

# Statistics for Applications

## Chapter 8: Bayesian Statistics

what we have been doing so far is frequentist approach.

# The Bayesian approach (1)

- ▶ So far, we have studied the frequentist approach of statistics.
- ▶ The frequentist approach:
  - ▶ Observe data
  - ▶ These data were generated randomly (by Nature, by measurements, by designing a survey, etc...)
  - ▶ We made assumptions on the generating process (e.g., i.i.d., Gaussian data, smooth density, linear regression function, etc...)
  - ▶ The generating process was associated to some object of interest (e.g., a parameter, a density, etc...)
  - ▶ This object was unknown but fixed and we wanted to find it: we either estimated it or tested a hypothesis about this object, etc...

- \* The first thing we had data, we observed some data and we assumed that this data is generated randomly, The reason we did that is because this will allow to leverage tools from probability.
- \* Then we made some assumption's on the data generating process, for eg, we assumed they are iid, sometimes we assume gaussian, etc..
- \* So this data, we are interested in some Parameter  $\theta$ , or  $\beta$  in regression model, so we have this unknown problem, unknown parameter we want to find so we want it either Estimate it, or test it, or may be find confidence interval on this object.

"The object was unknown but fixed" is different for Bayesian approach

## The Bayesian approach (2)

- ▶ Now, we still observe data, assumed to be randomly generated by some process. Under some assumptions (e.g., parametric distribution), this process is associated with some fixed object.
- ▶ We have a **prior belief** about it.
- ▶ Using the data, we want to update that belief and transform it into a **posterior belief**.

In the Bayesian approach we still assume we observe some random data, but the generating process is slightly different, it's sort of two layer process,

- ① That generate the Parameter
- ② Given this parameter generate the data

So what the 1<sup>st</sup> layer does, nobody really believes that there's some random process that's happening, about generating what's gonna to be the true expected number of people when they kiss.

\* But often we actually have Prior Belief

Ex: when we did least squares, we looked over for all of the vector's in all of the  $\mathbb{R}^P$  including the ones that are 50 million. So those are things we might be able to rule out, maybe we might be rule out that on a much smaller scale.

\* This Prior belief is gonna play hopefully less and less important role as we collect more and more data, but if we have smaller amount of data, we might want to be able to use this information, rather than just shooting in the dark.

\* So idea is to have this Prior Belief & then we want to update this Prior Belief what's called Posterior Belief, after we have seen some data.

\* So belief encompasses basically what you think and how strongly you think about it.

\* Ex: if we have some belief about parameter  $\Theta$ , may be my belief is telling me, where the  $\Theta$  should be and how strongly I believe means I have a very narrow region where  $\Theta$  could be.

# The Bayesian approach (3)

## Example

- ▶ Let  $p$  be the proportion of woman in the population.
- ▶ Sample  $n$  people randomly with replacement in the population and denote by  $X_1, \dots, X_n$  their gender (1 for woman, 0 otherwise).
- ▶ In the frequentist approach, we estimated  $p$  (using the MLE), we constructed some confidence interval for  $p$ , we did hypothesis testing (e.g.,  $H_0 : p = .5$  v.s.  $H_1 : p \neq .5$ ).
- ▶ Before analyzing the data, we may believe that  $p$  is likely to be close to  $1/2$ .
- ▶ The Bayesian approach is a tool to:
  1. include mathematically our prior belief in statistical procedures.
  2. update our prior belief using the data.

$P$ : Proportion of women in population

Collect data:  $X_1, X_2, \dots, X_n \sim \text{Ber}(P)$   $P \in (0,1)$

frequentist approach: Let's just estimate

the  $\hat{P} = \frac{1}{n} \sum X_i = \bar{X}_n$  & do some tests whether  $P=0.5$  or not etc.

But here this is the case where we definitely have prior belief of what  $P$  should be. We are pretty confident that  $P$  is not going to be 0.7. We actually believe  $P$  is extremely close to 0.5, maybe not exactly.

