

# Bayesian Methods for Biomedical Research — Part I: Bayesian theory

Boris Hejblum

<https://bayesee.borishejblum.science>

ISPED summer school  
at the University of Bordeaux  
June 3<sup>rd</sup>, 2024

# Course Presentation

# Introduce yourself



<https://www.menti.com/14uj38ttuj>

# Introduce yourself



<https://www.menti.com/14uj38ttuj>



<https://www.menti.com/8ztu8q7ke1>

# Bayesian vocabulary


- **paradigm**
- *a priori*
- *a posteriori*
- **elicitation**

# Course objectives

## I Familiarize oneself with the **Bayesian framework**:

- 1 understand and assess a Bayesian modeling strategy, and discuss its underlying assumptions
- 2 rigorously describe expert knowledge by a quantitative prior distribution

## II Study and perform Bayesian analyses in **biomedical applications**:

- 1 understand, discuss and reproduce a Bayesian (re-)estimation of a Relative Risk
- 2 perform a Bayesian regression using , applied to meta-analysis
- 3 put into perspective the results from a Bayesian analysis described in a scientific article

**NB** : this course is by no means exhaustive, and the curious reader will be referred to more complete works such as *The Bayesian Choice* by C Robert.

# Disclaimer

**Audience** is often **diverse**:

Students with *different backgrounds & different expertise* will get a **different experience** of this class

Some parts can feel *hard, frustrating* or even *not very relevant to you*.

**My goal:** *everyone* finds interesting ideas, concept and tools to learn.

For some, the important focus will be the *medical applications*, for others it will be the *programming tools*, or the new *philosophical framework*, or the *statistical tools*...

OK to feel a bit lost at first  
Things should make more sense as we progress !  
⇒ **Ask questions !**

# Motivational examples: diagnostic tests



## The obscure maths theorem that governs the reliability of Covid testing

There's been much debate about lateral flow tests - their accuracy depends on context and the theories of a 18th-century cleric

[Good, *J GEN INTERN MED* 2020]

Table 1 Estimates for Post-Test Probability of Acute COVID-19 Infection for Simulated Patient Scenarios

Clinical Scenarios	Pre-test probability (%)	PCR assay sensitivity (%)	Post-test probability of acute COVID-19 infection	
			Positive test (%)	Negative test (%)
Patient 1: high pre-test probability	70	70	100	41.2
		90	100	18.9
		70	100	73.0
Patient 2: low pre-test probability	5	90	100	47.4
		70	97.4	1.6
		90	97.9	0.5
	10	70	98.7	3.2
		90	99.0	1.1



Original Article

## Bayesian analysis of tests with unknown specificity and sensitivity

Andrew Gelman , Bob Carpenter

First published: 13 August 2020 | <https://doi.org/10.1111/rssc.12435> | Citations: 6



# Motivational examples: clinical trial design

Design

---

## Anti-Thrombotic Therapy to Ameliorate Complications of COVID-19 (ATTACC): Study design and methodology for an international, adaptive Bayesian randomized controlled trial

**Methods:** An international, open-label, adaptive randomized controlled trial. Using a Bayesian framework, the trial will declare results as soon as pre-specified posterior probabilities for superiority, futility, or harm are reached. The trial uses response-adaptive randomization to maximize the probability that patients will receive the more beneficial treatment approach, as treatment effect information accumulates within the trial. By leveraging a common data safety monitoring

[Houston *et al.*, *Clinical Trials*, 17(5):491-500, 2020]

**CLINICAL  
TRIALS**

*Clinical Trials*

1-10

© The Author(s) 2020

Article reuse guidelines:

[sagepub.com/journals-permissions](http://sagepub.com/journals-permissions)

DOI: 10.1177/1740774520943846

[journals.sagepub.com/home/ctj](http://journals.sagepub.com/home/ctj)

 SAGE

# Motivational examples: study/trial analyses

## The NEW ENGLAND JOURNAL of MEDICINE

ESTABLISHED IN 1812

APRIL 22, 2021

VOL. 384 NO. 16

### Interleukin-6 Receptor Antagonists in Critically Ill Patients with Covid-19

The REMAP-CAP Investigators\*

lumab group, and 0 (interquartile range, -1 to 15) in the control group. The median adjusted cumulative odds ratios were 1.64 (95% credible interval, 1.25 to 2.14) for tocilizumab and 1.76 (95% credible interval, 1.17 to 2.91) for sarilumab as compared with control, yielding posterior probabilities of superiority to control of more than 99.9% and of 99.5%, respectively. An analysis of 90-day survival showed improved survival in the pooled interleukin-6 receptor antagonist groups, yielding a hazard ratio for the comparison with the control group of 1.61 (95% credible interval, 1.25 to 2.08) and a posterior probability of superiority of more than 99.9%. All secondary analyses supported efficacy of these interleukin-6 receptor antagonists.

## ORIGINAL

### Dexamethasone 12 mg versus 6 mg for patients with COVID-19 and severe hypoxaemia: a pre-planned, secondary Bayesian analysis of the COVID STEROID 2 trial

## Abstract

**Purpose:** We compared dexamethasone 12 versus 6 mg daily for up to 10 days in patients with coronavirus disease 2019 (COVID-19) and severe hypoxaemia in the international, randomized, blinded COVID STEROID 2 trial. In the primary, conventional analyses, the predefined statistical significance thresholds were not reached. We conducted a pre-planned Bayesian analysis to facilitate probabilistic interpretation.

**Methods:** We analysed outcome data within 90 days in the intention-to-treat population (data available in 967 to 982 patients) using Bayesian models with various sensitivity analyses. Results are presented as median posterior probabilities with 95% credible intervals (CrIs) and probabilities of different effect sizes with 12 mg dexamethasone.

