

# Motivating example: Spelling correction

- **Problem:** Someone types '*radom*'.
- **Question:** What did they mean to type? Random?

## Ingredients:

- **Data**  $x$ : The observed word — *radom*
- **Parameter of interest**  $\theta$ : The correct word

## Comments: To solve this we need

- *background information* on which words are usually typed.
- an idea about how words are typically mistyped.

Example adapted from *Bayesian Data Analysis* by Gelman et al.

# Bayesian Idea

- **Data model:** Conditional on  $\theta$ , data  $x$  is distributed according to pf/pdf  $\pi(x)$ :

$$\pi(x|\theta) \propto L(\theta; x) \quad \leftarrow \text{the likelihood}$$

- **Prior:** Prior knowledge (i.e. *before* collecting data) about  $\theta$  is summaries by a pf/pdf,

$$\pi(\theta) \quad \leftarrow \text{the prior}$$

- **Posterior** : The updated knowledge about  $\theta$  *after* collecting data: The conditional distribution of  $\theta$  given data  $x$ :

$$\begin{aligned}\pi(\theta|x) &= \frac{\pi(x|\theta)\pi(\theta)}{\pi(x)} \\ &\propto \pi(x|\theta)\pi(\theta)\end{aligned}$$

$$\text{“posterior} = \text{likelihood} \times \text{prior”}$$

## Example: Prior

Without any other prior knowledge Google provides the following *prior probabilities* (for three candidate words):

$\theta$	$\pi(\theta)$
random	$7.60 \times 10^{-5}$
radon	$6.05 \times 10^{-6}$
radom	$3.12 \times 10^{-7}$

### Comments

- The relative high probability for the word *radom* is surprising!  
Name of Polish airshow and nickname for Polish hand gun...
- In the context of writing a scientific report these prior probabilities would look different...?

## Example: Likelihood

Google provides the following conditional probabilities *prior probabilities*:

$\theta$	$\pi(x = \text{'radom'} \theta)$
random	0.00193
radon	0.000143
radom	0.975

### Comments

- This is *not* a probability distribution!
- If one in fact intends to write 'radom' this actually happens in 97.5% of cases.
- If one intends to write either 'random' or 'radon' this is rarely misspelled 'radom'.

## Example: Posterior

Combining the prior and likelihood we obtain the *posterior probabilities*:

$$\pi(\theta|x) = \frac{\pi(x|\theta)\pi(\theta)}{\pi(x)} \propto \pi(x|\theta)\pi(\theta)$$

$\theta$	$\pi(x = \text{'radom'} \theta)\pi(\theta)$	$\pi(\theta x = \text{'radom'})$
random	$1.47 \times 10^{-7}$	0.325
radon	$8.65 \times 10^{-10}$	0.002
radom	$3.04 \times 10^{-7}$	0.673

### Conclusion

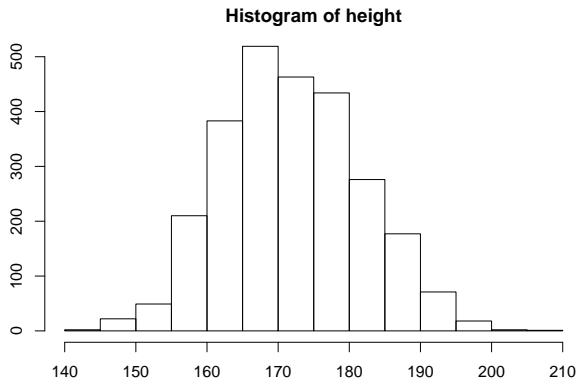
- With the given prior and likelihood the word 'radom' is twice as likely as 'random'.

### Criticism

- The posterior probability for 'radom' seems too high.
- Likelihood or prior to blame
- Likelihood is probably ok in this case
- Prior depends on context — and hence might be “wrong”.

