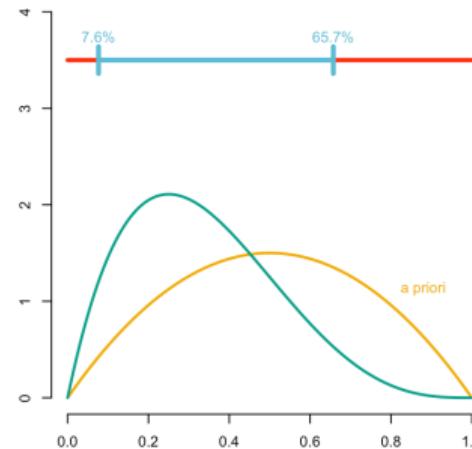
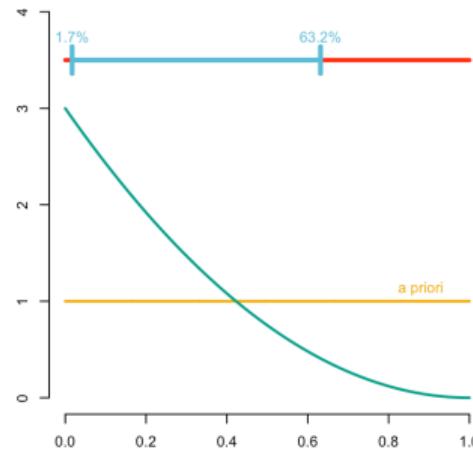


# Bayesianism, statistics and calculus IV

## ► Posterior distribution

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

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

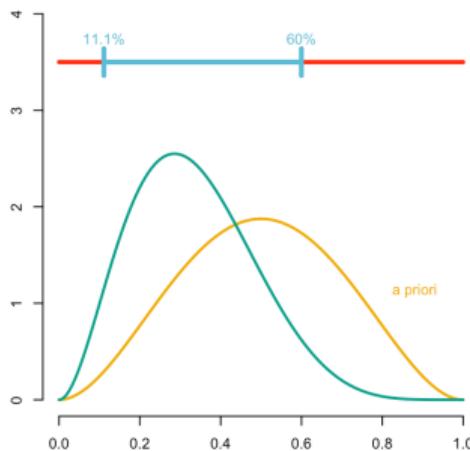
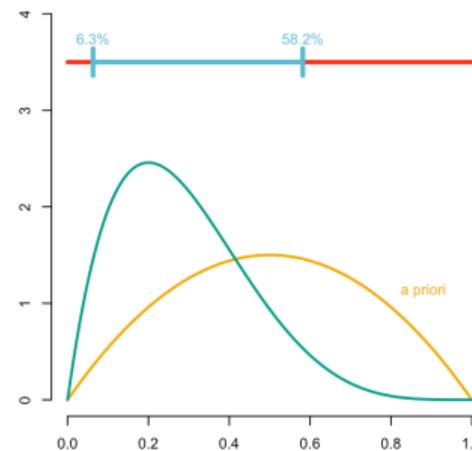
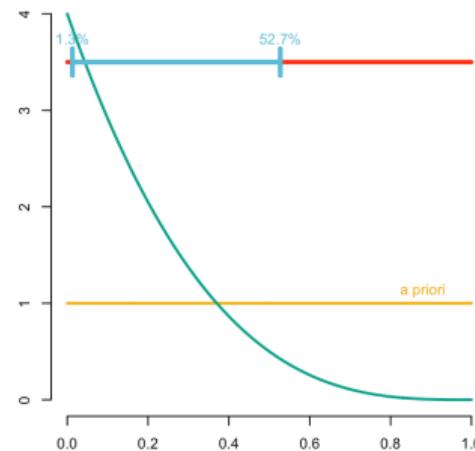


# Bayesianism, statistics and calculus V

## ► Posterior distribution

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

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

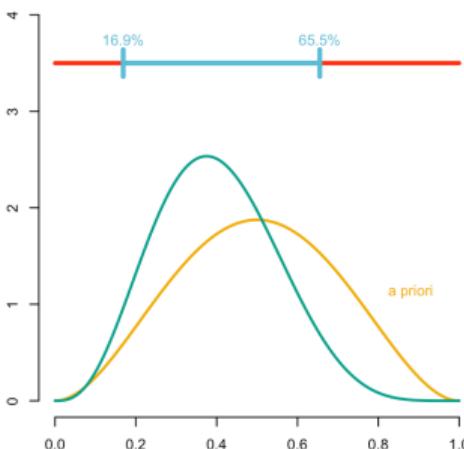
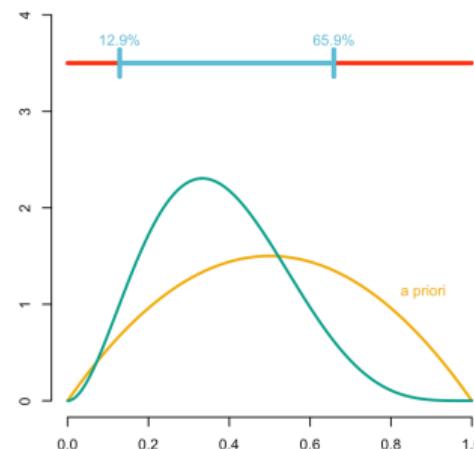
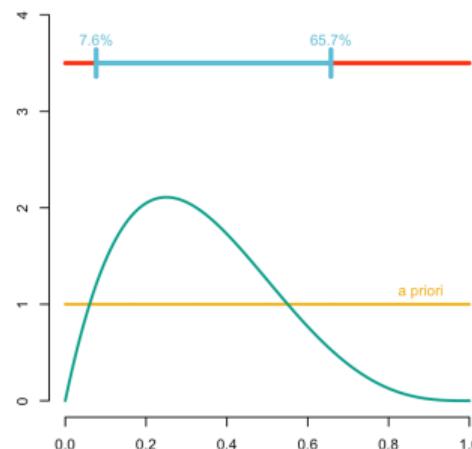


# Bayesianism, statistics and calculus VI

## ► Posterior distribution

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

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

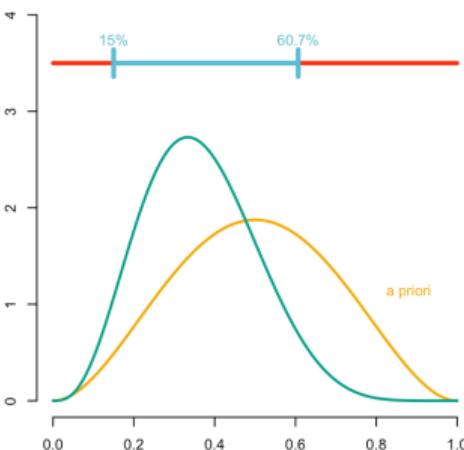
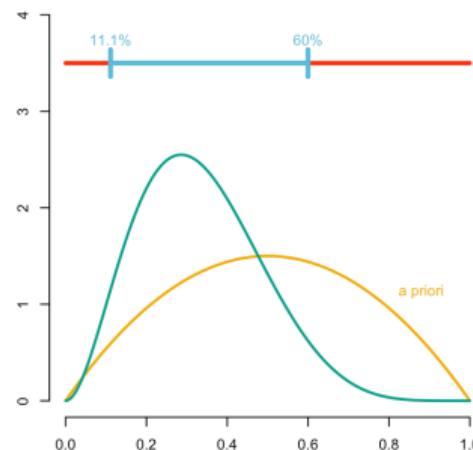
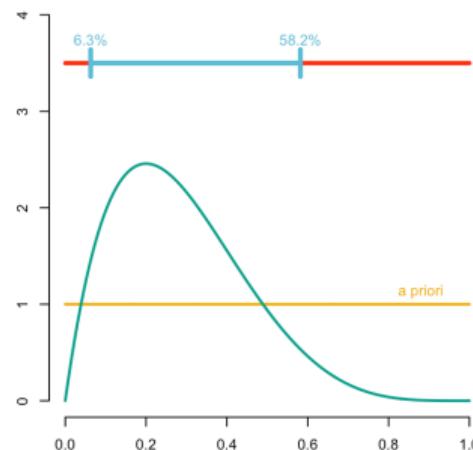


# Bayesianism, statistics and calculus VII

## ► Posterior distribution

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

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

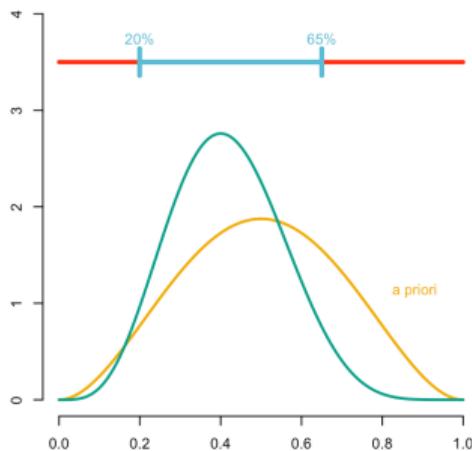
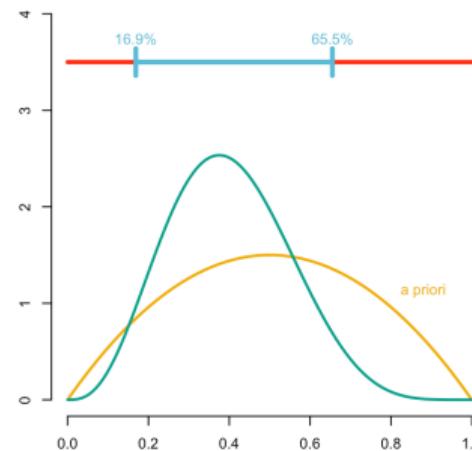
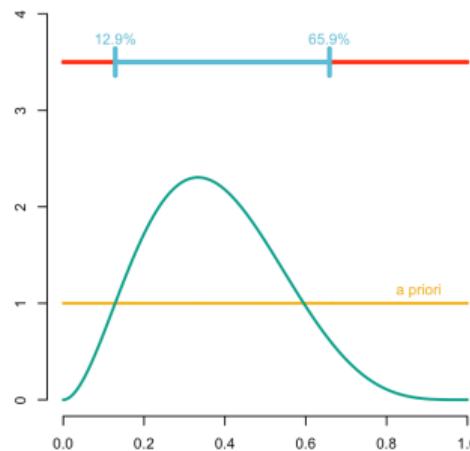


# Bayesianism, statistics and calculus VIII

## ► Posterior distribution

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

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

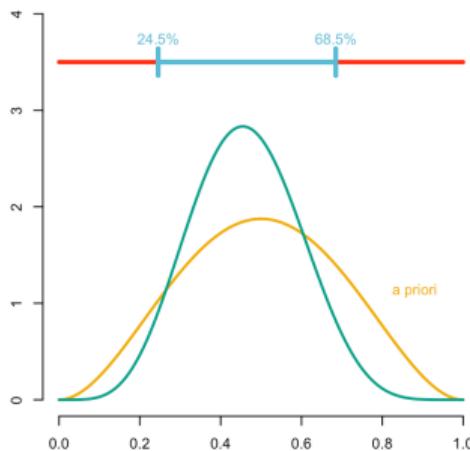
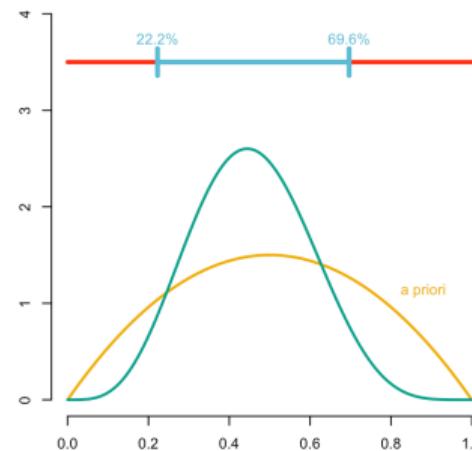
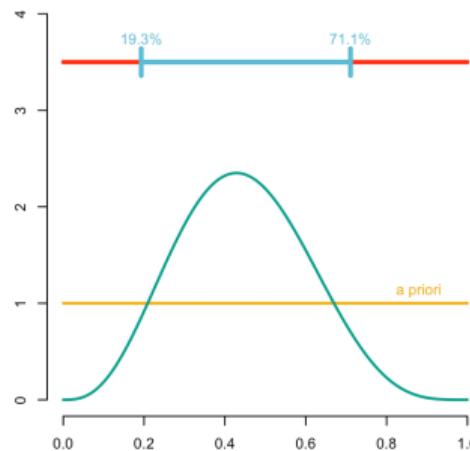


# Bayesianism, statistics and calculus IX

## ► Posterior distribution

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

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

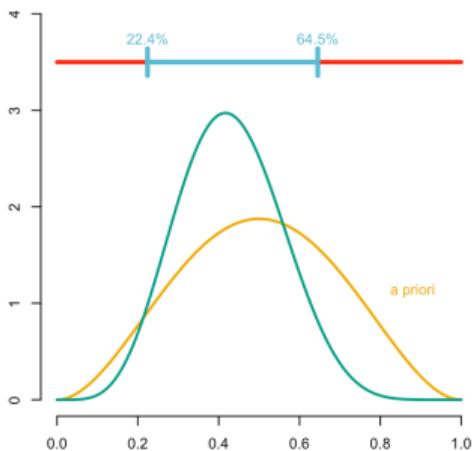
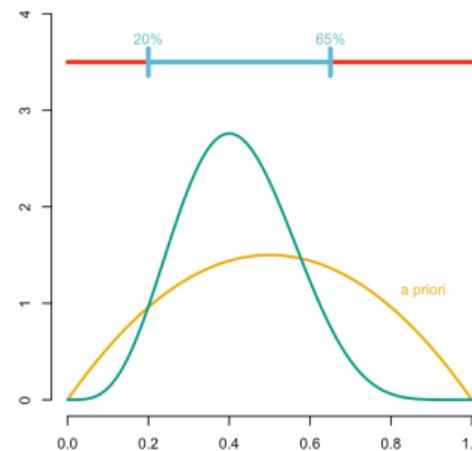
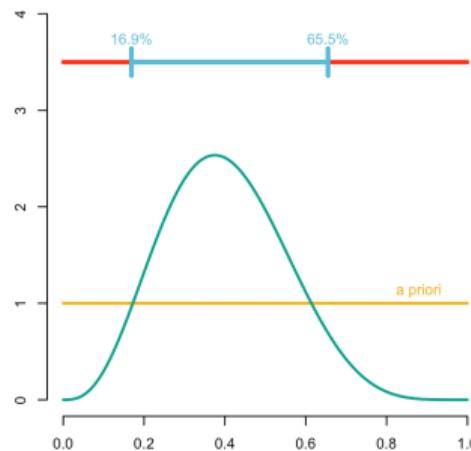