**Results:** The adjusted mean difference on days alive without life support at day 28 (primary outcome) was 1.3 days (95% CrI -0.3 to 2.9; 94.2% probability of benefit). Adjusted relative risks and probabilities of benefit on serious adverse reactions were 0.85 (0.63 to 1.16; 84.1%) and on mortality 0.87 (0.73 to 1.03; 94.8%) at day 28 and 0.88 (0.75 to 1.02; 95.1%) at day 90. Probabilities of benefit on days alive without life support and days alive out of hospital at day 90 were 85 and 95.7%, respectively. Results were largely consistent across sensitivity analyses, with relatively low probabilities of clinically important harm with 12 mg on all outcomes in all analyses.

## The NEW ENGLAND JOURNAL of MEDICINE

ESTABLISHED IN 1812

DECEMBER 31, 2020

VOL. 383 NO. 27

### Safety and Efficacy of the BNT162b2 mRNA Covid-19 Vaccine

Fernando P. Polack, M.D., Stephen J. Thomas, M.D., Nicholas Kitchin, M.D., Judith Absalon, M.D., Alejandra Gurtman, M.D., Stephen Lockhart, D.M., John L. Perez, M.D., Gonzalo Pérez Marc, M.D., Edson D. Moreira, M.D., Cristiano Zerbini, M.D., Ruth Bailey, B.Sc., Kena A. Swanson, Ph.D., Satrajit Roychowdhury, Ph.D., Kenneth Koury, Ph.D., Ping Li, Ph.D., Warren V. Kalina, Ph.D., David Cooper, Ph.D., Robert W. Frencik, Jr., M.D., Laura L. Hammit, M.D., Özlem Türeci, M.D., Haylene Nell, M.D., Auel Schaefer, M.D., Serhat Uenal, M.D., Dina B. Tresnan, D.V.M., Ph.D., Susan Mather, M.D., Philip R. Dormitzer, M.D., Ph.D., Ugur Sahin, M.D., Kathrin U. Jansen, Ph.D., and William C. Gruber, M.D., for the C4591001 Clinical Trial Group†

Table 2. Vaccine Efficacy against Covid-19 at Least 7 days after the Second Dose.\*

Efficacy End Point	BNT162b2		Placebo		Vaccine Efficacy, % (95% Credible Interval)‡	Posterior Probability (Vaccine Efficacy >30%)§
	No. of Cases	Surveillance Time (y)†	No. of Cases	Surveillance Time (y)†		
Covid-19 occurrence at least 7 days after the second dose in participants without evidence of infection	(N=18,198)		(N=18,225)		95.0 (90.3–97.6)	>0.9999
	8	2,214 (17,411)	162	2,222 (17,511)		
Covid-19 occurrence at least 7 days after the second dose in participants with and those without evidence of infection	(N=19,965)		(N=20,172)		94.6 (89.9–97.3)	>0.9999
	9	2,332 (18,559)	169	2,345 (18,708)		

\* The total population without baseline infection was 36,521; total population including those with and those without prior evidence of infection was 40,137.

† The surveillance time is the total time in 1000 person-years for the given end point across all participants within each group at risk for the end point. The time period for Covid-19 case accrual is from 7 days after the second dose to the end of the surveillance period.

‡ The credible interval for vaccine efficacy was calculated with the use of a beta-binomial model with prior beta (0.700102, 1) adjusted for the surveillance time.

§ Posterior probability was calculated with the use of a beta-binomial model with prior beta (0.700102, 1) adjusted for the surveillance time.

# Introduction

## Statistics:

- a **mathematical** science
- to **describe** what has happened and
- to assess what **may** happen in **the future**
- relies on the **observation** of natural phenomena in order to propose an interpretation, often through **probabilistic models**

## Statistics:

- a **mathematical** science
- to **describe** what has happened and
- to assess what **may** happen in **the future**
- relies on the **observation** of natural phenomena in order to propose an interpretation, often through **probabilistic models**

## Frequentist statistics:

- Neyman & Pearson
- **deterministic** view of the parameters
- **Maximum Likelihood Estimation**
- statistical **test theory** & **confidence interval**



# Bayes' theorem

Reverend Thomas Bayes posthumous article in 1763



$$\Pr(A|E) = \frac{\Pr(E|A) \Pr(A)}{\Pr(E|A) \Pr(A) + \Pr(E|\bar{A}) \Pr(\bar{A})} = \frac{\Pr(E|A) \Pr(A)}{\Pr(E)}$$

(conditional probability formula:  $\Pr(A|E) = \frac{\Pr(A \cap E)}{\Pr(E)}$ )

# Bayes' theorem

Reverend Thomas Bayes posthumous article in 1763



$$\Pr(A|E) = \frac{\Pr(E|A) \Pr(A)}{\Pr(E|A) \Pr(A) + \Pr(E|\bar{A}) \Pr(\bar{A})} = \frac{\Pr(E|A) \Pr(A)}{\Pr(E)}$$

(conditional probability formula:  $\Pr(A|E) = \frac{\Pr(A \cap E)}{\Pr(E)}$ )

## In practice:

Last time you visited the doctor, you got **tested for a rare disease**. Unluckily, the result was positive. . .

*Given the test result, what is the probability that I actually have this disease?*

(Medical tests are, after all, not perfectly accurate.)

→ *Seeing Theory*, Brown University

# Bayes theorem: exercise

In June 2022, about 0.33% of the French population was estimated to have COVID-19.

