Bayesian Statistics for Computer intensive methods sf2955

Timo Koski





- (a) Statistical Inference
- (b) Parametric statistical model
- (c) Bayes rule, Posterior Distributions
- (d) Bayes' Billiard Balls
- (e) Predictive Distributions
- (f) Asymptotic form of the posterior density





By learning/inference from data one often means the process of inferring a general law or principle from the observations of particular instances. The general law is a piece of knowledge about the mechanism of nature that generates the data.





The intended learning is done by use of 'MODELS', which serve as the language in which the constraints predicated on the data can be described. We deal here with parametric statistical models.



4 / 87



x is an observation of a random variable (X).

 $x \in f(x|\theta)$

 $f(x|\theta)$ is a probability density on R^{p} . $f(x|\theta)$ is a known function of x and θ .

 θ is an unknown parameter.





x can be of the form $x^{(n)} = (x_1, \ldots, x_n)$. x can be continuous or discrete variate or a mix thereof.

 θ is an unknown parameter $\in \Theta$ = a vector space of finite dimension. Hence we exclude, e.g., non-parametric statistics.





$$x \in f(x|\theta)$$

- x is distributed according to f,
- x is an observation from the distribution f.

An outcome x of a random variable (r.v.) X.





Parametric statistical model: Examples; Normal distribution

$$\theta = (\mu, \sigma^2) \in \Theta = R \times (0, \infty) .$$
$$f(x|\theta) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{\sigma^2}(x-\mu)^2}, -\infty < x < \infty.$$

We say that x is an observation from the normal distribution $N(\mu, \sigma^2)$.





Parametric statistical model: Examples; Bernoulli Distribution

Consider r.v. X with values 0, 1, $0 < \theta < 1$ and

$$\begin{array}{ccc} x = 1 & x = 0 \\ f(x \mid \theta) & \theta & 1 - \theta \end{array}$$

then we say x is distributed according to the Bernoulli distribution with the parameter θ .

$$X = x \in Be(\theta),$$



 $f(x|\theta)$ is a probabilistic mechanism of generating data, characterizes the behaviour of future observations conditional on θ .





Retrieve the parameters of the probabilistic generating mechanism using x. $f(x|\theta)$ is a probabilistic generating mechanism of data, characterizes the behaviour of future observations conditional on θ , but in inference the roles of x and θ are inverted.





Since

$$p(A \mid E) \cdot p(E) = p(E \mid A) \cdot p(A)$$

we have in a formal way Bayes' Rule of inversion

$$p(A \mid E) = \frac{p(E \mid A) \cdot p(A)}{p(E)}.$$

$$p(E) = p(E \mid A)p(A) + p(E \mid A^{c})p(A^{c}).$$

 A^c is the complement set of A.

12 / 87

Bayes' rule extended to continuous random variables:

$$g(y|x) = \frac{f(x|y) \cdot g(y)}{\int f(x|y) \cdot g(y) \, dy'}$$

Due to the standardization g(y|x) is a probability density; $g(y|x) \ge 0$, $\int g(y|x) \, dy = 1$.





$$\pi\left(\theta|x\right) = \frac{f\left(x\mid\theta\right)\cdot\pi\left(\theta\right)}{\int_{\Theta}f\left(x\mid\theta\right)\cdot\pi\left(\theta\right)\,d\theta}$$

Terminology for Bayes' Rule:

- $\pi(\theta)$: prior distribution on Θ .
- $\pi(\theta|x)$: posterior distribution on Θ .
- $m(x) = \int_{\Theta} f(x \mid \theta) \cdot \pi(\theta) d\theta$: marginal distribution of x.



14 / 87

Uncertainty about the unknown θ is modeled by a probability distribution $\pi(\theta)$, and $\pi(\theta|x)$ expresses the uncertainty about the unknown θ after the observation of x.

We use probability as tool for all parts of our analysis. This is called <u>coherence</u>.

Mathematically: the unknown θ becomes a random variable. (x,θ) will have a joint distribution.



15 / 87

$$\pi\left(\theta|x\right) = \frac{f\left(x \mid \theta\right) \cdot \pi\left(\theta\right)}{m(x)},$$

Terminology:

-

2

The notation

$$\int_{\Theta} f(x \mid \theta) \cdot \pi(\theta) \, d\theta$$

is imprecise by intent, as it can mean both a single integral and a multiple integral.