## Another example

Heights of some Copenhageners in 1995

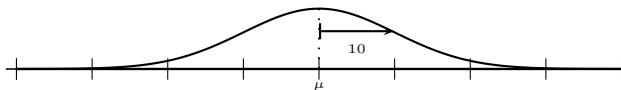


**Assume:** Heights are normal,  $X \sim \mathcal{N}(\mu, \tau)$ .

**For now:** Assume precision  $\tau$  known.

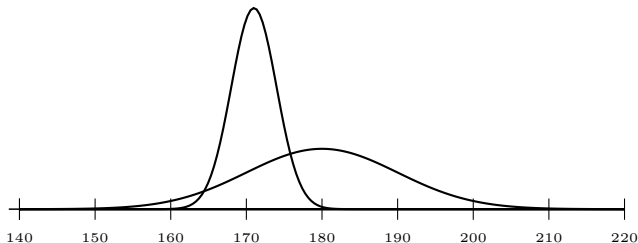
# Bayesian Idea: Illustration

**Data model:**  $X \sim \mathcal{N}(\mu, 0.01)$  (i.e. pop. sd = 10)



**Prior:** We believe that the population mean is most likely between 160 cm and 200 cm.  $\pi(\mu) = \mathcal{N}(180, 0.01)$ .

**Posterior:** After observing a number of heights ( $n = 10$ ,  $\bar{x} = 169$ ), our knowledge about  $\mu$  is updated. Summarised by the posterior.



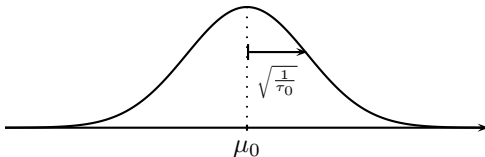
# Normal example: One (!) observation

Data model:  $X \sim \mathcal{N}(\mu, \tau)$

Assume: Precision  $\tau$  known.

Interest: The unknown mean  $\mu$ .

**Prior:** The prior for  $\mu$  is specified as  $\mu \sim \mathcal{N}(\mu_0, \tau_0)$ .





## Normal example: Data density

**Data:** One observation,  $X$ , from a normal distribution:

$$\begin{aligned}\pi(x|\mu) &= \sqrt{\frac{\tau}{2\pi}} \exp\left(-\frac{1}{2}\tau(x - \mu)^2\right) \\ &= \sqrt{\frac{\tau}{2\pi}} \exp\left(-\frac{1}{2}\tau x^2 - \frac{1}{2}\tau\mu^2 + \tau\mu x\right) \\ &\propto \exp\left(-\frac{1}{2}\tau x^2 + \tau\mu x\right)\end{aligned}$$

**Notice** the “pattern” inside the exponential.

## Normal example: Posterior density

$$\textit{Posterior} \propto \textit{Likelihood} \times \textit{Prior}$$

$$\begin{aligned}\pi(\mu|x) &\propto \pi(x|\mu)\pi(\mu) \\&= \sqrt{\frac{\tau}{2\pi}} \exp\left(-\frac{1}{2}\tau(x-\mu)^2\right) \sqrt{\frac{\tau_0}{2\pi}} \exp\left(-\frac{1}{2}\tau_0(\mu-\mu_0)^2\right) \\&\propto \exp\left(-\frac{1}{2}\tau\mu^2 + \tau x\mu - \frac{1}{2}\tau_0\mu^2 + \tau_0\mu\mu_0\right) \\&= \exp\left(-\frac{1}{2}(\tau + \tau_0)\mu^2 + (\tau x + \tau_0\mu_0)\mu\right) \\&\propto \mathcal{N}\left(\frac{\tau x + \tau_0\mu_0}{\tau + \tau_0}, \tau + \tau_0\right)\end{aligned}$$

**Notice:** Prior for  $\mu$  was normal, now the posterior for  $\mu$  is also normal!

## Normal example: Posterior mean & variance

The posterior:  $\pi(\mu|x) = \mathcal{N}\left(\frac{\tau x + \tau_0 \mu_0}{\tau + \tau_0}, \tau + \tau_0\right)$

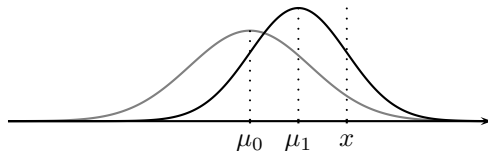
**Posterior expectation:**

$$\mathbb{E}[\mu|x] = \frac{\tau x + \tau_0 \mu_0}{\tau + \tau_0} = \frac{\tau}{\tau + \tau_0} x + \frac{\tau_0}{\tau + \tau_0} \mu_0 (= \mu_1).$$

*Weighted average* of prior mean and observation  $x$ .