Rapid tests have the following statistical properties:

- if someone has COVID-19, its test will come out positive 71% of the time
- if someone does not have the disease, its test will come out negative 98% of the time

*Given that someone got a positive result, what is his/her probability to truly have COVID-19 ?*



# Bayes theorem: exercise

In June 2022, about 0.33% of the French population was estimated to have COVID-19.

Rapid tests have the following statistical properties:

- if someone has COVID-19, its test will come out positive 71% of the time
- if someone does not have the disease, its test will come out negative 98% of the time

*Given that someone got a positive result, what is his/her probability to truly have COVID-19 ?*

$$\Pr(D=+) = 0.0033$$

$$\Pr(T=+|D=+) = 0.71$$

$$\Pr(T=-|D=-) = 0.98$$

# Bayes theorem: exercise

In June 2022, about 0.33% of the French population was estimated to have COVID-19.

Rapid tests have the following statistical properties:

- if someone has COVID-19, its test will come out positive 71% of the time
- if someone does not have the disease, its test will come out negative 98% of the time

*Given that someone got a positive result, what is his/her probability to truly have COVID-19 ?*

$$\Pr(D=+) = 0.0033$$

$$\Pr(T=+|D=+) = 0.71$$

$$\Pr(T=-|D=-) = 0.98$$

$$\Pr(D=+|T=+) = ?$$

# Bayes theorem: exercise

In June 2022, about 0.33% of the French population was estimated to have COVID-19.

Rapid tests have the following statistical properties:

- if someone has COVID-19, its test will come out positive 71% of the time
- if someone does not have the disease, its test will come out negative 98% of the time

*Given that someone got a positive result, what is his/her probability to truly have COVID-19 ?*

$$\Pr(D=+) = 0.0033 \quad \Pr(T=+|D=+) = 0.71 \quad \Pr(T=-|D=-) = 0.98$$

$$\begin{aligned} \Pr(D=+|T=+) &= \frac{\Pr(T=+|D=+)\Pr(D=+)}{\Pr(T=+)} \\ &= \frac{\Pr(T=+|D=+)\Pr(D=+)}{\Pr(T=+|D=+)\Pr(D=+) + \Pr(T=+|D=-)\Pr(D=-)} \\ &= \frac{\Pr(T=+|D=+)\Pr(D=+)}{\Pr(T=+|D=+)\Pr(D=+) + (1 - \Pr(T=-|D=-))(1 - \Pr(D=+))} \\ &= (0.71 \times 0.0033) / (0.71 \times 0.0033 + (1 - 0.98) \times (1 - 0.0033)) = 11\% \end{aligned}$$

# Continuous Bayes' theorem

- parametric (probabilistic) model  $f(y|\theta)$
- parameters  $\theta$
- probability distribution  $\pi$

Continuous Bayes' theorem:

$$p(\theta|y) = \frac{f(y|\theta)\pi(\theta)}{\int f(y|\theta)\pi(\theta) d\theta} = \frac{f(y|\theta)\pi(\theta)}{f(y)}$$

# Continuous Bayes' theorem

- parametric (probabilistic) model  $f(y|\theta)$
- parameters  $\theta$
- probability distribution  $\pi$

Continuous Bayes' theorem:

$$p(\theta|y) = \frac{f(y|\theta)\pi(\theta)}{\int f(y|\theta)\pi(\theta) d\theta} = \frac{f(y|\theta)\pi(\theta)}{f(y)}$$



Pierre-Simon de Laplace !

# Bayes philosophy

**Parameters are random variables !** – *no “true” value*

⇒ induces a marginal probability distribution  $\pi(\theta)$  on the parameters:  
the **prior** distribution

😊 allows to **formally** take into account hypotheses in the modeling

😞 necessarily introduces **subjectivity** into the analysis

# Bayesian vs. Frequentists: a historical note

- 1 **Bayes + Laplace**  $\Rightarrow$  development of statistics in the **18-19<sup>th</sup> centuries**
- 2 Galton & Pearson, then Fisher & Neymann  $\Rightarrow$  **frequentist** theory became dominant during the **20<sup>th</sup> century**
- 3 turn of the **21<sup>th</sup> century**: rise of the computer  $\Rightarrow$  **Bayes' comeback**



# Bayesian vs. Frequentists: an outdated debate

Fisher firmly rejected Bayesian reasoning

⇒ community split in 2 in the 20<sup>th</sup> century



# Bayesian vs. Frequentists: an outdated debate

Fisher firmly rejected Bayesian reasoning

⇒ community split in 2 in the 20<sup>th</sup> century

*To be, or not to be, Bayesian, that is no longer the question: it is a matter of wisely using the right tools when necessary*

Gilbert Saporta

# Bayesian modeling

## Refresher on frequentist modeling

- a series of *iid* (independent and identically distributed) random variables  $\mathbf{Y} = (Y_1, \dots, Y_n)$

## Refresher on frequentist modeling

- a series of *iid* (independent and identically distributed) random variables  $\mathbf{Y} = (Y_1, \dots, Y_n)$
- we observe a sample  $\mathbf{y} = (y_1, \dots, y_n)$



## Refresher on frequentist modeling

- a series of *iid* (independent and identically distributed) random variables  $\mathbf{Y} = (Y_1, \dots, Y_n)$
- we observe a sample  $\mathbf{y} = (y_1, \dots, y_n)$
- model their probability distribution as  $f(y|\theta)$ ,  $\theta \in \Theta$

This model assumes there is a “true” distribution of  $Y$  characterized by the “true” value of the parameter  $\theta^*$

$\hat{\theta} ?$

# Historical motivating example

## Laplace

What is the probability of birth of girls rather than boys ?

