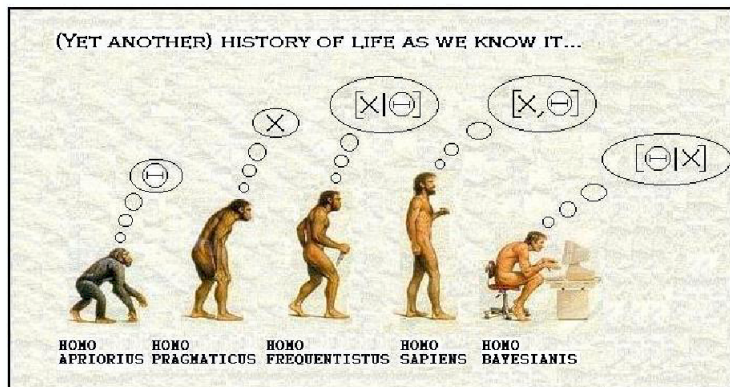


# Bayesian Thinking

Dr Ana Paula Palacios

# Objective

To provide an introduction to the Bayesian framework, highlighting its differences with the frequentist approach and focusing on estimation, modelisation and interpretation of the results.



# Contents

- Motivation
- Bayes Theorem
- Prior, Posterior and Conjugacy
- Beta-Binomial model
- Bayesian inference
- Predictive distribution
- Choosing a prior

# Thinking Bayesian

- Suppose that when you finish your MSc you plan to submit your dissertation for publication to the *Journal of the American Statistical Association* (JASA).
- *JASA Theory and Methods* receives approximately 700 original submissions a year and has an acceptance rate of about 10%.
- You submit your first paper and it is accepted! What is the probability that your next submission to this journal will be accepted? Take a minute to think on this and provide a numerical answer.

# Thinking Bayesian

- If your answer was greater than 10% and lower than 100%, then you are thinking Bayesian!
- From the frequentist point of view the probability that your next paper will be accepted is 100%, as the evidence suggest that your successful rate is 100%.
- But as the journal acceptance rate is only 10%, it seems reasonable to pick a number smaller than 100% and greater than 10%.

## Remark I

Bayesian thinking provides a way of formalising the process of learning from data to update beliefs in accord with recent notions of knowledge synthesis.

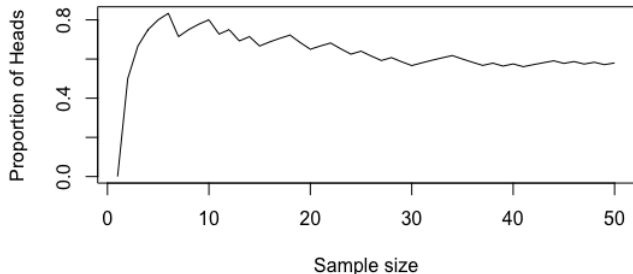
## Example: coin toss

- Suppose that we have a coin and we would like to estimate the probability of Heads ( $\theta$ )
- Suppose the following sequence of flips was observed: H T H H H H T H H H
- What do you think is the probability of Heads more “compatible” with this sample?
- From the frequentist point of view, the Maximum Likelihood Estimator (MLE) of the probability of Heads is  $\hat{\theta} = 8/10 = 80\%$
- But, we have a lot of experience flipping coins, and my experience says that the probability of Heads should be around 0.50. Why the MLE went wrong?
- Frequentist approach relies in large sample argument: they perform well in the long-run over some infinite number of repetitions of the experiment

### Probability as a Frequency

The probability of an event is equal to the long-term relative frequency of the event's occurrence when the same process is repeated many times.

## Example: coin toss



Sample Size	50	100	500	1000	5000
$\hat{\theta}$	0.5800	0.5900	0.5180	0.5070	0.4956

Frequentist approach works well when the experiment is replicable. But, not all experiments are replicable...

