# 04 - **Binomial model 01**

Bayesian Statistics
Spring 2022-2023

## Josep Fortiana

Matemàtiques - Informàtica UB

Monday, March 06, 2023

Estimating a probability

Which is the least informative prior?

# 04 - Binomial model 01

Estimating a probability

Which is the least informative prior?

# Bayesian Bernoulli model

Sample: $X = (X_1, \ldots, X_n)$ iid $\sim$ Ber$(\theta)$.

Estimate the probability $\theta \in \Theta = (0, 1)$.

*Prior distribution* for $\theta$: if no previous information, assume Unif$(0, 1)$:

$$p(\theta) = 1, \quad 0 < \theta < 1.$$

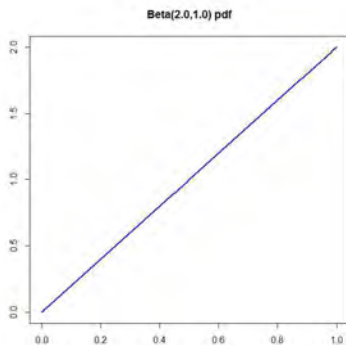*Non-Informative Prior (NIP).*

# A family of prior distributions

More generally: prior pdf of $\theta$ is Beta$(\alpha, \beta)$:

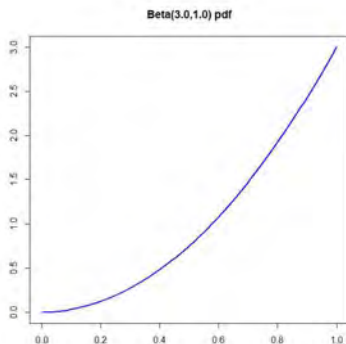$$p(t; \alpha, \beta) = \frac{1}{B(\alpha, \beta)} t^{\alpha-1} (1-t)^{\beta-1}, \quad 0 < t < 1,$$

where $B(\alpha, \beta), \alpha > 0, \beta > 0$, is the Beta function.

In particular, Beta$(1, 1) = $ Unif$(0, 1)$.

# Examples of Beta pdf's with $\alpha, \beta \geq 1$



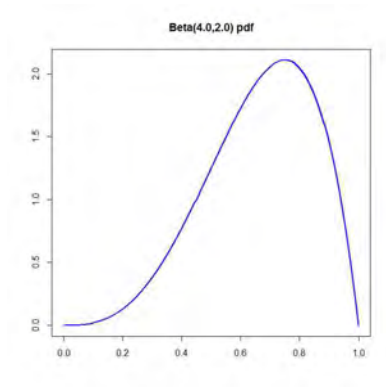Beta(2.0,1.0) pdf

# Examples of Beta pdf's with $\alpha, \beta \geq 1$



Beta(3.0,1.0) pdf

# Examples of Beta pdf's with $\alpha$, $\beta \geq 1$



Beta(4.0,1.0) pdf

# Examples of Beta pdf's with $\alpha, \beta \geq 1$



Beta(4.0,2.0) pdf

# Examples of Beta pdf's with $\alpha, \beta \geq 1$



Beta(4.0,3.0) pdf

# Examples of Beta pdf's with $\alpha, \beta \geq 1$



Beta(4.0,4.0) pdf

# Examples of Beta pdf's with $\alpha, \beta \geq 1$



Beta(3.0,4.0) pdf

# Examples of Beta pdf's with $\alpha, \beta \geq 1$



Beta(2.0,4.0) pdf

# Examples of Beta pdf's with $\alpha, \beta \geq 1$



Beta(1.0,4.0) pdf

# Examples of Beta pdf's with $\alpha, \beta < 1$



Beta(1.0,0.9) pdf

# Examples of Beta pdf's with $\alpha, \beta < 1$



Beta(1.0,0.8) pdf

# Examples of Beta pdf's with $\alpha, \beta < 1$



Beta(1.0,0.7) pdf

# Examples of Beta pdf's with $\alpha, \beta < 1$



Beta(0.9,0.7) pdf

# Examples of Beta pdf's with $\alpha, \beta < 1$



Beta(0.8,0.7) pdf

# Examples of Beta pdf's with $\alpha, \beta < 1$



Beta(0.7,0.7) pdf

# Likelihood

We observe $n$ values $X_i = x_i$, $1 \le i \le n$.

The *likelihood* is the joint pmf of $X = (X_1, \ldots, X_n)$, conditional to a given $\theta$, is:

$$p(x \mid \theta) = \theta^{n_1} (1 - \theta)^{n - n_1},$$

where $n_1 = \sum_{i=1}^{n} x_i$ is the absolute frequency of ones.