- \* So, we are going to want to integrate that knowledge, so we could integrate it in blunt manner by saying, like disregard the data and say that  $P = 0.5$ .
- \* But maybe that's just too much, so how do we do this trade off between adding the data & combining with this prior knowledge?

in many ways, in many instances,  
essentially what's gonna happen in this  
O.S is going to act like one new obser-  
vation. So if we have five observations,  
this will be sixth observation, which  
will play a role.

\* if we have million observations, we are  
going to have one million and one observa-  
tion's, and it's not going to play so much  
of a role.

\* we have 5 observations  $\Rightarrow$  5 + 1 (prior belief)  
 $= 6$  observations

if we have 1 million observations  $\Rightarrow$  1 million + 1 (prior  
belief)  
(Not gonna play big role)

- ① The Bayesian approach is a tool to one to  
include mathematically our prior belief into  
statistical procedures,
- ② update this prior belief into a posterior  
belief using the data.

How do we do this?

there is two layer's.

- ① one is where we draw the parameter at random
- ② once we have the parameter, condition this parameter and draw our data.

No body believed this actually happening that Nature is just rolling dice, for us and choosing parameter's at random. The idea that parameter coming from some random distribution actually capture very well that this idea that how you would encompass your prior.

# The Bayesian approach (4)

## Example (continued)

- ▶ Our prior belief about  $p$  can be quantified:
- ▶ E.g., we are 90% sure that  $p$  is between .4 and .6, 95% that it is between .3 and .8, etc...
- ▶ Hence, we can model our prior belief using a distribution for  $p$ , *as if*  $p$  was random.
- ▶ In reality, the true parameter is not random ! However, the Bayesian approach is a way of modeling our belief about the parameter by doing **as if** it was random.
- ▶ E.g.,  $p \sim \mathcal{B}(a, a)$  (*Beta distribution*) for some  $a > 0$ .
- ▶ This distribution is called the *prior distribution*.

Ex:

there are many way's

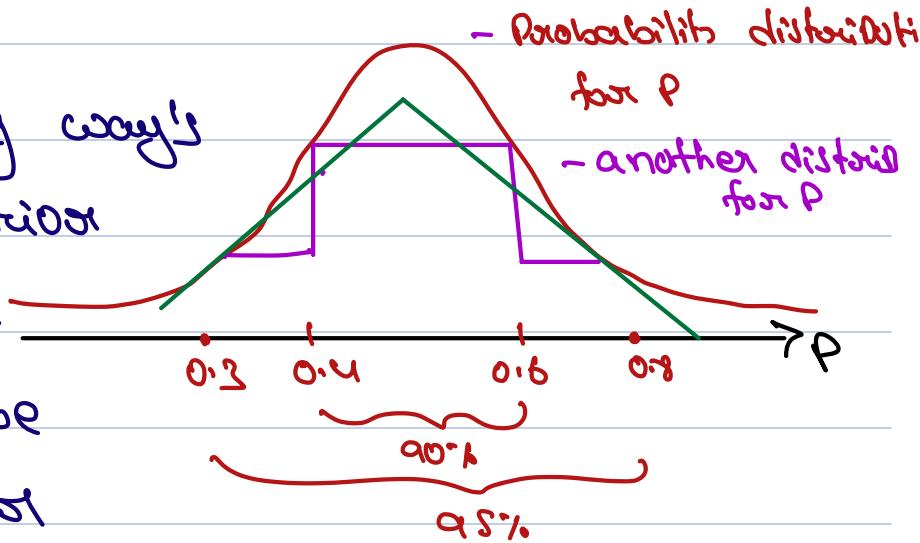
I can assume my prior

Belief. Some of them

are definitely to be

mathematically most

Convenient than others. Hopefully we are gonna have things that we can parameterize very well.



$$\theta \in \mathbb{C}$$

Gaussian

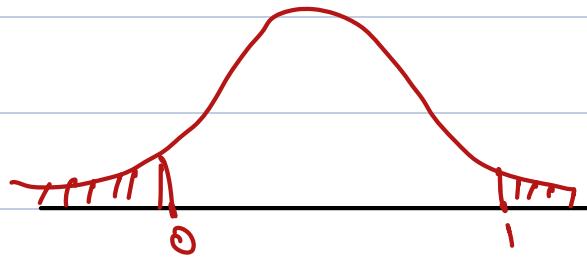
$$\mathbb{R} \quad [0, 1]$$

Bernoulli

$$\in \mathbb{P}$$

$$(0, \infty)$$

Exponential



cannot have Gaussian

as it has Prob outside

$$[0, 1]$$

Beta distribution:

X ~ Beta( $\alpha, \beta$ )

$$f(x) = \begin{cases} \frac{1}{C} x^{\alpha-1} \cdot (1-x)^{\beta-1} & \text{if } x \in (0, 1) \\ 0 & \text{otherwise} \end{cases}$$

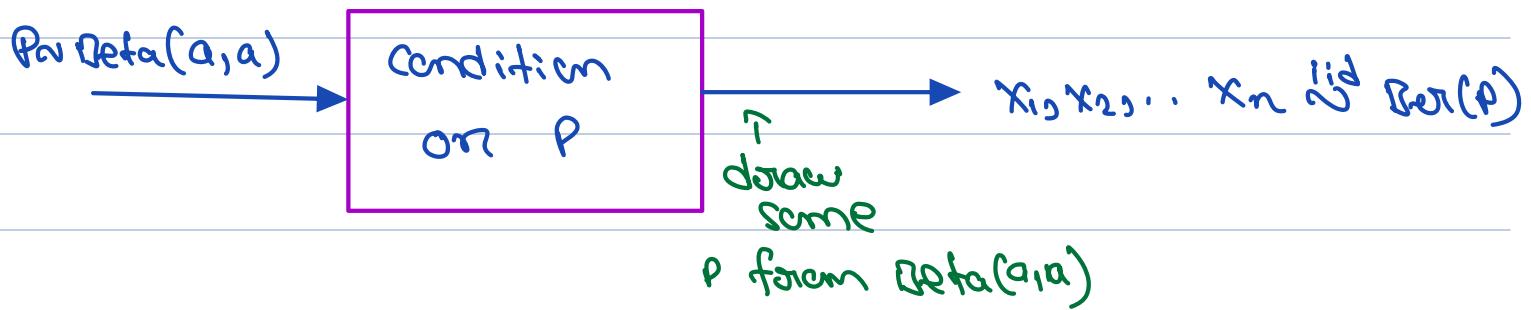
# The Bayesian approach (5)

## Example (continued)

- ▶ In our statistical experiment,  $X_1, \dots, X_n$  are assumed to be i.i.d. Bernoulli r.v. with parameter  $p$  **conditionally on**  $p$ .
- ▶ After observing the available sample  $X_1, \dots, X_n$ , we can update our belief about  $p$  by taking its distribution conditionally on the data.
- ▶ The distribution of  $p$  conditionally on the data is called the *posterior distribution*.
- ▶ Here, the posterior distribution is