A Bayesian parametric statistical model consists of

• a parametric model

 $x \in f(x|\theta)$

a prior distribution

 $\theta \in \pi(\theta)$

The quantity of interest

 $\theta | x \in \pi \left(\theta | x \right)$





Any function $\pi(\cdot)$ such that

 $\pi \left(\theta \right) \geq \mathbf{0,}$

and

$$\int_{\Theta}\pi\left(\theta\right) d\theta=1\text{,}$$

can serve as a prior distribution.



19 / 87

But even functions with the properties

$$\pi\left(heta
ight) \geq$$
 0,

and

$$\int_{\Theta}\pi\left(\theta\right) d\theta=\infty\text{,}$$

are also invoked as priors, and are called improper priors.



$$X_i \mid M = m \in N(m, \sigma_0^2), M \in N(\mu, s^2). x^{(n)} = (x_1, \dots, x_n)$$
 a sample of I.I.D. $X_i, \overline{x} = \frac{1}{n} \sum_{i=1}^n x_i.$

$$M \mid (X_1,\ldots,X_n) \in \mathbb{N}\left(\frac{n\overline{x}/\sigma_0^2 + \mu/s^2}{n/\sigma_0^2 + 1/s^2}, \frac{1}{n/\sigma_0^2 + 1/s^2}\right)$$

i.e., $\pi(m|x^{(n)})$ is the density of this normal distribution. Here μ and s^2 are hyperparameters.



21 / 87



•
$$P(a(x) \le \theta \le b(x)) = \int_{a(x)}^{b(x)} \pi(\theta|x) d\theta$$

• $\underbrace{P(a(x) \le \theta \le b(x))}_{a(x)}$

This is a probability, not a degree of confidence





Take the example above

$$N\left(\frac{n\overline{x}/\sigma_0^2 + \mu/s^2}{n/\sigma_0^2 + 1/s^2}, \frac{1}{n/\sigma_0^2 + 1/s^2}\right)$$

Let $s \to \infty$ (the prior becomes improper). Then

$$N\left(\frac{n\overline{x}/\sigma_0^2 + \mu/s^2}{n/\sigma_0^2 + 1/s^2}, \frac{1}{n/\sigma_0^2 + 1/s^2}\right) \to N\left(\overline{x}, \frac{\sigma_0^2}{n}\right)$$

But then the Bayesian confidence interval $P(a(x) \le \theta \le b(x)) = 0.95$ becomes the familiar $\overline{x} \pm \lambda_{0.025} \frac{\sigma_0}{\sqrt{n}}$, and the common mistaken but natural interpretation of the confidence interval is correct !

23 / 87

- inference is based on the observed *x*, not on an unobserved sample space.
- $\pi(\theta|x)$ is the only quantity evaluated for inference about θ .

On the other hand , evaluation of $\pi(\theta|x)$ is not in general possible by explicit means of integral calculus. This is where statistical inference needs Markov chain Monte Carlo (McMC).



24 / 87

The distribution $f(x \mid \theta)$ regarded as a function of θ is known as the *likelihood function*

$$I(x;\theta) = f(x \mid \theta).$$

The likelihood function $I(\theta|x)$ thus compares the plausibilities of different parameter values for given x.





$$ll(x;\theta) = -\log l(x;\theta).$$

is called the log likelihood function.





The information brought by x about θ is entirely contained in the likelihood function $I(\theta \mid x)$.

If x_1 and x_2 are two observations depending on the same parameter θ such that there exists a constant c such that

$$l_1(x_1;\theta) = c \cdot l_2(x_2;\theta)$$

for every θ , then they bring the same information about θ and must lead to identical inferences.

27 / 87

Conditions for the likelihood principle:

- inference should be about the same parameter set
- θ should include every unknown factor in the model.





$$\pi \left(\theta | x \right) = \frac{f\left(x \mid \theta \right) \cdot \pi \left(\theta \right)}{\int_{\Theta} f\left(x \mid \theta \right) \cdot \pi \left(\theta \right) d\theta},$$
$$= \frac{I\left(x; \theta \right) \cdot \pi \left(\theta \right)}{\int_{\Theta} I\left(x \mid \theta \right) \cdot \pi \left(\theta \right) d\theta}$$

Hence Likelihood Principle is satisfied by Bayesian inference. There are ways of implementing the likelihood principle: MLE and MAP \Rightarrow



29 / 87

The maximum likelihood estimate MLE, $\widehat{\theta}_{ML}$ of $\theta,$ is defined by

 $\widehat{\theta}_{\mathrm{ML}} = \operatorname{argmax}_{\theta \in \Theta} f\left(x \mid \theta \right)$

 $= \operatorname{argmin}_{\theta \in \Theta} II\left(x; \theta\right)$

MLE is a parameter value that gives the observed x the highest possible probability.