# Bayesianism, statistics and calculus X

## ► Posterior distribution

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

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

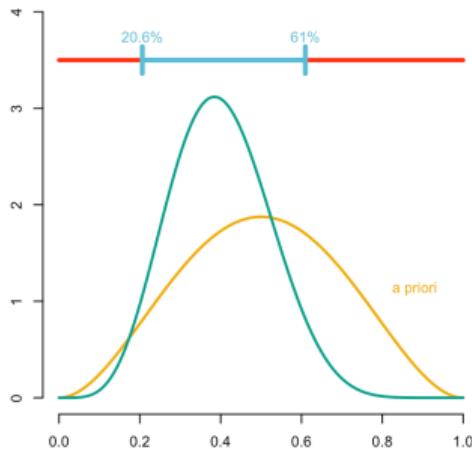
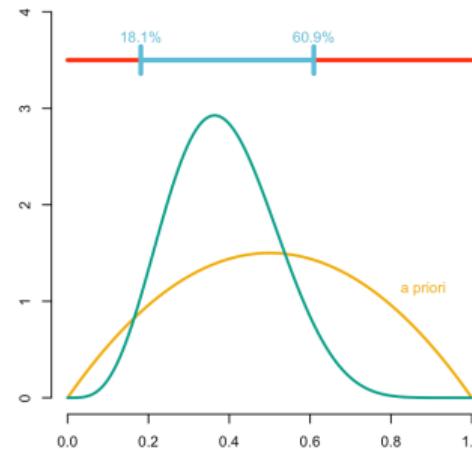
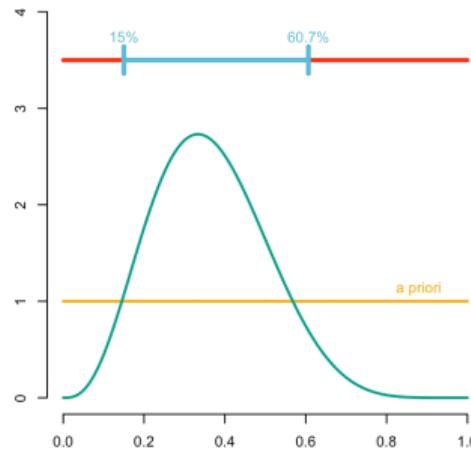


# Bayesianism, statistics and calculus XI

## ► Posterior distribution

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

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

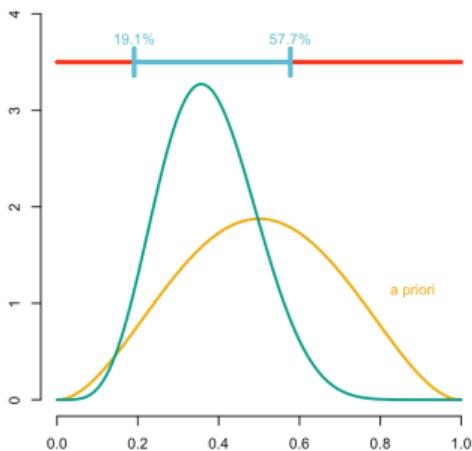
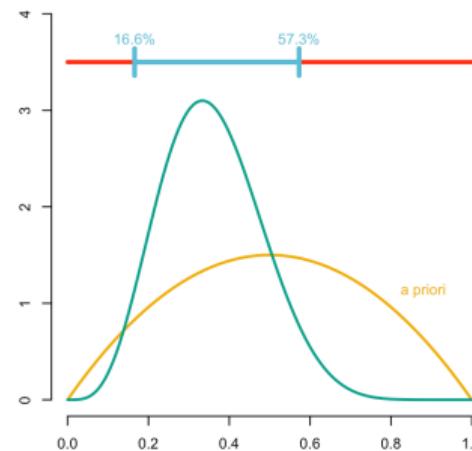
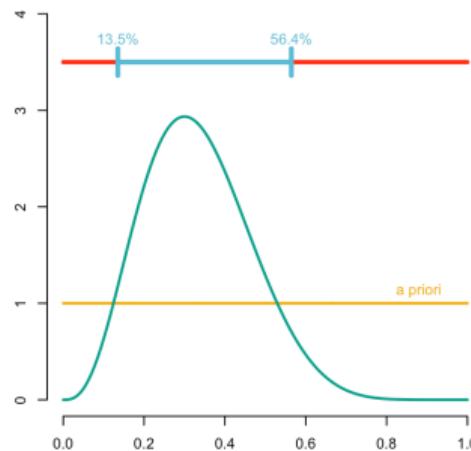


# Bayesianism, statistics and calculus XII

## ► Posterior distribution

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

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

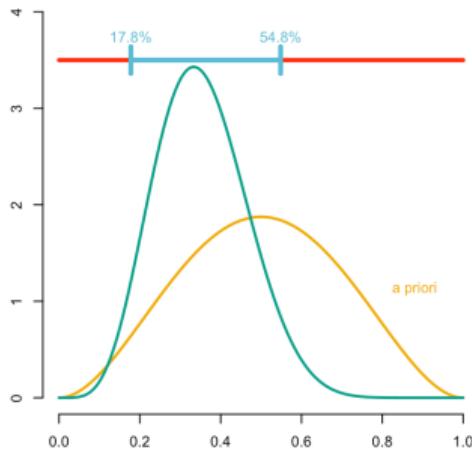
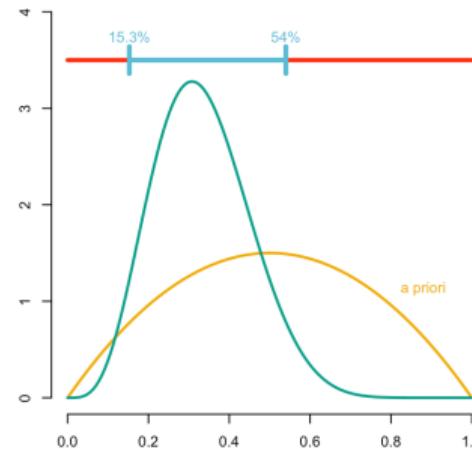
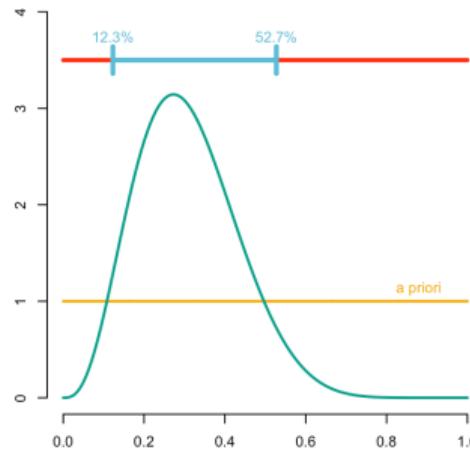


# Bayesianism, statistics and calculus XIII

## ► Posterior distribution

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

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

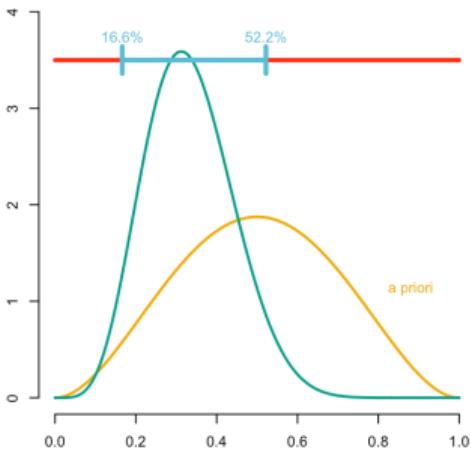
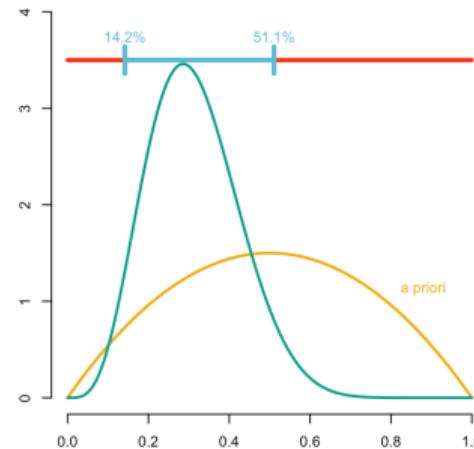
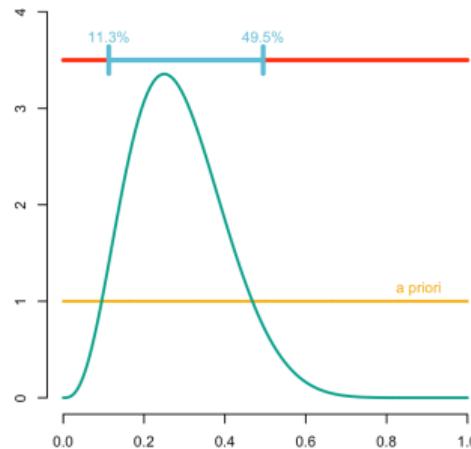


# Bayesianism, statistics and calculus XIV

## ► Posterior distribution

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

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

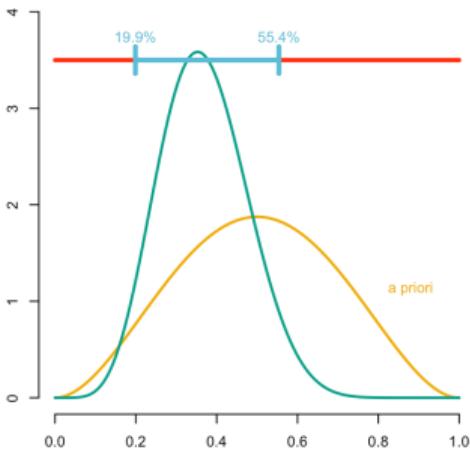
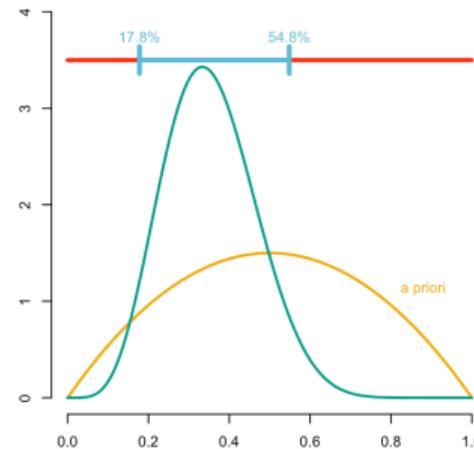
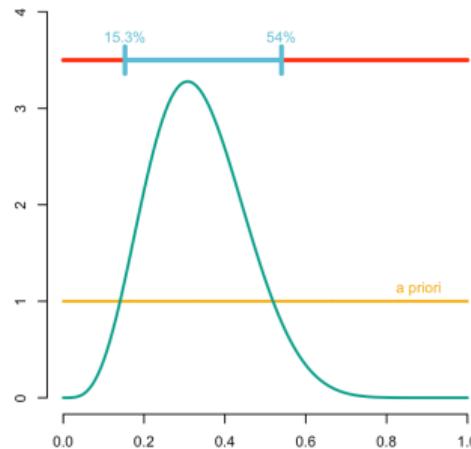


# Bayesianism, statistics and calculus XV

## ► Posterior distribution

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

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

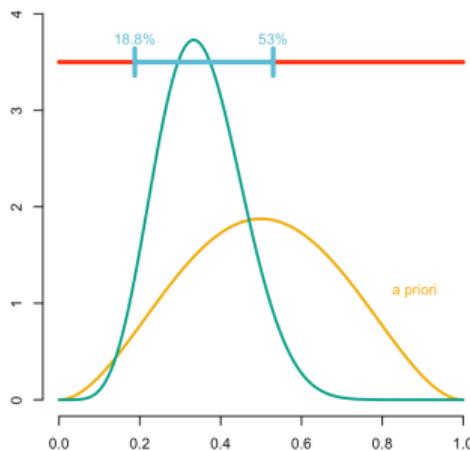
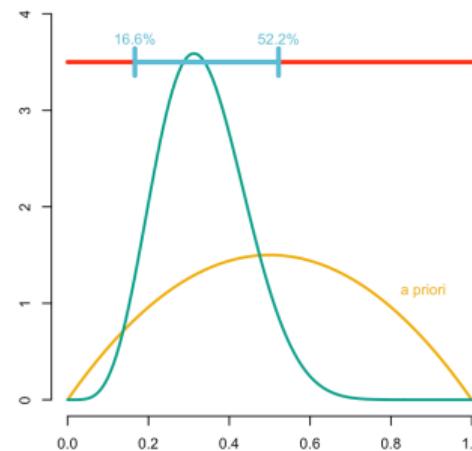
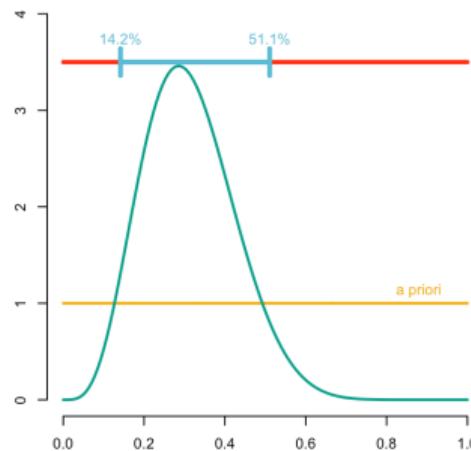


# Bayesianism, statistics and calculus XVI

## ► Posterior distribution

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

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

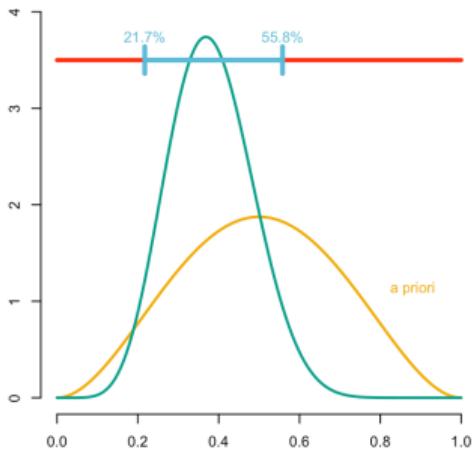
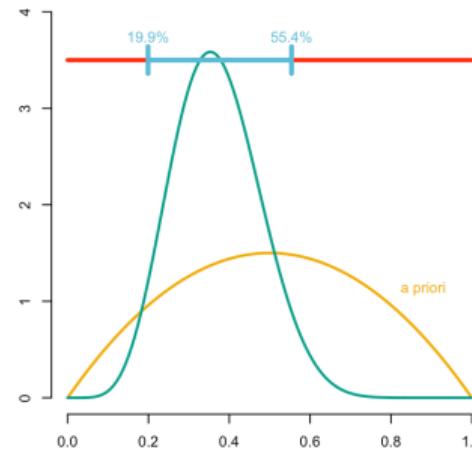
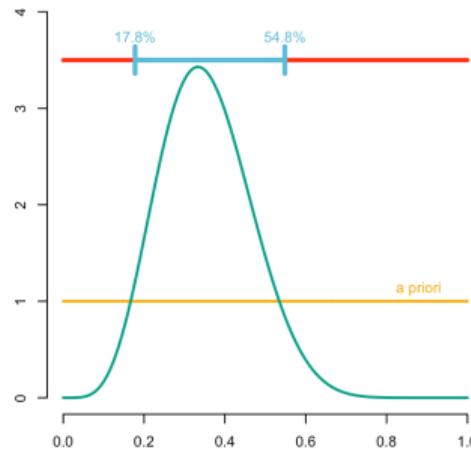