**Posterior variance:**

$$\mathbb{V}\text{ar}[\mu|x] = \frac{1}{\tau + \tau_0} (= \frac{1}{\tau_1})$$



# Posterior as prior — or updating believes

**General setup:** We are interested in parameter  $\theta$ .

- Data model:  $\pi(x|\theta)$
- Prior:  $\pi(\theta)$
- Data: First observation  $x_1 \sim \pi(x_1|\theta)$
- Posterior:  $\pi(\theta|x_1) \propto \pi(x|\theta)\pi(\theta)$

Assume we have a second independent observation  $x_2 \sim \pi(x_2|\theta)$ .

**Posterior:**

$$\begin{aligned}\pi(\theta|x_1, x_2) &\propto \pi(x_1, x_2|\theta)\pi(\theta) \\ &= \pi(x_1|\theta)\pi(x_2|\theta)\pi(\theta) \\ &\propto \underbrace{\pi(x_2|\theta)}_{\text{likelihood}} \underbrace{\pi(\theta|x_1)}_{\text{prior}}\end{aligned}$$

**Notice:** The posterior after observing  $x_1$  is the prior before observing  $x_2$ .

# Independent normal case

Posterior mean and precision after one observation  $x_1$ :

$$\mu_1 = \frac{x_1\tau + \mu_0\tau_0}{\tau + \tau_0} \quad \text{and} \quad \tau_1 = \tau + \tau_0.$$

Next,  $\mu_1$  and  $\tau_1$  are prior mean and precision before observing  $x_2$ .

Hence, posterior mean and precision after observing (independent)  $x_1$  and  $x_2$  are

$$\begin{aligned} \mu_2 &= \mathbb{E}[\mu_2|x_1, x_2] = \frac{x_2\tau + \mu_1\tau_1}{\tau + \tau_1} \\ &= \frac{x_2\tau + x_1\tau + \mu_0\tau_0}{\tau + \tau + \tau_0} = \frac{(x_1 + x_2)\tau + \mu_0\tau_0}{2\tau + \tau_0} \\ \tau_2 &= 2\tau + \tau_0 \end{aligned}$$

This can easily be generalised.

# Many independent normal observations

Assume...

- $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \tau)$ .
- $\tau$  is known.
- Prior  $\pi(\mu) = \mathcal{N}(\mu_0, \tau_0)$ .

The posterior is then

$$\pi(\mu|x_1, x_2, \dots, x_n) = \mathcal{N}\left(\frac{\tau \sum_i x_i + \tau_0 \mu_0}{n\tau + \tau_0}, n\tau + \tau_0\right)$$

# Posterior mean: Sanity check

The posterior is

$$\pi(\mu|x_1, x_2, \dots, x_n) = \mathcal{N}\left(\frac{\tau \sum_i x_i + \tau_0 \mu_0}{n\tau + \tau_0}, n\tau + \tau_0\right)$$

Does the posterior mean seem “sane”?

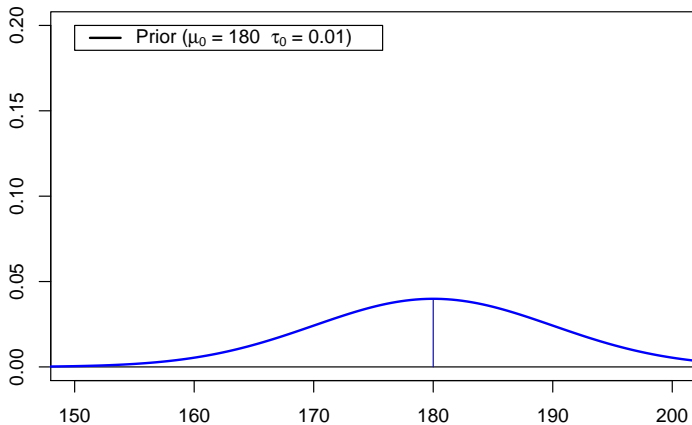
$$\begin{aligned}\mathbb{E}[\mu|x_1, \dots, x_n] = \mu_n &= \frac{\tau \sum_i x_i + \tau_0 \mu_0}{n\tau + \tau_0} \\ &= \frac{\tau n \frac{1}{n} \sum_i x_i + \tau_0 \mu_0}{n\tau + \tau_0} \\ &= \frac{n\tau}{n\tau + \tau_0} \bar{x} + \frac{\tau_0}{n\tau + \tau_0} \mu_0\end{aligned}$$

Weighted average of sample average  $\bar{x}$  and prior mean  $\mu_0$ .

