## Statistical inference (MVE155/MSG200)

Bayesian inference

<ロト <回ト < 国ト < 国ト < 国ト < 国ト < 国 > 日 > 1/18

Frequentistic approach: Data x are generated from some population distribution  $f(x|\theta)$ , where  $\theta$  is an unknown (constant) parameter.

Bayesian approach:

- ▶ Parameters of interest are treated as random variables and generated from some prior distribution  $g(\theta)$ .
- Given  $\theta$ , data has the distribution or likelihood  $f(x|\theta)$ .
- ► Parameters are estimated by finding the posterior distribution  $h(\theta|x)$ .

We have two events A and B, where  $P(A) \neq 0$  and  $P(B) \neq 0$ . The Bayes theorem says that

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

Also, for random variables X and Y with density (or probability mass) functions  $f_X$  and  $f_Y$ , respectively, and  $f_X(x) \neq 0$ ,  $f_Y(y) \neq 0$ ,

$$f_{X|Y}(x|y) = rac{f_{Y|X}(y|x)f_X(x)}{f_Y(y)}.$$

## Posterior distribution

Given a prior distribution  $g(\theta)$  and likelihood  $f(x|\theta)$ , the posterior distribution  $h(\theta|x)$  can be computed by using the Bayes theorem:

$$h( heta|x) = rac{f(x| heta)g( heta)}{\phi(x)},$$

where

$$\phi(x) = \int f(x|\theta)g(\theta) \, d\theta$$
 or  $\phi(x) = \sum P(X = x|\theta)g(\theta)$ 

depending on whether X is continuous or discrete. This gives that

posterior  $\propto$  likelihood  $\times$  prior.

- ► We choose the prior.
- If we do not have any prior information on the parameter(s), we can choose uninformative, uniform priors.
- If we have some prior information, we can take it into account when choosing the prior.
- ▶ The prior should be chosen before the data are collected.

A sample  $x_1, ..., x_n$  from a normal distribution with known variance  $\sigma^2$ .

We choose  $N(\mu_0, \sigma_0)$  as the prior distribution  $g(\theta)$  for the mean  $\theta$  and the likelihood

$$f(x_1,...,x_n|\theta) = \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2}\sum_{i=1}^n (x_i-\theta)^2\right).$$

The posterior distribution  $h(\theta|x) \propto f(x_1, ..., x_n|\theta)g(\theta)$  is also normal.

Let the data be generated from a parametric model having the likelihood  $f(x|\theta)$  and let us have a parametric family of prior distributions  $\mathcal{G}$ .

Then  $\mathcal{G}$  is called a family of conjugated priors for the likelihood function  $f(x|\theta)$  if for any prior  $g(\theta) \in \mathcal{G}$ , the posterior

 $h(\theta|x) \propto f(x|\theta)g(\theta)$ 

also belongs to  $\mathcal{G}$ .

Normal distributions are conjugated priors for normal distributions when estimating the mean.

Model for the data	$\theta$	Prior	Posterior
$N(\mu, \sigma)$	$\mu$	$N(\mu_0, \sigma_0)$	$N(\gamma_n\mu_0 + (1-\gamma_N)\bar{x}, \sigma_0\sqrt{\gamma_n})$
Bin(n, p)	р	Beta( <i>a</i> , <i>b</i> )	Beta(a+x, b+n-x)
$Pois(\mu)$	$\mu$	$Gam(\alpha_0,\lambda_0)$	$Gam(\alpha_0 + n\bar{x}, \lambda_0 + n)$
$\textsf{Gam}(lpha,\lambda)$	$\lambda$	$Gam(\alpha_0,\lambda_0)$	$Gam(\alpha_0 + \alpha n, \lambda_0 + n\bar{x})$

Above, 
$$\gamma_n = \frac{\sigma^2}{\sigma^2 + n\sigma_0^2}$$
.

Note that as n increases the posterior becomes less effected by the prior.

Data:  $X \sim Bin(n, p)$ , where  $X = X_1 + ... + X_n$  and each  $X_i \sim Bin(1, p)$ .

Task: Estimate p using a Beta(a, b) prior, which has the density function

$$g(p) = rac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} p^{a-1} (1-p)^{b-1}, \quad 0$$

where a > 0 and  $\beta > 0$ ,

## Beta distribution

-

0

0.0 0.2 0.4

 $\sim$ ~ Beta(1,1) Beta(0.5,0.5) Beta(1,3) 9 9 φ ŝ ŝ ŝ 4 ž ž ž ĉ ŝ c N 2 N ~ ~ -0 0 0 0.0 0.2 0.4 0.6 0.8 1.0 0.0 0.2 0.4 0.6 0.8 1.0 0.0 0.2 0.4 0.6 0.8 1.0 x х х ~ ~ Beta(12,40) -Beta(3,0.5) Beta(3,9) 9 9 9 ŝ ß ŝ -4 4 ž ž ž ĉ c ŝ N 2 N

0.0 0.2 0.4 0.6 0.8 1.0

х

~

0

0.0

0.2 0.4 0.6 0.8 1.0

<ロ> <回<sub>2</sub> < 回<sub>2</sub> < 回 > < 回 >

<del>.</del>

0

0.6 0.8 1.0

x

≣ ∽ < ⊂ 10 / 18 A point estimate *a* for the paramter  $\theta$  is chosen by minimizing the posterior risk (given the data)