⇒ **observations:** births observed in Paris between 1745 and 1770  
(241,945 girls & 251,527 boys)

When a child is born, is it equally likely to be a girl or a boy ?

# Three building blocks

- 1 the question
- 2 the sampling model
- 3 the prior



# Three building blocks

## 1 the question

The first step in building a model is always to identify the question you want to answer

## 2 the sampling model

## 3 the prior

# Three building blocks

## 1 the question

The first step in building a model is always to identify the question you want to answer

## 2 the sampling model

Which **observations** are available to inform our response to this ?  
How can they be **described**?

## 3 the prior

# Three building blocks

## 1 the question

The first step in building a model is always to identify the question you want to answer

## 2 the sampling model

Which **observations** are available to inform our response to this ?  
How can they be **described**?

## 3 the prior

A probability distribution on the parameters  $\theta$  of the sampling model

# The sampling model

$\mathbf{y}$ : the observations available

⇒ (parametric) **probabilistic model** underlying their **generation**:

$$Y_i \stackrel{iid}{\sim} f(y|\theta)$$

# The *prior* distribution

In Bayesian modeling, compared to frequentist modeling, we add a **probability distribution** on the **parameters**  $\theta$

$$\theta \sim \pi(\theta)$$

$$Y_i|\theta \stackrel{iid}{\sim} f(y|\theta)$$

$\theta$  will thus be treated like a random variable,  
but which is never observed !

# Back to Laplace's historical example

- 1 The question
- 2 Sampling model
- 3 *prior*

# Back to Laplace's historical example

- 1 The question  
...
- 2 Sampling model  
...
- 3 *prior*  
...

# Back to Laplace's historical example

## 1 The question

When a child is born, is it equally likely to be a girl or a boy ?

## 2 Sampling model

...

## 3 *prior*

...



# Back to Laplace's historical example

## 1 The question

When a child is born, is it equally likely to be a girl or a boy ?

## 2 Sampling model

Bernoulli's law for  $Y_i = 1$  if the new born  $i$  is a girl, 0 if it is a boy:

$$Y_i \sim \text{Bernoulli}(\theta) \quad \theta \in [0, 1]$$

## 3 *prior*

...

# Back to Laplace's historical example

## 1 The question

When a child is born, is it equally likely to be a girl or a boy ?

## 2 Sampling model

Bernoulli's law for  $Y_i = 1$  if the new born  $i$  is a girl, 0 if it is a boy:

$$Y_i \sim \text{Bernoulli}(\theta) \quad \theta \in [0, 1]$$

## 3 *prior*

A uniform prior on  $\theta$  (the probability that a newborn would be a girl rather than a boy):

$$\theta \sim \mathcal{U}_{[0,1]}$$

# Posterior distribution

Purpose of a Bayesian modeling: **infer the *posterior*** distribution of the **parameters**

- **Posterior**: the law of  $\theta$  conditionally on the observations  $p(\theta|\mathbf{y})$

# Posterior distribution

Purpose of a Bayesian modeling: **infer the *posterior*** distribution of the **parameters**

- **Posterior**: the law of  $\theta$  conditionally on the observations  $p(\theta|\mathbf{y})$

Bayes' theorem:

$$p(\theta|\mathbf{y}) = \frac{f(\mathbf{y}|\theta)\pi(\theta)}{f(\mathbf{y})}$$

# Posterior distribution

Purpose of a Bayesian modeling: **infer the *posterior*** distribution of the **parameters**

- **Posterior**: the law of  $\theta$  conditionally on the observations  $p(\theta|\mathbf{y})$

Bayes' theorem:

$$p(\theta|\mathbf{y}) = \frac{f(\mathbf{y}|\theta)\pi(\theta)}{f(\mathbf{y})}$$

Posterior is calculated from:

- 1 the sampling model  $f(\mathbf{y}|\theta)$  – which yields the likelihood  $f(\mathbf{y}|\theta)$  for all observations
- 2 the *prior*  $\pi(\theta)$

# Application to the historical example

- 1 the likelihood
- 2 the prior
- 3 the posterior

# Application to the historical example

## ① the likelihood

...

## ② the prior

...

## ③ the posterior

...

# Application to the historical example

## 1 the likelihood

$$f(\mathbf{y}|\theta) = \prod_{i=1}^n \theta^{y_i} (1-\theta)^{(1-y_i)} = \theta^S (1-\theta)^{n-S} \quad \text{where } S = \sum_{i=1}^n y_i$$

## 2 the prior

...

## 3 the posterior

...



# Application to the historical example

## 1 the likelihood

$$f(\mathbf{y}|\theta) = \prod_{i=1}^n \theta^{y_i} (1-\theta)^{(1-y_i)} = \theta^S (1-\theta)^{n-S} \quad \text{where } S = \sum_{i=1}^n y_i$$

## 2 the prior

Uniform:  $\pi(\theta) = 1$

## 3 the posterior

...

# Application to the historical example

## 1 the likelihood

$$f(\mathbf{y}|\theta) = \prod_{i=1}^n \theta^{y_i} (1-\theta)^{(1-y_i)} = \theta^S (1-\theta)^{n-S} \quad \text{where } S = \sum_{i=1}^n y_i$$

## 2 the prior

Uniform:  $\pi(\theta) = 1$

## 3 the posterior

$$p(\theta|\mathbf{y}) = \frac{\theta^S (1-\theta)^{n-S}}{f(\mathbf{y})} = p(\theta|\mathbf{y}) = \binom{n}{S} (n+1) \theta^S (1-\theta)^{n-S}$$

# Application to the historical example

## 1 the likelihood

$$f(\mathbf{y}|\theta) = \prod_{i=1}^n \theta^{y_i} (1-\theta)^{(1-y_i)} = \theta^S (1-\theta)^{n-S} \quad \text{where } S = \sum_{i=1}^n y_i$$

## 2 the prior

Uniform:  $\pi(\theta) = 1$

## 3 the posterior

$$p(\theta|\mathbf{y}) = \frac{\theta^S (1-\theta)^{n-S}}{f(\mathbf{y})} = p(\theta|\mathbf{y}) = \binom{n}{S} (n+1) \theta^S (1-\theta)^{n-S}$$

To answer the question of interest, we can then compute: ...