For  $n$  large we have  $\mu_n \approx \bar{x}$ . Choice of  $\mu_0$  of little importance.

Precision  $\tau_n = n\tau + \tau_0$ . Knowledge is ever more precise.

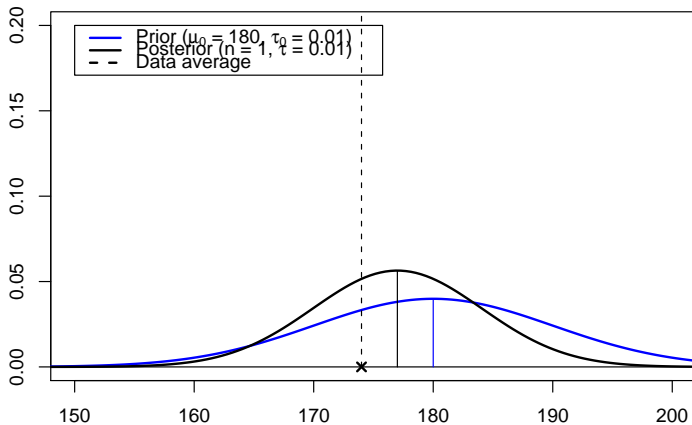
# Heights in Copenhagen: Prior





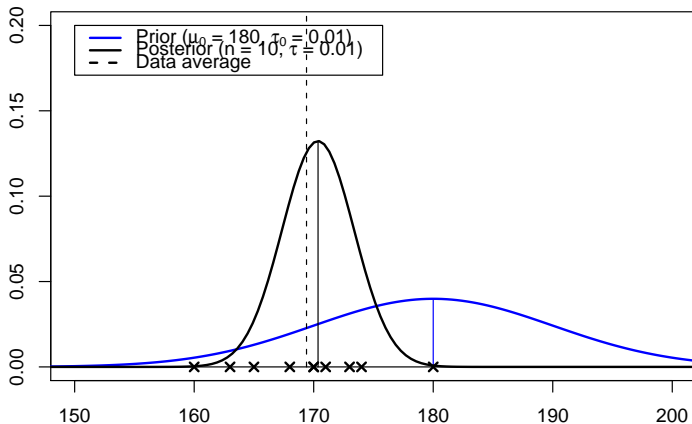
# Heights in Copenhagen: Posterior

One observation



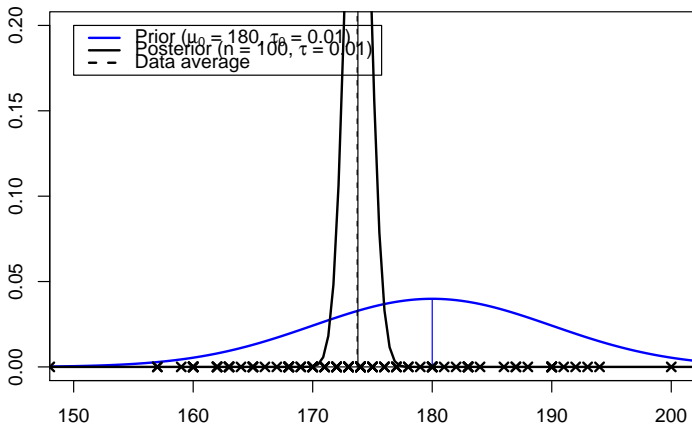
# Heights in Copenhagen: Posterior

Ten observations



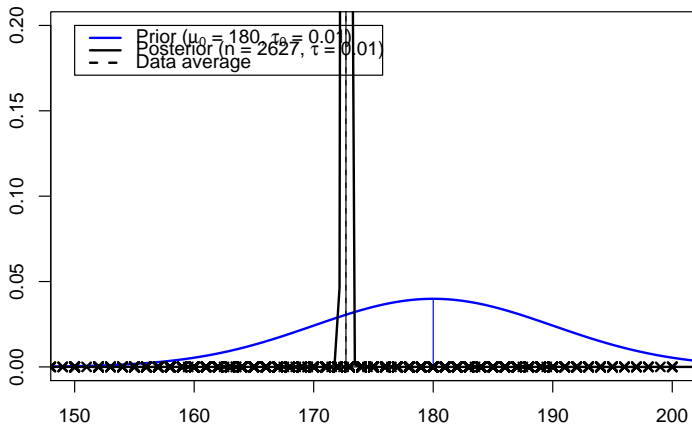
# Heights in Copenhagen: Posterior

100 observations



# Heights in Copenhagen: Posterior

2627 observations



# How to summaries the posterior $\pi(\theta|x)$ ?

The posterior is usually summaries using one or more of the following:

- Plot of posterior density  $\pi(\theta|x)$ . See previous slides.
- Posterior mean and variance/precision.
- Central Posterior Interval (CPI). See next slide.
- Maximum A Posteriori (MAP) estimate

$$MAP(\theta) = \underset{\theta}{\operatorname{argmax}} \pi(\theta|x).$$

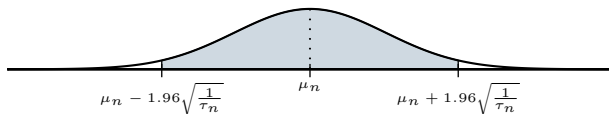
# Central Posterior Interval

- The CPI is an interval estimate.
- Also known as credibility interval.
- A 95% CPI for a parameter  $\theta$  is the shortest (connected) interval which contains  $\theta$  with 95% posterior probability.
- In case of the normal example, we have

$$P\left(\mu_n - 1.96\sqrt{\frac{1}{\tau_n}} \leq \mu \leq \mu_n + 1.96\sqrt{\frac{1}{\tau_n}}\right) = 0.95$$

