

# BAYESIAN INFERENCE

PHD COURSE IN STATISTICAL INFERENCE

MATTIAS VILLANI

**DEPARTMENT OF STATISTICS**

**STOCKHOLM UNIVERSITY**

**AND**

**DEPARTMENT OF COMPUTER AND INFORMATION SCIENCE**

**LINKÖPING UNIVERSITY**

- **Introduction to subjective probability and Bayesian inference**
- **Prediction, Decisions, Exponential family**
- **Bayesian large sample theory and posterior approximation**
- **Bayesian computations**
- **Bayesian model comparison**

- **Subjective probability**
- **The Bayesics**
- **Bayes for the Normal model**
- **Priors**

# THE LIKELIHOOD FUNCTION

- The likelihood function is  
the probability of the observed data  
considered as a function of the parameter.
- Likelihood function is **NOT** a probability distribution for  $\theta$ .
- Statements like  $\Pr(\theta > c)$  makes no sense.
- Unless ...

# UNCERTAINTY AND SUBJECTIVE PROBABILITY

- $\Pr(\theta < 0.6 | \text{data})$  only makes sense if  $\theta$  is random.
- But  $\theta$  may be a fixed natural constant?
- **Bayesian: doesn't matter if  $\theta$  is fixed or random.**
- Do **You** know the value of  $\theta$  or not?
- $p(\theta)$  reflects Your knowledge/**uncertainty** about  $\theta$ .
- **Subjective probability.**
- The statement  $\Pr(10\text{th decimal of } \pi = 9) = 0.1$  makes sense.

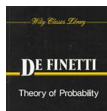


# UNCERTAINTY AND SUBJECTIVE PROBABILITY

*“The only relevant thing is uncertainty - the extent of our knowledge and ignorance. The actual fact of whether or not the events considered are in some sense determined, or known by other people, and so on, is of no consequence.” - Bruno de Finetti*

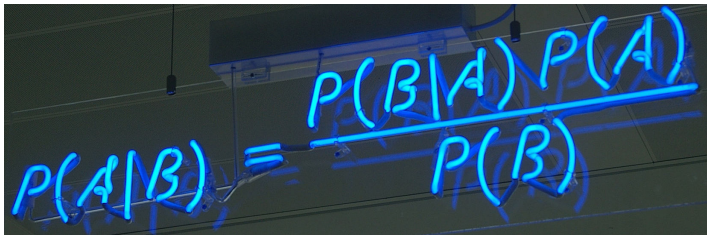
*“**Probability does not exist**” - Bruno de Finetti in the introduction to his classic book *A Theory of Probability**

- Subjective probability applies also to **non-repeatable experiments**.
- Subjective probabilities must satisfy the usual **axioms for probabilities**. Dutch book arguments. Axiomatic theory.



- **Bayesian learning** about a model parameter  $\theta$ :
  - state your **prior** knowledge as a probability distribution  $p(\theta)$ .
  - collect **data**  $\mathbf{x}$  and form the **likelihood** function  $p(\mathbf{x}|\theta)$ .
  - **combine** prior knowledge  $p(\theta)$  with data information  $p(\mathbf{x}|\theta)$ .
- **How to combine** the two sources of information?

## Bayes' theorem


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

- How to **update** from **prior**  $p(\theta)$  to **posterior**  $p(\theta|Data)$ ?
- **Bayes' theorem** for events  $A$  and  $B$

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}.$$

- Bayes' Theorem for a model parameter  $\theta$

$$p(\theta|Data) = \frac{p(Data|\theta)p(\theta)}{p(Data)}.$$

- It is the prior  $p(\theta)$  that takes us from  $p(Data|\theta)$  to  $p(\theta|Data)$ .
- A probability distribution for  $\theta$  is extremely useful.  
**Predictions. Decision making.**



# GREAT THEOREMS MAKE GREAT TATTOOS

- Bayes theorem

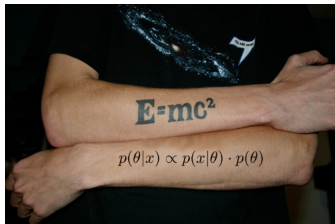
$$p(\theta|Data) = \frac{p(Data|\theta)p(\theta)}{p(Data)}$$

