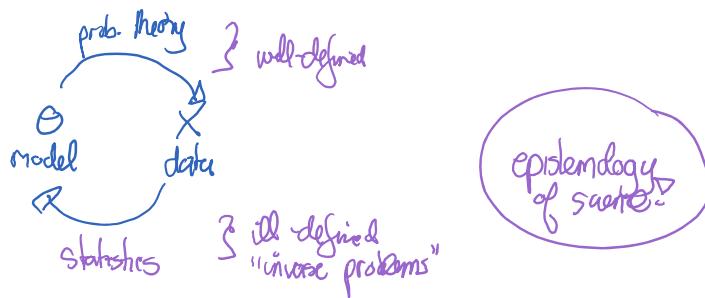


today: statistician frequentist vs. Bayesian

### Statistical concepts

cartoon



example: model  $\Lambda$  indep coin-flips

prob. theory  $\rightarrow$  prob.  $k$  heads in a row

statistics: I have observed  $k$  heads,  $\Lambda-k$  tails, what is  $\Theta$ ?

### Frequentist vs. Bayesian

Semantic of prob.: meaning of a prob?

a) (traditional) frequentist semantic

$P\{X=x\}$  represents the limiting frequency of observing  $X=x$

If I could repeat  $\infty$  # of i.i.d. experiments

b) Bayesian (subjective) semantic

$P\{X=x\}$  encodes an agent "belief" that  $X=x$

laws of prob. characterizes a "rational" way to combine "beliefs"  
and "evidence" [observations]

[has motivation in terms of gambling, utility/decision theory, etc...]

operationally:

Bayesian approach:  $\otimes$  very simple philosophically

treat all uncertain quantities as R.V.

i.e. encode all knowledge about the system ("beliefs")  
as a "prior" on probabilistic models

and then

use law of prob. (and Bayes rule) to  
get updated beliefs and answer?

justification for frequentist semantic

for discrete R.V.  $X$ , suppose  $P\{X=x\} = \theta$

$$\Rightarrow P\{X \neq x\} = 1-\theta$$

$$B \stackrel{\Delta}{=} \mathbb{1}_{\{X=x\}} \quad \Rightarrow B \sim \text{Bern}(\theta) \text{ R.V.}$$

$\mathbb{1}_A(u) = \begin{cases} 1 & \text{if } u \in A \\ 0 & \text{o.w.} \end{cases}$

repeat i.i.d. ie.  $B_i \stackrel{i.i.d.}{\sim} \text{Bern}(\theta)$

by L.L.N.  
law of large  
numbers

$$\frac{1}{n} \sum B_i \xrightarrow{a.s.} E[B_i] = \theta$$

↑ Limiting frequency

by C.L.T.  $\sqrt{n} \left( \frac{1}{n} \sum B_i - \theta \right) \xrightarrow{d} N(0, \theta(1-\theta))$

$\sim N \text{Bin}(n, \theta)$

### Coin flips - Bayesian approach

biased coin flips  $\downarrow$  unknown  $\Rightarrow$  model it as R.V.

we believe  $X \sim \text{Bin}(n, \theta)$   $\Rightarrow$  need a  $p(\theta)$  "prior distribution"

$$\Omega_{\theta} = [0, 1]$$

suppose we observe  $X=x$  (result of  $n$  coin flips)

then we can "update" our belief about  $\theta$  by using Bayes rule

$$p(\theta=\theta | X=x) = \frac{p(X=x | \theta=\theta) p(\theta=\theta)}{p(X=x)}$$

posterior belief

prior belief

normalization

marginal likelihood

observation model / likelihood

[ note:  $p(x|\theta) \rightarrow \text{pmf}$      $p(x, \theta)$  is a "mixed distribution" ]

$p(\theta) \rightarrow \text{pdf}$

example: suppose  $p(\theta)$  is uniform on  $[0, 1]$  "no specific preference"

$p(\theta|x) \propto p(x|\theta) p(\theta)$

"proportional to"  $p(\theta|x) \propto \theta^x (1-\theta)^{n-x} \cdot \mathbb{1}_{[0,1]}(\theta)$

up to a constant in  $\theta$

Scaling:  $\int_0^1 \theta^x (1-\theta)^{n-x} d\theta = B(x+1, n-x+1)$

Scaling:  $\int_0^1 \theta^{\alpha-1} (1-\theta)^{n-x} d\theta = B(x+1, n-x+1)$

normalization constant  $\int_0^1 p(\theta|x) d\theta = 1$

$B(a,b) \triangleq \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$

$\Gamma(a) = \int_0^\infty u^{a-1} e^{-u} du$

beta-fct.      gamma-fct.

here  $p(\theta|x)$  is called a "beta distribution"

$$B(\theta|\alpha, \beta) \triangleq \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)} \mathbb{1}_{[0,1]}(\theta)$$

parameters

• uniform distribution:  $B(\theta|1,1)$

posterior density:  $B(\theta|x+1, n-x+1)$  "prior counts"

exercise to the reader: if use  $B(\alpha_0, \beta_0)$  as prior

posterior will be  $B(\theta|x+\alpha_0, n-x+\beta_0)$

15h32

② posterior  $p(\theta|X=x)$  contains all the info from data  $x$  that we need  
observation to answer questions about  $\theta$

e.g. question: what is prob. of head ( $F=1$ ) on the next flip?

