

# Introduction to the Bayesian framework in Biometrics

## Part I: Bayesian theory

Boris Hejblum

PhD Training course  
*Digital Public Health* program, University of Bordeaux

# 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 and discuss assumptions and methodological choices in biomedical literature using Bayesian methods, including "under the hood" estimation machinery
- ② understand, discuss and reproduce a Bayesian estimation of a relative risk
- ③ understand and perform a Bayesian meta-analysis using 
- ④ understand continuous reassessment method for phase I/II dose trials, along with the associated decision-rule

**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.

# Motivational examples: diagnostic tests

[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 90 90	70 90 90	100 100 100	41.2 18.9 73.0
Patient 2: low pre-test probability	5 10	70 70	97.4 98.7	1.6 3.2
			97.9 99.0	0.5 1.1



## 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



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

**CLINICAL  
TRIALS**

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

Clinical Trials  
I–10  
© The Author(s) 2020  
Article reuse guidelines:  
[sagepub.com/journals-permissions](http://sagepub.com/journals-permissions)  
DOI: [10.1177/1740774520943846](https://doi.org/10.1177/1740774520943846)  
[journals.sagepub.com/home/ctj](http://journals.sagepub.com/home/ctj)



**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]

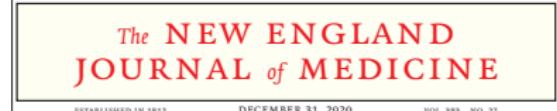
# Motivational examples: study/trial analyses



## 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.



## 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., Satrajita Roychoudhury, Ph.D., Kenneth Couris, Ph.D., Ping Li, Ph.D., Warren V. Kalina, Ph.D., David Cooper, Ph.D., Robert W. French, Jr., M.D., Laura L. Hammitt, M.D., Özlem Türeci, M.D., Haylene Nell, M.D., Axel Schaefer, M.D., Serhat Ünal, M.D., Dina B. Tresnan, D.V.M., Ph.D., Susan Mather, M.D., Philip R. Dormitzer, M.D., Ph.D., Uğur Şahin, 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 (N=18,196)	Surveillance Time (n)†	No. of Cases (N=18,325)	Surveillance Time (n)†		
Covid-19 occurrence at least 7 days after the second dose in participants without evidence of infection	8	2,214 (17,411)	162	2,222 (17,511)	95.0 (90.3–97.4)	>0.9999
Covid-19 occurrence at least 7 days after the second dose in participants with and without evidence of infection	9	2,332 (18,559)	169	2,345 (18,708)	94.6 (89.9–97.3)	>0.9999

\* The total population without baseline infection was 36,523; 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.)

## Bayes theorem: exercise

As of May 11<sup>th</sup>, about 7% of the French population was estimated to have had COVID-19. A medical test has the following 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

As of May 11<sup>th</sup>, about 7% of the French population was estimated to have had COVID-19. A medical test has the following 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 ?*

$$\Pr(D=+) = 0.07 \quad \Pr(T=+|D=+) = 0.71 \quad \Pr(T=-|D=-) = 0.98$$

## Bayes theorem: exercise

As of May 11<sup>th</sup>, about 7% of the French population was estimated to have had COVID-19. A medical test has the following 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 ?*

$$\Pr(D=+) = 0.07 \quad \Pr(T=+|D=+) = 0.71 \quad \Pr(T=-|D=-) = 0.98$$

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

## Bayes theorem: exercise

As of May 11<sup>th</sup>, about 7% of the French population was estimated to have had COVID-19. A medical test has the following 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 ?*

$$\Pr(D=+) = 0.07 \quad \Pr(T=+|D=+) = 0.71 \quad \Pr(T=-|D=-) = 0.98$$

$$\begin{aligned}
 \Pr(D=+|T=+) &= \frac{\Pr(T=+|D=+)\Pr(D=+)}{\Pr(T=+)} \\
 &= \frac{\Pr(T=+|D=+)\Pr(M=+)}{\Pr(T=+|D=+)\Pr(D=+) + \Pr(T=+|D=-)\Pr(D=-)} \\
 &= \frac{\Pr(T=+|D=+)\Pr(M=+)}{\Pr(T=+|D=+)\Pr(D=+) + (1 - \Pr(T=+|D=-))(1 - \Pr(D=+))} \\
 &= (0.71 \times 0.07) / (0.71 \times 0.07 + (1 - 0.98) \times (1 - 0.07)) = 0.73
 \end{aligned}$$

## Continuous Bayes' theorem

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

### Continuous Bayes' theorem:

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

## Continuous Bayes' theorem

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

## Continuous Bayes' theorem:

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



remember 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** ⇒ development of statistics in the **18-19<sup>th</sup> centuries**
- ② Galton & Pearson, then Fisher & Neymann ⇒ **frequentist** theory became dominant during the **20<sup>th</sup> century**
- ③ turn of the **21<sup>th</sup> century**: rise of the computer  
⇒ **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  $f(y|\theta)$ ,  $\theta \in \Theta$

# 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 ?

## Construction of a Bayesian model

## 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

$y$ : 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  $p(\theta|y)$

# Posterior distribution

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

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

**Bayes' theorem:**

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

# Posterior distribution

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

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

**Bayes' theorem:**

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

Posterior is calculated from:

- ① the sampling model  $f(y|\theta)$  – which yields the likelihood  $f(y|\theta)$  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