- All you need to know:

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

or

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



## ■ Model

$$x_1, \dots, x_n | \theta, \sigma^2 \stackrel{iid}{\sim} N(\theta, \sigma^2).$$

## ■ Prior

$$p(\theta) \propto c \text{ (a constant)}$$

## ■ Likelihood

$$p(x_1, \dots, x_n | \theta, \sigma^2) \propto \exp \left[ -\frac{1}{2(\sigma^2/n)} (\theta - \bar{x})^2 \right]$$

## ■ Posterior

$$\theta | x_1, \dots, x_n \sim N(\bar{x}, \sigma^2/n)$$

## ■ Prior

$$\theta \sim N(\mu_0, \tau_0^2)$$

## ■ Posterior

$$\begin{aligned} p(\theta|x_1, \dots, x_n) &\propto p(x_1, \dots, x_n|\theta, \sigma^2)p(\theta) \\ &\propto N(\theta|\mu_n, \tau_n^2), \end{aligned}$$

where

$$\begin{aligned} \frac{1}{\tau_n^2} &= \frac{n}{\sigma^2} + \frac{1}{\tau_0^2}, \\ \mu_n &= w\bar{x} + (1-w)\mu_0, \end{aligned}$$

and

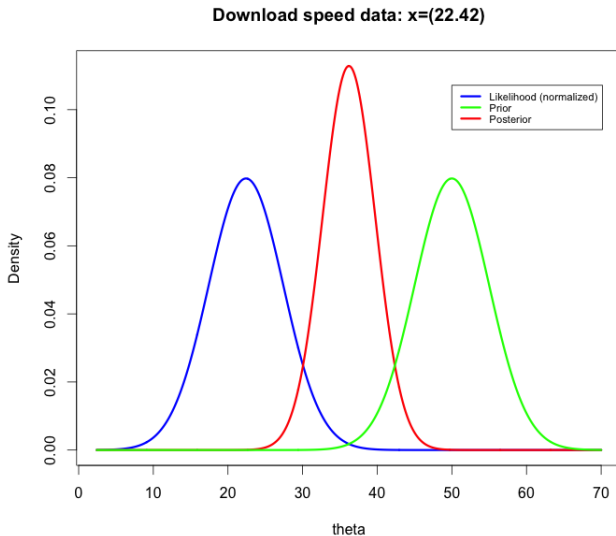
$$w = \frac{\frac{n}{\sigma^2}}{\frac{n}{\sigma^2} + \frac{1}{\tau_0^2}}.$$

■ Proof: complete the squares in the exponential.

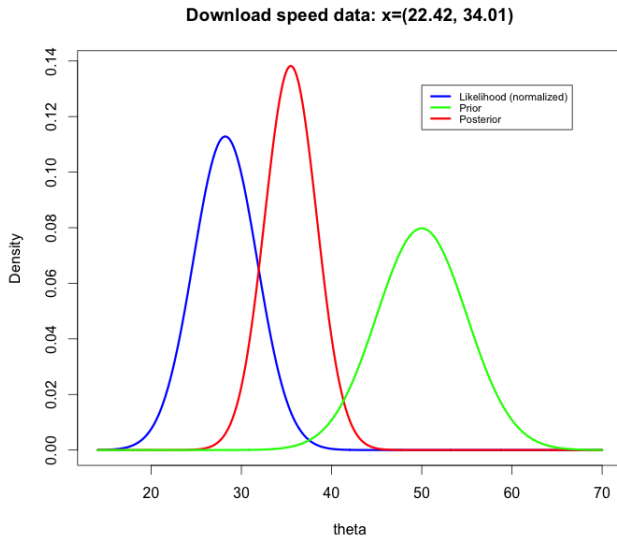
## EXAMPLE: AM I REALLY GETTING MY 50MBIT/SEC?

- **Data:**  $x = (22.42, 34.01, 35.04, 38.74, 25.15)$  Mbit/sec.
- **Model:**  $X_1, \dots, X_5 \sim N(\theta, \sigma^2)$ .
- Assume  $\sigma = 5$  (measurements can vary  $\pm 10$  MBit with 95% probability)
- My **prior:**  $\theta \sim N(50, 5^2)$ .

