## Frequentist

① Probability is consistent with long term relative frequency.

② Parameters are fixed unknown constants, not random quantities.

③ Statistical procedures should have well-defined long-run properties.

## Bayesian

① Probability is a degree of belief that an event will occur.

② Can make probabilistic statements about parameters even though they are fixed constants.

③ Parameter inference is done by computing a probability distribution of $\theta$, point estimates are computed after the fact.

## Example of a Bayesian Method

Goal  make a statement about some unknown parameter $\theta$.

① Choose a prior distribution $f(\theta)$ which captures initial belief about $\theta$.

② Choose a statistical model $f(\vec{x}|\theta)$ which captures beliefs about data $\vec{x}$ given the parameter $\theta$: conditional distribution

　　　Note: Write $f(\vec{x}|\theta)$ instead of $f(\vec{x}; \theta)$.

③ Observe $X_1, ..., X_n$, update our belief about $\theta$ in the form: $f(\theta|\vec{x})$ ← the posterior distribution

The posterior is just a conditional distribution:

$$f(\theta \mid \vec{x}) = \frac{f(\theta, \vec{x})}{f(\vec{x})} = \frac{f(\vec{x} \mid \theta) f(\theta)}{\int f(\theta, \vec{x}) \, d\theta}$$

$$= \frac{f(\vec{x} \mid \theta) f(\theta)}{\int f(\vec{x} \mid \theta) f(\theta) \, d\theta} \Bigg\} \text{ Bayes Theorem.}$$

In the case of $n$ IID observations $X_1, ..., X_n$,

$$f(\vec{x} \mid \theta) = \prod_{i=1}^{n} f(x_i \mid \theta) = L(\theta),$$

and therefore $f(\theta \mid \vec{x}) = \underbrace{\frac{f(\vec{x} \mid \theta) f(\theta)}{\int f(\vec{x} \mid \theta) f(\theta) \, d\theta}}_{\text{a constant } C = C(\vec{x})} = \frac{L(\theta) f(\theta)}{C}$

$\underbrace{\text{function of } \theta}$

The constant $C$ is such that $\int f(\theta \mid \vec{x}) \, d\theta = 1$

$$= \int \frac{L(\theta) f(\theta)}{C} \, d\theta.$$

Often for this reason, we will write

$$\underbrace{f(\theta \mid \vec{x})}_{\text{posterior}} \propto \underbrace{L(\theta)}_{\text{model}} \underbrace{f(\theta)}_{\text{prior}}.$$

To generate a point estimate compute a functional of the posterior $f(\theta \mid \vec{x})$:

$$\underbrace{\bar{\theta}}_{\text{Bayesian estimator}} = E(\theta \mid \vec{x}) = \int \theta \, f(\theta \mid \vec{x}) \, d\theta = \frac{\int \theta \, L(\theta) f(\theta) \, d\theta}{\int L(\theta) f(\theta) \, d\theta}$$

2

Posterior $\alpha$ interval   (_not_ a confidence interval)

Find   a, b   such that

$$P(\theta | \vec{x} \in (a,b)) = 1 - \alpha = \int_a^b f(\theta | \vec{x}) \, d\theta.$$

Example   $X_1, ..., X_n \sim$ Bernoulli($p$)   random variables

Our   model :   $f(\vec{x} | p) = \mathcal{L}(p)$

$$= \prod_{i=1}^{n} p^{X_i} (1-p)^{1-X_i}$$

prior   :   $p \sim$ Uniform $(0,1)$

$$f(p) = 1 \quad \text{on} \quad (0,1).$$

posterior   $f(p | \vec{x}) \propto \mathcal{L}(p) \, f(p)$

$$= \prod_{i=1}^{n} p^{X_i} (1-p)^{1-X_i}$$

$$= p^s (1-p)^{n-s} \qquad\qquad S = \sum_{i=1}^{n} X_i.$$

think of this as a function of $p$, not the data $X_i$ anymore

$$= p^{(s+1)-1} (1-p)^{(n-s+1)-1}$$

Now identify which family of probability distribution
$\mathcal{L}(p) \, f(p)$   belongs to.