# Bayesianism, statistics and calculus XVII

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

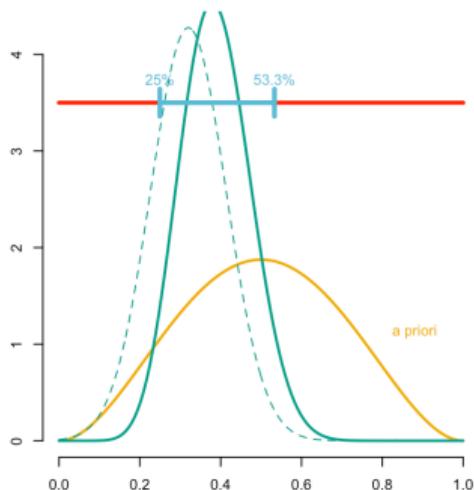
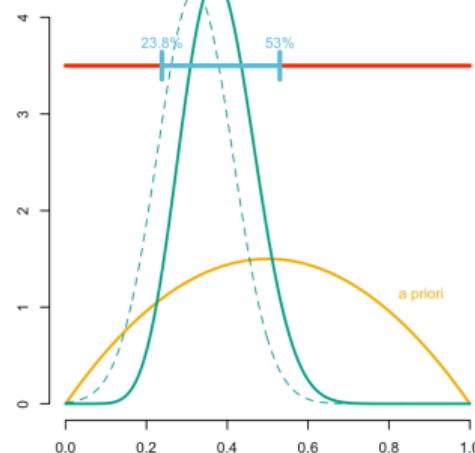
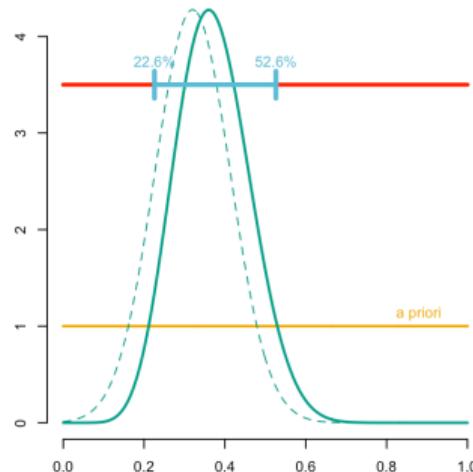


# Bayesianism, statistics and calculus XVIII

## ► Posterior distribution

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

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



freakonometrics

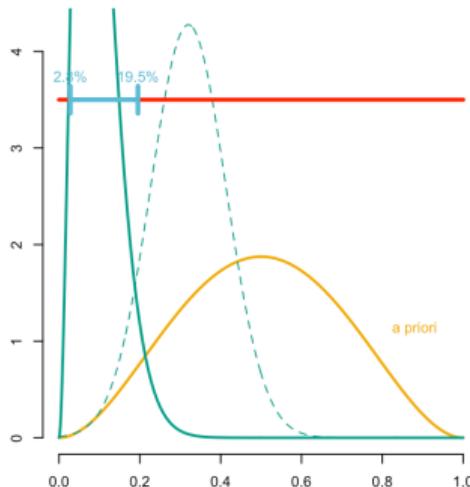
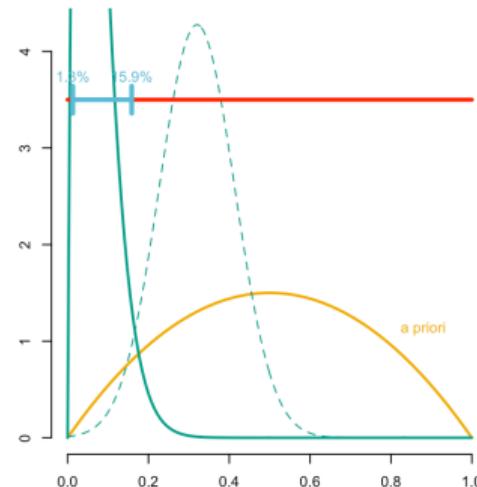
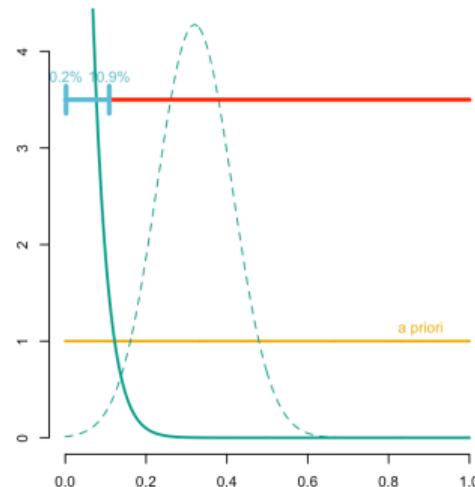
freakonometrics.hypotheses.org

# Bayesianism, statistics and calculus XIX

## ► Posterior distribution

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

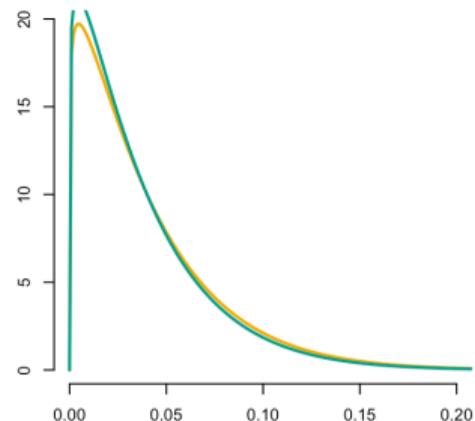
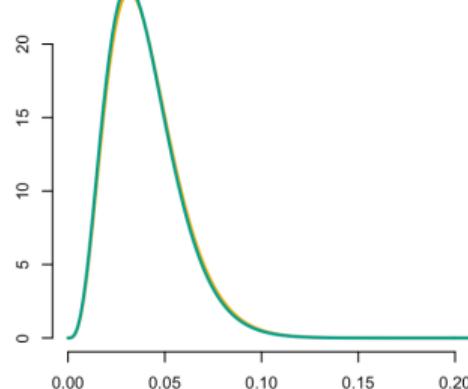
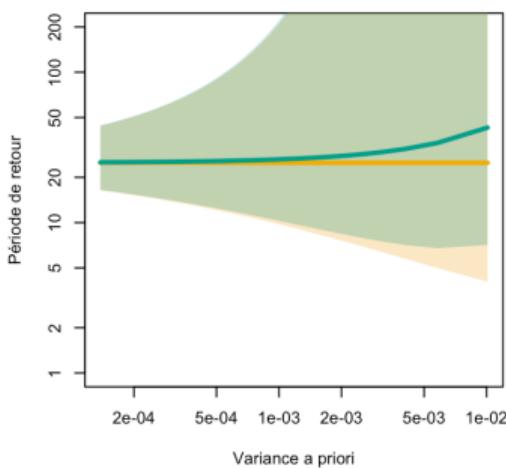
Bayesian approach,  $\hat{\theta}|\mathbf{x} \sim \text{Beta}(\alpha_0, \beta_0 + n)$ , since  $\sum_{i=1}^n x_i = 0$



# Bayesianism, statistics and calculus XX

## ► Posterior distribution

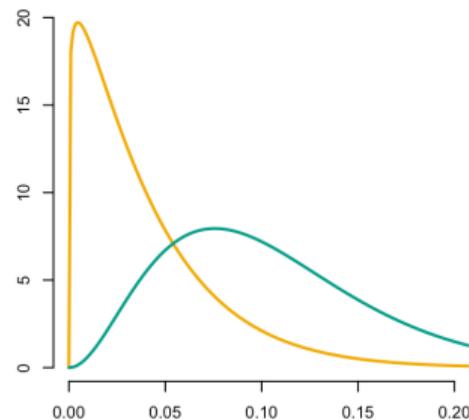
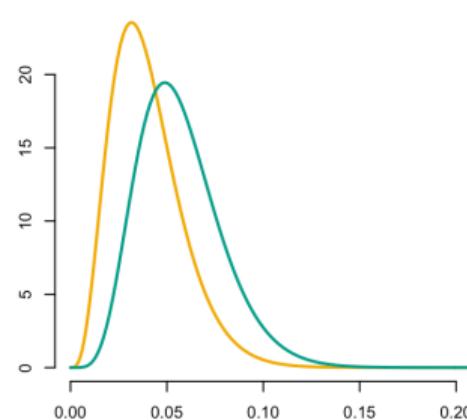
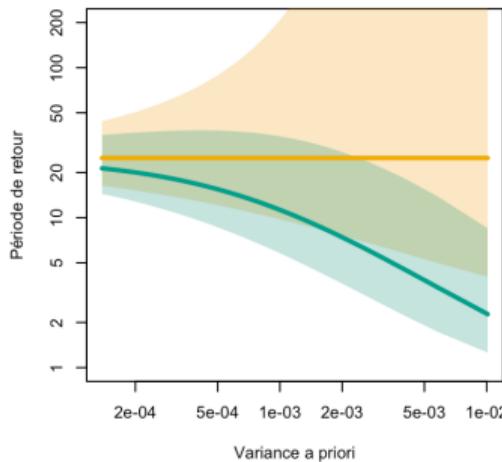
Ministère de l'intérieur (2019) “A single threshold for qualifying a geotechnical drought as abnormal: a return period greater than or equal to 25 years ” (probabilité 1/25) (probability 1/25) No drought has been observed over 2 years ( $\{0, 0\}$ ), what happens to our belief about the return period?



# Bayesianism, statistics and calculus XXI

## ► Posterior distribution

As a comparison, if we have observed two major droughts ( $\{1, 1\}$ ), our beliefs a posteriori are very influenced by these unexpected events



# Bayesianism, statistics and calculus XXII

## ► 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"

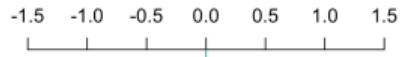
For the confidence interval, we look for  $[\hat{a}_{\mathcal{D}}, \hat{b}_{\mathcal{D}}]$  such that  $P[\theta \in [\hat{a}_{\mathcal{D}}, \hat{b}_{\mathcal{D}}]] \geq 95\%$ .

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 XXIII

## ► "confidence interval"

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

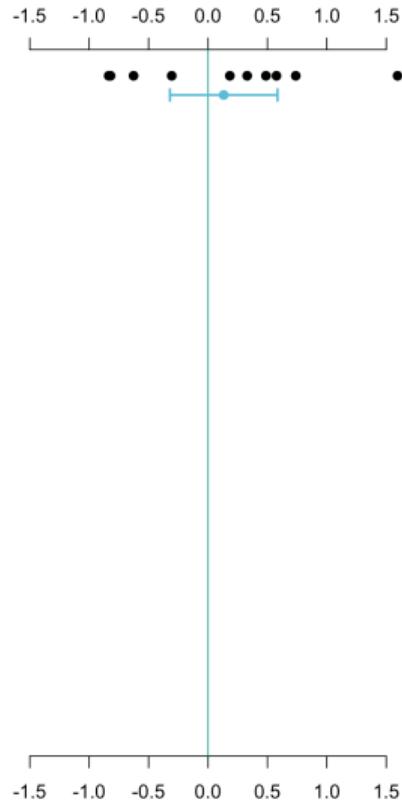


## Bayesianism, statistics and calculus XXIV

### ► "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 XXV

## ► "confidence interval"

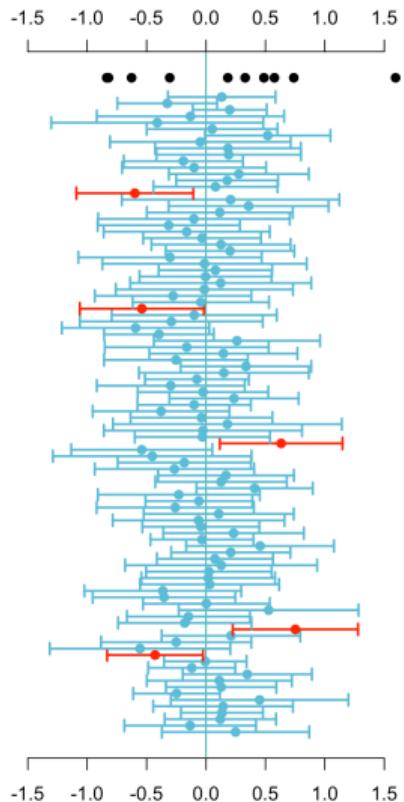
Suppose  $\mathcal{D} = \{x_1, \dots, x_n\}$ ,  $X_i \sim \mathcal{N}(\theta, \sigma^2)$   
(here  $\theta = 0$ )

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  $[a, b]$ .

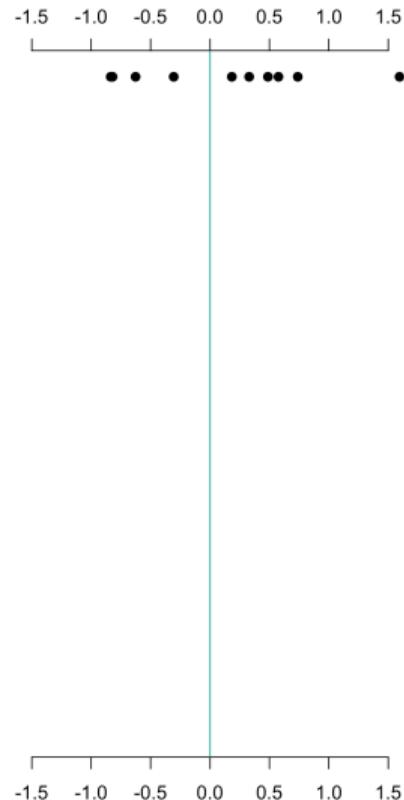


## Bayesianism, statistics and calculus XXVI

### ► "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  $\theta$

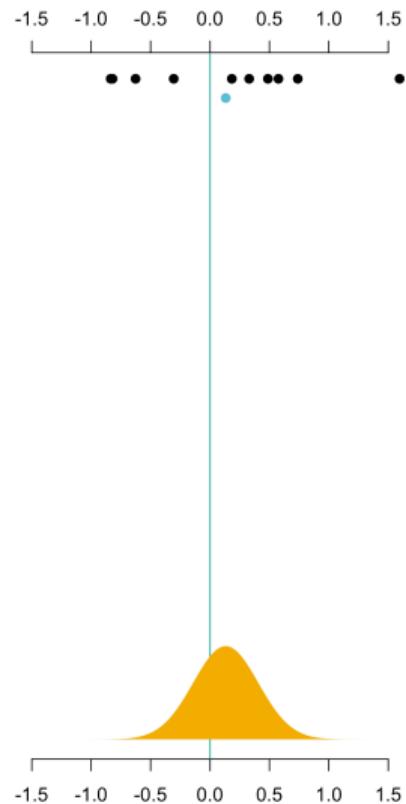


## Bayesianism, statistics and calculus XXVII

### ► "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  $\theta$   
and  $\pi(\cdot|\mathcal{D})$  is the posterior distribution  
(potentially complicated)



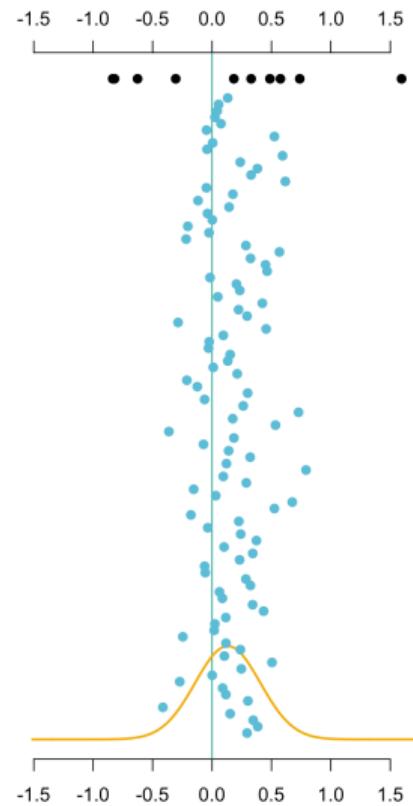
## Bayesianism, statistics and calculus XXVIII