# Probability as uncertainty

- The Bayesian approach will include our prior beliefs.
- For example, I believe that the probability of Heads is quite likely to be around 0.50, or some value near that (0.52, 0.48) and it is less likely to be far from that (0.10, 0.90)

## Bayesian Probability

Probabilities represent how certain we are about the truth of statements. These statements/beliefs can refer to the past, the present, or the future. In the most general sense, a probability is a number attached to a statement. That number specifies how likely it is that the statement is true.



# Frequentists vs Bayesian

The basic philosophical difference between the frequentists and Bayesian paradigms is that

- Bayesians treat an unknown parameters  $\theta$  as *random* and use probability to quantify their uncertainty about it.
- In contrast, frequentists treat  $\theta$  as unknown but *fixed*, and probabilities are interpreted as long-run relative frequencies.

# Maximum Likelihood Estimator

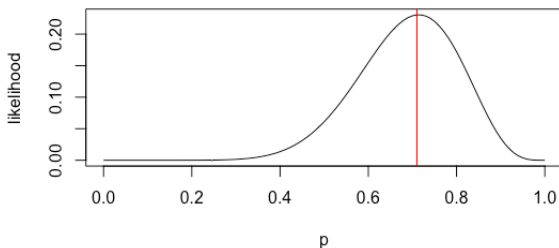
- Decide a model for the data (the likelihood function)
- For the coin example, the number of heads we get when flipping a coin  $n$  times follows a binomial distribution  $X \sim B(n, p)$ :

$$f(x|p) = \binom{n}{x} p^x (1-p)^{n-x}$$

where  $x = 0, 1, \dots, n$  and  $0 \leq p \leq 1$  is the probability of Heads.

- Find the value of  $p$  that maximises the likelihood function for the observed data
- Suppose we observe 10 Heads and 4 Tails ( $x = 10$  and  $n = 14$ ).

# Maximum Likelihood Estimator



$$\hat{p}_{MLE} = \max_p L(x|p) = \max_p \left\{ \binom{n}{x} p^x (1-p)^{n-x} \right\}$$

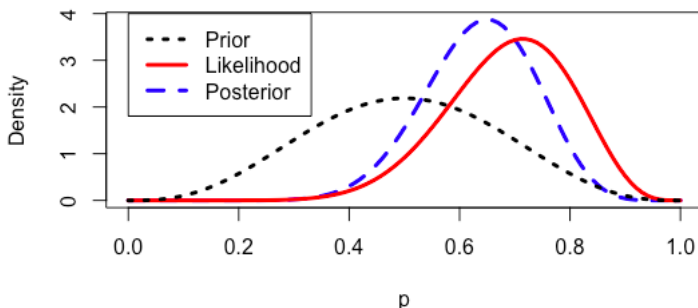
$$\hat{p}_{MLE} = \max_p \ln L(x|p) = \max_p \left\{ \ln \frac{n!}{x!(n-x)!} + x \ln p + (n-x) \ln(1-p) \right\}$$

$$\hat{p}_{MLE} = \frac{x}{n}$$

In our example:  $\hat{p} = 10/14 \approx 0.71$

# The Bayesian Theorem

- We are interested in the values of the unknown parameter  $\theta$
- The uncertainty about the parameter can be modelled through a probability distribution,  $p(\theta)$ , called the **prior distribution**. It will normally represents our *a priori* beliefs about the unknown parameter
- We have some relevant data, suppose we have  $n$  observations  $x = (x_1, \dots, x_n)$  which have a probability distribution that depends on the unknown parameter value:  $p(x|\theta)$  (the **likelihood**).
- After observing the data, we update the beliefs about the unknown parameter.



# The Bayesian Theorem

- The belief about  $\theta$  are updated by applying the **Bayes Theorem** to random variables

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)}$$

$$p(\theta|x) \propto p(x|\theta)p(\theta)$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

where  $p(x) = \int_{\theta} P(x|\theta)p(\theta)d\theta$ . So,  $p(x)$  is the normalising constant, and its distribution does not depend on the value of  $\theta$ .