The maximum a posterior estimate MAP $\hat{\theta}_{MAP}$ of θ is defined by

$$\widehat{\theta}_{\mathrm{MAP}} = \mathrm{argmax}_{\theta \in \Theta} \pi\left(\theta \mid x\right)$$





A family \mathcal{F} of probability distributions on Θ is said to be **conjugate** or **closed under sampling** for a likelihood function

$$I(x;\theta) = f(x \mid \theta).$$

if for every $\pi \in \mathcal{F}$, the posterior distribution $\pi(\theta|x)$ also belongs to $\pi \in \mathcal{F}$.

Above we have seen an example with normal prior density on the mean.

32 / 87

An intuitive way of understanding conjugate priors is that with conjugate priors the prior knowledge can be translated into equivalent sample information. See, e.g.,

$$N\left(\frac{n\overline{x}/\sigma_{0}^{2}+\mu/s^{2}}{n/\sigma_{0}^{2}+1/s^{2}},\frac{1}{n/\sigma_{0}^{2}+1/s^{2}}\right).$$

Next we reconsider another problem with a conjugate family of priors.



33 / 87

Bayes' Billiard Ball



The square table is made in such a way that, if the White or Orange ball be thrown on it, the probability that it rest on any part of the plane is the same. First, the White ball is thrown, and suppose it rests on line og; then the Orange ball is strown it times. We define event H as any event the Orange ball rests between **A** and **0**, and not-M as its resting between **0** and **B**. What this means is this:

The first throw of the White ball determines the value of probability x (i.e. the probability of an unknown event) from a uniform distribution between 0 and 1; and then a series of trials, with probability x of success (i.e., M) is generated (this provides the data on which to infer the cornect value of x).

Then, thanks to the geometrical representation of the problem in the Figure, we can obtain the solution to the initial problem, by calculating integration. Although we have omitted mathematical formulas, the preceding is the central idea.



34 / 87

・ロト ・同ト ・ヨト ・ヨ

A billiard ball W is rolled on a line of length one, with a uniform probability of stopping anywhere. It stops at p, not disclosed to us. A second ball O is rolled n times under the same assumptions and X denotes the number of times O stops to the left of W. Given X = x, what inference can we make on p? (In the figure above $x \leftrightarrow p$.)



35 / 87

We let **P** be a random variable, whose values are denoted by p, $0 \le p \le 1$. Parametric statistical model for rolls of Bayes' Billiard Ball *O*: Conditional on **P** = p, the rolls are outcomes of I.I.D Be(p) R.V's.




Hence for x = 0, 1, 2, ..., n,

$$f(x|p) = P(X = x | \mathbf{P} = p)$$
$$= {n \choose x} p^{x} \cdot (1-p)^{n-x},$$

(the Binomial distribution)





Bayes' rule

$$\pi\left(p \mid x\right) = \frac{f\left(x \mid p\right) \cdot \pi\left(p\right)}{\int_{0}^{1} f\left(x \mid p\right) \cdot \pi\left(p\right) dp}, 0 \le p \le 1$$

and zero elsewhere. The marginal distribution of x is

$$m(x) = \int_0^1 f(x \mid p) \cdot \pi(p) \, dp.$$



The posterior $\pi(p \mid x)$ expresses our updated uncertainty of the 'true' position of W given the data X = x.





One way to get further from here is to use an explicit expression for $\pi(p)$. There are many possible choices (some more systematic choices outlined below), some have straightforward analytical advantages. Laplace assumed that $p \in R(0, 1)$. i.e.,

$$\pi\left(oldsymbol{p}
ight) = \left\{ egin{array}{cc} 1 & 0\leq oldsymbol{p}\leq 1 \ 0 & {
m elsewhere,} \end{array}
ight.$$



40 / 87

The marginal distribution of x: uniform prior

$$m(x) = \int_0^1 f(x \mid p) \cdot \pi(p) \, dp$$
$$= \binom{n}{x} \int_0^1 p^x \cdot (1-p)^{n-x} \, dp.$$

We use the Beta integral:





$$\int_0^1 p^{\alpha-1}(1-p)^{\beta-1}dp = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}.$$

Recall also that $\Gamma(x+1) = x!$, if x is a positive integer. $\alpha = \beta = 1$ gives the distribution R(0,1). We set

$$B(\alpha,\beta) := \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}.$$



42 / 87

$$\pi(p) = \begin{cases} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} & 0$$

is a probability density $Beta(\alpha, \beta)$. $\alpha > 0$ and $\beta > 0$ are hyperparameters.