A function of the *sufficient statistic, $n_1$*.

# Marginal pmf of $X$

$$p(x) \;=\; \int_{\Theta} p(x \mid \theta)\, p(\theta)\, d\theta$$

$$=\; \int_0^1 \frac{1}{B(\alpha,\beta)}\, t^{\alpha+n_1-1}\,(1-t)^{\beta+n-n_1-1}\, dt$$

$$=\; \frac{1}{B(\alpha,\beta)}\, B(\alpha+n_1,\, \beta+n-n_1).$$

# Prior predictive pdf

$p(x)$ is also called Prior predictive pmf of $X$.

# Prior predictive pdf

$p(x)$ is also called Prior predictive pmf of $X$.

Why?

# Prior predictive pdf

$p(x)$ is also called Prior predictive pmf of $X$.

Why?

$p(x)$ averages $p(x \mid \theta)$ over all possible $\theta$, each with a relative weight *proportional to the prior $p(\theta)$.*

# The Beta-Binomial distribution

For real numbers $\alpha, \beta > 0$, and integer $n > 0$, the pmf:

$$p(k; n, \alpha, \beta) = \binom{n}{k} \cdot \frac{B(\alpha + k, \beta + n - k)}{B(\alpha, \beta)},$$

defines the *Beta-binomial distribution,*

r.v. with support on the set of

nonnnegative integers $k$ such that $0 \leq k \leq n$.

# Moments of the Beta-Binomial distribution

For a r.v. $Y \sim$ Beta-Binom$(n, \alpha, \beta)$

$$E(Y) = n \cdot \frac{\alpha}{\alpha + \beta},$$

$$var(Y) = n \cdot \frac{\alpha\,\beta\,(\alpha + \beta + n)}{(\alpha + \beta)^2\,(\alpha + \beta + 1)}.$$

# Examples of Beta-Binomial pmf's



Beta-Binomial(10,1.0,1.0) pmf

# Examples of Beta-Binomial pmf's



Beta-Binomial(10,2.0,1.0) pmf

# Examples of Beta-Binomial pmf's



Beta-Binomial(10,3.0,1.0) pmf

# Examples of Beta-Binomial pmf's



Beta-Binomial(10,4.0,1.0) pmf

# Examples of Beta-Binomial pmf's



Beta-Binomial(10,4.0,2.0) pmf

# Examples of Beta-Binomial pmf's



Beta-Binomial(10,4.0,3.0) pmf

# Examples of Beta-Binomial pmf's



Beta-Binomial(10,4.0,4.0) pmf

# Examples of Beta-Binomial pmf's

# *Posterior* pdf of $\theta$

Bayes' formula $\rightarrow \quad p(\theta \mid x)$

$$= \frac{p(x \mid \theta)\, p(\theta)}{f(x)}$$

$$= \frac{1}{B(\alpha + n_1, \beta + n - n_1)}\, \theta^{\alpha + n_1 - 1} \left(1 - \theta\right)^{\beta + n - n_1 - 1}.$$

# A conjugate family

The resulting pdf is another Beta distribution,

$$\text{Beta}(\alpha + n_1, \beta + n - n_1).$$

> The pair
>
> Bernoulli likelihood / Beta prior
>
> is a conjugate pair.

# Posterior expectation of $\theta$

$$\mathsf{E}[\theta \mid X = x] \;=\; \frac{\alpha + n_1}{\alpha + \beta + n}.$$

Can be written as a convex combination

$$\mathsf{E}[\theta \mid X = x] \;=\; \lambda \cdot \frac{n_1}{n} + (1 - \lambda) \cdot \frac{\alpha}{\alpha + \beta},$$

where $\lambda = \dfrac{n}{\alpha + \beta + n}$.

# Posterior expectation of $\theta$

$$\frac{n_1}{n} \quad = \quad \text{empirical probability.}$$

$$\frac{\alpha}{\alpha + \beta} \quad = \quad \text{prior expectation.}$$

Think of prior expectation as the result of a

previous experiment, $\alpha$ sucesses

out of $\alpha + \beta$ realizations.

# Posterior expectation of $\theta$

The coefficient in the convex combination:

$$\lambda = \frac{n}{\alpha + \beta + n}$$

is the ratio of sizes,

actually observed sample

*vs.* a previous *"virtual"* sample.

# Posterior predictive distribution

The Posterior predictive distribution for a new observation $\tilde{x}$, given the observed $x$, is the average of the pmf $p(x \mid \theta)$ over all possible values of $\theta$, where now relative weights of $\theta$ are given by the posterior pdf.

We integrate with respect to $\theta$, the product of the pmf $\mathrm{Binom}(n, \theta)$ times the posterior pdf $\mathrm{Beta}(\alpha + x, \beta + n - x)$.