- The inference about  $\theta$  is made is based on the distribution of  $\theta$  conditional on the data,  $p(\theta|x)$ , called the **posterior distribution**

# Proper posterior distribution

- To perform valid inference, the posterior distribution must be proper. For the continuous case that is  $\int_{\theta} p(\theta|x) d\theta = 1$ .
- Often we do not need to calculate the normalising constant because we can recognize the form of  $p(x|\theta)p(\theta)$  as a probability distribution that we know.
- Sometimes we can derive analytically the posterior distribution
- Other times we can obtain random draws from the posterior distribution of the parameter by using Markov Chain Monte Carlo methods (Gibbs sampling, Metropolis Hasting algorithm, etc.).
- Sometimes the posterior distribution belongs to the same probability distribution family as the prior distribution

## Conjugacy

If the posterior distribution belongs to the same probability distribution family as the prior distribution, the prior and posterior are then called **conjugate distributions**, and the prior is called a conjugate prior for the likelihood.

# The Beta-Binomial model

## Example (Bayesian Data Analysis, Gelman et al., 2003)

The proportion of births that are female has long been a topic of interest both scientifically and to the lay public. Two hundred years ago it was established that the proportion of female birth in European populations was less than 0.5, while the currently accepted value of the proportion in very large European-race populations is 0.485

- Let  $x$  be the number of girls in  $n$  recorded births. The data can be modelled with the binomial distribution, that is  $X \sim B(n, \theta)$ , where  $\theta$  is the probability of a female birth. Therefore, the likelihood is of the form

$$p(x|\theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$$

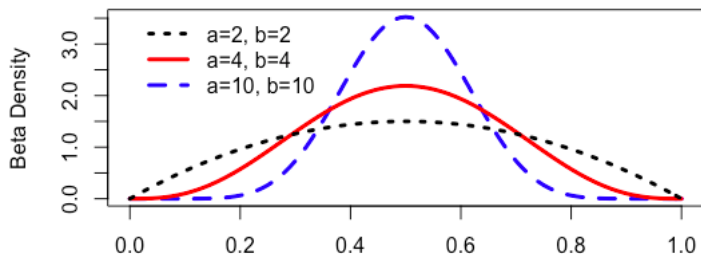
- What about the prior? Well, we know that  $\theta$  is a probability and therefore  $0 \leq \theta \leq 1$ . It is also likely to take a value near 0.50, and less likely to be far away from 0.5. A good candidate will be the Beta distribution

# The Beta-Binomial model

$$p(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \quad \text{where } \theta \in [0, 1]$$

$$E(\theta) = \frac{\alpha}{\alpha + \beta} \quad \text{and} \quad V(\theta) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

where  $\alpha > 0$  and  $\beta > 0$  are called hyperparameters





# The Beta-Binomial model

The posterior density for  $\theta$  is then

$$\begin{aligned}p(\theta|x) &\propto p(x|\theta)p(\theta) \\&\propto \binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\&\propto \theta^x (1-\theta)^{n-x} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\&\propto \theta^{\alpha+x-1} (1-\theta)^{\beta+n-x-1}\end{aligned}$$

Note that this looks like a Beta function with updated parameters, that is

$$P(\theta|x) = \text{Beta}(\alpha + x, \beta + n - x)$$

- As the posterior distribution follows the same parametric form as the prior distribution this is a case of conjugacy.
- Note that  $\alpha$  and  $\beta$  can be seen as the prior number of success (female birth in our example) and failures (males births)

# More cases of conjugacy

Likelihood	Conjugate Prior
Binomial	Beta
Geometric	Beta
Multinomial	Dirichlet
Poisson	Gamma
Normal ( $\sigma$ known)	Normal
Normal ( $\mu$ known)	Inverse Chi-Square
Normal ( $\mu$ known)	Gamma
Exponential	Gamma