$$\mathcal{B} \left( a + \sum_{i=1}^n X_i, a + n - \sum_{i=1}^n X_i \right).$$

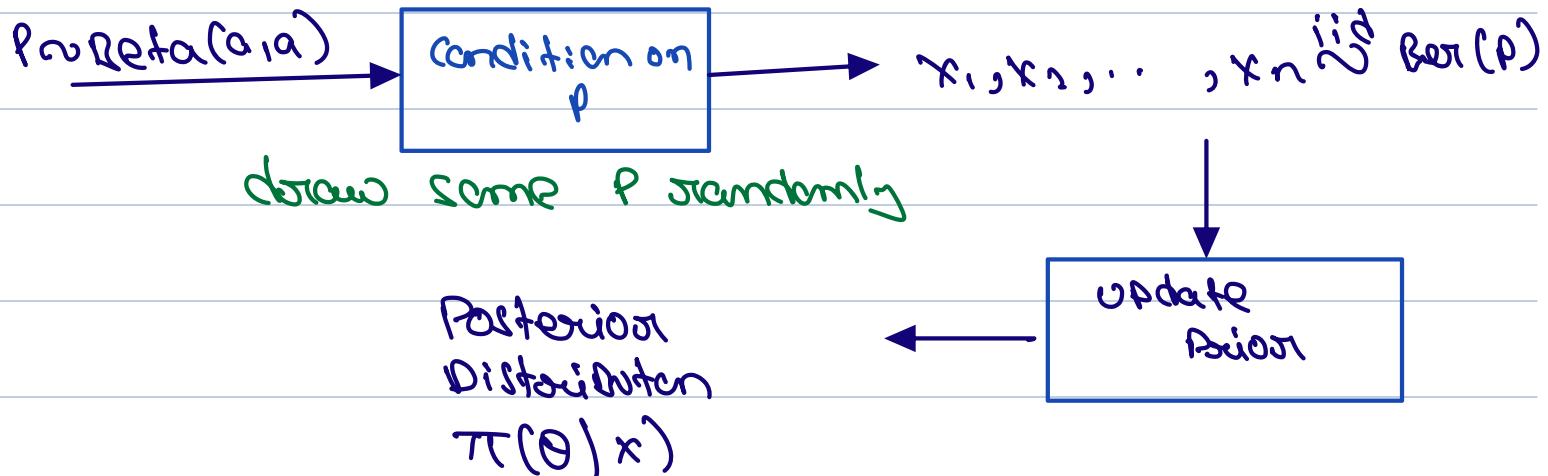
## Flow Chart:



$P$  is randomly drawn from  $\text{Beta}(a, a)$

- \* So we first draw Parameter  $P \sim \text{Beta}(a, a)$  (assume  $p=0.52$ ) and then we just flip those independent biased coin's with this particular  $P=0.52$  (2 layer process)

- \* This is just a thought process. it's not anything that actually happens in general. This is our way of thinking about how the data was generated.



In this case we will see the posterior distribution is still Beta

Beta is Conjugate: meaning I put Beta as Prior and I get Beta as Posterior.

$$P(\theta|x) \sim \text{Beta}(\alpha + \sum x_i, \alpha + n - \sum x_i)$$

so how do we get this posterior distribution given the prior? How do we update this?

This is called Bayesian Statistics.

Bayes Formula:

$$P = \text{data}(x_1, x_2, \dots, x_n)$$

$$A = \text{Parameter } (\theta)$$

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

$\downarrow \quad \downarrow \quad \downarrow$

$$\underbrace{P(\theta)}_{\text{dist of parameter } \theta} \cdot \underbrace{P(A)}_{\text{dist of data given } \theta}$$

$P(B) \leftarrow P_x(x)$

dist of parameter  $\theta$ ,

given I have seen the

data  $P(\theta|x_1, x_2, \dots, x_n)$   
of  $x$

distribution of data given

that I know what my

parameter  $\theta$  is

$P_{x|\theta}(x_1, x_2, \dots, x_n | \theta)$

$P(A) = P(\theta)$  = Prior distribution

$P(B) = P_x(x_1, x_2, \dots, x_n)$  = distribution of data itself.

# The Bayes rule and the posterior distribution (1)

- ▶ Consider a probability distribution on a parameter space  $\Theta$  with some pdf  $\pi(\cdot)$ : the *prior distribution*.
- ▶ Let  $X_1, \dots, X_n$  be a sample of  $n$  random variables.
- ▶ Denote by  $p_n(\cdot|\theta)$  the joint pdf of  $X_1, \dots, X_n$  conditionally on  $\theta$ , where  $\theta \sim \pi$ .
- ▶ Usually, one assumes that  $X_1, \dots, X_n$  are i.i.d. conditionally on  $\theta$ .
- ▶ The conditional distribution of  $\theta$  given  $X_1, \dots, X_n$  is called the *posterior distribution*. Denote by  $\pi(\cdot|X_1, \dots, X_n)$  its pdf.

