

# Appendix

Christophe Ambroise

2021-2022

# Snooker Frequentist example I

Let consider a snooker table, and a ball. The ball  $A$  is launched perpendicularly to a reference edge and stop at a distance  $l$  from this billiard edge.

A second ball  $B$  is thrown  $n$  times and we denote by  $X$  the number of times that  $B$  stops at a distance  $l'$  of the edge such that  $l' > l$ .

We try to estimate the proportion  $p$  the number of times  $u$   $l' > l$  knowing that  $X = x$ . We assume that  $X \sim \mathcal{B}(n, p)$  (binomial distribution).

In this case we have a sample of size 1 and the likelihood of parameter  $p$  writes:

$$\ell(p; x) = C_n^x p^x (1 - p)^{n-x},$$

and by canceling the first derivative of the log-likelihood, we obtain

$$\hat{p}_{MV} = \arg \max_p \ell(p; x) = \frac{x}{n}$$

# Snooker Frequentist example II

If many launches have been performed, the estimate of  $p$  will be satisfactory. On the other hand, if only one throw is observed, we find:

- $x = 0$  gives  $\hat{p} = 0$ ;
- $x = 1$ , which gives  $\hat{p} = 1$ .

In both cases the estimate appears intuitively of very poor quality.

# Bayesian Statistics in a nutshell I

The maximum likelihood framework produces point estimates

## Reverend T. Bayes (1701-1761)

Two years after the death of the Reverend T. Bayes (1701-1761), a friend of this one, published his *essay in view of solving the doctrine of chances* (Bayes 1763).

## Parameters are random variables

In this little booklet, which is the source of the inference modern Bayesian statistic,

- Parameters are no longer treated as deterministic quantities but random as are observations.
- The dual role of the parameters  $\theta$  and the observations  $x$  is described thanks to conditioning by Bayes' theorem:

For a distribution (called *Prior*)  $\pi$  on the parameter  $\theta$ , and a observation  $x$  of density  $f(x|\theta)$ , the distribution of  $\theta$  conditionally on  $x$

$$\pi(\theta|x) = \frac{f(x|\theta)\pi(\theta)}{\int f(x|\theta)\pi(\theta)d\theta}.$$

The main innovation of the Bayesian statistical model is the law  $\pi$  on the model parameters.

Thus, in the context of a Bayesian statistical approach three functions must be specified:

- the law of the observations, the so-called likelihood  $f(x|\theta)$ ;
- the prior distribution on the parameters,  $\pi(\theta)$ ;
- the cost  $C$  associated with the decision  $\delta$  for parameters  $\theta$ .

Cost is a numerical measure of the quality of a decision.

We call the Bayes estimator associated with a prior distribution  $\pi$  and a cost  $C$ , any estimator  $\delta^\pi$  which, given an observation vector  $x$ , minimizes the cost a posteriori

$$\rho(\pi, \delta|x) = E^\pi[C(\theta, \delta)|x] = \int_{\theta} C(\theta, \delta) \pi(\theta|x) d\theta.$$

$$E_{\theta|x}[(\theta - \delta)^2] = E[\theta^2 - 2\delta\theta + \delta^2]$$

$$\hat{\theta} = \arg\min_{\delta} \rho(\pi, \delta|x) \quad \frac{\partial \rho(\pi, \delta|x)}{\partial \delta} = -2E[\theta|x] \stackrel{!}{=} 0$$

$$\delta = E_{\theta|x}[\theta]$$

When the prior and the posterior have the same form, we say that the prior is a conjugate prior for the corresponding likelihood. Conjugate priors are widely used because they simplify computation, and are easy to interpret

The posterior  $p(\theta|x)$  summarizes everything we know about the unknown quantities  $\theta$ .



We can easily compute a point estimate of an unknown quantity by computing the posterior mean, median or mode.

However, the posterior mode, aka the MAP estimate, is the most popular choice:

- it reduces to an optimization problem, for which efficient algorithms often exist.
- MAP estimation can be interpreted in non-Bayesian terms, by thinking of the log prior as a regularizer (see Lasso e.g.)