43 / 87



The marginal distribution of *x*: uniform prior

$$m(x) = \int_0^1 f(x \mid p) \cdot \pi(p) \, dp$$
$$= \binom{n}{x} \frac{x!(n-x)!}{(n+1)!}$$



The marginal distribution of x, $p \in U(0, 1)$

$$m(x) = \int_0^1 f(x \mid p) \cdot dp = \binom{n}{x} \frac{x!(n-x)!}{(n+1)!}$$
$$= \frac{n!}{x!(n-x)!} \frac{x!(n-x)!}{(n+1)!} = \frac{1}{(n+1)}$$

There is an interpretation of Bayes' work claiming that the problem really attacked and solved by Bayes was: What should $\pi(p)$ be so that

$$\int_0^1 f(x \mid p) \cdot \pi(p) dp = \frac{1}{(n+1)}$$

holds for the Billiard Balls.

April 29, 2009

45 / 87

The Posterior Density for *n* rolls of Bayes' Orange Ball

$$\pi \left(p \mid x \right) = \frac{\binom{n}{x} p^{x} \cdot (1-p)^{n-x}}{m(x)}$$
$$= \begin{cases} \frac{(n+1)!}{x!(n-x)!} \cdot p^{k} \left(1-p\right)^{n-k} & 0 \le p \le 1\\ 0 & \text{elsewhere.} \end{cases}$$

This is again a Beta density, i.e., we have used a conjugate family of priors.



46 / 87

The Posterior Density for *n* rolls of Bayes' Orange Ball

$$\frac{(n+1)!}{x!(n-x)!} = \frac{\Gamma(n+2)}{\Gamma(x+1)\Gamma(n-x+1)} = \frac{1}{B(x+1,n-x+1)}.$$



47 / 87



The Posterior Density for *n* rolls of Bayes' Orange Ball



The Posterior Density with *n* rolls of Bayes Ball, $p \in Beta(\alpha, \beta)$

$$\pi\left(p \mid x\right) = \begin{cases} \frac{1}{B(x+\alpha, n-x+\beta)} \cdot p^{x+\alpha-1} \left(1-p\right)^{\beta+n-x-1} & 0 \le p \le 1\\ 0 & \text{elsewhere.} \end{cases}$$

This is the Beta density $Beta(\alpha + x, \beta + n - x)$.



49 / 87

The example of Bayes' Billiard Balls W and O is discontinued for the moment, but will be reconsidered later.





The probability distribution $\phi(x, y)$ is a joint distribution of (x, y). We assume a parametric model $(x, y, \theta) \in \phi(x, y, \theta)$ so that:

$$\phi(x,y) = \int \phi(x,y,\theta) \, d\theta.$$



The predictive distribution g(y|x) of y conditional on x, denoted by g(y|x) is

$$g(y|x) = \frac{\phi(x,y)}{m(x)} = \frac{\int \phi(x,y,\theta) \, d\theta}{m(x)}$$
$$\frac{\phi(x,y)}{m(x)} = \frac{\int \phi(x,y,\theta) \, d\theta}{m(x)} = \frac{\int g(y|x,\theta) \, \phi(x,\theta) \, d\theta}{m(x)}$$
$$= \frac{\int g(y|x,\theta) \, f(x|\theta) \, \pi(\theta) \, d\theta}{m(x)}$$



Predictive Distribution (2)

i.e.,

$$g(y|x) = \frac{\int g(y|x,\theta) f(x|\theta) \pi(\theta) d\theta}{m(x)}$$
$$= \frac{\int g(y|x,\theta) f(x|\theta) \pi(\theta) d\theta}{\int f(x|\theta) \cdot \pi(\theta) d\theta}$$
$$= \int g(y|x,\theta) \frac{f(x|\theta) \pi(\theta)}{\int f(x|\theta) \cdot \pi(\theta) d\theta} d\theta$$