$x_1, x_2, \dots, x_n$  iid

Give Parameter  $\theta$

Conditioned on random Parameter, which  
was a fixed Number

$x_1, x_2, \dots, x_n \sim P_n(x_1, x_2, \dots, x_n | \theta)$   
pdf or pmf

Ex:  $x_1, x_2, \dots, x_n$  iid  $N(0, 1)$

$$\Rightarrow P_n(x_1, x_2, \dots, x_n | \theta) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum (x_i - \theta)^2}$$

$\Rightarrow \theta$  has a prior distribution  $\Pi(\cdot)$  pdf or pmf

Ex:  $\Pi(\theta) = C \cdot \theta^{a-1} (1-\theta)^{a-1}$ ,  $\theta \sim \text{Beta}(a, a)$

Posterior Distribution:

$$\Pi(\theta | x_1, x_2, \dots, x_n) = \frac{P_n(x_1, x_2, \dots, x_n | \theta) \Pi(\theta)}{\int_{\Theta} P_n(x_1, x_2, \dots, x_n | \theta) \Pi(\theta) d\theta}$$

↓  
marginal distribution

weighted version of likelihood

Posterior dist

$$\pi(\theta | x_1, x_2, \dots, x_n) =$$

Likelihood

$$P_n(x_1, x_2, \dots, x_n | \theta)$$

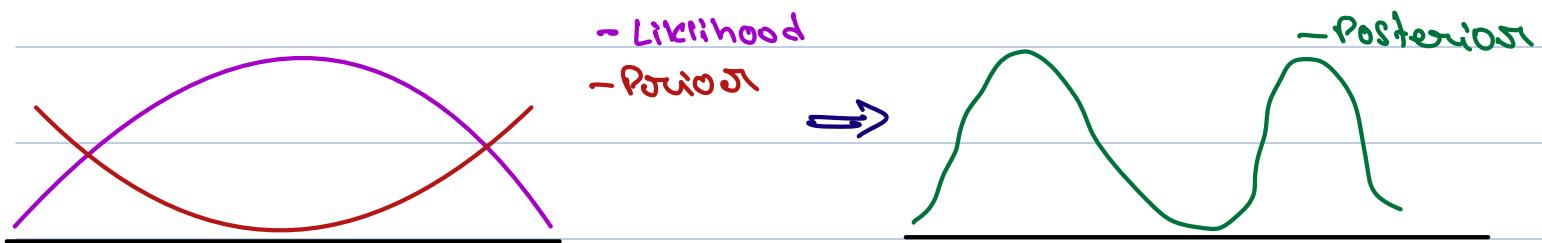
Prior dist

$$\pi(\theta)$$

Constant that  
does not depend  
on  $\theta$

$$\int_{\Theta} P_n(x_1, x_2, \dots, x_n | \theta) \pi(\theta) d\theta$$

\* Important to remember: MLE & Bayesian are really not that different. because our posterior is really just like your likelihood times something that's just putting some weight's on the  $\theta$ 's. depending on where the  $\theta$  should be.



Likelihood

Prior dist

$$P_n(x_1, x_2, \dots, x_n | \theta)$$

$$\pi(\theta)$$

$\Rightarrow$  weighted version of likelihood.

- weighting the likelihood using my prior belief on  $\theta$ .  
- if we follow MLE principle, if we maximize the weighted likelihood we would get Maximum A Posterior Estimator  $\hat{\theta}_{MAP}$

MAP  $\rightarrow$  Maximum a Posterior

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{arg\,max}} \pi(\theta | x_1, x_2, \dots, x_n)$$

The beauty of Bayesian Statistics is we don't have to take any number in particular (like maximum, mean, median etc). We have an entire posterior distribution.

## The Bayes rule and the posterior distribution (2)

- ▶ Bayes' formula states that:

$$\pi(\theta|X_1, \dots, X_n) \propto \pi(\theta)p_n(X_1, \dots, X_n|\theta), \quad \forall \theta \in \Theta.$$

- ▶ The constant does not depend on  $\theta$ :

$$\pi(\theta|X_1, \dots, X_n) = \frac{\pi(\theta)p_n(X_1, \dots, X_n|\theta)}{\int_{\Theta} p_n(X_1, \dots, X_n|t) d\pi(t)}, \quad \forall \theta \in \Theta.$$

# The Bayes rule and the posterior distribution (3)

**In the previous example:**

- ▶  $\pi(p) \propto p^{a-1}(1-p)^{a-1}, p \in (0, 1).$

- ▶ Given  $p, X_1, \dots, X_n \stackrel{i.i.d.}{\sim} Ber(p)$ , so

$$p_n(X_1, \dots, X_n | \theta) = p^{\sum_{i=1}^n X_i} (1-p)^{n - \sum_{i=1}^n X_i}.$$