# Application to the historical example

## 1 the likelihood

$$f(\mathbf{y}|\theta) = \prod_{i=1}^n \theta^{y_i} (1-\theta)^{(1-y_i)} = \theta^S (1-\theta)^{n-S} \quad \text{where } S = \sum_{i=1}^n y_i$$

## 2 the prior

Uniform:  $\pi(\theta) = 1$

## 3 the posterior

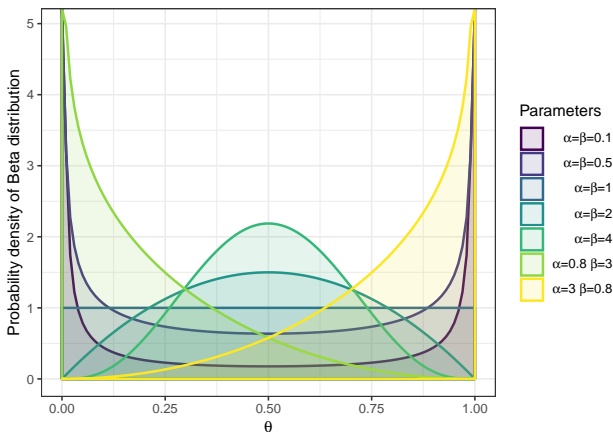
$$p(\theta|\mathbf{y}) = \frac{\theta^S (1-\theta)^{n-S}}{f(\mathbf{y})} = p(\theta|\mathbf{y}) = \binom{n}{S} (n+1) \theta^S (1-\theta)^{n-S}$$

To answer the question of interest, we can then compute:

$$P(\theta \geq 0.5|\mathbf{y}) = \int_{0.5}^1 p(\theta|\mathbf{y}) = \binom{n}{S} (n+1) \int_{0.5}^1 \theta^S (1-\theta)^{n-S} d\theta \approx 1.15 \cdot 10^{-42}$$

# The Beta distribution

$$f(\theta) = \frac{(\alpha + \beta - 1)!}{(\alpha - 1)! (\beta - 1)!} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \text{ for } \alpha > 0 \text{ and } \beta > 0$$



Examples of various parametrizations for the Beta distribution

**Construction of a Bayesian model**

# Conjugacy of the Beta distribution

**Beta *prior*:**  $\pi = \text{Beta}(\alpha, \beta)$

# Conjugacy of the Beta distribution

**Beta prior:**  $\pi = \text{Beta}(\alpha, \beta)$

**Corresponding posterior:**  $p(\theta|\mathbf{y}) \propto \theta^{\alpha+S-1}(1-\theta)^{\beta+(n-S)-1}$

...

The  $\propto$  symbol means: “proportional to”



# Conjugacy of the Beta distribution

**Beta prior:**  $\pi = \text{Beta}(\alpha, \beta)$

**Corresponding posterior:**  $p(\theta|\mathbf{y}) \propto \theta^{\alpha+S-1}(1-\theta)^{\beta+(n-S)-1}$

$\Rightarrow \theta|\mathbf{y} \sim \text{Beta}(\alpha + S, \beta + (n - S))$

The  $\propto$  symbol means: “proportional to”

# Conjugacy of the Beta distribution

**Beta prior:**  $\pi = \text{Beta}(\alpha, \beta)$

**Corresponding posterior:**  $p(\theta|\mathbf{y}) \propto \theta^{\alpha+S-1} (1-\theta)^{\beta+(n-S)-1}$

$\Rightarrow \theta|\mathbf{y} \sim \text{Beta}(\alpha + S, \beta + (n - S))$

This is called a **conjugated distribution** because the **posterior** and the **prior** belong to the **same parametric family**

The  $\propto$  symbol means: “proportional to”

Impact of the *prior* choice

Interpretation of the <i>prior</i>	Parameters of the Beta distribution	$P(\theta \geq 0.5 \mathbf{y})$
#boys > #girls	$\alpha = 0.1, \beta = 3$	$1.08 \cdot 10^{-42}$
#boys < #girls	$\alpha = 3, \beta = 0.1$	$1.19 \cdot 10^{-42}$
#boys = #girls	$\alpha = 4, \beta = 4$	$1.15 \cdot 10^{-42}$
#boys $\neq$ #girls	$\alpha = 0.1, \beta = 0.1$	$1.15 \cdot 10^{-42}$
non-informative	$\alpha = 1, \beta = 1$	$1.15 \cdot 10^{-42}$

For 493,472 newborns including 241,945 girls

Impact of the *prior* choice

Interpretation of the <i>prior</i>	Parameters of the Beta distribution	$P(\theta \geq 0.5 \mathbf{y})$
#boys > #girls	$\alpha = 0.1, \beta = 3$	$1.08 \cdot 10^{-42}$
#boys < #girls	$\alpha = 3, \beta = 0.1$	$1.19 \cdot 10^{-42}$
#boys = #girls	$\alpha = 4, \beta = 4$	$1.15 \cdot 10^{-42}$
#boys $\neq$ #girls	$\alpha = 0.1, \beta = 0.1$	$1.15 \cdot 10^{-42}$
non-informative	$\alpha = 1, \beta = 1$	$1.15 \cdot 10^{-42}$