### ► "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  $\theta$   
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 XXIX

## ► "credibility interval"

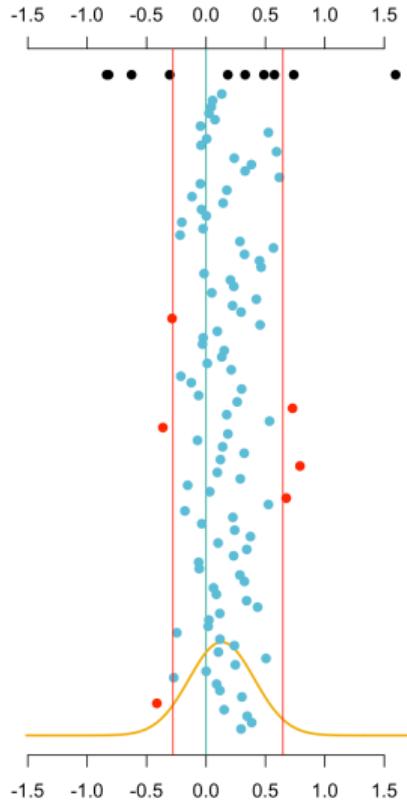
Suppose  $\mathcal{D} = \{x_1, \dots, x_n\}$ ,  $X_i \sim \mathcal{N}(\theta, \sigma^2)$

Consider some prior distribution  $\pi(\cdot)$  for  $\theta$   
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 XXX

## ► "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  $\theta$   
and  $\pi(\cdot|\mathcal{D})$  is the posterior distribution  
(potentially complicated)

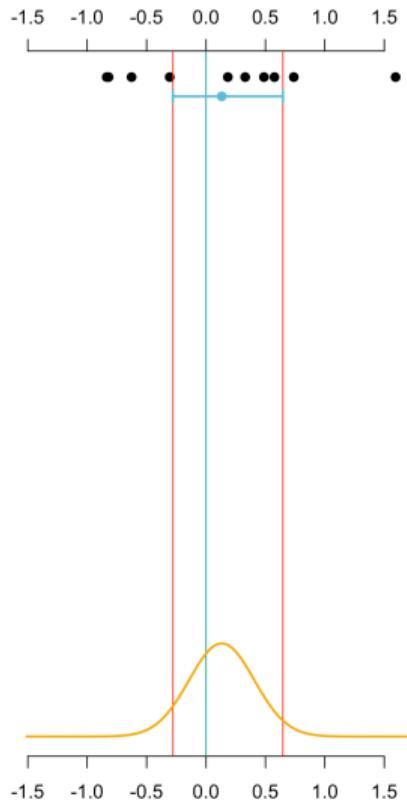
Suppose we generate  $\tilde{\theta}_1, \dots, \tilde{\theta}_k$  given  $\pi(\cdot|\mathcal{D})$ .

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

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



## Bayesianism, statistics and calculus XXXI

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  $\pi$  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 XXXII

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

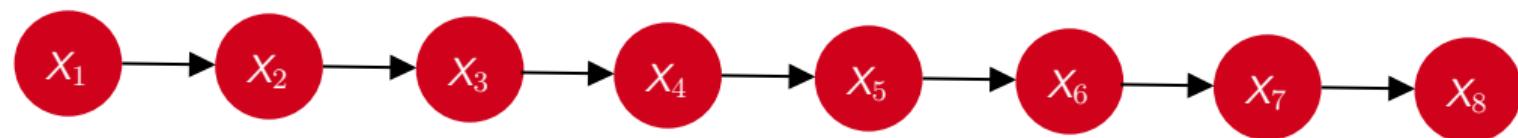
$$F_n = \frac{n}{n + \alpha} \widehat{F}_n + \frac{\alpha}{n + \alpha} F_0, \text{ where } \widehat{F}_n(x) = \frac{1}{n} \sum_{j=1}^n \mathbf{1}(x \leq x_j)$$

# Bayes and Markov property I

## ► Markov property

This property allows to simplify the writing (and the calculation) of the posterior distribution

$$\mathbb{P}[X_{t+1} = x_{t+1} | X_t = x_t, X_{t-1} = x_{t-1}, \dots] = \mathbb{P}[X_{t+1} = x_{t+1} | X_t = x_t]$$



As a reminder, under some technical assumptions, the transition kernel  $p(x_{t+1}|x_t)$  converges ( $t \rightarrow \infty$ ) to a stationary measure  $p^*(x)$ .

If  $x_t \in \mathcal{X}$  of finite cardinal,  $p(\cdot|\cdot)$  reads in a (stochastic) matrix  $P$ .

$$\mathbb{P}[X_{t+k} = j | X_t = i] = [P^k]_{ij} \text{ (Chapman Kolmogorov)}$$

# Bayes and Markov property II

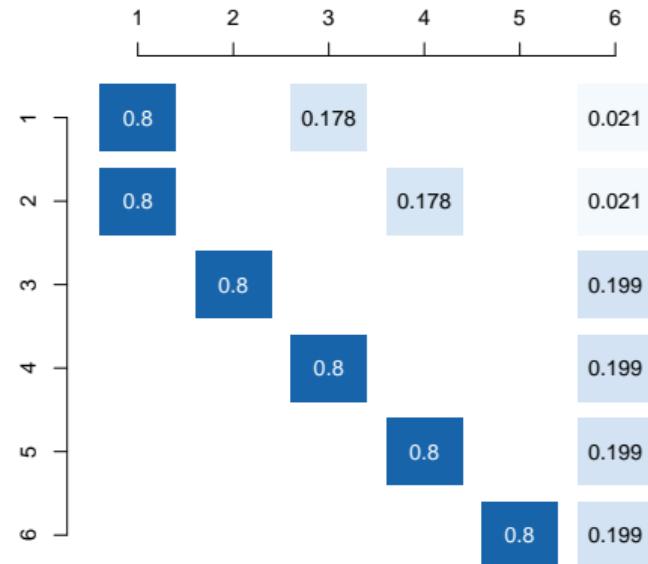
Example bonus-malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+1 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $\mathbb{P}(N = 0) = 20\%$ .

# Bayes and Markov property III

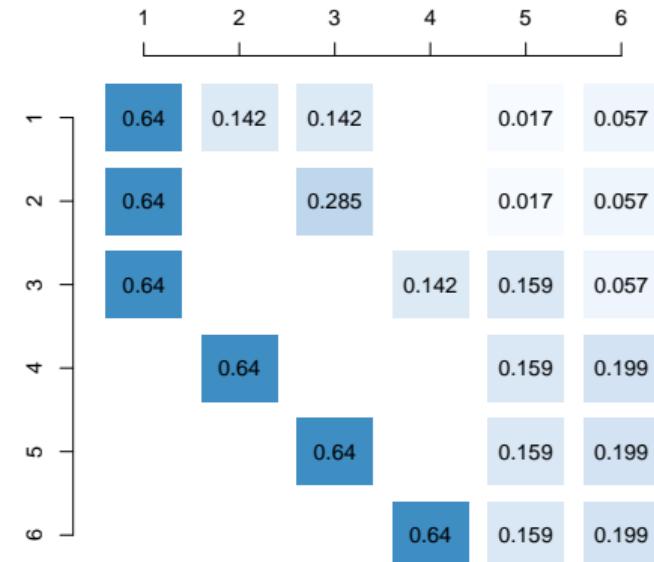
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+2 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $\mathbb{P}(N = 0) = 20\%$ .

# Bayes and Markov property IV

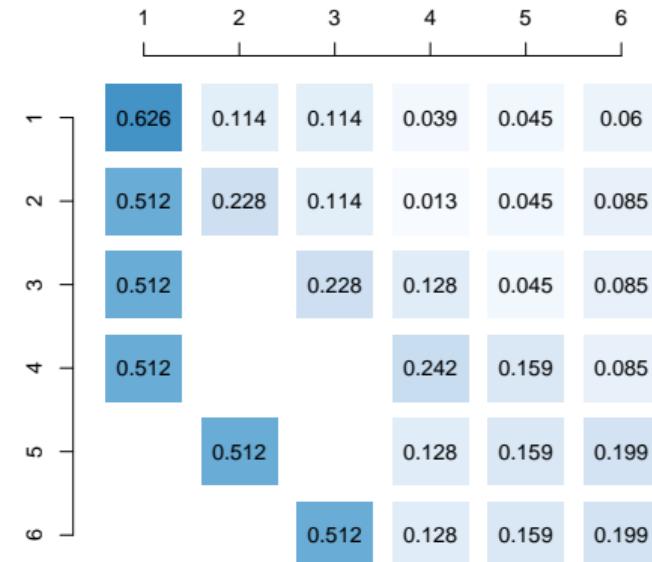
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+3 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $P(N = 0) = 20\%$ .

# Bayes and Markov property V

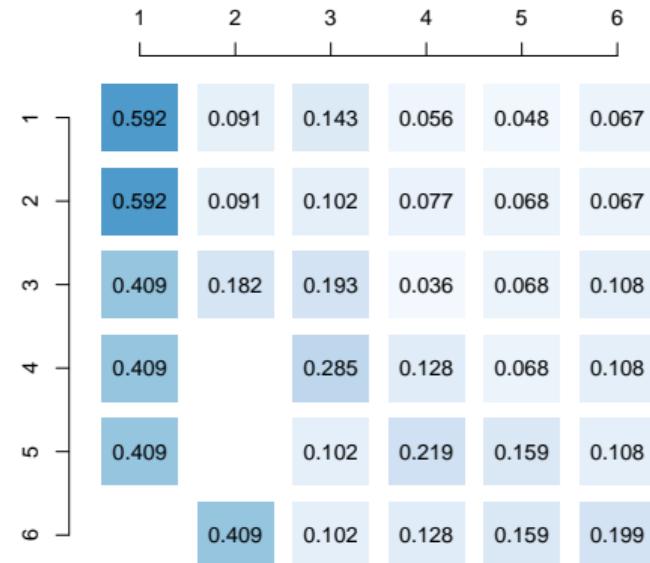
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+4 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $P(N = 0) = 20\%$ .

# Bayes and Markov property VI

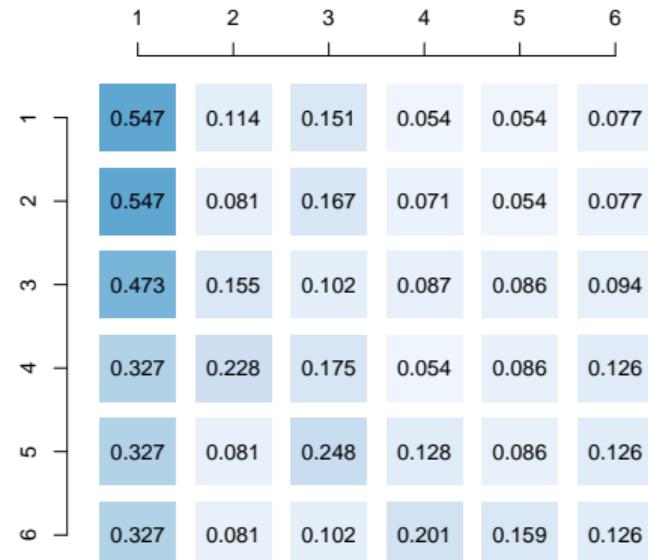
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+5 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $P(N = 0) = 20\%$ .

# Bayes and Markov property VII

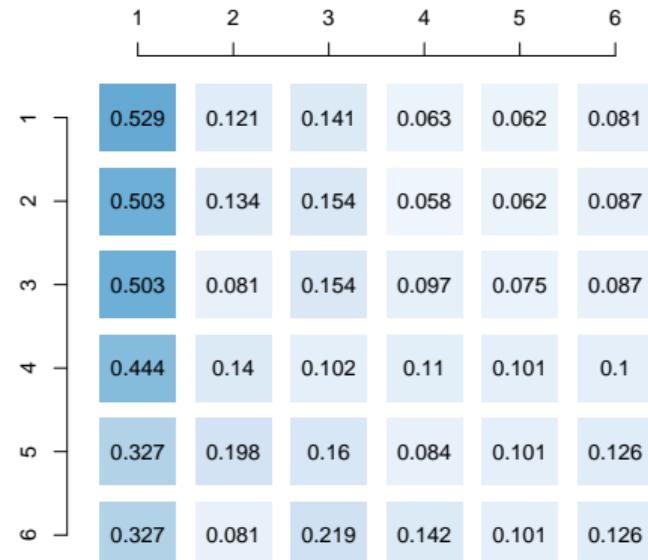
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+6 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $P(N = 0) = 20\%$ .

# Bayes and Markov property VIII

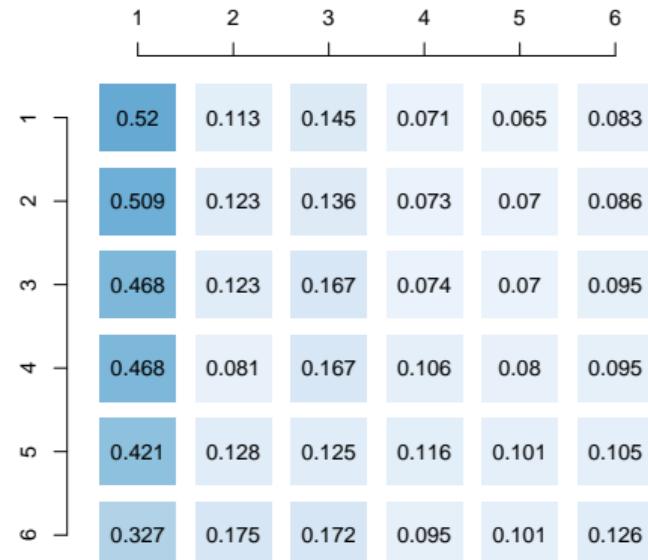
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+7 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $P(N = 0) = 20\%$ .

# Bayes and Markov property IX

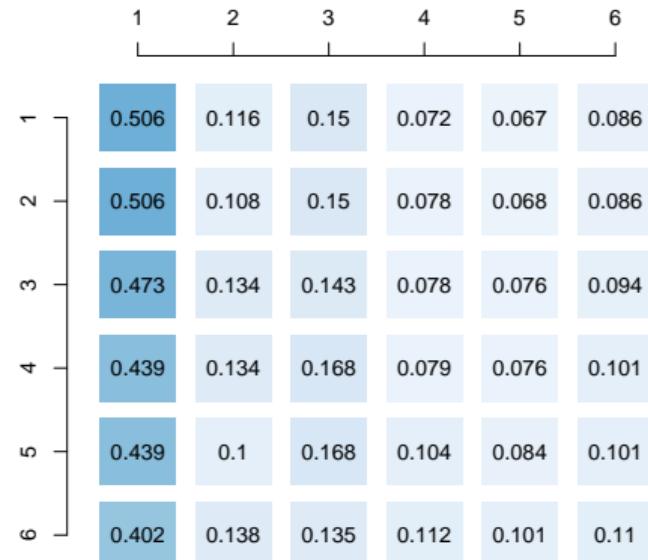
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+8 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $P(N = 0) = 20\%$ .