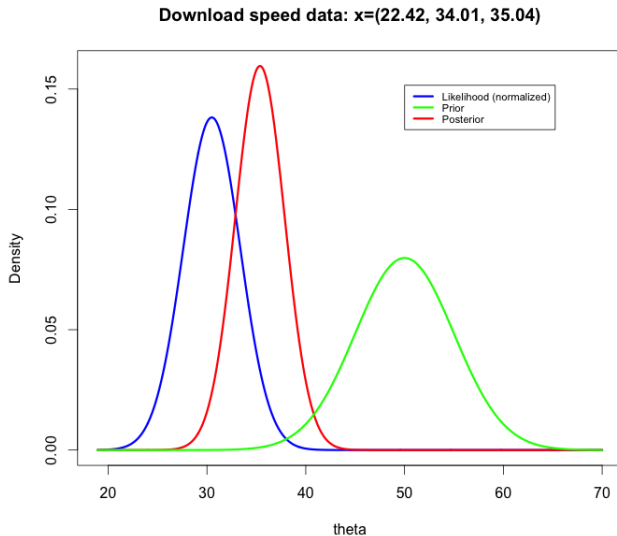
# DOWNLOAD SPEED N=1



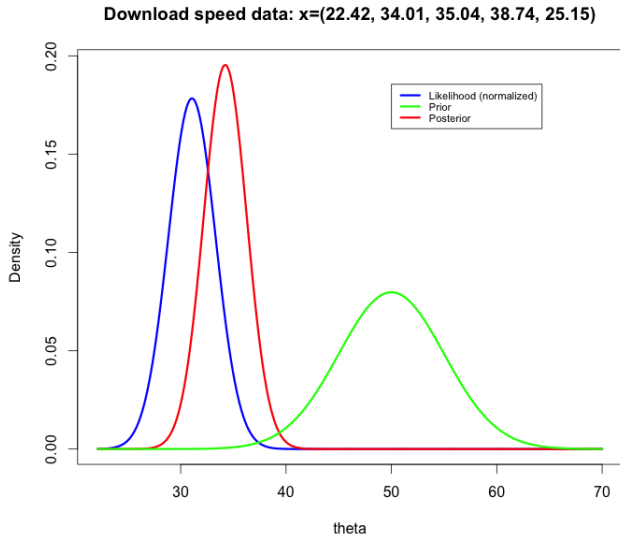
# DOWNLOAD SPEED N=2



# DOWNLOAD SPEED $N=3$



# DOWNLOAD SPEED N=5





- Models with **multiple parameters**  $\theta_1, \theta_2, \dots$
- Examples:  $x_i \stackrel{iid}{\sim} N(\theta, \sigma^2)$ ; multiple regression ...
- **Joint posterior distribution**

$$p(\theta_1, \theta_2, \dots, \theta_p | y) \propto p(y | \theta_1, \theta_2, \dots, \theta_p) p(\theta_1, \theta_2, \dots, \theta_p).$$

$$p(\theta | y) \propto p(y | \theta) p(\theta).$$

- **Marginalize** out parameter of no direct interest (**nuisance**).
- Example:  $\theta = (\theta_1, \theta_2)'$ . **Marginal posterior** of  $\theta_1$

$$p(\theta_1 | y) = \int p(\theta_1, \theta_2 | y) d\theta_2 = \int p(\theta_1 | \theta_2, y) p(\theta_2 | y) d\theta_2.$$

# NORMAL MODEL - NORMAL PRIOR

## ■ Model

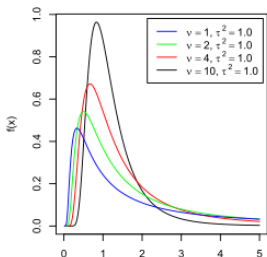
$$y_1, \dots, y_n | \theta, \sigma^2 \stackrel{iid}{\sim} N(\theta, \sigma^2)$$