Recall:   Beta($\alpha, \beta$) density :   $f(p; \alpha, \beta) = \dfrac{\Gamma(\alpha, \beta)}{\Gamma(\alpha) \, \Gamma(\beta)} p^{a-1} (1-p)^{\beta-1}$

$\Rightarrow f(p | \vec{x}) = \dfrac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \, \Gamma(\beta)} p^{a-1} (1-p)^{\beta-1}$   with   $\alpha = 1 + \Sigma X_i$

$\beta = 1 + n - \Sigma X_i$

$\Rightarrow p | \vec{x} \sim$ Beta $(\alpha, \beta)$.

3

Bayes estimator:

$$\bar{p} = E(p \mid \bar{x}) = \underbrace{\frac{s+1}{n+2}}_{\text{mean of } \beta(\alpha, \beta)} = \frac{\alpha}{\alpha + \beta}.$$

Graphically

Uniform prior



$0 \qquad 1$

Beta



$0 \quad \bar{p} \qquad 1$

Part 2   If we use the prior $p \approx \text{Beta}(\alpha, \beta)$, then repeating the calculation gives

$$p \mid \bar{x} \sim \text{Beta}(\alpha + s, \; \beta + n - s).$$

(Note: that $p \sim \text{Unifor}(0,1)$ is just $p \sim \text{Beta}(1,1)$)

This is a case when the prior family equals the posterior family.

We say that "the prior is conjugate with respect to the model".

$$f(p \mid \bar{x}) \; \alpha \; \overbrace{f(\bar{x} \mid p)}^{\text{model}} \; f(p)$$

same family

# Functions of parameters

- Recall from MLE that if $\hat{p}$ is the MLE estimate for $p$, then the MLE estimate for $\tau = g(p)$ was just $\hat{\tau} = g(\hat{p})$.

- Furthermore if $Y = g(X)$, then

$$F(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y) = \int_{g(x) \leq y} f(x)\, dx$$

and $f(y) = F'(y)$.

We can use these ideas for Bayesian inference as well.

Let $\tau = g(\theta)$.

$$\text{Bayes says}: \quad f(\theta | \vec{x}) \propto L(\theta)\, f(\theta).$$

Posterior CDF for $\tau = g(\theta) \mid \vec{x}$

$$= H(\tau | \vec{x})$$
$$= \mathbb{P}(g(\theta) \leq \tau \mid \vec{x})$$
$$= \int_{g(\theta) \leq \tau} f(\theta | \vec{x})\, d\theta$$

Then the posterior $h(\tau | \vec{x}) = \frac{d}{d\tau} H(\tau | \vec{x})$.

**Example** $X_i \sim \text{Bernoulli}(p)$, prior $f(p) = 1$.

$$\Rightarrow \quad p | \vec{x} \sim \text{Beta}(s+1, n-s+1). \qquad \text{Let } \psi = \log\left(\frac{p}{1-p}\right).$$
$$p \in (0,1) \Rightarrow \psi \in (-\infty, \infty).$$

$$\Rightarrow H(\psi | \vec{x}) = \mathbb{P}\left(\log \frac{p}{1-p} \leq \psi \mid \vec{x}\right)$$
$$= \mathbb{P}\left(p \leq \frac{e^\psi}{1+e^\psi} \mid \vec{x}\right)$$

[5]

$$= \int_0^{e^\gamma/1+e^\gamma} f(p \mid \vec{x}) \, dp$$

$$= \int_0^{e^\gamma/1+e^\gamma} \frac{\Gamma(n+2)}{\Gamma(s+1)\Gamma(n-s+1)} p^s (1-p)^{n-s} \, dp$$