For 493,472 newborns including 241,945 girls

Interpretation of the <i>prior</i>	Parameters of the Beta distribution	$P(\theta \geq 0.5 \mathbf{y})$
#boys > #girls	$\alpha = 0.1, \beta = 3$	0.39
#boys < #girls	$\alpha = 3, \beta = 0.1$	0.52
#boys = #girls	$\alpha = 4, \beta = 4$	0.46
#boys $\neq$ #girls	$\alpha = 0.1, \beta = 0.1$	0.45
non-informative	$\alpha = 1, \beta = 1$	0.45

For 20 newborns including 9 girls

Construction of a Bayesian model

# Impact of the *prior* choice for 20 observed births

# Priors: pros & cons

Having a *prior* distribution:

- 😊 brings **flexibility**
- 😊 allows to incorporate **external knowledge**
- 😞 adds intrinsic **subjectivity**

⇒ choice (or elicitation) of a *prior* distribution is sensitive !

# Prior properties

- 1 *posterior* support must be included in the support of the *prior*:  
if  $\pi(\theta) = 0$ , then  $p(\theta|\mathbf{y}) = 0$
- 2 independence of the different parameters *a priori*

# Prior Elicitation

**Strategies to communicate** with non-statistical experts

⇒ transform their **knowledge** into *prior distribution*

- **histogram method**: experts give weights to ranges of values  
⚠ might give a zero *prior* for plausible parameter values
- choose a **parametric family** of distributions  $p(\theta|\eta)$  in **agreement with what the experts think** (e.g. for quantiles or moments)  
(solves the support problem but the parametric family has a big impact)
- elicit *priors* from the **literature**
- ...



# SHELF: a tool for prior elicitation from expert knowledge

Your turn !



**Practicals:** exercise 1

# The quest for non-informative *priors*

Sometimes, one has **no prior knowledge whatsoever**  
Which *prior* distribution to use ?



# The quest for non-informative *priors*

Sometimes, one has **no prior knowledge whatsoever**  
⇒ the Uniform distribution, a **non-informative prior** ?

# The quest for non-informative *priors*

Sometimes, one has **no prior knowledge whatsoever**

⇒ the Uniform distribution, a **non-informative prior** ?

2 major difficulties:

- 1 **Improper distributions**  $\int_{\Theta} \pi(\theta) d\theta = \infty$
- 2 **Non-invariant distributions**

# The quest for non-informative *priors*

Sometimes, one has **no prior knowledge whatsoever**

⇒ the Uniform distribution, a **non-informative prior** ?

2 major difficulties:

- 1 **Improper distributions**  $\int_{\Theta} \pi(\theta) d\theta = \infty$
- 2 **Non-invariant distributions**

*Other solutions ?*

# Jeffreys' priors

A **weakly informative prior** invariant through re-parameterization

- unidimensional Jeffreys' prior:

$$\pi(\theta) \propto \sqrt{I(\theta)} \quad \text{where } I \text{ is Fisher's information matrix}$$

- multidimensional Jeffreys' prior:

$$\pi(\theta) \propto \sqrt{|I(\theta)|}$$

In practice, parameter are considered independent *a priori*

# Hyper-priors & hierarchical models

**Hierarchical levels:**

①  $\pi(\theta)$

②  $f(\mathbf{y}|\theta)$

# Hyper-priors & hierarchical models

## Hierarchical levels:

①  $\eta \sim h(\eta)$

②  $\pi(\theta|\eta)$

③  $f(\mathbf{y}|\theta)$



# Hyper-priors & hierarchical models

**Hierarchical levels:**

①  $\eta \sim h(\eta)$

②  $\pi(\theta|\eta)$

③  $f(\mathbf{y}|\theta)$

$$p(\theta|\mathbf{y}) = \frac{f(\mathbf{y}|\theta)\pi(\theta)}{f(\mathbf{y})} = \frac{\int f(\mathbf{y}|\theta, \eta)\pi(\theta|\eta)h(\eta)d\eta}{f(\mathbf{y})}$$

# Hyper-priors & hierarchical models

**Hierarchical levels:**

①  $\eta \sim h(\eta)$

②  $\pi(\theta|\eta)$

③  $f(\mathbf{y}|\theta)$

$$p(\theta|\mathbf{y}) = \frac{f(\mathbf{y}|\theta)\pi(\theta)}{f(\mathbf{y})} = \frac{\int f(\mathbf{y}|\theta, \eta)\pi(\theta|\eta)h(\eta)d\eta}{f(\mathbf{y})} = \frac{f(\mathbf{y}|\theta)\int \pi(\theta|\eta)h(\eta)d\eta}{f(\mathbf{y})}$$

**NB:** 3 hierarchical levels  $\Leftrightarrow$  two levels with *prior*:  $\pi(\theta) = \int \pi(\theta|\eta)h(\eta)d\eta$

# Hyper-priors & hierarchical models

**Hierarchical levels:**

①  $\eta \sim h(\eta)$

②  $\pi(\theta|\eta)$

③  $f(\mathbf{y}|\theta)$

$$p(\theta|\mathbf{y}) = \frac{f(\mathbf{y}|\theta)\pi(\theta)}{f(\mathbf{y})} = \frac{\int f(\mathbf{y}|\theta, \eta)\pi(\theta|\eta)h(\eta)d\eta}{f(\mathbf{y})} = \frac{f(\mathbf{y}|\theta)\int \pi(\theta|\eta)h(\eta)d\eta}{f(\mathbf{y})}$$

**NB:** 3 hierarchical levels  $\Leftrightarrow$  two levels with *prior*:  $\pi(\theta) = \int \pi(\theta|\eta)h(\eta)d\eta$

$\Rightarrow$  can **ease modeling** and **elicitation** of the *prior*...