The predictive distribution g(y|x) is thus

$$g(y|x) = \int g(y|x,\theta) \pi(\theta \mid x) d\theta.$$





If y and x are conditionally independent given θ , then

$$g(y|x,\theta) = g(y|\theta)$$

and

$$g(y|x) = \int g(y|\theta) \pi(\theta \mid x) d\theta.$$



Introduce order (e.g., in time). $x = x^{(n)} = (x_1, \dots, x_n) y = x_{n+1}$.

$$g\left(x_{n+1}|x^{(n)}\right) = \int g\left(x_{n+1}|x^{(n)},\theta\right) \pi\left(\theta \mid x^{(n)}\right) d\theta.$$

We assume conditional independence

$$\phi(x, y|\theta) = \phi\left(x^{(n)}, x_{n+1}|\theta\right) = (f(x_1|\theta) \cdots f(x_n|\theta)) \cdot f(x_{n+1}|\theta)$$
$$= \phi\left(x^{(n)}|\theta\right) \cdot f(x_{n+1}|\theta)$$



Predictive Distribution (6)

$$\phi\left(x^{(n)}, x_{n+1} \mid \theta\right) = \phi\left(x^{(n)} \mid \theta\right) \cdot f\left(x_{n+1} \mid \theta\right)$$

Then

$$g\left(x_{n+1}|x^{(n)},\theta\right) = \frac{\phi\left(x^{(n)},x_{n+1}\mid\theta\right)}{\phi\left(x^{(n)}\mid\theta\right)} = f\left(x_{n+1}\mid\theta\right)$$

 and

$$g\left(x_{n+1}|x^{(n)}\right) = \int f\left(x_{n+1}|\theta\right) \pi\left(\theta \mid x^{(n)}\right) d\theta$$



$$g\left(x_{n+1}|x^{(n)}\right) = \int f\left(x_{n+1}|\theta\right) \pi\left(\theta \mid x^{(n)}\right) d\theta$$

We often need to update $g(x_{n+1}|x^{(n)})$ for new observations. We need clearly only to update $\pi(\theta | x^{(n)})$.





Up-date of Posterior Distribution (1)

$$\pi \left(\theta \mid x^{(n+1)} \right) = \frac{f\left(x^{(n+1)} \mid \theta\right) \pi\left(\theta\right)}{\int f\left(x^{(n+1)} \mid \theta\right) \cdot \pi\left(\theta\right) d\theta}$$
$$= \frac{f\left(x_{n+1} \mid \theta\right) f\left(x_1 \mid \theta\right) \cdots f\left(x_n \mid \theta\right) \pi\left(\theta\right)}{\int f\left(x_{n+1} \mid \theta\right) f\left(x_1 \mid \theta\right) \cdots f\left(x_n \mid \theta\right) \cdot \pi\left(\theta\right) d\theta}$$
$$= \frac{f\left(x_{n+1} \mid \theta\right) \frac{f(x_1 \mid \theta) \cdots f(x_n \mid \theta) \pi(\theta)}{m(x^{(n)})}}{\int f\left(x_{n+1} \mid \theta\right) \frac{f(x_1 \mid \theta) \cdots f(x_n \mid \theta) \pi(\theta)}{m(x^{(n)})} d\theta}$$
$$= \frac{f\left(x_{n+1} \mid \theta\right) \pi\left(\theta \mid x^{(n)}\right)}{\int f\left(x_{n+1} \mid \theta\right) \pi\left(\theta \mid x^{(n)}\right) d\theta}$$



Timo Koski ()

April 29, 2009 59 / 87

$$\pi\left(\theta \mid x^{(n+1)}\right) = \frac{f\left(x_{n+1}\mid\theta\right)\pi\left(\theta\mid x^{(n)}\right)}{\int f\left(x_{n+1}\mid\theta\right)\pi\left(\theta\mid x^{(n)}\right)\,d\theta}$$

Hence, under the assumptions made, we can update posterior distribution in a sequential manner. Or, we can use the posterior of $\pi(\theta|x^{(n)})$ as a new prior for computing $\pi(\theta | x^{(n+1)})$.





The fundamental task of statistical inference is to pass from one set of observations x to express an opinion about another, as yet unobserved set y.