To compute $h(\gamma \mid \vec{x})$, compute $\frac{d}{d\gamma} H(\gamma \mid \vec{x})$:

$$h(\gamma \mid \vec{x}) = \frac{d}{d\gamma} H(\gamma \mid \vec{x})$$

$$= \frac{d}{d\gamma} \left( \int_0^{e^\gamma/1+e^\gamma} \frac{\Gamma(n+2)}{\Gamma(s+1)\Gamma(n-s+1)} p^s (1-p)^{n-s} \, dp \right)$$

$$= \frac{\Gamma(n+2)}{\Gamma(s+1)\Gamma(n-s+1)} e^\gamma \left( \frac{e^\gamma}{1+e^\gamma} \right)^s \left( \frac{1}{1+e^\gamma} \right)^{n-s}$$

## Types of Priors

The big question in Bayesian analysis is <u>where does</u> the prior come from in a general problem.

<u>Subjective prior</u>   vs.   <u>Non-informative prior</u>

e.g. $\left( \begin{array}{c} \includegraphics \\ f(p) \end{array} , \begin{array}{c} \text{normal} \\ \text{distribution} \end{array} \right)$        e.g. $\left( \begin{array}{c} \includegraphics \end{array} , \begin{array}{c} \text{flat, uniform} \\ \text{pdf} \end{array} \right)$.

<u>Improper prior</u>: Let $x \mid \mu \sim N(\mu, \sigma^2)$, $\sigma^2$ is known.

Set a flat prior on $\mu$: $f(\mu) = c > 0$.

$$\Rightarrow \int_{-\infty}^{\infty} f(\mu) \, d\mu = \infty \qquad \leftarrow \text{not a prob. density.} \quad \boxed{6}$$

This is known as an improper prior, but Bayes can formally be carried out:

$$f(\mu \mid x) = \frac{L(\mu) \cdot f(\mu)}{\int L(\mu) \cdot f(\mu) \, d\mu}$$

$$= \frac{L(\mu) \cdot \cancel{c}}{\int L(\mu) \cdot \cancel{c} \, d\mu} = \frac{L(\mu)}{\underbrace{\int L(\mu) \, d\mu}}$$

$$\longrightarrow \int_{-\infty}^{\infty} \frac{e^{-(x-\mu)^2/2\sigma^2}}{\sqrt{2\pi}\,\sigma} \, d\mu$$

$$= 1$$

$$\Rightarrow \mu \mid x \sim N(x, \sigma^2)$$

In general, improper priors are <u>not</u> a problem so long as $f(\bar{x} \mid \theta)$ decays fast enough as a function of $\theta$.

<u>Flat priors are not transformation invariant</u>

Logically, if we know nothing about a parameter $p$, the we should also know nothing about $\psi = \log\left(\frac{p}{1-p}\right)$, but $f(p) = 1 \Rightarrow f(\psi) \neq 1$. Contradiction?

My interpretation: flat priors do <u>not</u> mean non-informative.

# Jeffrey's Prior

Set $f(\theta) \propto \sqrt{I(\theta)}$

         ↰ Fisher information

The Jeffrey's Prior is <u>transformation invariant</u>.

If $\tau = g(\theta)$, then what does the prior for $\tau$ look like?

$$I(\theta) = -\mathbb{E}\left( \frac{\partial^2}{\partial \theta^2} \log f(\vec{x}; \theta) \right)$$

                        → statistical model.

$$= -\mathbb{E}\left( \left(\frac{d\tau}{d\theta}\right)^2 \frac{\partial^2}{\partial \tau^2} \log f(\vec{x}; \tau) \right)$$

$$= -\left(\frac{d\tau}{d\theta}\right)^2 \mathbb{E}\left( \frac{\partial^2}{\partial \tau^2} \log f(\vec{x}; \tau) \right)$$

$$= \left(\frac{d\tau}{d\theta}\right)^2 I(\tau)$$

$$\Rightarrow \sqrt{I(\theta)} = \sqrt{I(\tau)} \left| \frac{d\tau}{d\theta} \right|$$

This the usual change of variables formula.