- ▶ Hence,

$$\pi(\theta | X_1, \dots, X_n) \propto p^{a-1 + \sum_{i=1}^n X_i} (1-p)^{a-1 + n - \sum_{i=1}^n X_i}.$$

- ▶ The posterior distribution is

$$\mathcal{B}\left(a + \sum_{i=1}^n X_i, a + n - \sum_{i=1}^n X_i\right).$$

Prior  $\pi(\theta) \sim \text{Beta}(\alpha, \alpha)$

$\Rightarrow \pi(p) \propto p^{\alpha-1} (1-p)^{\alpha-1}$  (Density of  $p$ )

Likelihood:  $P_n(x_1, x_2, \dots, x_n | p) = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i}$

Posterior distribution:  $\propto p^{\alpha-1} (1-p)^{\alpha-1} p^{\sum x_i} (1-p)^{n - \sum x_i}$

$\propto p^{\alpha + \sum x_i - 1} (1-p)^{\alpha + n - \sum x_i - 1}$

$\pi(\theta | x_1, x_2, \dots, x_n) \sim \text{Beta}(\alpha + \sum x_i, \alpha + n - \sum x_i)$

\* The Prior  $\pi(\theta)$  in the Posterior formula

is to weigh some theta's more than other's, depending on their Prior Belief

\* if our Prior belief does not want to put any preference towards some theta's than to other's, what do we do?  $\pi(\theta) = 1$

\* if  $\pi(\theta) = 1$ , constant and does not depend on  $\theta$ , that would mean that we are not preferring anything, and we are looking at the Likelihood.

$\Rightarrow$  we are looking at the likelihood, but not of a function that we are trying to maximize, but it is a function that we normalize in such a way its actually a distribution.

# Non informative priors (1)

- ▶ Idea: In case of ignorance, or of lack of prior information, one may want to use a prior that is as little informative as possible.
- ▶ Good candidate:  $\pi(\theta) \propto 1$ , i.e., constant pdf on  $\Theta$ .
- ▶ If  $\Theta$  is bounded, this is the uniform prior on  $\Theta$ .
- ▶ If  $\Theta$  is unbounded, this does not define a proper pdf on  $\Theta$  !
- ▶ An improper prior on  $\Theta$  is a measurable, nonnegative function  $\pi(\cdot)$  defined on  $\Theta$  that is not integrable.
- ▶ In general, one can still define a posterior distribution using an improper prior, using Bayes' formula.

## Non informative priors (2)

**Examples:**

- If  $p \sim \mathcal{U}(0, 1)$  and given  $p, X_1, \dots, X_n \stackrel{i.i.d.}{\sim} Ber(p)$ :

$$\pi(p|X_1, \dots, X_n) \propto p^{\sum_{i=1}^n X_i} (1-p)^{n-\sum_{i=1}^n X_i},$$

i.e., the posterior distribution is

$$\mathcal{B}\left(1 + \sum_{i=1}^n X_i, 1 + n - \sum_{i=1}^n X_i\right).$$

- If  $\pi(\theta) = 1, \forall \theta \in \mathbb{R}$  and given  $\theta, X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \mathcal{N}(\theta, 1)$ :

$$\pi(\theta|X_1, \dots, X_n) \propto \exp\left(-\frac{1}{2} \sum_{i=1}^n (X_i - \theta)^2\right),$$

i.e., the posterior distribution is

$$\mathcal{N}\left(\bar{X}_n, \frac{1}{n}\right).$$

Prior:  $\theta \sim \text{Unif}(0,1)$

Likelihood:  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(\theta)$

$$\Rightarrow P_n(X_1, X_2, \dots, X_n | \theta) .$$

$$= \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

Posterior:

$$\pi(\theta | x_1, x_2, \dots, x_n) = \frac{\theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \cdot 1}{\int_0^1 \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} d\theta}$$

$$\pi(\theta) = \text{Beta}\left(1 + \sum_{i=1}^n x_i, 1 + n - \sum_{i=1}^n x_i\right)$$

---

$$\hat{\theta}_{MLE} = \bar{x} \quad (\text{MLE})$$

## Non informative priors (3)

- ▶ *Jeffreys prior:*

$$\pi_J(\theta) \propto \sqrt{\det I(\theta)},$$

where  $I(\theta)$  is the Fisher information matrix of the statistical model associated with  $X_1, \dots, X_n$  in the frequentist approach (provided it exists).