## ■ Conjugate prior

$$\theta | \sigma^2 \sim N\left(\mu_0, \frac{\sigma^2}{\kappa_0}\right)$$

$$\sigma^2 \sim \text{Inv-}\chi^2(\nu_0, \sigma_0^2)$$

## ■ Scaled inverse $\chi^2$ distribution



## ■ Posterior

$$\theta | \mathbf{y}, \sigma^2 \sim N \left( \mu_n, \frac{\sigma^2}{\kappa_n} \right)$$
$$\sigma^2 | \mathbf{y} \sim \text{Inv-}\chi^2(\nu_n, \sigma_n^2).$$

where

$$\begin{aligned}\mu_n &= \frac{\kappa_0}{\kappa_0 + n} \mu_0 + \frac{n}{\kappa_0 + n} \bar{y} \\ \kappa_n &= \kappa_0 + n \\ \nu_n &= \nu_0 + n \\ \nu_n \sigma_n^2 &= \nu_0 \sigma_0^2 + (n - 1) s^2 + \frac{\kappa_0 n}{\kappa_0 + n} (\bar{y} - \mu_0)^2.\end{aligned}$$

## ■ Marginal posterior

$$\theta | \mathbf{y} \sim t_{\nu_n} (\mu_n, \sigma_n^2 / \kappa_n)$$

## ■ Binomial sampling:

$$s|\theta \stackrel{iid}{\sim} \text{Bin}(n, \theta), \text{ } n \text{ fixed.}$$

## ■ Prior

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

## ■ Posterior

$$\begin{aligned} p(\theta|s) &\propto p(s|\theta)p(\theta) \\ &= \binom{n}{s} \theta^s (1-\theta)^{n-s} \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &\propto \theta^{\alpha+s-1} (1-\theta)^{\beta+n-s-1} \end{aligned}$$

$$\text{so, } \theta|s \sim \text{Beta}(\alpha + s, \beta + n - s).$$

## ■ Negative binomial sampling:

$$n|\theta \stackrel{iid}{\sim} \text{Negbin}(s, \theta), s \text{ fixed.}$$

## ■ Posterior

$$\begin{aligned} p(\theta|n) &\propto p(n|\theta)p(\theta) \\ &= \binom{n-1}{s-1} \theta^s (1-\theta)^f \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \\ &\propto \theta^{\alpha+s-1} (1-\theta)^{\beta+f-1} \end{aligned}$$

so again,  $\theta|s \sim \text{Beta}(\alpha + s, \beta + f)$ .

## ■ Same posterior regardless of how the data was obtained.

## ■ Bayesian inference respects the **likelihood principle**:

$$p(\theta|\mathbf{x}) = \frac{c \cdot p(\mathbf{x}|\theta)p(\theta)}{\int c \cdot p(\mathbf{x}|\theta)p(\theta)d\theta} = \frac{p(\mathbf{x}|\theta)p(\theta)}{\int p(\mathbf{x}|\theta)p(\theta)d\theta}$$

for any  $c > 0$ .

- The prior should be determined (**elicited**) by an **expert**. Typically, expert  $\neq$  statistician.
- Elicit the prior on **a quantity that the expert knows well**. Convert afterwards.
- **Ask probabilistic questions** to the expert:
  - $E(\theta) = ?$
  - $SD(\theta) = ?$
  - $Pr(\theta < c) = ?$
  - $Pr(y > c) = ?$
- **Show some consequences** of the elicited prior to the expert.
- Beware of **psychological effects**, such as anchoring.

- **Autoregressive process** of order  $p$

$$y_t = \mu + \phi_1(y_{t-1} - \mu) + \dots + \phi_p(y_{t-p} - \mu) + \varepsilon_t, \varepsilon_t \stackrel{iid}{\sim} N(0, \sigma^2)$$

- Informative prior on the unconditional mean:  $\mu \sim N(\mu_0, \tau_0^2)$ .
- “Noninformative” prior on  $\sigma^2$ :  $p(\sigma^2) \propto 1/\sigma^2$
- Assume  $\phi_i \sim N(\mu_i, \psi_i)$ ,  $i = 1, \dots, p$  are independent a priori.
- Prior on  $\phi = (\phi_1, \dots, \phi_p)$  centered on persistent AR(1) process:  
 $\mu_1 = 0.8, \mu_2 = \dots = \mu_p = 0$
- $\text{Var}(\phi_i) = \frac{c}{i^\lambda}$ . “Longer” lags are more likely to be zero a priori.