Hence, a 95% CPI for  $\mu$  is

$$95\% \text{ CPI: } \mu_n \pm 1.96\sqrt{\frac{1}{\tau_n}}.$$



## CPI compared to confidence interval

The classical 95% confidence interval for  $\mu$  is

$$95\% \text{ CI: } \bar{x} \pm 1.96 \sqrt{\frac{1}{n\tau}}.$$

For CPI: Assume the prior precision is  $\tau_0 = 0$ , i.e. infinite variance. Then  $\mu_n = \bar{x}$  and  $\tau_n = n\tau$ , i.e.

$$95\% \text{ CPI: } \bar{x} \pm 1.96 \sqrt{\frac{1}{n\tau}}.$$

Same interval. Different interpretations.

# Conjugate priors

In the normal example: Both prior and posterior were normal! Very convenient!

We say that the normal distribution is conjugate.

## Definition: Conjugate priors

Let  $\pi(x|\theta)$  be the data model.

A class  $\Pi$  of prior distributions for  $\theta$  is said to be conjugate for  $\pi(x|\theta)$  if

$$\pi(\theta|x) \propto \pi(x|\theta)\pi(\theta) \in \Pi$$

whenever  $\pi(\theta) \in \Pi$ . I.e. prior and posterior are in the same class of distributions.

**Notice:**  $\Pi$  should be a class of “natural” distributions for this to be useful.



# Improper priors

If we have no prior knowledge we may be tempted to use a “flat” prior, i.e.

$$\pi(\theta) \propto k$$

If  $\theta \in \mathbf{R}$  this is an example of an improper prior, as

$$\int_{-\infty}^{\infty} \pi(\theta) d\theta = \int_{-\infty}^{\infty} k = \infty.$$