# Bayes and Markov property X

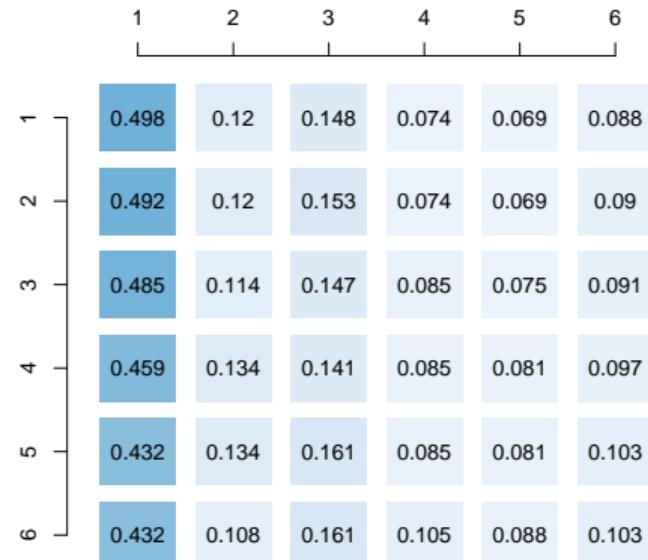
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+9 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $P(N = 0) = 20\%$ .

# Bayes and Markov property XI

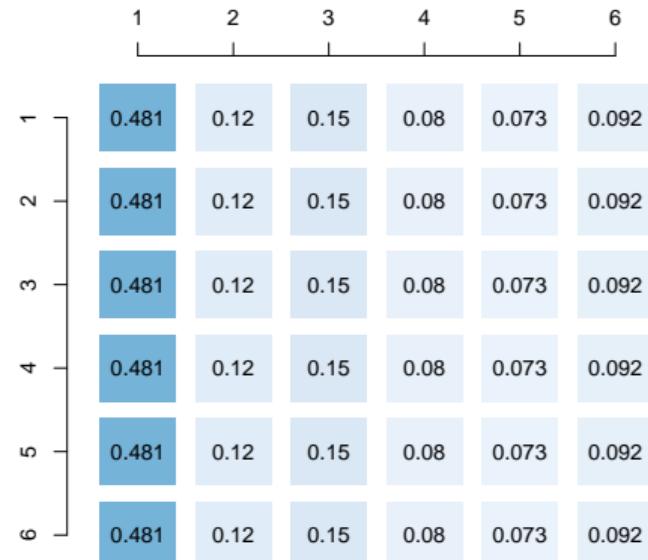
Example bonus malus schemes Lemaire (1995),

## HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	$\geq 2$
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.



t+100 vs. t

If claims frequency is  $N \sim \mathcal{P}(0.225)$ ,  
 $\mathbb{P}(N = 0) = 20\%$ .

## Bayes and Markov property XII

### ► Expected values and MCMC

#### Law of large numbers

if  $X_1, \dots, X_n, \dots$  i.i.d. with law  $p^*$ ,  $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{a.s.} \mathbb{E}_{p^*}(X) = \int x dp^*(x)$

Ergodic theorem (if  $p(\cdot|\cdot)$  has invariant distribution  $p^*$ )

if  $X_1, \dots, X_t, X_{t+1}, \dots$  is generated from  $p(\cdot|\cdot)$ ,  $\frac{1}{n} \sum_{t=t_0+1}^{t_0+n} X_t \xrightarrow{a.s.} \mathbb{E}_{p^*}(X) = \int x dp^*(x)$

where  $(X_t)$  is generated from  $p(\cdot|\cdot)$  using either d'Hasting-Metropolis or Gibbs sampler, Andrieu et al. (2003) or Kruschke (2014).

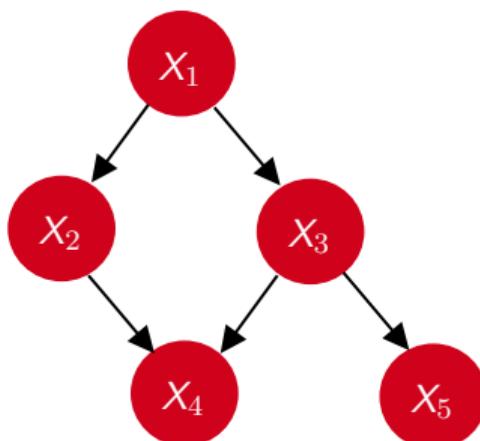
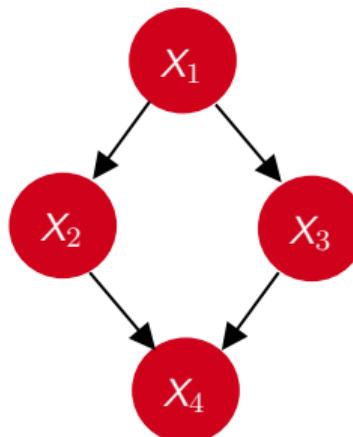
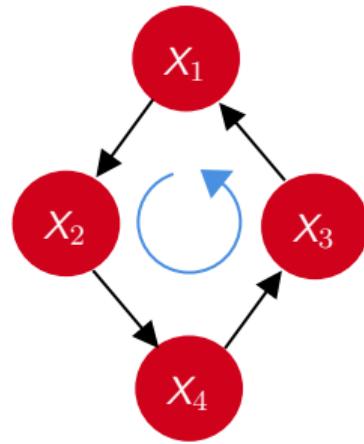
# Bayes and Markov property XIII

Using Markov property

$$\mathbb{P}(\mathbf{x}) = \prod_{i=2}^p \mathbb{P}(x_i|x_{i-1}) \cdot \mathbb{P}(x_1)$$

That can be extended on a DAG for the  $p$  variables.

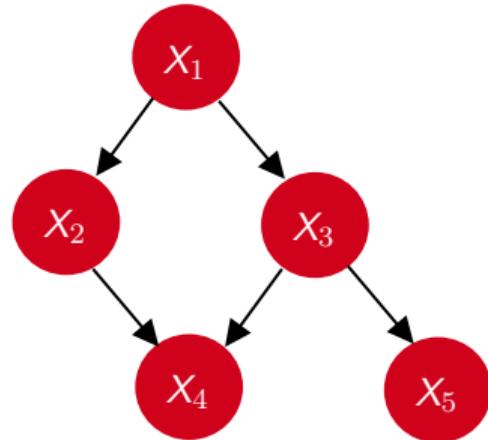
- Directed acyclic graph (DAG)



# Bayes and Markov property XIV

## ► Bayesian Network

A couple  $\{G, \mathbb{P}\}$  is a Bayesian network, if  $G = \{V, E\}$  is a DAG and if it satisfies the Markov property : each variable  $X$  in  $V$  is independent from its non-descendants, in  $G$ , conditional on its parents,



$$\mathbb{P}(\mathbf{x}) = \prod_{i=1}^p \mathbb{P}(x_i | \mathbf{x}_{\text{parents}_i})$$

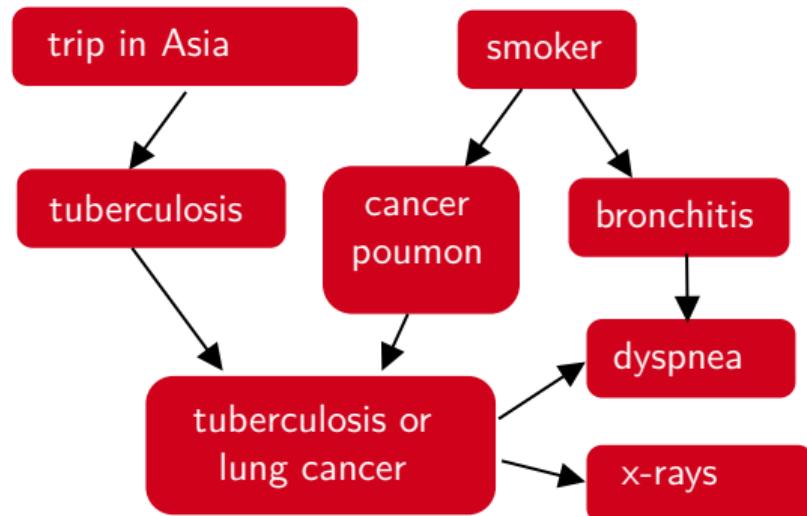
$$\left\{ \begin{array}{l} X_2 \perp\!\!\!\perp \{X_3, X_4\} \mid X_1 \\ X_3 \perp\!\!\!\perp X_2 \mid X_1 \\ X_4 \perp\!\!\!\perp \{X_1, X_5\} \mid \{X_2, X_3\} \\ X_5 \perp\!\!\!\perp \{X_1, X_2, X_4\} \mid X_3 \end{array} \right.$$

$$\mathbb{P}(\mathbf{x}) = \mathbb{P}(x_5 | x_3) \mathbb{P}(x_4 | x_2, x_3) \mathbb{P}(x_3 | x_1) \mathbb{P}(x_2 | x_1) \mathbb{P}(x_1)$$

# Bayes and Markov property XV

## ► Bayesian Network and Medical Diagnostics

via Lauritzen and Spiegelhalter (1988) and Højsgaard et al. (2012)



We have network (DAG)  
and conditional probabilities

# Bayesianism and statistical learning I

Econometrics is based on a probabilistic model, unlike most machine learning approaches, see [Charpentier et al. \(2018\)](#)

- ▶ in SVMs, the distance to the separation line is used as a score which can then be interpreted as a probability - [Platt scaling, Platt et al. \(1999\)](#) or [isotonic regression Zadrozny and Elkan \(2001, 2002\)](#) (see also [Niculescu-Mizil and Caruana \(2005\)](#) "good probabilities")
- ▶ GLM models (under additional conditions) satisfy the [autocalibration](#) property, [Denuit et al. \(2021\)](#), not machine learning models, i.e.

$$\mathbb{E}[Y|\hat{Y} = y] = y, \quad \forall y$$

[Lichtenstein et al. \(1977\)](#), [Dawid \(1982\)](#) or [Oakes \(1985\)](#), [Gneiting et al. \(2007\)](#)

## Bayesianism and statistical learning II

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

Very close to what exists to quantify uncertainty in weather models,

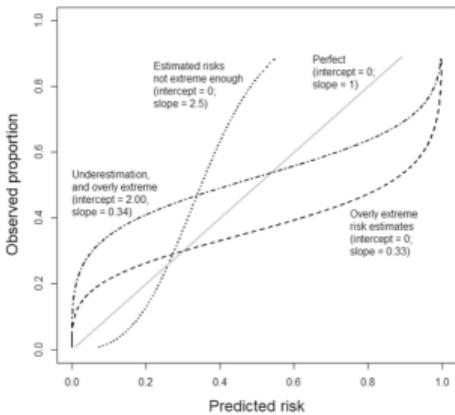
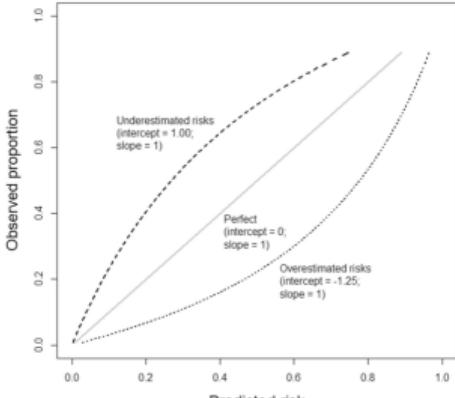
“*Suppose that a forecaster sequentially assigns probabilities to events. He is well calibrated if, for example, of those events to which he assigns a probability 30 percent, the long-run proportion that actually occurs turns out to be 30 percent*”, [Dawid \(1982\)](#) ou “*we desire that the estimated class probabilities are reflective of the true underlying probability of the sample*”, [Kuhn et al. \(2013\)](#)

# Bayesianism and statistical learning III

As explained in Van Calster et al. (2019), "*among patients with an estimated risk of 20%, we expect 20 in 100 to have or to develop the event*" ,

- ▶ if 40 out of 100 in this group are found to have the disease, the risk is **underestimated**
- ▶ If we observe that in this group, 10 out of 100 have the disease, we have **overestimated** the risk.

Hosmer-Lemeshow test (Hosmer Jr et al. (2013)) for the logistic model.



## Bayesianism and statistical learning IV

- Ridge estimate, Hoerl and Kennard (1970) (linear model)

We look for  $\hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta\|_2^2 \right\}$ , "equivalent" to the constrained optimization problem  $\underset{\beta \in \mathbb{R}^p: \|\beta\|_2 \leq c}{\operatorname{argmin}} \left\{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \right\}$ .

Consider

$$\begin{cases} \mathbf{y} = \mathbf{X}\beta + \varepsilon \text{ or } \mathbf{y} | \mathbf{X}, \beta \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbb{I}) \\ \beta \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbb{I}) \text{ posterior} \end{cases}$$

Maximum a posteriori (MAP) satisfies

$$\hat{\beta}_{MAP} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \frac{\sigma^2}{\tau^2} \|\beta\|_2^2 \right\}$$

## Bayesianism and statistical learning V

- LASSO estimate, Tibshirani (1996) (linear regression)

We look for  $\hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta\|_1 \right\}$ , "equivalent" (Gill et al. (2019)) to the constrained optimization problem  $\underset{\beta \in \mathbb{R}^p: \|\beta\|_1 \leq c}{\operatorname{argmin}} \left\{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \right\}$ .

Consider (Tibshirani (1996) and Park and Casella (2008))

$$\begin{cases} \mathbf{y} = \mathbf{X}\beta + \varepsilon \text{ ou } \mathbf{y} | \mathbf{X}, \beta \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbb{I}) \\ \beta \sim \mathcal{L}(\tau) \text{ posterior, i.e. } \pi(\beta) = (\tau/2)^p \exp[-\tau \|\beta\|_1] \end{cases}$$

Maximum a posteriori (MAP) satisfies

$$\hat{\beta}_{MAP} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \sigma^2 \tau \|\beta\|_1 \right\}$$

# Bayesianism and statistical learning VI

Tibshirani (1996) suggested that Lasso estimates can be interpreted as posterior mode estimates when the regression parameters have independent and identical Laplace (i.e., double-exponential) priors

- Neural nets

Rumelhart et al. (1985), Rumelhart et al. (1986) Hertz et al. (1991) and Buntine and Weigend (1991) proposed to formalize back-propagation in a Bayesian context, taken up by MacKay (1992) and Neal (1992).

State of the art in Neal (2012), more than 25 years ago (or more recently Neal (2012) Theodoridis (2015), Gal and Ghahramani (2016) and Goulet et al. (2021))

# Bayesianism as a learning process I

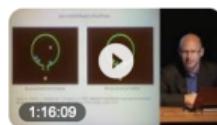
Old topic, see  
Shepard (1987) or Tenenbaum (1998).

*"How does abstract knowledge guide learning and reasoning from sparse data? How does the mind get so much from so little?",  
Tenenbaum et al. (2011)*

Discussed in Dehaene (2012),

[www.youtube.com > watch](http://www.youtube.com/watch)

la révolution Bayésienne... (1) - Stanislas Dehaene (2011-2012)



Enseignement 2011-2012 : Le cerveau statisticien : la révolution Bayésienne en sciences cognitives Cours du ma...