We have above accomplished this in terms of the predictive distribution

$$g(y|x) = \int f(y|\theta) \pi(\theta \mid x) d\theta,$$

which shows the role of learning about parameters in accomplishing this task



61 / 87

 $f(x \mid \theta) \leftrightarrow N(\mu, \sigma^2)(x)$, the mean μ and variance σ^2 are unknown, i.e., $\theta = (\mu, \sigma^2)$. The (improper) prior density is taken as

$$\pi\left(\theta\right) \propto d\mu \frac{1}{\sigma} d\sigma$$

Let $x = \underline{x}_n = (x^{(1)}, \dots, x^{(n)})$, $x^{(i)} \in \mathbb{N}(\mu, \sigma^2)$. We have the estimates $\widehat{\mu} = \frac{1}{n} \sum_{i=1}^n x^{(i)}$, and $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x^{(i)} - \widehat{\mu})^2$.



62 / 87

An Example of Bayesian prediction $f(x \mid \theta) \leftrightarrow N(\mu, \sigma^2)$ $\pi(\theta) \propto d\mu \frac{1}{\sigma} d\sigma$

The predictive density is

$$g\left(x^{(n+1)}|\underline{x}_{n}\right) = \sqrt{\frac{n}{(n^{2}-1)\pi}} \cdot \frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{1}{2}(n-1)s\right)} \cdot \left(1 + \frac{n\left(x^{(n+1)}-\widehat{\mu}\right)^{2}}{(n^{2}-1)s^{2}}\right)^{-n/2}$$



The predictive density

$$\sqrt{\frac{n}{(n^2-1)\pi}} \cdot \frac{\Gamma\left(\frac{n}{2}\right)}{\Gamma\left(\frac{1}{2}(n-1)s\right)} \cdot \left(1 + \frac{n\left(x^{(n+1)} - \widehat{\mu}\right)^2}{(n^2-1)s^2}\right)^{-n/2}$$

is known in Bayesian statistics as 't-like distribution'¹.

¹J.M. Dickey (1968): Three Multidimensional-Integral Identities with Bayesian Applications. *The Annals of Mathematical Statistics*, 39, pp. 1615–1628. (E)





Predictive Distribution for the next m rolls of Bayes' Orange Ball

The predictive distribution of y positions of O left of W in m additional rolls is

$$g(y|x) = \int_0^1 f(y|p) \pi(p \mid x) dp$$
$$= \binom{m}{y} \int_0^1 p^y \cdot (1-p)^{m-y} \pi(p \mid x) dp$$

 $y=0,1,\ldots,m$



65 / 87

Predictive Distribution for the next m rolls of Bayes' Ball

$$\int_{0}^{1} p^{y} \cdot (1-p)^{m-y} \pi (p \mid x) dp =$$

$$\int_{0}^{1} p^{y} \cdot (1-p)^{m-y} \frac{1}{B(x+1, n-x+1)} \cdot p^{x} (1-p)^{n-x} dp =$$

$$= \frac{1}{B(x+1, n-x+1)} \int_{0}^{1} p^{y+x} \cdot (1-p)^{m+n-y-x} dp =$$

$$= \frac{B(y+x+1, m+n-x-y+1)}{B(x+1, n-x+1)}.$$

by the Beta integral.

April 29, 2009

66 / 87

Predictive Distribution for the next m rolls of Bayes' Ball

$$g(y|x;m) = \binom{m}{y} \int_0^1 p^y \cdot (1-p)^{m-y} \pi(p \mid x) \, dp$$
$$= \binom{m}{y} \frac{B(y+x+1,m+n-x-y+1)}{B(x+1,n-x+1)}.$$
$$= \frac{m!}{(m-y)!y!} \frac{\Gamma(n+2)\Gamma(y+x+1)\Gamma(m+n-x-y+1)}{\Gamma(x+1)\Gamma(n-x+1)\Gamma(m+n+2)}.$$



67 / 87

Predictive Distribution for y = 1 in the next m = 1 roll of Bayes' Ball

$$g(1|x;1) = \frac{\Gamma(x+2)\Gamma(n-x+1)\Gamma(n+2)}{\Gamma(x+1)\Gamma(n-x+1)\Gamma(n+3)}$$
$$= \frac{(x+1)!(n-x)!(n+1)!}{x!(n-x)!(n+2)!} = \frac{x+1}{n+2}$$

A famous predictive probability, known as Laplace's rule of succession,

$$\frac{x+1}{n+2} = \frac{x+1}{(n+1)+1}$$

The prior knowledge can be translated into equivalent sample information.