# DIFFERENT TYPES OF PRIOR INFORMATION

- Real **expert information**. Combo of previous studies and experience.
- **Vague prior** information.
- **Reporting priors**. Easy to understand the information they contain.
- **Smoothness priors**. Regularization. Shrinkage. Big thing in modern statistics/machine learning.



## ■ Observed information

$$J_{\theta, \mathbf{x}} = - \frac{\partial^2 \ln p(\mathbf{x}|\theta)}{\partial \theta^2} \Big|_{\theta=\hat{\theta}}$$

## ■ Fisher information

$$I_{\theta} = E_{\mathbf{x}|\theta} (J_{\theta, \mathbf{x}})$$

## ■ A common non-informative prior is **Jeffreys' prior**

$$p(\theta) = |I_{\theta}|^{1/2}.$$

- **Invariant** to 1:1 transformations of  $\theta$ .
- Often non-conjugate.
- Often problematic in multiparameter settings.

## JEFFREYS' PRIOR FOR BERNOULLI SAMPLING

$$x_1, \dots, x_n | \theta \stackrel{iid}{\sim} \text{Bern}(\theta).$$

$$\ln p(\mathbf{x}|\theta) = s \ln \theta + f \ln(1 - \theta)$$

$$\frac{d \ln p(\mathbf{x}|\theta)}{d\theta} = \frac{s}{\theta} - \frac{f}{(1 - \theta)}$$

$$\frac{d^2 \ln p(\mathbf{x}|\theta)}{d\theta^2} = -\frac{s}{\theta^2} - \frac{f}{(1 - \theta)^2}$$

$$I(\theta) = \frac{E_{\mathbf{x}|\theta}(s)}{\theta^2} + \frac{E_{\mathbf{x}|\theta}(f)}{(1 - \theta)^2} = \frac{n\theta}{\theta^2} + \frac{n(1 - \theta)}{(1 - \theta)^2} = \frac{n}{\theta(1 - \theta)}$$

Thus, the Jeffreys' prior is

$$p(\theta) = |I(\theta)|^{1/2} \propto \theta^{-1/2}(1 - \theta)^{-1/2} \propto \text{Beta}(1/2, 1/2).$$

# JEFFREYS' PRIOR FOR NEGATIVE BINOMIAL SAMPLING

- Jeffreys' prior:

$$n|\theta \stackrel{iid}{\sim} \text{NegBin}(s, \theta).$$

$$\ln p(\mathbf{x}|\theta) = \ln \binom{n-1}{s-1} + s \ln \theta + f \ln(1-\theta)$$

$$\frac{d^2 \ln p(\mathbf{x}|\theta)}{d\theta^2} = -\frac{s}{\theta^2} - \frac{f}{(1-\theta)^2}$$

$$I(\theta) = \frac{s}{\theta^2} + \frac{E_{n|\theta}(n-s)}{(1-\theta)^2} = \frac{s}{\theta^2} + \frac{s/\theta - s}{(1-\theta)^2} = \frac{s}{\theta^2(1-\theta)}$$

- Thus, the Jeffreys' prior is

$$p(\theta) = |I(\theta)|^{1/2} \propto \theta^{-1}(1-\theta)^{-1/2} \propto \text{Beta}(\theta|0, 1/2).$$

- Jeffreys' prior is **improper**, but the posterior is proper:  
 $\theta|n \sim \text{Beta}(s, f + 1/2)$  which is proper since  $s \geq 1$ .
- Jeffreys' prior **violates the likelihood principle** because  $I(\theta)$  is sampling-based.

- **Maximum entropy prior:** choose prior with maximum entropy (most uncertain). Problematic for continuous parameters.
- **Reference prior:** Choose the prior that maximizes the expected value of perfect information about  $\theta$ .
- Under the conditions that guarantees asymptotic normality of posterior: Reference = Jeffreys.