# Hyperprior in the historical example

Historical example of birth sex with a Beta *prior*

⇒ two Gamma hyper-priors for  $\alpha$  and  $\beta$  (conjugated):

$$\alpha \sim \text{Gamma}(4, 0.5)$$

$$\beta \sim \text{Gamma}(4, 0.5)$$

$$\theta | \alpha, \beta \sim \text{Beta}(\alpha, \beta)$$

$$Y_i | \theta \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$$

# Empirical Bayes

Eliciting the *prior* according to its empirical marginal distribution

⇒ estimate the *prior* from the data

- 1 hyper-parameters
- 2 estimate them through frequentist methods (e.g. MLE) by  $\hat{\eta}$
- 3 plug-in estimates into the *prior*
- 4 ⇒ *posterior*:  $p(\theta|\mathbf{y}, \hat{\eta})$

# Empirical Bayes

Eliciting the *prior* according to its empirical marginal distribution

⇒ estimate the *prior* from the data

- 1 hyper-parameters
- 2 estimate them through frequentist methods (e.g. MLE) by  $\hat{\eta}$
- 3 plug-in estimates into the *prior*
- 4 ⇒ *posterior*:  $p(\theta|\mathbf{y}, \hat{\eta})$

- Combines Bayesian and frequentist frameworks
- Concentrated *posterior*: ↘ variance – but ↗ bias (data used twice ⇒ shrinkage around the average!)
- Approximate a fully Bayesian approach

# Sequential Bayes

Bayes' theorem can be used sequentially:

$$p(\theta|\mathbf{y}) \propto f(\mathbf{y}|\theta)\pi(\theta)$$

If  $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2)$ , then:

$$p(\theta|\mathbf{y}) \propto f(\mathbf{y}_2|\theta)f(\mathbf{y}_1|\theta)\pi(\theta) \propto f(\mathbf{y}_2|\theta)p(\theta|\mathbf{y}_1)$$

⇒ *posterior* distribution updates as new observations are acquired/available (*online updates*)

# Sequential Bayes in the historical example

Let's imagine that we start by observing 20 births  $\mathbf{y}_{1:20}$  at the start of 1745, including 9 girls, and that we have a uniform *prior* on  $\theta$  :

$$\theta | \mathbf{y}_{1:20} \sim \dots$$



## Sequential Bayes in the historical example

Let's imagine that we start by observing 20 births  $\mathbf{y}_{1:20}$  at the start of 1745, including 9 girls, and that we have a uniform *prior* on  $\theta$  :

$$\theta | \mathbf{y}_{1:20} \sim \text{Beta}(10, 12)$$

Then we observe  $\mathbf{y}_{21:493472}$  the remaining 493 452 births between 1745 and 1770, including 241 936 girls, and we then uses this  $\text{Beta}(10, 12)$  *prior* for  $\theta$  :

$$\theta | \mathbf{y}_{1:20}, \mathbf{y}_{21:493472} \sim \dots$$

## Sequential Bayes in the historical example

Let's imagine that we start by observing 20 births  $\mathbf{y}_{1:20}$  at the start of 1745, including 9 girls, and that we have a uniform *prior* on  $\theta$  :

$$\theta | \mathbf{y}_{1:20} \sim \text{Beta}(10, 12)$$

Then we observe  $\mathbf{y}_{21:493472}$  the remaining 493 452 births between 1745 and 1770, including 241 936 girls, and we then uses this  $\text{Beta}(10, 12)$  *prior* for  $\theta$  :

$$\begin{aligned} \theta | \mathbf{y}_{1:20}, \mathbf{y}_{21:493472} &\sim \text{Beta}(10 + 241\,936, 12 + 251\,516) \\ &\sim \text{Beta}(241\,946, 251\,528) \end{aligned}$$

We get the same *posterior* distribution as with all the observations taken together at once

# Bayesian inference

# Bayesian Inference

Bayesian modeling  $\Rightarrow$  *posterior* distribution:

- all of the information on  $\theta$ , **conditionally to both the model and the data**

# Bayesian Inference

Bayesian modeling  $\Rightarrow$  *posterior* distribution:

- all of the information on  $\theta$ , **conditionally to both the model and the data**

*Summary* of this *posterior* distribution ?

- center
- spread
- ...

# Decision theory

Context: estimating an unknown parameter  $\theta$

Decision: choice of an “optimal” point estimator  $\hat{\theta}$

**cost function**: quantify the penalty associated with the choice of a particular  $\hat{\theta}$

⇒ minimize the cost function to choose the optimal  $\hat{\theta}$

a large number of cost functions are available: each one yields a different point estimator based on its own minimum rule

# Point estimates

- **Posterior mean:**  $\mu_P = \mathbb{E}(\theta|\mathbf{y}) = \mathbb{E}_{\theta|\mathbf{y}}(\theta)$   
not always easy because it assumes the calculation of an integral...  
⇒ minimize the quadratic error cost
- **Maximum A Posteriori (MAP):**  
easy(er) to compute: just a simple maximization of the *posterior*  
 $f(\mathbf{y}|\theta)\pi(\theta)$
- **Posterior median:** the median of  $p(\theta|\mathbf{y})$   
⇒ minimize the absolute error cost

⚠ the Bayesian approach gives a full characterization of the *posterior* distribution that goes beyond point estimation