April 29, 2009

68 / 87

The Maximum Likelihood Estimate of p in Bayes' Billiard

$$\begin{aligned} \widehat{p}_{\mathrm{ML}} &= \mathrm{argmin}_{0 \le p \le 1} II(p \mid x) \\ &= \mathrm{argmin}_{0 \le p \le 1} \left[-\log \left(\begin{array}{c} n \\ x \end{array} \right) - x \log p - (n - x) \log (1 - p) \right]. \\ &= \mathrm{argmin}_{0 \le p \le 1} \left(-x \log p - (n - x) \log (1 - p) \right). \\ &\Rightarrow \\ &\widehat{p}_{\mathrm{ML}} = \frac{x}{n} \end{aligned}$$

If you observed x = 0, would you belive in the estimate $\hat{p} = 0$ for all future purposes ?

The predictive probability found above

$$\frac{x+1}{n+2} = \frac{x+1}{(n+1)+1}$$

is a maximum likelihood estimate of p when n + 1 rolls of the ball O and the first roll of the ball W are included in the data.



70 / 87



- Assessment (by Questionnaries)
- Conjugate prior
- Non-informative or reference prior
 - Laplace's prior
 - Jeffreys' prior
- Maximum entropy prior





(One form of) Bayesian statistics relies upon a **personalistic theory of probability** for quantification of prior knowledge. In such a theory

- probability measures the confidence that a particular individual (assessor) has in the truth of a particular proposition
- no attempt is made to specify which assessments are correct
- personal probabilities should satisfy certain postulates of coherence.




R.L.Winkler (work published in

 Robert L. Winkler: The Assessment of Prior Distributions in Bayesian Analysis Journal of the American Statistical Association, Vol. 62, No. 319.

(Sep., 1967), pp. 776-800.)

devises questionnaires (or interviews) to elicit information to write down a prior distribution. Students of Univ. of Chigago were asked to, e.g., assess the uncertainty about the probability of a randomly chosen student of Univ. of Chigago being Roman Catholic using a probability distribution. The assessment was done by four different methods, like giving fractiles, making bets, assessing impact of additional data, drawing graphs. One interesting finding is that the assessments by the same person using different methods may be conflicting.



The priors in Winkler's study are not diffuse: the students of Univ. of Chigago have, since they have been around, an idea about the number of Roman Catholics at the campus of of Univ. of Chigago.



74 / 87



The interviews by Winkler were mathematically speaking all concerned with assessing the prior of θ in a Bernoulli Be (θ) – I.I.D. process. Winkler claims a sensitivity analysis (loc.cit p. 791) showing that the prior distributions assessed by the interviews yielded posterior distributions that were 'only little' different (by a test of goodness-of-fit) from those obtained from Beta densities on θ .





• A. O'Hagan: Eliciting Expert Beliefs in Substantial Practical Applications. *The Statistician*, 47, pp. 21–35, 1998.

Not only priors are elicited in

 R.L. Keeney & D. von Winterfeldt: Eliciting Probabilities in Complex Technical Problems. *IEEE Transactions on Engineering Management*, 38, pp.191–201, 1991.



76 / 87



But this line of study can evolve rapidly to a topic of research in psychology or (economic) behaviour, c.f.,

- C-A. S. Stael von Holstein: Assessment and Evaluation of Subjective Probability Distributions. 1970, Stockholm School of Economics.
- A. G. Wilson: Cognitive factors affecting subjective probability assessments.
 ISDS Discussion Paper # 94-02,

http://www.isds.duke.edu/

so it feels safe to leave the matter at rest here.



$$x_i| heta \in f(x| heta)$$
 , I.I.D. ,

or independent, identically, distributed conditional on θ

$$x^{(n)} = (x_1, x_2, \ldots, x_n) \in \mathcal{X}^n$$

 $f(x|\theta)$ is a probability density on R^{p} . $f(x|\theta)$ is a known function of x and θ . θ is an unknown parameter $\in \Theta$ = a vector space of finite dimension.