YouTube · Sciences de la vie - Collège de France · Il y a 1 semaine

Le cerveau statisticien : la révolution Bayésienne en sciences cognitives

Présentation

10 janvier 2012 ~ 09:30 ~

Cours

Introduction au raisonnement Bayésien et à ses applications

Stanislas Dehaene

17 janvier 2012 ~ 09:30 ~

Cours

Les mécanismes Bayésiens de l'induction chez l'enfant

Stanislas Dehaene

24 janvier 2012 ~ 09:30 ~

Cours

Les illusions visuelles : des inférences optimales ?

Stanislas Dehaene

31 janvier 2012 ~ 09:30 ~

Cours

Combinaison de contraintes et sélection d'un percept unique

Stanislas Dehaene

07 février 2012 ~ 09:30 ~

Cours

La prise de décision Bayésienne

Stanislas Dehaene

14 février 2012 ~ 09:30 ~

Cours

L'implémentation neuronale des mécanismes Bayésiens

Stanislas Dehaene

21 février 2012 ~ 09:30 ~

Cours

Le cerveau vu comme un système prédictif

Stanislas Dehaene

## Bayesianism as a learning process II

The simplifications managed by the brain are known since a long time, [Goodman \(1955\)](#).

We have an urn containing 100 balls, a person draws a blue ball, what can we say ?  
A priori not much... except if in the past, we observed that all the urns always contained balls of the same color. A single observation can then be very informative  
Allows to learn how to learn, [Kemp and Tenenbaum \(2008\)](#), [Kemp et al. \(2010\)](#), [Tenenbaum et al. \(2011\)](#)

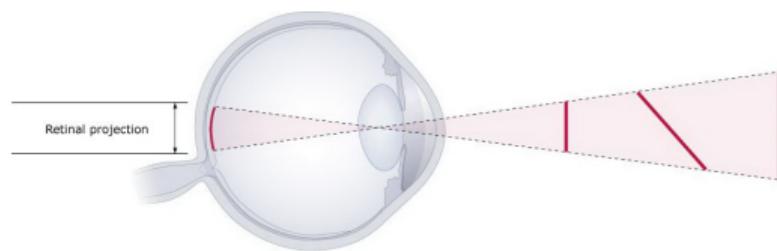
Language learning, [Stolcke \(1994\)](#), [Watanabe and Chien \(2015\)](#), [Duh \(2018\)](#) or [Murawaki \(2019\)](#).

Since [Shepard \(1992\)](#), many experiences on vision

## Bayesianism as a learning process III

Von Helmholtz (1867) defined “unbewusste Schluss”, or unconscious inference.

The view is constructed (more or less) as a projection, but (see linear algebra course) projections are not invertible: several images could have the same projection. Our brain looks for the most likely image

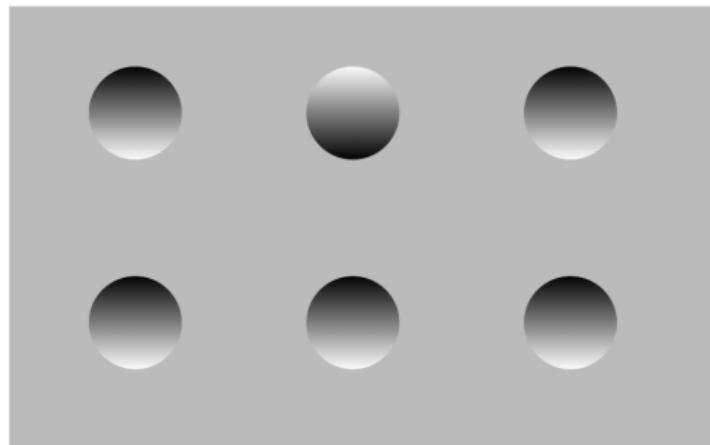


Sensory inputs are always ambiguous, so our perceptual system must select, among an infinite number of possible solutions, the one that is most plausible, Ernst and Banks (2002).

On vision as a Bayesian learning process Yuille and Kersten (2006), Clark (2013)  
Moreno-Bote et al. (2011)

## Bayesianism as a learning process IV

Classic example on "biases" of image perception, for example the [forms](#).

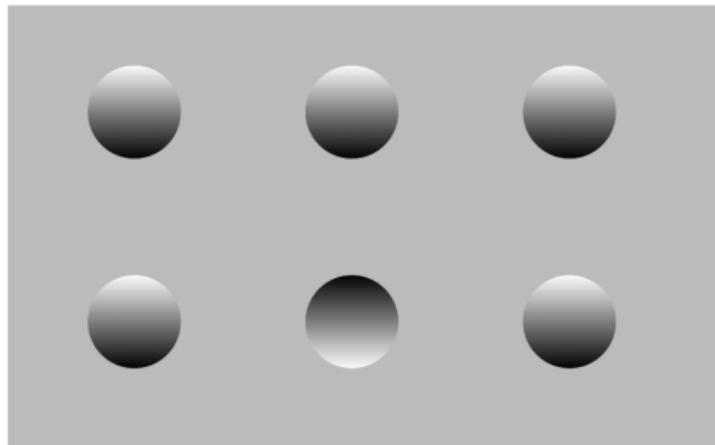


Consider the image above, what do we see?

Classically, we see 5 "holes" and 1 "bump"

## Bayesianism as a learning process V

Classic example on "biases" of image perception, for example the [forms](#).

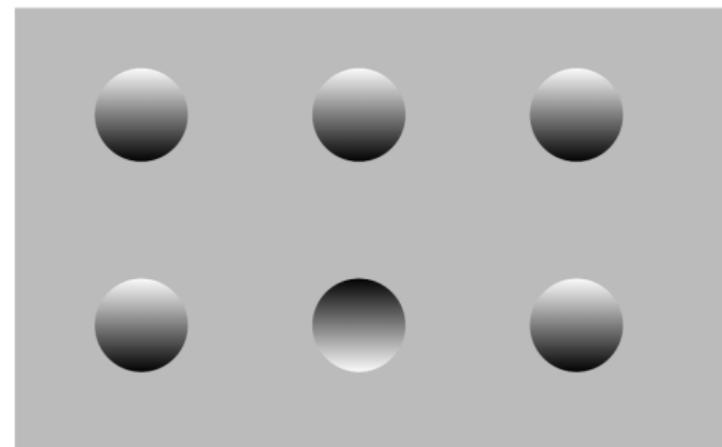
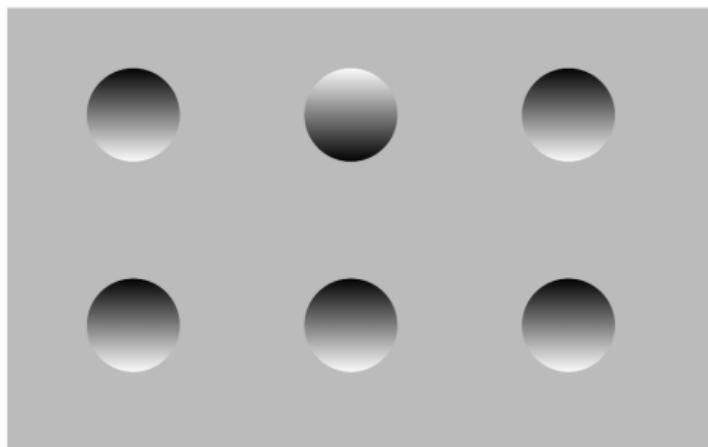


Consider the picture above, what do you see ?

Classically, 5 "bumps" et 1 "hole"

## Bayesianism as a learning process VI

Classic example on "biases" of image perception, for example the [forms](#).



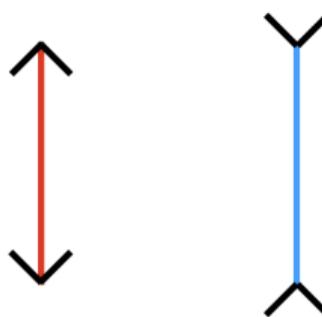
It is however the same figure (having undergone a rotation of 180°. (grey rectangle with 6 disks with a black/white gradient)). Ambiguous problem, [Ramachandran \(1988\)](#).

**Note:** our eye makes an inference about the light source (comes from above, without any other information - a priori assumption) to infer the shape.

## Bayesianism as a learning process VII

Classic example on "biases" of image perception, for example the lengths

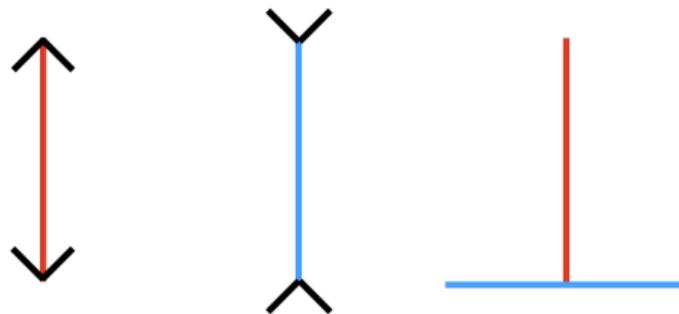
Among red and blue lines,  
which one is the longest?



## Bayesianism as a learning process VIII

Classic example on "biases" of image perception, for example the lengths

Among red and blue lines,  
which one is the longest?



As mentioned by Dehaene (2012), “*Bayesian inference gives a good account of perception processes: given ambiguous inputs, our brain reconstructs the most likely interpretation.*

# Bayesianism as a learning process IX

Classic example on "biases" of image perception, for example the [lengths](#)

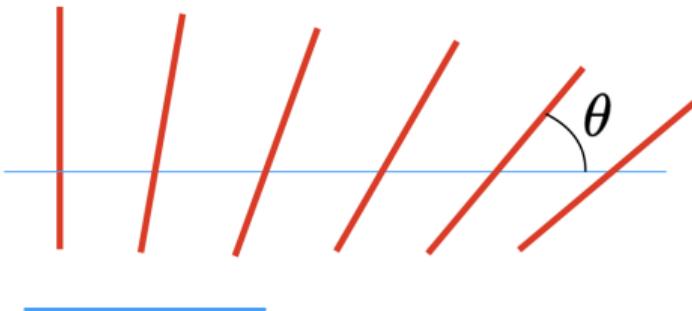
Among **red** and **blue** lines,  
which one is the longest?



Generally, all strokes **red** are seen as larger than the stroke **blue**.

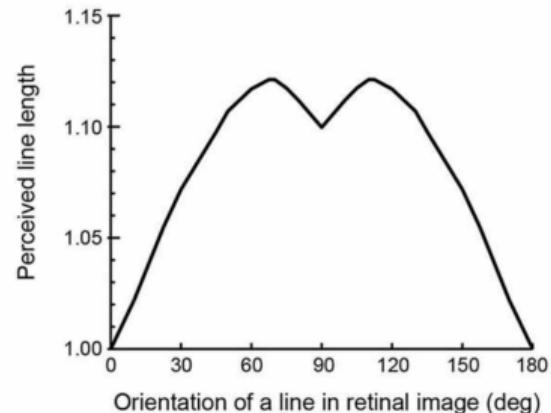
# Bayesianism as a learning process X

Which of the lines red and blue  
is larger?



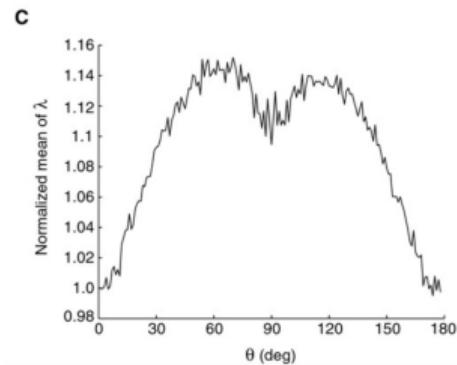
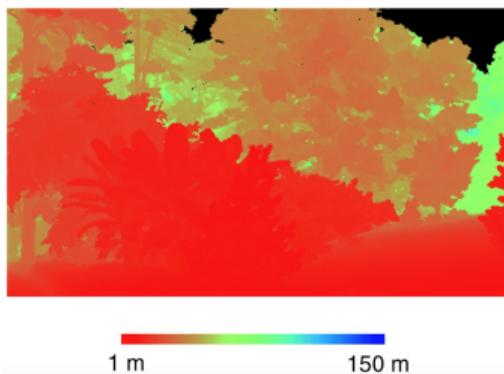
Several studies on the perception of the size of an object,  
according to its orientation (angle  $\theta$ )

Shipley et al. (1949), Pollock and Chapanis (1952), Cormack and Cormack (1974) and Purves et al. (2008) noted that the vertical line appears 10% larger than the horizontal line.



## Bayesianism as a learning process XI

The deformation made by the brain corresponds to a priori distributions that can be observed on images in nature, [Howe and Purves \(2002\)](#), [Purves \(2009\)](#), [Girshick et al. \(2011\)](#) or [Purves et al. \(2011\)](#) (based on (real) distances measured, by laser telemetry and compared to the measurement on the retina)



In other words, our retina has learned to correct the perceived distances according to the angle of inclination, in an everyday environment (3d), but continues to reproduce it for a drawing on a sheet (2d).

# Bayesianism as a learning process XII

One can also learn from Ensemble methods and by aggregation of opinions. For example, guess the weight of a cow, Cornwall, England, 1906, Galton (1907).

787 participants,  $x_1, \dots, x_n$ .

Unique prediction  $x_j$  v.s average  $\bar{x}$ ,

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

where  $t$  is the truth ("ambiguity decomposition").

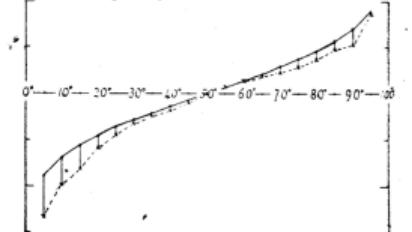
"Bayesian methods are sometimes proposed as mathematical aggregations of expert judgements", Hanea et al. (2021)

Distribution of the estimates of the dressed weight of a particular living ox, made by 787 different persons.

Degrees of the length of Array $o = 100$	Estimates in lbs.	* Centiles		Excess of Observed over Normal
		Observed deviates from 1198 lbs.	Normal p.e = .37	
5	1074	-133	-90	+43
10	1109	-98	-70	+28
15	1126	-81	-57	+24
20	1148	-59	-40	+13
25	1162	-45	-37	+8
30	1174	-33	-29	+4
35	1181	-26	-21	+5
40	1188	-19	-14	+5
45	1197	-10	-7	+3
50	1207	0	0	0
55	1214	+7	+7	0
60	1219	+12	+14	-2
65	1225	+18	+21	-3
70	1230	+23	+29	-6
75	1230	+29	+37	-8
80	1243	+36	+46	-10
85	1254	+47	+57	-10
90	1267	+52	+70	-18
95	1293	+86	+90	-4

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

Diagram, from the tabular values.



The continuous line is the normal curve with p.e = .37.  
The broken line is drawn from the observations.  
The lines connecting them show the differences between the observed and the normal.

## Bayesianism as a learning process XIII

*"I have approximate answers and possible beliefs and different degrees of certainty about different things"*, Feynman (2005)

*"Diversity and independence are important because the best collective decisions are the product of disagreement and contest, not consensus or compromise"*, Surowiecki (2005)

Merrick (2008), Karvetski et al. (2013) on model aggregation  $m_1, \dots, m_k$ ,

$$m(\mathbf{x}) = \sum_{i=1}^k \theta_i m_i(\mathbf{x}, \alpha_i)$$

with weights  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$  in the simplex  $\mathcal{S}_k$ . We assume a prior Dirichlet distribution.