11

## 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

• • •

## Construction of a Bayesian model

### 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 calculate: ...

# 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

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

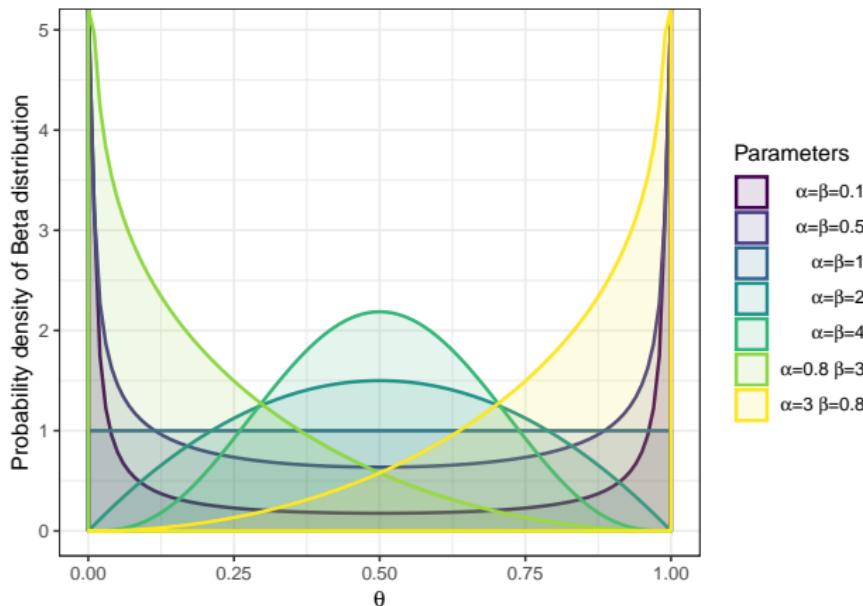
To answer the question of interest, we can then calculate:

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

## Construction of a Bayesian model

## 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

Introduction to Bayesian statistics  
oooooooo

## Construction of a Bayesian model

Bayesian modeling  
oooooooooooo●oooooooooooo

Bayesian Inference  
oooooooooooo

Conclusion  
ooo

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

...

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

# Conjugacy of the Beta distribution

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

**Corresponding posterior:**  $p(\theta|y) \propto \theta^{\alpha+S-1} (1-\theta)^{\beta+(n-S)-1}$   
 $\Rightarrow \theta|y \sim \text{Beta}(\alpha + S, \beta + (n - S))$

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

# Conjugacy of the Beta distribution

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

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

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

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

## Construction of a Bayesian model

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 ≠ #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 ≠ #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 – continued

# 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  $p(\theta|y) = 0$
  
- 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**
- ...

# 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**
- 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**
- 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, parameters 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:

$$① \quad \eta \sim h(\eta)$$

$$② \quad \pi(\theta|\eta)$$

$$③ \quad f(\mathbf{y}|\theta)$$

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

**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:

$$① \quad \eta \sim h(\eta)$$

$$② \quad \pi(\theta|\eta)$$

$$③ \quad f(\mathbf{y}|\theta)$$

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

**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*:  $p(\theta|y, \hat{\eta})$

# 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*:  $p(\theta|y, \hat{\eta})$

- Combines Bayesian and frequentist frameworks
- Concentrated *posterior* ( $\searrow$  variance) but  $\nearrow$  bias (data used twice !)
- Approximate a fully Bayesian approach

# Empirical Bayes: example

For a distribution Beta( $\alpha, \beta$ ):

- $\frac{\alpha}{\alpha+\beta}$  is the mean
- $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$  is the variance

Thanks to the **method of moments** we get:  $\hat{\alpha}_M = 0,020$  et  
 $\hat{\beta}_M = 0,021$

since  $\hat{\theta} = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = 0,49$  et  $\widehat{Var}(\theta) = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = 0,24$

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

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

# 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  $y_{1:20}$  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  $y_{1:20}$  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  $y_{21:493472}$  the remaining 493 452 births between 1745 and 1770, including 241 936 girls, and we then uses this 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  $y_{1:20}$  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  $y_{21:493472}$  the remaining 493 452 births between 1745 and 1770, including 241 936 girls, and we then uses this Beta(10, 12) *prior* for  $\theta$ :

$$\begin{aligned}\theta | \mathbf{y}_{1:20}, \mathbf{y}_{21:493472} &\sim \text{Beta}(10 + 241936, 12 + 251516) \\ &\sim \text{Beta}(241946, 251528)\end{aligned}$$

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
- ...

## Point estimates

## 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

## Point estimates

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

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

## Point estimates

## 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 ?

...

# 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

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

**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  $\mathbf{y}$

## Posterior odds

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

# Concentration de la loi *a posteriori*

## Doob's convergence theorem

*Posteriori* distribution concentrate on the “real” value of the  $\theta^*$  parameter when  $n \rightarrow \infty$  :

$$p(\theta | \mathbf{y}_n) \xrightarrow{\mathcal{L}} \delta_{\theta^*}$$

→ *Seeing Theory*, Brown University  
<http://students.brown.edu/seeing-theory/bayesian-inference/index.html#section3>

# 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

# Illustration on the historical example

→ *Seeing Theory*, Brown University

# Conclusion

# Essential concepts

## ① Bayesian modeling:

$$\theta \sim \pi(\theta) \quad \text{the } prior$$

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

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

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

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

$$p(\theta|y) \propto f(y|\theta)\pi(\theta)$$

## 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 ?