Problematic, but ok, if posterior is proper, i.e. if

$$\int \pi(\theta|x) d\theta = \int \pi(x|\theta)\pi(\theta) d\theta < \infty.$$

**Notice:** If  $\pi(\theta) \propto 1$  then MAP estimator = Maximum likelihood estimator.

## Normal example: Unknown precision, known mean

- **Data model:**  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \tau)$ :

$$\pi(x|\tau) = \left(\frac{\tau}{2\pi}\right)^{\frac{n}{2}} \exp\left(-\frac{1}{2}\tau \sum_{i=1}^n (x_i - \mu)^2\right)$$

- **Prior:** Gamma distribution:  $\pi(\tau) = \text{Gamma}(\alpha, \beta)$

$$\pi(\tau) = \frac{1}{\beta^\alpha \Gamma(\alpha)} \tau^{\alpha-1} \exp\left(-\frac{\tau}{\beta}\right)$$

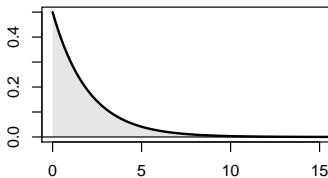
Shape parameter  $\alpha$  and scale parameter  $\beta$ .

Properties of the gamma distribution:

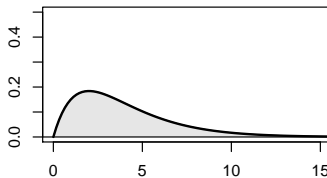
$$\mathbb{E}[\tau] = \alpha\beta \quad \text{Var}[\tau] = \alpha\beta^2.$$

# Gamma distribution

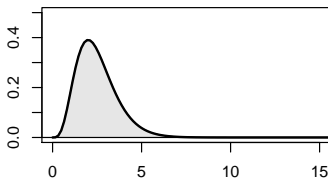
$\alpha=1, \beta=2$



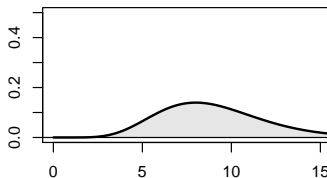
$\alpha=2, \beta=2$



$\alpha=5, \beta=0.5$



$\alpha=9, \beta=1$



# Normal example: Posterior precision

- **Data model:**  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \tau)$ :
- **Prior:**  $\pi(\tau) = \text{Gamma}(\alpha, \beta)$
- **Posterior:**

$$\pi(\tau|x) = \text{Gamma} \left( \frac{n}{2} + \alpha, \left\{ \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 + \frac{1}{\beta} \right\}^{-1} \right)$$

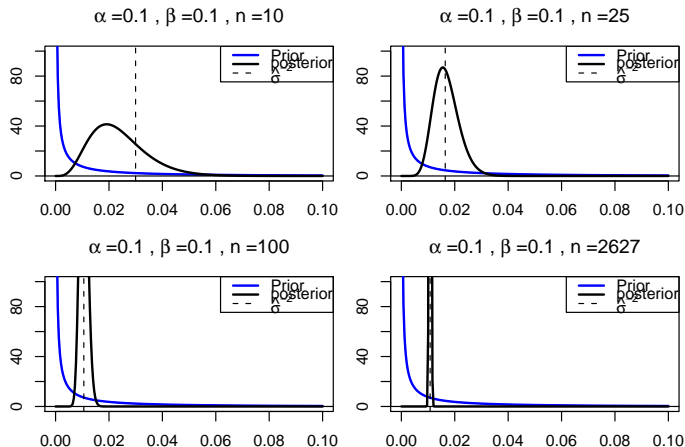
Posterior mean and variance

$$\mathbb{E}[\tau|x] = \frac{\frac{n}{2} + \alpha}{\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 + \frac{1}{\beta}} \quad \mathbb{V}\text{ar}[\tau|x] = \frac{\frac{n}{2} + \alpha}{\left( \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 + \frac{1}{\beta} \right)^2}$$

For large  $n$  we have

$$\mathbb{E}[\tau|x] \approx \frac{1}{\hat{\sigma}^2} \quad \text{where } \hat{\sigma}^2 \text{ is the usual ML variance estimate.}$$

# Known mean: Priors and posteriors



# Binomial example

- **Data model:** Binomial,  $X \sim B(n, p)$ ,  $n$  known.

$$\pi(x|p) = \binom{n}{x} p^x (1-p)^{n-x} \quad 0 \leq p \leq 1.$$

- **Prior:** Beta distribution,

$$\pi(p) = \text{Be}(\alpha, \beta) \quad , \alpha, \beta > 0.$$

Where we have to specify  $\alpha$  and  $\beta$ .