See also Mongin (1995, 2001), inspired by Karni et al. (1983).

## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

$$\alpha_0 = 1, \beta_0 = 1 \quad 0.1331$$

$$\alpha_0 = 1, \beta_0 = 1 \quad 0.1282$$



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

1  $\alpha_1 = 2, \beta_1 = 1$  0.7369

$\alpha_1 = 1, \beta_1 = 1$  0.8081



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad \alpha_2 = 2, \beta_2 = 1 \quad 0.8835$$

$$0 \quad \alpha_2 = 1, \beta_2 = 2 \quad 0.3092$$



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

1    1  $\alpha_3 = 3, \beta_3 = 1$     0.9407

0     $\alpha_3 = 1, \beta_3 = 2$     0.8079



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

1	1	0	$\alpha_4 = 3$	$\beta_4 = 2$	0.6529
0	0	1	$\alpha_4 = 1$	$\beta_4 = 2$	0.5452



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

1	1	0	1	$\alpha_5 = 4$	$\beta_5 = 2$	0.5835
0				$\alpha_5 = 1$	$\beta_5 = 2$	0.8632



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad 1 \ 0 \ 1 \quad \alpha_6 = 4, \beta_6 = 2 \quad 0.7858$$

$$0 \quad \boxed{1} \quad \alpha_6 = 2, \beta_6 = 2 \quad 0.4509$$



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

1	1	0	1	1	$\alpha_7 = 5, \beta_7 = 2$	0.8871
0			1		$\alpha_7 = 2, \beta_7 = 2$	0.1643



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

1	1	0	1	1	1	$\alpha_8 = 6$	$\beta_8 = 2$	0.8052
0			1			$\alpha_8 = 2$	$\beta_8 = 2$	0.9383



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

1	1	0	1	1	1	$\alpha_9 = 6$	$\beta_9 = 2$	0.5769
0		1		0	1	$\alpha_9 = 2$	$\beta_9 = 3$	0.6047



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

$$1 \quad 1 \quad 0 \quad 1 \quad 1 \quad 1 \quad \alpha_{10} = 6, \beta_{10} = 2 \quad 0.692$$

$$0 \quad 1 \quad 0 \quad 0 \quad \alpha_{10} = 2, \beta_{10} = 4 \quad 0.5244$$



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

1	1	0	1	1	1	0	1	$\alpha_{15} = 9$	$\beta_{15} = 4$	0.5462
0		1	0					$\alpha_{15} = 2$	$\beta_{15} = 4$	0.2837



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$

$$\begin{array}{ccccccccc} 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 0 & & & & & & & & \\ & 0 & 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{array} \boxed{1} \quad \alpha_{20} = 11, \beta_{20} = 6 \quad 0.5201$$
$$0 \quad 1 \quad 0 \quad 0 \quad 1 \quad \alpha_{20} = 3, \beta_{20} = 4 \quad 0.2459$$



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$



## Bayesianism as a learning process

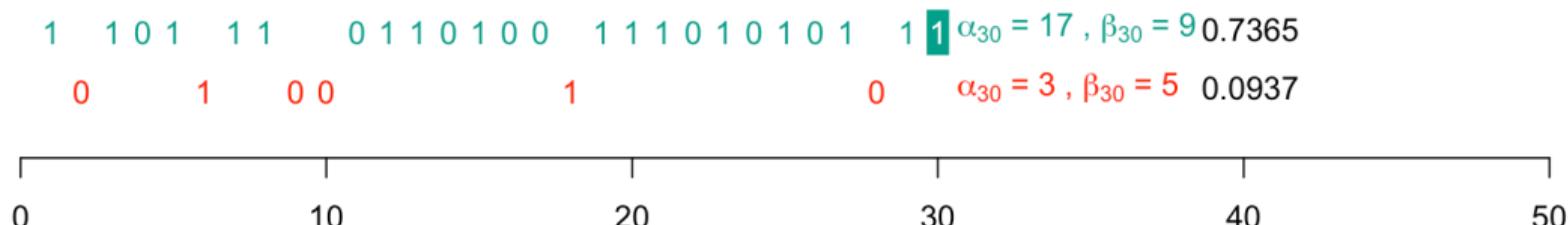
Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname*{argmin}_{k=1,\dots,K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
  - ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$



## Bayesianism as a learning process

Thompson sampling (or posterior sampling and probability matching), by [Thompson \(1933, 1935\)](#), and Beta-Bernoulli bandits.

We have to choose among  $K$  alternatives, that yield  $\mathbf{X} = (X_1, \dots, X_K)$ , with  $X_k \sim \mathcal{B}(\theta_k)$ .

Assume (prior)  $\theta_k \sim \text{Beta}(\alpha_k, \beta_k)$ . At time  $t$ , draw  $K$  Beta variables (independents)  $B_k \sim \text{Beta}(\alpha_k, \beta_k)$ , and select  $k^* = \operatorname{argmin}_{k=1, \dots, K} \{B_k\}$ .

Consider updating  $(\alpha_{k^*}, \beta_{k^*}) \leftarrow (\alpha_{k^*} + x_{k^*}, \beta_{k^*} + (1 - x_{k^*}))$ ,

- ▶ simulated data, i.i.d.,  $X_1 \sim \mathcal{B}(72\%)$
- ▶ simulated data, i.i.d.,  $X_2 \sim \mathcal{B}(24\%)$



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



$$\alpha_0 = 1, \beta_0 = 1 \quad 0.4267$$

$$\alpha_0 = 1, \beta_0 = 1 \quad 0.8151$$



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



$$\alpha_1 = 1, \beta_1 = 1 \quad 0.4473$$

$$1 \quad \alpha_1 = 2, \beta_1 = 1 \quad 0.6376$$



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



$$\alpha_2 = 1, \beta_2 = 1 \quad 0.5947$$

$$1 \boxed{0} \alpha_2 = 2, \beta_2 = 2 \quad 0.6936$$



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



$$\alpha_3 = 1, \beta_3 = 1 \quad 0.5552$$

$$1 \ 0 \ 0 \quad \alpha_3 = 2, \beta_3 = 3 \quad 0.8841$$



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



$$\alpha_4 = 1, \beta_4 = 1 \quad 0.1631$$

$$1 \ 0 \ 0 \ 1 \quad \alpha_4 = 3, \beta_4 = 3 \quad 0.5958$$



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



$$\alpha_5 = 1, \beta_5 = 1 \quad 0.8508$$

1 0 0 1 0  $\alpha_5 = 3, \beta_5 = 4 \quad 0.3475$



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



$$\begin{array}{l} \boxed{1} \quad \alpha_6 = 2, \beta_6 = 1 \quad 0.8951 \\ 10010 \quad \alpha_6 = 3, \beta_6 = 4 \quad 0.8137 \end{array}$$



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

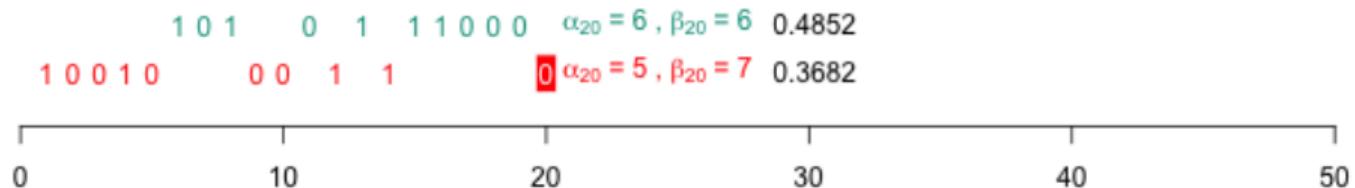
- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

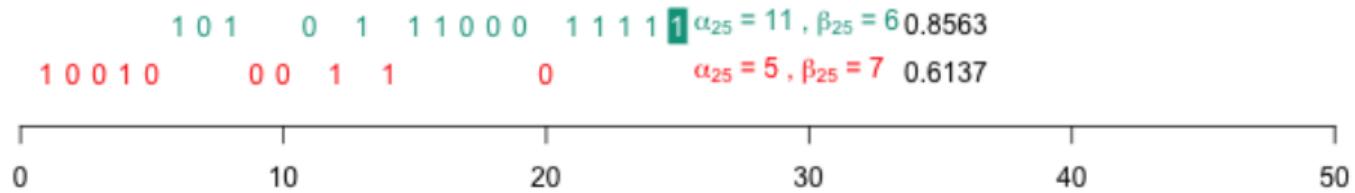
- ▶ strategy 1 : always switch the door
- ▶ strategy 2 : never switch the door



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



# Bayesianism as a learning process

We can use that approach in the context of Monty Hall

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



## "Conclusion" or wrap-up

- ▶ the Bayesian approach is interesting to describe beliefs in front of uncertain events, in particular if the events will occur only once
- ▶ Bayesian computation can be interpreted as a belief update or as an inverse problem
- ▶ is very strongly linked to causal graphs
- ▶ allows to take into account expert opinions, and proposes an ensemble method modeling describes both human and machine learning



## "Conclusion" or wrap-up

MODIFIED BAYES' THEOREM:

$$P(H|X) = P(H) \times \left( 1 + P(C) \times \left( \frac{P(X|H)}{P(X)} - 1 \right) \right)$$

H: HYPOTHESIS

X: OBSERVATION

P(H): PRIOR PROBABILITY THAT H IS TRUE

P(X): PRIOR PROBABILITY OF OBSERVING X

P(C): PROBABILITY THAT YOU'RE USING  
BAYESIAN STATISTICS CORRECTLY

(via <https://xkcd.com/2059/>)

# Références I

- Adjemian, S. and Pelgrin, F. (2008). Un regard bayésien sur les modèles dynamiques de la macroéconomie. *Economie prévision*, (2):127–152.
- Andrieu, C., De Freitas, N., Doucet, A., and Jordan, M. I. (2003). An introduction to mcmc for machine learning. *Machine learning*, 50(1):5–43.
- Anscombe, F. J., Aumann, R. J., et al. (1963). A definition of subjective probability. *Annals of mathematical statistics*, 34(1):199–205.
- Bailey, A. L. (1950). *Credibility Procedures: Laplace's generalization of Bayes' Rule and the combination of collateral knowledge with observed data*. New York State Insurance Department,.
- Barclay, S. et al. (1977). Handbook for decisions analysis.
- Baron, K. and Lange, J. (2006). *Parimutuel applications in finance: new markets for new risks*. Springer.
- 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.

## Références II

- Bell, E. T. (1945). *The development of mathematics*. Courier Corporation.
- Berliner, L. M., Levine, R. A., and Shea, D. J. (2000). Bayesian climate change assessment. *Journal of Climate*, 13(21):3805–3820.
- 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.
- Buehler, R. J. (1976). Coherent preferences. *The Annals of Statistics*, 4(6):1051–1064.
- Bühlmann, H. (1967). Experience rating and credibility. *ASTIN Bulletin: The Journal of the IAA*, 4(3):199–207.
- Buntine, W. L. and Weigend, A. S. (1991). Bayesian back-propagation. *Complex Systems*, 5.
- Cardano, G. (1564). *Liber de ludo aleae*. Franco Angeli.

## Références III

- Charpentier, A., Flachaire, E., and Ly, A. (2018). Econometrics and machine learning. *Economie et Statistique*, 505(1):147–169.
- Chen, Y. and Pennock, D. M. (2010). Designing markets for prediction. *AI Magazine*, 31(4):42–52.
- Chipman, H., George, E., and McCulloch, R. (2006). Bayesian ensemble learning. *Advances in neural information processing systems*, 19.
- Choquet, G. (1954). Theory of capacities. In *Annales de l'Institut Fourier*, volume 5, pages 131–295.
- Clark, A. (2013). Whatever next? predictive brains, situated agents, and the future of cognitive science. *Behavioral and brain sciences*, 36(3):181–204.
- Cormack, E. O. and Cormack, R. H. (1974). Stimulus configuration and line orientation in the horizontal-vertical illusion. *Perception & Psychophysics*, 16(2):208–212.
- Cournot, A. A. (1843). *Exposition de la théorie des chances et des probabilités*. Hachette.
- Cox, R. T. (1946). Probability, frequency and reasonable expectation. *American journal of physics*, 14(1):1–13.
- Dawid, A. P. (1982). The well-calibrated bayesian. *Journal of the American Statistical Association*, 77(379):605–610.

## Références IV

- De Finetti, B. (1931). *Probabilismo: saggio critico sulla teoria delle probabilità e sul valore della scienza*. Francesco Perrella.
- De Finetti, B. (1937). La prévision: ses lois logiques, ses sources subjectives. In *Annales de l'institut Henri Poincaré*, volume 7, pages 1–68.
- Dehaene, S. (2012). *Le cerveau statisticien : la révolution Bayésienne en sciences cognitives*. Collège de France.
- Denuit, M., Charpentier, A., and Trufin, J. (2021). Autocalibration and tweedie-dominance for insurance pricing with machine learning. *Insurance: Mathematics & Economics*.
- Donnat, C., Miolane, N., Bunbury, F., and Kreindler, J. (2020). A bayesian hierarchical network for combining heterogeneous data sources in medical diagnoses. In *Machine Learning for Health*, pages 53–84. PMLR.
- Drouet, I. (2016). *Le bayésianisme aujourd’hui. Fondements et pratiques*. Éditions Matériologiques.
- Duh, K. (2018). Bayesian Analysis in Natural Language Processing. *Computational Linguistics*, 44(1):187–189.
- Eisenberg, E. and Gale, D. (1959). Consensus of subjective probabilities: The pari-mutuel method. *The Annals of Mathematical Statistics*, 30(1):165–168.

## Références V

- Ernst, M. O. and Banks, M. S. (2002). Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415(6870):429–433.
- 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.
- Feynman, R. P. (2005). *The pleasure of finding things out: The best short works of Richard P. Feynman*. Basic Books.
- Gal, Y. and Ghahramani, Z. (2016). Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR.
- Galton, F. (1907). Vox populi (the wisdom of crowds). *Nature*, 75(7):450–451.
- Gelman, A. and Hill, J. (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge university press.

## Références VI

- Ghavamzadeh, M., Mannor, S., Pineau, J., Tamar, A., et al. (2015). Bayesian reinforcement learning: A survey. *Foundations and Trends® in Machine Learning*, 8(5-6):359–483.
- Ghosh, J. K. and Ramamoorthi, R. V. (2003). *Bayesian Nonparametrics*. Springer Verlag.
- Gigerenzer, G. and Edwards, A. (2003). Simple tools for understanding risks: from innumeracy to insight. *Bmj*, 327(7417):741–744.
- Gigerenzer, G. and Hoffrage, U. (1995). How to improve bayesian reasoning without instruction: frequency formats. *Psychological review*, 102(4):684.
- Gill, P. E., Murray, W., and Wright, M. H. (2019). *Practical optimization*. SIAM.
- Girshick, A. R., Landy, M. S., and Simoncelli, E. P. (2011). Cardinal rules: visual orientation perception reflects knowledge of environmental statistics. *Nature neuroscience*, 14(7):926–932.
- Gneiting, T., Balabdaoui, F., and Raftery, A. E. (2007). Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(2):243–268.
- Gneiting, T. and Raftery, A. E. (2005). Weather forecasting with ensemble methods. *Science*, 310(5746):248–249.
- Good, I. J. (1966). Speculations concerning the first ultraintelligent machine. In *Advances in computers*, volume 6, pages 31–88. Elsevier.