# Posterior predictive distribution

The result is again a Beta-Binomial distribution:

$$p(\tilde{x}) = \frac{1}{B(\alpha + x, \beta + n - x)}$$

$$\times\, B(\alpha + x + \tilde{x}, \beta + n - x + \tilde{n} - \tilde{x}) \binom{\tilde{n}}{\tilde{x}}.$$

*[To allow for the case when the new observation $\tilde{x}$ comes from a different number $\tilde{n}$ of Bernoulli experiment repetitions, $\tilde{x} \sim \text{Binom}(\tilde{n}, \theta)$.]*

# Summary: Beta-Binomial (Bernoulli) model

▶ Prior distribution of $\theta$: A Beta pdf,

▶ Prior predictive of $x$: A Beta-Binomial pdf,

▶ Posterior of $\theta$, given $x$: A Beta pdf,

▶ Posterior predictive of $\tilde{x}$, given $x$: A Beta-Binomial pdf.

# 04 - Binomial model 01

Estimating a probability

## Which is the least informative prior?

# How does choice of prior reflect on the posterior?

With a Bernoulli likelihood, it is not obvious that Unif$(0, 1)$ is "the" Non-Informative Prior (NIP).

Beta priors, plus improper Beta distributions of the form:

$$p(\theta) \propto \theta^{\alpha-1} \cdot (1-\theta)^{\beta-1}, \qquad \alpha, \beta \in \mathbb{R}.$$

*Zhu, Mu; Lu, Arthur Y. (2004), The Counter-Intuitive Non-informative Prior for the Bernoulli Family, Journal of Statistics Education, 12 (2).*

# Useful formulas (1)

With a Beta$(\alpha, \beta)$ prior pdf, the marginal pmf of $x$ is a Beta-binomial:

$$p(x) = \frac{1}{B(\alpha, \beta)} B(\alpha + n_1, \beta + n - n_1),$$

where $n_1 = \sum_{i=1}^{n} x_i$.

# Useful formulas (2)

The expectation and variance of $U \sim \text{Beta}(\alpha, \beta)$ are:

$$\mathsf{E}(U) = \frac{\alpha}{\alpha + \beta},$$

$$\mathsf{var}(U) = \frac{\alpha \beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)}.$$

# Useful formulas (3)

The posterior pdf of $\theta$, given $x$:

$$p(\theta \mid x) \;=\; \frac{p(x \mid \theta) \cdot p(\theta)}{p(x)}$$

$$=\; \frac{1}{B(\alpha + n_1, \beta + n - n_1)} \, \theta^{\alpha + n_1 - 1} \left(1 - \theta\right)^{\beta + n - n_1 - 1},$$

is a $\text{Beta}(\alpha + n_1, \beta + n - n_1)$ distribution.

# Posterior expectation and variance

For the posterior pdf, a $\text{Beta}(\alpha + n_1, \beta + n - n_1)$,

$$\mathsf{E}(\theta \mid x) \;=\; \frac{\alpha + n_1}{\alpha + \beta + n},$$

$$\mathrm{var}(\theta \mid x) \;=\; \frac{(\alpha + n_1)(\beta + n - n_1)}{(\alpha + \beta + n)^2 (\alpha + \beta + n + 1)}.$$

# NIP 1: The uniform law

$$p_1(\theta) \sim \text{Unif}[0, 1] = \text{Beta}(1, 1).$$

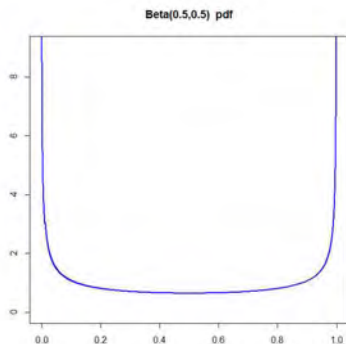$$\mathsf{E}(\theta \mid x) \quad = \quad \frac{n_1 + 1}{n + 2},$$

$$\text{var}(\theta \mid x) \quad = \quad \frac{(n_1 + 1)(n - n_1 + 1)}{(n + 2)^2 (n + 3)}.$$

# NIP 2: Jeffreys' prior

$$p_2(\theta) \sim \text{Beta}(1/2, 1/2).$$

Drawback is, its appearance is not "non-informative": probability concentrates near 0 and 1.