78 / 87

Let us assume that $f(x|\theta)$ is a density with a scalar parameter (for simplicity of notation), and that $f(x|\theta)$ is some $k \ge 2$ times differentiable in θ . We let $\hat{\theta}_{ML}$ be the maximum likelihood estimate of θ . We expand the log likelihood function around $\hat{\theta}_{ML}$

$$\log f\left(x^{(n)}|\theta\right) =$$

$$\log f\left(x^{(n)}|\widehat{\theta}_{ML}\right) + \left(\theta - \widehat{\theta}_{ML}\right)\frac{d}{d\theta}\log f\left(x^{(n)}|\widehat{\theta}_{ML}\right) \\ + \frac{1}{2}\left(\theta - \widehat{\theta}_{ML}\right)^{2}\frac{d^{2}}{d\theta^{2}}\log f\left(x^{(n)}|\widehat{\theta}_{ML}\right) + R_{n}\left(\theta\right)$$



But here $\hat{\theta}_{ML}$ is a solution of the equation

$$\frac{d}{d\theta}\log f\left(x^{(n)}|\widehat{\theta}_{ML}\right)=0$$

Hence

$$\log f\left(x^{(n)}|\theta\right) = \log f\left(x^{(n)}|\widehat{\theta}_{ML}\right) + \frac{1}{2}\left(\theta - \widehat{\theta}_{ML}\right)^2 \frac{d^2}{d\theta^2}\log f\left(x^{(n)}|\widehat{\theta}_{ML}\right) + R_n\left(\theta\right)$$



We have by assumption of I.I.D. data

$$\frac{d^2}{d\theta^2}\log f\left(x^{(n)}|\theta\right) = \sum_{l=1}^n \frac{d^2}{d\theta^2}\log f\left(x_l|\theta\right)$$

We set $Y_l = \frac{d^2}{d\theta^2} \log f(x_l|\theta)$. Then the Law of Large Numbers says that

$$\frac{1}{n}\sum_{l=1}^{n}Y_{l}\rightarrow E\left[Y\right],n\rightarrow\infty$$

where

$$E[Y] = \int_{\mathcal{X}} \frac{d^2}{d\theta^2} \log f(x|\theta) f(x|\theta) dx$$



The integral

$$I\left(\theta\right) = -\int_{\mathcal{X}} \frac{d^{2}}{d\theta^{2}} \log f\left(x|\theta\right) f\left(x \mid \theta\right) dx$$

is called Fisher information.



82 / 87



Then we may feel inclined to believe that

$$\frac{d^2}{d\theta^2}\log f\left(x^{(n)}|\widehat{\theta}_{ML}\right) = \sum_{l=1}^n \frac{d^2}{d\theta^2}\log f\left(x_l|\widehat{\theta}_{ML}\right) \approx -n \cdot I\left(\widehat{\theta}_{ML}\right)$$

Note that even $\hat{\theta}_{ML}$ depends on *n*.



83 / 87



This gives

$$\log f\left(x^{(n)}|\theta\right) \approx \log f\left(x^{(n)}|\widehat{\theta}_{ML}\right) - \frac{1}{2}\left(\theta - \widehat{\theta}_{ML}\right)^2 n \cdot I\left(\widehat{\theta}_{ML}\right)$$

The first term does not involve θ .





Then

$$f\left(x^{(n)}|\theta\right) \approx e^{-\frac{n}{2}\left(\theta - \widehat{\theta}_{ML}\right)^2 \cdot I\left(\widehat{\theta}_{ML}\right)}$$

The interpretation of the relation is that the likelihood function can be for large *n* be approximated by a normal density for which the mean is $\hat{\theta}_{ML}$ and the variance is $\frac{1}{nI(\hat{\theta}_{ML})}$.



85 / 87

Let $I(\theta)$ be the Fisher information of a parametric model. Take the prior density as

$$\pi(\theta) \stackrel{\text{def}}{=} \frac{\sqrt{I(\theta)}}{\int \sqrt{I(\theta)} d\theta},$$

assuming the integral exists. This choice of prior is known as Jeffreys' prior. This prior is invariant to monotonous transformations of θ .



It turns out that Jeffreys' prior for a binomial likelihood is obtained by $\alpha=1/2$ and $\beta=1/2$ in

$$\pi(p) = \begin{cases} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} & 0$$

and with $\Gamma(1/2+1/2)=1,\ \Gamma(1/2)=\sqrt{\pi},$

$$\pi_{Jeffreys}(p) = \begin{cases} \frac{1}{\pi} p^{-1/2} (1-p)^{-1/2} & 0$$