# MAP on the historical example

**Maximum A Posteriori** on the historical example of feminine birth in Paris with a uniform prior:

$$p(\theta|\mathbf{y}) = \binom{n}{S} (n+1)\theta^S(1-\theta)^{n-S}$$

with  $n = 493,472$  et  $S = 241,945$

$$\hat{\theta}_{MAP} = \frac{S}{n} = 0.4902912$$



## Posterior mean on the historical example

**Posterior mean** on the historical example of feminine birth in Paris with a uniform prior:

$$p(\theta|\mathbf{y}) = \binom{n}{S} (n+1)\theta^S(1-\theta)^{n-S}$$

with  $n = 493,472$  et  $S = 241,945$

$$E(\theta|\mathbf{y}) = \int_0^1 \theta p(\theta|\mathbf{y}) d\theta$$

$$\tilde{\theta} = \binom{n}{S} (n+1) \frac{S+1}{\binom{n}{S} (n+1)(n+2)} = \frac{S+1}{n+2} = 0.4902913$$

# Confidence Interval reminder

What is the interpretation of a frequentist confidence interval at a 95% level ?

...

⇒ Socrative: <https://b.socrative.com/login/student/>  
Room: BAYESMED2024

# Confidence Interval reminder

What is the interpretation of a frequentist confidence interval at a 95% level ?

*95% of the intervals computed on all possible samples (all those that could have been observed) contain the true value  $\theta$*

**Warning:** one cannot interpret a realization of a confidence interval in probabilistic terms ! It is a common mistake. . .

# Credibility interval

The **credibility interval** is interpreted much more naturally than the confidence interval:

It is an interval that has a 95% chance of containing  $\theta$  (for a 95% level, obviously)

Defined as an interval with a high *posterior* probability of occurrence.

For example, a **95% credibility interval** is an interval  $[t_{inf}, t_{sup}]$  such

$$\text{that } \int_{t_{inf}}^{t_{sup}} p(\theta|\mathbf{y}) d\theta = 0.95$$

**NB:** usually interested in the shortest possible 95% credibility interval (also called Highest Density Interval).

# Bayes Factor

**Bayes Factor:** marginal likelihood ratio between two hypotheses

$$BF_{10} = \frac{f(\mathbf{y}|H_1)}{f(\mathbf{y}|H_0)}$$

⇒ favored support for either hypothesis from the observed data  $\mathbf{y}$

# Bayes Factor

**Bayes Factor:** marginal likelihood ratio between two hypotheses

$$BF_{10} = \frac{f(\mathbf{y}|H_1)}{f(\mathbf{y}|H_0)}$$

⇒ favored support for either hypothesis from the observed data  $\mathbf{y}$

<i>BF</i> value	Interpretation
$BF < 1$	Negative (favors $H_0$ )
$1 \leq BF < 10^{1/2}$	Barely worth mentioning
$10^{1/2} \leq BF < 10$	Substantial
$10 \leq BF < 10^{3/2}$	Strong
$10^{3/2} \leq BF < 100$	Very strong
$100 \leq BF$	Decisive

# Bayes Factor

**Bayes Factor:** marginal likelihood ratio between two hypotheses

$$BF_{10} = \frac{f(\mathbf{y}|H_1)}{f(\mathbf{y}|H_0)}$$

⇒ favored support for either hypothesis from the observed data  $\mathbf{y}$

<i>BF</i> value	Interpretation
$BF < 1$	Negative (favors $H_0$ )
$1 \leq BF < 10^{1/2}$	Barely worth mentioning
$10^{1/2} \leq BF < 10$	Substantial
$10 \leq BF < 10^{3/2}$	Strong
$10^{3/2} \leq BF < 100$	Very strong
$100 \leq BF$	Decisive

**Posterior odds:**  $\frac{p(H_1|\mathbf{y})}{p(H_0|\mathbf{y})} = BF_{10} \times \frac{p(H_1)}{p(H_0)}$

Asymptotics

# Concentration of the posterior

## Doob's convergence



# Normal approximation

**Bernstein-von Mises Theorem (or Bayesian central-limit theorem):**  
For a large  $n$  the *posterior* can be approximated by a normal distribution.

$$p(\theta|\mathbf{y}) \approx \mathcal{N}(\hat{\theta}, I(\hat{\theta})^{-1})$$

## Consequences:

- Bayesian methods and frequentist procedures based on maximum likelihood give, for large enough  $n$ , very close results
- the *posterior* can be computed as a normal whose mean and variance we can calculate simply using the MAP

# Conclusion

# Essential concepts

## 1 Bayesian modeling:

$\theta \sim \pi(\theta)$  the *prior*

$Y_i|\theta \stackrel{iid}{\sim} f(y|\theta)$  sampling model

## 2 Bayes' formula: $p(\theta|\mathbf{y}) = \frac{f(\mathbf{y}|\theta)\pi(\theta)}{f(\mathbf{y})}$

with  $p(\theta|\mathbf{y})$  the *posterior*,  $f(\mathbf{y}|\theta)$  the likelihood (inherited from the sampling model),  $\pi(\theta)$  the *prior* and  $f(\mathbf{y}) = \int f(\mathbf{y}|\theta)\pi(\theta)$  is the marginal distribution of the data, i.e. the normalizing constant (with respect to  $\theta$ )

## 3 The *posterior* distribution is given by:

$$p(\theta|\mathbf{y}) \propto f(\mathbf{y}|\theta)\pi(\theta)$$

## 4 *Posterior* mean, MAP, and credibility intervals

# Practical use

The Bayesian framework is (just) another statistical tool for data analysis

Particularly **useful when:**

- few observations only are available
- there is important knowledge *a priori*

Like any statistical method, Bayesian analysis has advantages and disadvantages that will be more or less important depending on the application considered.

# Questions ?