 $R(a|x) = \mathbb{E}(I(\Theta, a)|x)$ 

which is computed by using the posterior distribution, i.e.

$${\it R}({\it a}|{\it x}) = \int {\it I}({\scriptstyle heta},{\it a}) {\it h}({\scriptstyle heta}|{\it x}) \, d{\scriptstyle heta}, \qquad \left( {
m or} \, \sum_{\scriptstyle heta} {\it I}({\scriptstyle heta},{\it a}) {\it h}({\scriptstyle heta}|{\it x}) 
ight)$$

where *I* is a loss function, for example

► zero-loss function: *l*(θ, a) = 1<sub>θ≠a</sub> (maximum a posteriori (map), the value that maximizes the posterior, posterior mode)

▶ squared loss:  $I(\theta, a) = (\theta - a)^2$  (posterior mean)

A parameter is a random variable  $\Theta$  having the (posterior) distribution  $h(\theta|x)$  and we can compute  $100(1-\alpha)$ % credibility intervals for  $\Theta$ . They are of the form

 $J_{\theta} = (b_1(x), b_2(x))$ 

such that

 $P(b_1(x) < \Theta < b_2(x)|x) = 1 - \alpha.$ 

Consider the case of two simple hypotheses

 $H_0: \theta = \theta_0$  versus  $H_1: \theta = \theta_1$ .

Likelihood functions connected to these hypotheses are  $f(x|\theta_0)$ and  $f(x|\theta_1)$  and priors  $P(H_0) = \pi_0$  and  $P(H_1) = \pi_1 = 1 - \pi_0$ . The rejection region  $\mathcal{R}$  and whether to reject the null hypothesis is decided based on the cost function:

	Decision	<i>H</i> <sub>0</sub> true	$H_1$ true
$x \notin \mathcal{R}$	Do not reject $H_0$	0	cost <sub>1</sub>
$x \in \mathcal{R}$	Reject <i>H</i> 0	cost <sub>0</sub>	0

Here,

- cost<sub>0</sub> is the cost for the type I error
- cost<sub>1</sub> the cost for the type II error

The rejection region is chosen by minimizing the average cost (weighted mean of  $cost_0$  and  $cost_1$ )

```
\operatorname{cost}_0 \pi_0 \operatorname{P}(X \in \mathcal{R}|H_0) + \operatorname{cost}_1 \pi_1 \operatorname{P}(X \notin \mathcal{R}|H_1).
```

This leads to rejecting  $H_0$  if

 $\frac{f(x|\theta_0)}{f(x|\theta_1)} < \frac{\operatorname{cost}_1 \pi_1}{\operatorname{cost}_0 \pi_0},$ 

where  $\pi_0/\pi_1$  is called the prior odds and  $cost_1/cost_0$  the cost ratio. Equivalently,  $H_0$  is rejected if

$$\frac{h(\theta_0|x)}{h(\theta_1|x)} < \frac{\text{cost}_1}{\text{cost}_0}.$$

The person N, who is charged for rape, is a male of age 37 living in the area not very far from the crime scene. The jury has to decide whether the person is innocent ( $H_0$ : N is innocent) or guilty ( $H_1$ : N is guilty).

There are three conditionally independent pieces of evidence:

- E1: a DNA match
- ► E2: defendant N is not recognised by the victim
- E3: an alibi supported by the N's girlfriend.

## Example (compendium)

The reliability of E1-E3 was quantified as

▶  $P(E1|H_0) = 1/200,000,000 \text{ and } P(E1|H_1) = 1$ → very strong evidence for  $H_1$ ,  $P(E1|H_0)/P(E1|H_1) = 1/200,000,000$ 

▶  $P(E2|H_0) = 0.9$  and  $P(E2|H_1) = 0.1$ → strong evidence for  $H_0$ ,  $P(E2|H_0)/P(E2|H_1) = 9$ 

•  $P(E3|H_0) = 0.5$  and  $P(E3|H_1) = 0.25$ 

 $\rightarrow$  some evidence for  $H_0$ ,  $P(E3|H_0)/P(E3|H_1) = 2$ 

The non-informative prior probability

 $\pi_1 = P(H_1) = 1/200,000$ 

was used (thinking about the number of males who theoretically could have committed the crime without any evidence taken into account). Posterior odds become

 $\frac{\mathrm{P}(H_0|E_1, E_2, E_3)}{\mathrm{P}(H_1|E_1, E_2, E_3)} = \frac{\mathrm{P}(E_1|H_0)\mathrm{P}(E_2|H_0)\mathrm{P}(E_3|H_0)\pi_0}{\mathrm{P}(E_1|H_1)\mathrm{P}(E_2|H_1)\mathrm{P}(E_3|H_1)\pi_1} = 0.018.$ 

The person N would be found guilty if the cost values assigned by the jury were such that

 $\frac{\text{cost}_1}{\text{cost}_0} = \frac{\text{cost for unpunished crime}}{\text{cost for punishing an innocent}} > 0.018.$