## Références VII

- Goodman, N. (1955). Axiomatic measurement of simplicity. *The Journal of Philosophy*, 52(24):709–722.
- Goulet, J.-A., Nguyen, L. H., and Amiri, S. (2021). Tractable approximate gaussian inference for bayesian neural networks. *J. Mach. Learn. Res.*, 22:251–1.
- Hájek, A. (2002). Interpretations of probability. *Stanford Encyclopedia of Philosophy*.
- Hájek, A. (2009). Dutch book arguments. In *The Handbook of Rational and Social Choice*. Oxford University Press.
- Hanea, A., Wilkinson, D. P., McBride, M., Lyon, A., van Ravenzwaaij, D., Singleton Thorn, F., Gray, C., Mandel, D. R., Willcox, A., Gould, E., et al. (2021). Mathematically aggregating experts predictions of possible futures. *PLoS one*, 16(9):e0256919.
- Hasselmann, K. (1998). Conventional and bayesian approach to climate-change detection and attribution. *Quarterly Journal of the Royal Meteorological Society*, 124(552):2541–2565.
- Hertz, J., Krogh, A., and Palmer, R. G. (1991). *Introduction to the theory of neural computation*. CRC Press.
- Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67.

## Références VIII

- Højsgaard, S., Edwards, D., and Lauritzen, S. (2012). *Graphical models with R*. Springer Verlag.
- Hosmer Jr, D. W., Lemeshow, S., and Sturdivant, R. X. (2013). *Applied logistic regression*, volume 398. John Wiley & Sons.
- Howe, C. Q. and Purves, D. (2002). Range image statistics can explain the anomalous perception of length. *Proceedings of the National Academy of Sciences*, 99(20):13184–13188.
- Hunt, I. and Mostyn, J. (2020). Probability reasoning in judicial fact-finding. *The international Journal of evidence & Proof*, 24(1):75–94.
- Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physical review*, 106(4):620.
- Jaynes, E. T. (1988). How does the brain do plausible reasoning? In *Maximum-entropy and Bayesian methods in science and engineering*, pages 1–24. Springer.
- Jaynes, E. T. (2003). *Probability theory: The logic of science*. Cambridge university press.
- Jeffrey, R. (1965). *The logic of decision*. University of Chicago press.
- Jeffrey, R. (2004). *Subjective probability: The real thing*. Cambridge University Press.
- 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.

## Références IX

- Kahneman, D. and Tversky, A. (1972). Subjective probability: A judgment of representativeness. *Cognitive psychology*, 3(3):430–454.
- Karni, E., Schmeidler, D., and Vind, K. (1983). On state dependent preferences and subjective probabilities. *Econometrica*, pages 1021–1031.
- Karvetski, C. W., Olson, K. C., Mandel, D. R., and Twardy, C. R. (2013). Probabilistic coherence weighting for optimizing expert forecasts. *Decision Analysis*, 10(4):305–326.
- Kause, A., Bruine de Bruin, W., Persson, J., Thorén, H., Olsson, L., Wallin, A., Dessai, S., and Vareman, N. (2022). Confidence levels and likelihood terms in ipcc reports: a survey of experts from different scientific disciplines. *Climatic Change*, 173(1):1–18.
- Kemeny, J. G. (1955). Fair bets and inductive probabilities1. *The Journal of Symbolic Logic*, 20(3):263–273.
- Kemp, C. and Tenenbaum, J. B. (2008). The discovery of structural form. *Proceedings of the National Academy of Sciences*, 105(31):10687–10692.
- Kemp, C., Tenenbaum, J. B., Niyogi, S., and Griffiths, T. L. (2010). A probabilistic model of theory formation. *Cognition*, 114(2):165–196.
- Klugman, S. A. (1991). *Bayesian statistics in actuarial science: with emphasis on credibility*, volume 15. Springer Science & Business Media.

## Références X

- Kolmogorov, A. (1933). Grundbegriffe der wahrscheinlichkeitsrechnung.
- Kremer, W. (2014). Do doctors understand test results. *BBC World Service*.
- Kruschke, J. (2014). *Doing Bayesian data analysis: A tutorial with R, JAGS, and Stan*. Academic Press.
- Kuhn, M., Johnson, K., et al. (2013). *Applied predictive modeling*, volume 26. Springer.
- Laplace, P. S. (1774). Mémoire sur la probabilité de causes par les événements. *Mémoire de l'académie royale des sciences*.
- Lauritzen, S. L. and Spiegelhalter, D. J. (1988). Local computations with probabilities on graphical structures and their application to expert systems. *Journal of the Royal Statistical Society: Series B (Methodological)*, 50(2):157–194.
- Lebesgue, H. (1918). Remarques sur les théories de la mesure et de l'intégration. *Annales scientifiques de l'École Normale Supérieure*, 35:191–250.
- Lehman, R. S. (1955). On confirmation and rational betting. *The Journal of Symbolic Logic*, 20(3):251–262.
- Lemaire, J. (1995). *Bonus-malus systems in automobile insurance*, volume 19. Springer science & business media.

## Références XI

- Lichtenstein, S., Fischhoff, B., and Phillips, L. D. (1977). Calibration of probabilities: The state of the art. *Decision making and change in human affairs*, pages 275–324.
- Lindley, D. V. (2013). *Understanding uncertainty*. John Wiley & Sons.
- Lindley, D. V., Tversky, A., and Brown, R. V. (1979). On the reconciliation of probability assessments. *Journal of the Royal Statistical Society: Series A (General)*, 142(2):146–162.
- Longley-Cook, L. H. (1962). An introduction to credibility theory. Casualty Actuarial Society.
- MacKay, D. J. (1992). A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472.
- Martin, T. (2009). La probabilité, un concept pluriel. *Pour la science*, (385):46–50.
- Mastrandrea, M. D., Field, C. B., Stocker, T. F., Edenhofer, O., Ebi, K. L., Frame, D. J., Held, H., Kriegler, E., Mach, K. J., Matschoss, P. R., et al. (2010). Guidance note for lead authors of the ipcc fifth assessment report on consistent treatment of uncertainties.
- McGrayne, S. B. (2011). *The Theory That Would Not Die: How Bayes' Rule Cracked the Enigma Code, Hunted Down Russian Submarines, & Emerged Triumphant from Two Centuries of C.* Yale University Press.
- Merrick, J. R. (2008). Getting the right mix of experts. *Decision Analysis*, 5(1):43–52.

## Références XII

- Ministère de l'intérieur (2019). Procédure de reconnaissance de l'état de catastrophe naturelle - révision des critères permettant de caractériser l'intensité des épisodes de sécheresses-réhydrations des sols à l'origine des mouvements de terrains différentiels. Technical report.
- Mongin, P. (1995). Consistent bayesian aggregation. *Journal of Economic Theory*, 66(2):313–351.
- Mongin, P. (2001). The paradox of the bayesian experts. In *Foundations of Bayesianism*, pages 309–338. Springer.
- Moreno-Bote, R., Knill, D. C., and Pouget, A. (2011). Bayesian sampling in visual perception. *Proceedings of the National Academy of Sciences*, 108(30):12491–12496.
- Murawaki, Y. (2019). Bayesian Learning of Latent Representations of Language Structures. *Computational Linguistics*, 45(2):199–228.
- Murphy, A. H. and Epstein, E. S. (1967). Verification of probabilistic predictions: A brief review. *Journal of Applied Meteorology and Climatology*, 6(5):748–755.
- Neal, R. M. (1992). Bayesian training of backpropagation networks by the hybrid monte carlo method. Technical report, Citeseer.
- Neal, R. M. (2012). *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media.

## Références XIII

- Neyman, J. (1977). Frequentist probability and frequentist statistics. *Synthese*, pages 97–131.
- Niculescu-Mizil, A. and Caruana, R. (2005). Predicting good probabilities with supervised learning. In *Proceedings of the 22nd international conference on Machine learning*, pages 625–632.
- Oakes, D. (1985). Self-calibrating priors do not exist. *Journal of the American Statistical Association*, 80(390):339–339.
- Orbanz, P. and Teh, Y. W. (2010). Bayesian nonparametric models. *Encyclopedia of machine learning*, 1.
- Park, T. and Casella, G. (2008). The bayesian lasso. *Journal of the American Statistical Association*, 103(482):681–686.
- Pettigrew, R. (2016). *Accuracy and the Laws of Credence*. Oxford University Press.
- Pettigrew, R. (2020). *Dutch book arguments*. Cambridge University Press.
- Pherson, K. H. and Pherson, R. H. (2012). *Critical thinking for strategic intelligence*. CQ Press.
- Platt, J. et al. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3):61–74.
- Pollock, W. T. and Chapanis, A. (1952). The apparent length of a line as a function of its inclination. *Quarterly Journal of Experimental Psychology*, 4(4):170–178.

## Références XIV

- Pólya, G. (1958). *Les mathématiques et le raisonnement plausible*. Paris, Gauthier-Villars.
- Popper, K. R. (1955). Two autonomous axiom systems for the calculus of probabilities. *The British Journal for the Philosophy of Science*, 6(21):51–57.
- Popper, K. R. (1959). The propensity interpretation of probability. *The British journal for the philosophy of science*, 10(37):25–42.
- Purves, D. (2009). Vision. *Handbook of neuroscience for the behavioral sciences*.
- Purves, D., Wojtach, W. T., and Howe, C. (2008). Visual illusions: an empirical explanation. *Scholarpedia*, 3(6):3706.
- Purves, D., Wojtach, W. T., and Lotto, R. B. (2011). Understanding vision in wholly empirical terms. *Proceedings of the National Academy of Sciences*, 108(supplement\_3):15588–15595.
- Ramachandran, V. S. (1988). Perceiving shape from shading. *Scientific American*, 259(2):76–83.
- Ramsey, F. P. (1926). *Truth and probability*. Cambridge University Press.
- Roberts, H. V. (1968). On the meaning of the probability of rain. In *first national conference on statistical meteorology*.
- Rougier, J. (2007). Probabilistic inference for future climate using an ensemble of climate model evaluations. *Climatic Change*, 81(3):247–264.

## Références XV

- Rougier, J. and Crucifix, M. (2018). Uncertainty in climate science and climate policy. In *Climate Modelling*, pages 361–380. Springer.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1985). Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *nature*, 323(6088):533–536.
- Sadegh-Zadeh, K. (1980). Bayesian diagnostics: A bibliography part 1. *Metamedicine*, 1(1):107–124.
- Saini, A. (2011). A formula for justice. *The Guardian*, October 2nd.
- 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. (1972). *The foundations of statistics*. Courier Corporation.
- Shepard, R. N. (1987). Toward a universal law of generalization for psychological science. *Science*, 237(4820):1317–1323.
- Shepard, R. N. (1992). *L’Oeil qui pense: visions, illusions, perceptions*. Editions du Seuil.
- Shipley, W. C., Nann, B. M., and Penfield, M. J. (1949). The apparent length of tilted lines. *Journal of experimental psychology*, 39(4):548.

## Références XVI

- Silver, N. (2012). *The signal and the noise: Why so many predictions fail—but some don't*. Penguin.
- Skyrms, B. (1987). Dynamic coherence and probability kinematics. *Philosophy of science*, 54(1):1–20.
- Spiegelhalter, D. J., Dawid, A. P., Lauritzen, S. L., and Cowell, R. G. (1993). Bayesian analysis in expert systems. *Statistical science*, pages 219–247.
- Stoerk, T., Wagner, G., and Ward, R. E. (2020). Policy brief – recommendations for improving the treatment of risk and uncertainty in economic estimates of climate impacts in the sixth intergovernmental panel on climate change assessment report. *Review of Environmental Economics and Policy*.
- Stolcke, A. (1994). *Bayesian learning of probabilistic language models*. University of California, Berkeley.
- Surowiecki, J. (2005). *The wisdom of crowds*. Anchor.
- Tebaldi, C., Smith, R. L., Nychka, D., and Mearns, L. O. (2005). Quantifying uncertainty in projections of regional climate change: A bayesian approach to the analysis of multimodel ensembles. *Journal of Climate*, 18(10):1524–1540.
- Teller, P. (1973). Conditionalization and observation. *Synthese*, 26(2):218–258.
- Tenenbaum, J. (1998). Bayesian modeling of human concept learning. *Advances in neural information processing systems*, 11.

## Références XVII

- Tenenbaum, J. B., Kemp, C., Griffiths, T. L., and Goodman, N. D. (2011). How to grow a mind: Statistics, structure, and abstraction. *science*, 331(6022):1279–1285.
- Theodoridis, S. (2015). *Machine learning: a Bayesian and optimization perspective*. Academic press.
- Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294.
- Thompson, W. R. (1935). On the theory of apportionment. *American Journal of Mathematics*, 57(2):450–456.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288.
- Titelbaum, M. G. (2022a). *Fundamentals of Bayesian Epistemology 1: Introducing Credences*. Oxford University Press.
- Titelbaum, M. G. (2022b). *Fundamentals of Bayesian Epistemology 2: Arguments, Challenges, Alternatives*. Oxford University Press.
- Van Calster, B., McLernon, D. J., Van Smeden, M., Wynants, L., and Steyerberg, E. W. (2019). Calibration: the achilles heel of predictive analytics. *BMC medicine*, 17(1):1–7.

## Références XVIII

- van de Schoot, R., Depaoli, S., King, R., Kramer, B., Märkens, K., Tadesse, M. G., Vannucci, M., Gelman, A., Veen, D., Willemsen, J., et al. (2021). Bayesian statistics and modelling. *Nature Reviews Methods Primers*, 1(1):1–26.
- Vogel, H., Appelbaum, S., Haller, H., and Ostermann, T. (2022). The interpretation of verbal probabilities: A systematic literature review and meta-analysis. *German Medical Data Sciences 2022–Future Medicine: More Precise, More Integrative, More Sustainable!*, pages 9–16.
- Von Helmholtz, H. (1867). *Handbuch der physiologischen Optik: mit 213 in den Text eingedruckten Holzschnitten und 11 Tafeln*, volume 9. Voss.
- von Mises, R. (1928). *Wahrscheinlichkeit Statistik und Wahrheit*. Springer.
- von Mises, R. (1939). *Probability, statistics and truth*. Macmillan.
- Vul, E. and Pashler, H. (2008). Measuring the crowd within: Probabilistic representations within individuals. *Psychological Science*, 19(7):645–647.
- Watanabe, S. and Chien, J.-T. (2015). *Bayesian speech and language processing*. Cambridge University Press.
- Whitney, A. W. (1918). Theory of experience rating. *Proceedings of the Casualty Actuarial Society*, 4.
- Williamson, J. (2004). *Bayesian nets and causality: philosophical and computational foundations*. Oxford University Press.

## Références XIX

- Yuille, A. and Kersten, D. (2006). Vision as bayesian inference: analysis by synthesis? *Trends in cognitive sciences*, 10(7):301–308.
- Zadrozny, B. and Elkan, C. (2001). Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. In *ICML*, volume 1, pages 609–616. Citeseer.
- Zadrozny, B. and Elkan, C. (2002). Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 694–699.