as a frequentist  $P(F=1 | \text{data}) = \hat{\theta}$  relation to mean "estimate"

$$\begin{aligned} \text{as a Bayesian } P(F=1 | X=x) &= \int_{\theta} p(F=1, \theta | x) d\theta \\ &= \int_{\theta} \underbrace{p(F=1 | \theta=x, X=x)}_{=\hat{\theta} \text{ (by our model)}} \underbrace{p(\theta | x=x)}_{\text{posterior}} d\theta \\ &\approx \int_{\theta} \hat{\theta} p(\theta | x=x) d\theta = \mathbb{E}[\theta | X=x] \end{aligned}$$

"posterior mean" of  $\theta$

\* a meaningful "Bayesian" estimator of  $\theta$

$$\hat{\theta}_{\text{Bayes}}(x) \triangleq \mathbb{E}[\theta | X=x] \quad (\text{posterior mean})$$

relation:  $\hat{\theta}$ : observation  $\rightarrow \hat{\theta}_{\text{Bayes}}$

Our coin example:  $p(\theta|x) = \text{Beta}(\theta|\alpha=x+1, \beta=n-x+1)$

mean of a beta R.V.  $\frac{\alpha}{\alpha+\beta}$

thus  $\hat{\theta}_{\text{Bayes}}(x) = \mathbb{E}[\theta | x] = \frac{x+1}{n+2}$

$$\text{thus } \hat{\theta}_{\text{Bayes}}(x) = \mathbb{E}[G|x] = \frac{x+1}{n+2}$$

here, biased estimator  $\mathbb{E}_X[\hat{\theta}(x)] \neq G$

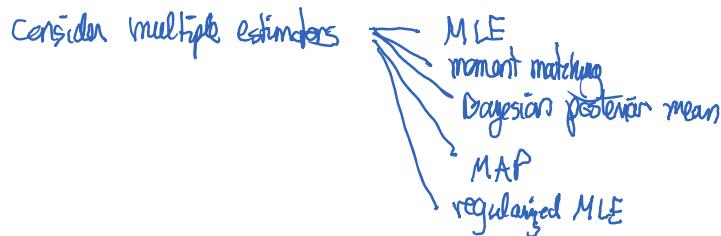
$$= \mathbb{E}\left[\frac{x+1}{n+2}\right] = \frac{\mathbb{E}x+1}{n+2} = \frac{nG+1}{n+2} \neq G$$

but asymptotically unbiased  $\xrightarrow{n \rightarrow \infty} G$

compare & contrast with  $\hat{\theta}_{\text{MLE}}(x) = \frac{x}{n}$  [unbiased  $\mathbb{E}[\frac{x}{n}] = \frac{\mathbb{E}x}{n} = G$ ]

To summarize:

- as a Bayesian: get a posterior + use law of probabilities
- in "frequentist statistics"



and then analyze the statistical properties of estimator:

- biased?
- variance?
- consistent?
- frequentist risk?

### Maximum Likelihood principle

Setup: given a parametric family  $p(x; \theta)$  for  $\theta \in \Theta$

we want to estimate/learn  $\theta$  from  $x$

$$\hat{\theta}_{\text{MLE}}(x) \triangleq \underset{\theta \in \Theta}{\operatorname{argmax}} \quad p(x; \theta)$$

$L \triangleq L(\theta)$

"likelihood function" of  $\theta$

$\hat{\theta}_{\text{MLE}}(x)$  maximizes  
 $p(x; \cdot)$

MLE example I: binomial:

n coin-flips  $\sum x = 0:n$

$$X \sim \text{Bin}(n, \theta) \quad p(x; \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

trick: to maximize  $\log L(\theta)$  instead of  $L(\theta)$

trick: to maximize  $\log L(\theta)$  instead of  $L(\theta)$

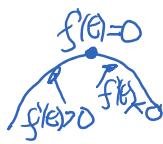
$$\stackrel{!}{=} l(\theta) \quad \boxed{\text{log likelihood}}$$

justification:  $\log(\cdot)$  is strictly increasing

$$\text{i.e. } a < b \Leftrightarrow \log a < \log b \quad (\forall a, b > 0)$$

$$\Rightarrow \underset{\theta \in \Theta}{\operatorname{argmax}} \log p(x; \theta) = \underset{\theta \in \Theta}{\operatorname{argmax}} p(x; \theta)$$

$$\log p(x; \theta) = \underbrace{\log \binom{n}{x}}_{\text{constant w.r.t. } \theta} + x \log \theta + (n-x) \log (1-\theta) \approx l(\theta)$$



look for  $\theta$  st.  $\frac{\partial l}{\partial \theta} = 0$

$$\text{want: } \frac{\partial \ell}{\partial \theta} = \frac{x}{\theta} - \frac{(n-x)}{1-\theta} = 0$$

$$x(1-\theta) = \theta(n-x), \\ x - x\theta = n\theta - \theta^2 x$$

$$\text{hence } \boxed{\hat{\theta}_{MLE}(x) = \frac{x}{n}}$$

$$\boxed{\theta^* = \frac{x}{n}}$$

Used often  
as solution in optimization