# Fonctions Gamma & Beta

$$y = (x-1)!$$

La fonction Gamma est "version" continue de

$$\Gamma(x) = (x-1)!$$

$$\Gamma(x) = \int_0^{+\infty} x^{x-1} e^{-x} dx$$

$$\text{si } x \in \mathbb{N}^+ \text{ alors } \Gamma(x) = (x-1)!$$

Montrons  
que  $\Gamma(x+1) = x \Gamma(x)$

$$\Gamma(x+1) = \int_0^{\infty} x^x e^{-x} dx$$

$$u = e^{-x} \quad u' = -e^{-x}$$
$$v = x^x \quad v' = x^{x-1}$$

$$= [uv]_0^{\infty} - \int_0^{\infty} u v'$$

$$= [-e^{-x} x^x]_0^{\infty} + \int_0^{\infty} e^{-x} x^{x-1} dx$$

$$= \underbrace{\lim_{x \rightarrow \infty} -e^{-x} x^x}_0 + x \underbrace{\int_0^{\infty} e^{-x} x^{x-1} dx}_{\Gamma(x)}$$

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$$

Montreons que  $B(\alpha, \beta) = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)}$

$$\begin{aligned} \Gamma(\alpha) \Gamma(\beta) &= \int_{u=0}^{\infty} e^{-u} u^{\alpha-1} du \int_{v=0}^{\infty} e^{-v} v^{\beta-1} dv \\ &= \int_{u=0}^{\infty} \int_{v=0}^{\infty} e^{-u-v} u^{\alpha-1} v^{\beta-1} du dv \end{aligned}$$

$u = zt \quad v = z(1-t) \quad \text{avec } t \in [0, 1] \\ z \in \mathbb{R}^+$

$u = tz \quad v = z(1-t)$

$$F: (z, t) \rightarrow (u, v)$$

$$J_F = \begin{pmatrix} \partial u / \partial t & \partial u / \partial z \\ \partial v / \partial t & \partial v / \partial z \end{pmatrix} = \begin{pmatrix} z & t \\ -z & 1-t \end{pmatrix}$$

$$|\det(J_F)| =$$

$$|z(1-t) + zt| =$$

$$\int_{y=0}^{\infty} \int_{t=0}^1 e^{-y} \left(\frac{t}{y}\right)^{\alpha-1} (y(1-t))^{\beta-1} y \, dt \, dy =$$

$$\int_{y=0}^{\infty} e^{-y} y^{\alpha+\beta-1} dy \underbrace{\int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt}_{B(\alpha, \beta)}$$

$\Gamma(\alpha)\Gamma(\beta) = \Gamma(\alpha+\beta)$

Distribution de probabilité Beta

$\alpha, \beta > 0$

$$p(\alpha | d, \beta) = \frac{\alpha^{d-1} (1-\alpha)^{\beta-1} \mathbb{1}_{\alpha \in [0,1]}}{B(\alpha, \beta)}$$

Supposons  $X \sim B(\alpha, \beta)$

$$\begin{aligned}
 E[X] &= \frac{1}{B(\alpha, \beta)} \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx \\
 &= \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)} = \\
 &= \frac{\Gamma(\alpha+1) \cancel{\Gamma(\beta)}}{\Gamma(\alpha+1+\beta)} \times \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \cancel{\Gamma(\beta)}} \\
 &= \frac{\alpha \cancel{\Gamma(\alpha)} \times \Gamma(\alpha+\beta)}{(\alpha+\beta) \cancel{\Gamma(\alpha+\beta)} \cancel{\Gamma(\alpha)}} \\
 &= \frac{\alpha}{\alpha+\beta}
 \end{aligned}$$

$$X \sim \mathcal{B}(n, p)$$

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

$$p \sim U[0, 1] \quad P(p) = \Delta(p \in [0, 1])$$

$$P(a < p < b, X=x) = \int_a^b \binom{n}{x} p^x (1-p)^{n-x} dp$$

AVFC

$$0 \leq a \leq b \leq 1$$

$$P(X=x) = \int_0^1 \binom{n}{x} p^x (1-p)^{n-x} dp = \binom{n}{x} B(x+1, n-x+1)$$

$$P(a < p < b | X=x) = \frac{\int_a^b \binom{n}{x} p^x (1-p)^{n-x} dp}{\binom{n}{x} B(x+1, n-x+1)}$$

$$p | X=x \sim \mathcal{B}(x+1, n-x+1)$$

$$E[p | X=x] = \frac{x+1}{x+1 + n-x+1} = \frac{x+1}{n+2}$$

# Example of Bayesian inference I

Let us take the snooker example again and look for Bayesian approach

- the law on the observations is a binomial,  $X \sim \mathcal{B}(n, p)$ ;  ~~$\mathcal{B}(n, \sqrt{p})$~~   $\mathcal{B}(n, p)$
- the ball can equally probably stop any distance from the edge. Hence  $p \sim U_{[0,1]}$ ;
- let consider a quadratic cost:  $C(p, \delta) = (p - \delta)^2$ .