- ▶ In the previous examples:

- ▶ Ex. 1:  $\pi_J(p) \propto \frac{1}{\sqrt{p(1-p)}}, p \in (0, 1)$ : the prior is  $\mathcal{B}(1/2, 1/2)$ .
- ▶ Ex. 2:  $\pi_J(\theta) \propto 1, \theta \in \mathbb{R}$  is an improper prior.

## Non informative priors (4)

- ▶ Jeffreys prior satisfies a reparametrization invariance principle:  
If  $\eta$  is a reparametrization of  $\theta$  (i.e.,  $\eta = \phi(\theta)$  for some one-to-one map  $\phi$ ), then the pdf  $\tilde{\pi}(\cdot)$  of  $\eta$  satisfies:

$$\tilde{\pi}(\eta) \propto \sqrt{\det \tilde{I}(\eta)},$$

where  $\tilde{I}(\eta)$  is the Fisher information of the statistical model parametrized by  $\eta$  instead of  $\theta$ .

# Bayesian confidence regions

- ▶ For  $\alpha \in (0, 1)$ , a Bayesian confidence region with level  $\alpha$  is a random subset  $\mathcal{R}$  of the parameter space  $\Theta$ , which depends on the sample  $X_1, \dots, X_n$ , such that:

$$\mathbb{P}[\theta \in \mathcal{R} | X_1, \dots, X_n] = 1 - \alpha.$$

- ▶ Note that  $\mathcal{R}$  depends on the prior  $\pi(\cdot)$ .
- ▶ "Bayesian confidence region" and "confidence interval" are two **distinct** notions.

# Bayesian estimation (1)

- ▶ The Bayesian framework can also be used to estimate the true underlying parameter (hence, in a frequentist approach).
- ▶ In this case, the prior distribution does not reflect a prior belief: It is just an artificial tool used in order to define a new class of estimators.
- ▶ **Back to the frequentist approach:** The sample  $X_1, \dots, X_n$  is associated with a statistical model  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$ .
- ▶ Define a distribution (that can be improper) with pdf  $\pi$  on the parameter space  $\Theta$ .
- ▶ Compute the posterior pdf  $\pi(\cdot | X_1, \dots, X_n)$  associated with  $\pi$ , seen as a prior distribution.

## Bayesian estimation (2)

- ▶ *Bayes estimator:*

$$\hat{\theta}^{(\pi)} = \int_{\Theta} \theta \, d\pi(\theta | X_1, \dots, X_n) :$$

This is the *posterior mean*.

- ▶ The Bayesian estimator depends on the choice of the prior distribution  $\pi$  (hence the superscript  $\pi$ ).

## Bayesian estimation (3)

- ▶ In the previous examples:
  - ▶ Ex. 1 with prior  $\mathcal{B}(a, a)$  ( $a > 0$ ):

$$\hat{p}^{(\pi)} = \frac{a + \sum_{i=1}^n X_i}{2a + n} = \frac{a/n + \bar{X}_n}{2a/n + 1}.$$

In particular, for  $a = 1/2$  (Jeffreys prior),

$$\hat{p}^{(\pi_J)} = \frac{1/(2n) + \bar{X}_n}{1/n + 1}.$$

- ▶ Ex. 2:  $\hat{\theta}^{(\pi_J)} = \bar{X}_n$ .
- ▶ In each of these examples, the Bayes estimator is consistent and asymptotically normal.
- ▶ In general, the asymptotic properties of the Bayes estimator do not depend on the choice of the prior.

MIT OpenCourseWare

<https://ocw.mit.edu>

## 18.650 / 18.6501 Statistics for Applications

Fall 2016

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.