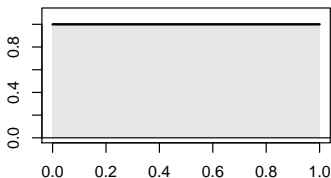
The Beta distribution has density

$$\pi(p) = \begin{cases} \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} p^{\alpha-1} (1-p)^{\beta-1} & \text{for } 0 \leq p \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

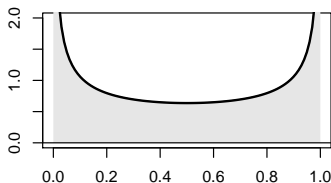
If  $\alpha = \beta = 1$  then  $\pi(p) = 1$  for  $0 \leq p \leq 1$ , i.e. uniform.

# Beta distribution: Examples

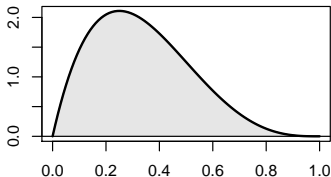
$\alpha=1, \beta=1$



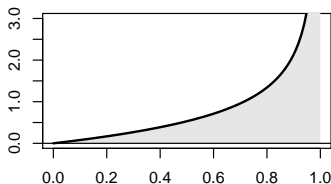
$\alpha=0.5, \beta=0.5$



$\alpha=2, \beta=4$



$\alpha=2, \beta=0.5$



$$\text{Mean: } \mathbb{E}[p] = \frac{\alpha}{\alpha + \beta}$$

$$\text{Variance: } \mathbb{V}\text{ar}[p] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

## Binomial example — cont.

- **Data model:**  $X \sim B(n, p)$
- **Prior:**  $\pi(p) = Be(\alpha, \beta)$ , that is

$$\pi(p) = \begin{cases} \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} p^{\alpha-1} (1-p)^{\beta-1} & \text{for } 0 \leq p \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

- **Posterior:**

$$\begin{aligned} \pi(p|x) &\propto \pi(x|p)\pi(p) \\ &= \binom{n}{x} p^x (1-p)^{n-x} \cdot \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} p^{\alpha-1} (1-p)^{\beta-1} \\ &\propto p^{x+\alpha-1} (1-p)^{n-x+\beta-1} \\ &= Be(x+\alpha, n-x+\beta). \end{aligned}$$



# Posterior mean & Variance

Posterior

$$\pi(p|x) = Be(x + \alpha, n - n + \beta).$$

Posterior mean

$$\mathbb{E}[p|x] = \frac{x + \alpha}{(x + \alpha) + (n - x + \beta)} = \frac{x + \alpha}{\alpha + \beta + n}$$

If  $x, n \gg \alpha, \beta$  then  $\mathbb{E}[p|x] \approx \frac{x}{n}$ .

Posterior variance

$$\begin{aligned}\mathbb{V}\text{ar}[p|x] &= \frac{(x + \alpha)(n - x + \beta)}{(x + \alpha + n - x + \beta)^2(x + \alpha + n - x + \beta + 1)} \\ &= \frac{(x + \alpha)(n - x + \beta)}{(\alpha + \beta + n)^2(\alpha + \beta + n + 1)} = O\left(\frac{1}{n}\right)\end{aligned}$$

## Example: Placenta Previa (PP)

- **Question:** Is the sex ratio different for PP births compared to normal births?
- **Prior knowledge:** 51.5% of new-borns are boys.
- **Data:** Of  $n = 980$  cases of PP  $x=543$  were boys ( $543/980=55.4\%$ ).
- **Data model:**  $X \sim B(n, p)$ .
- **Prior:**  $\pi(p) = Be(\alpha, \beta)$ .
- **Posterior:**

$$\begin{aligned}\pi(p|x) &= Be(x + \alpha, n - x + \beta) \\ &= Be(543 + \alpha, 437 + \beta)\end{aligned}$$

How to choose  $\alpha$  and  $\beta$ , and what difference does it make?

# Placenta Previa: Beta priors and posteriors