In that case,

$$\begin{aligned}\pi(p|X=x) &= \frac{C_n^x p^x (1-p)^{n-x} \mathbb{I}_{\{p \in [0,1]\}}}{\int_0^1 C_n^x p^x (1-p)^{n-x} dp} \\ &= \frac{p^x (1-p)^{n-x} \mathbb{I}_{p \text{ in } [0,1]}}{\int_0^1 p^x (1-p)^{n-x} dp}.\end{aligned}$$

## Example of Bayesian inference II

The posterior distribution is therefore a beta distribution,  $\mathcal{Be}(x + 1, n - x + 1)$ . It is easy to show that the Bayes estimator associated with a distribution  $\pi$  and a quadratic cost is the posterior mean

$$\delta^\pi(x) = E^\pi[p|x] = \int p \cdot \pi(p|x) dp.$$

The expectation of a random variable  $X$  following a beta distribution,  $\mathcal{Be}(\alpha, \beta)$  is given by

$$E[X] = \frac{\alpha}{\alpha + \beta}.$$

The Bayes estimator associated with the quadratic cost is therefore written

$$\delta^\pi(x) = \frac{x + 1}{n + 2}.$$



## Example of Bayesian inference III

If many launches have been performed, the estimate of  $p$  by this Bayesian procedure will be very close to the estimator of the maximum of likelihood. On the other hand, if only one throw is observed, we find:

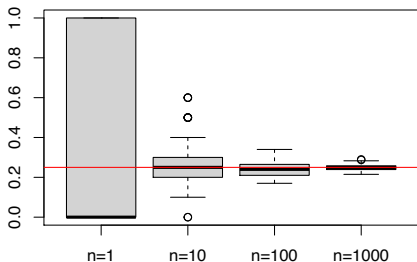
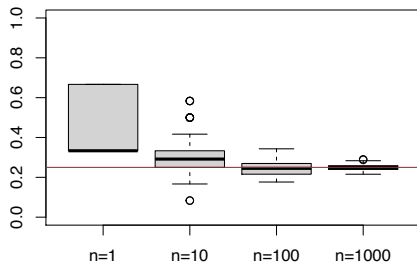
- $x = 0$  gives  $\hat{p} = \frac{1}{3}$ ;
- $x = 1$  gives  $\hat{p} = \frac{2}{3}$ .

Both of these results seem reasonable.

Note that by taking a cost which is equal to 0 if the decision is correct and 1 otherwise (cost 0-1), the Bayes estimator is, in this case, the same as that obtained by the maximum likelihood method.

Note that the Bayes estimators are justified for a size of finite sample, unlike estimators of the maximum of likelihood which only have asymptotic properties.

Comparing Bayesian and Frequentist estimator for  $p = \frac{1}{4}$  and  $n \in \{1, 10, 100, 1000\}$ .

**Frequentist****Bayesian**

# Dirichlet distribution $Dir(\alpha)$ I

- Peter Gustav Lejeune Dirichlet (13 février 1805 , Düren – · 5 mai 1859 , Göttingen)
- continuous multivariate probability distribution parameterized by a vector  $\alpha$  of positive reals.
- used as prior distributions in Bayesian statistics,
- conjugate prior of the categorical distribution and multinomial distribution.

The Dirichlet distribution of order  $K \geq 2$  with parameters  $\alpha_1, \dots, \alpha_K > 0$  has a probability density function

$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i-1}$$

where the normalizing constant is the multivariate beta function.

# Dirichlet distribution $Dir(\alpha)$ II

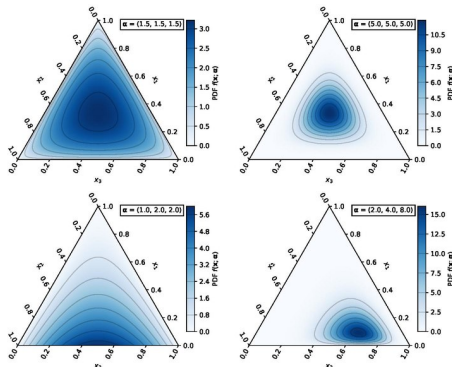


Figure 1: Dirichlet distributions

Denoting  $\alpha_0 = \sum_{i=1}^K \alpha_i$ , we have

$$E[X_i] = \frac{\alpha_i}{\alpha_0},$$

$$\text{Var}[X_i] = \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)}.$$

$$\text{Cov}[X_i, X_j] = \frac{-\alpha_i\alpha_j}{\alpha_0^2(\alpha_0 + 1)}.$$