# Probability density function of Jeffreys' prior



Beta(0.5,0.5) pdf

# NIP 2: Jeffreys' prior

With Jeffreys' prior,

$$E(\theta \mid x) = \frac{n_1 + 1/2}{n + 1},$$

$$var(\theta \mid x) = \frac{(n_1 + 1/2)(n - n_1 + 1/2)}{(n + 1)^2(n + 2)}.$$

# The Beta$(c, c)$ subfamily

For the Beta subfamily with $\alpha = \beta = c$, where both Jeffreys' and uniform belong:

$$E(\theta \mid x) = \frac{n_1 + c}{n + 2c},$$

$$\mathrm{var}(\theta \mid x) = \frac{(n_1 + c)(n - n_1 + c)}{(n + 2c)^2(n + 2c + 1)}.$$

# The Beta($c, c$) subfamily

A Beta($c, c$) prior is equivalent to adding $2c$ virtual observations to the sample, $c$ zeros and $c$ ones.

Writing: $N = n + 2c$, $N_1 = n_1 + c$,

$$\mathsf{E}(\theta \mid x) = \frac{N_1}{N}, \qquad \mathsf{var}(\theta \mid x) = \frac{N_1 (N - N_1)}{N^2 (N + 1)}.$$

# Comparing Jeffreys' and uniform prior

Jeffreys' prior is less influential than the uniform,

It meddles less in the experiment, contributing only one *virtual observation*, evenly distributed between 0 and 1,

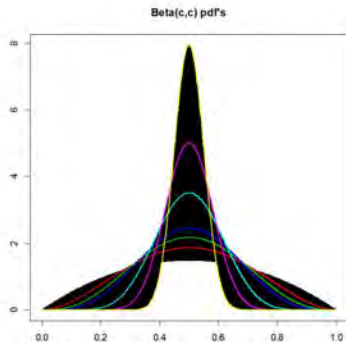The uniform adds two *virtual observations*, one of each.

# Other Beta($c$, $c$) priors

Within this subfamily,

What happens with a very large or a very small $c$?

# Other Beta($c, c$) priors

For $c = 2, 3, 4, 5, 10, 20, 50$,

# Other Beta($c$, $c$) priors

If $c \to \infty$, the Beta($c$, $c$) law tends to a degenerate (constant) distribution, with:

$$P\{\theta = 1/2\} = 1.$$

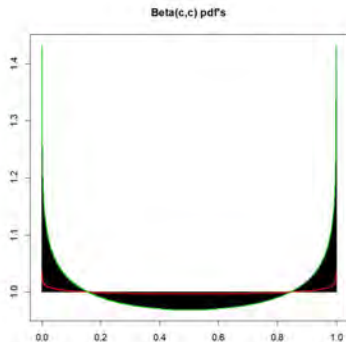Then the posterior is this same degenerate law.

# Other Beta($c$, $c$) priors

In agreement with the interpretation above, this is the *dogmatic estimator.*

The *a priori* information is so strong that it overrules any experimental evidence.

# Other Beta($c, c$) priors

For $c = 1, 0.995, 0.95,$

# Other Beta$(c, c)$ priors

In the opposite direction, if $c \to 0$, the less influential prior should be the limit $c = 0$,

$$p(\theta) \propto \theta^{-1} \cdot (1 - \theta)^{-1}, \quad \theta \in (0, 1),$$

for which,

$$E(\theta \mid x) = \frac{n_1}{n} = f_1, \quad \text{relative frequency of ones},$$

The classical ML estimator.

# Haldane's prior

This Beta$(0, 0)$ pdf can be derived by applying the change of variable formula to the (improper) uniform law:

$$p(\eta) = 1, \quad \eta \in (-\infty, \infty),$$

for the log-odds ratio $\eta = \log\left(\frac{\theta}{1-\theta}\right)$, the natural Bernoulli parameter (as a regular exponential family).

# Other Beta($c, c$) priors

For $c = 0$,

$$\text{var}(\theta \mid x) = \frac{n_1\,(n - n_1)}{n^2\,(n + 1)} = \frac{1}{n + 1} f_1(1 - f_1).$$

Smaller than $\text{var}_\theta(f_1) = \dfrac{1}{n}\,\theta\,(1 - \theta)$, the CR bound. !?

# Other Beta($c, c$) priors

Not a contradiction,

the variance of an estimator $\hat{\theta}(x)$ and

the posterior variance of the parameter $\theta$ itself

are entirely different concepts.

# Other Beta($c$, $c$) priors

The $c \to 0$ limit, Beta($0$, $0$), is the discrete law:

$$P[\theta = 0] = P[\theta = 1] = 1/2,$$

In a sense, the opposite case to setting $P = 1$ at $\theta = 0.5$: now there is a maximum indeterminacy between the two extreme possible $\theta$ values.

# Summary

Jeffreys' prior should appear as reasonably

non informative, the *aurea mediocritas*

between both "radical"priors.