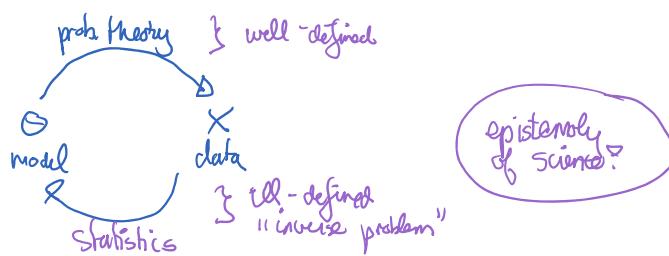


today: statistics frequentist vs. Bayesian

### Statistical concepts

Cartoon



example: model n indep. coin-flips

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

statistics: I have observed k heads, n-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 most 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:  $\textcircled{R}$  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 laws of prob. (and Bayes rule) to  
get updated beliefs and answer?

Justification for frequentist semantic:

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

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

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

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

repeat i.i.d experiments  
by L.L.N.  
(law of large numbers)  
by CLT.  $\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n b_i - \theta \right) \xrightarrow{d} N(0, \theta(1-\theta))$   
 $b_i \sim \text{Bern}(\theta)$   
 $\frac{1}{n} \sum_{i=1}^n b_i \xrightarrow{\text{a.s.}} E[B_i] = \theta$   
↳ limiting frequency

### Coin flips - Bayesian approach

biased coin flips  $\xrightarrow{\text{unknown}}$  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

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

↓  
 posterior belief      ↓  
 observation model      "marginal likelihood"  
 ↓  
 prior belief      ↓  
 normalization

[ 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"

$$p(\theta|x) \propto \frac{p(x|\theta) p(\theta)}{p(x)}$$

↓  
 "proportional to"  
 ↓  
 up to a constant

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

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

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

↓  
 Beta fct.

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

↓  
 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  $B(\theta|x+1, n-x+1)$  as "prior count!"

Exercise to the reader: If use  $B(\alpha_0, \beta_0)$  as prior

16NO7

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

⊕ posterior  $p(\theta | \underset{\text{observation}}{X=x})$  contains all the info from data  $x$  that we need to answer queries about  $\theta$

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

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

as a Bayesian  $P(F=1 | X=x) = \int_{\theta} P(F=1, \theta | X=x) d\theta$

$$\approx \int_{\theta} p(F=1 | \theta, x=x) p(\theta | x=x) d\theta$$

product rule  
by own model

$$= \int_{\theta} \theta p(\theta | x=x) d\theta = \mathbb{E}[\theta | x=x]$$

"Posterior mean of  $\theta$ "

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

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

notation:  $\hat{\theta}$ : observation  $\rightarrow$  ⊕

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

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

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

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

$$= \mathbb{E}\left[\frac{X+1}{n+2}\right] = \frac{\mathbb{E}X+1}{n+2} = \frac{n\theta+1}{n+2} \neq \theta$$

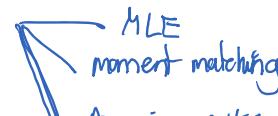
but asymptotically unbiased

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

to summarize:

- as a Bayesian: get a posterior + use laws of probabilities
- in "frequent statistics":

consider multiple estimators



consider multiple estimators

MLE  
moment matching

Bayesian posterior means

MAP  
regularized MLE

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  $x$

$\hat{\theta}_{MLE}(x)$  maximizes

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

$\hat{L}(\theta) \triangleq L(\theta)$   
"likelihood function" of  $\theta$

$p(x; \cdot)$

### 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  $L(\theta)$   
 $\stackrel{\triangle}{=} l(\theta)$  log likelihood

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

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

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

$$\log p(x; \theta) = \log \binom{n}{x} + x \log \theta + (n-x) \log(1-\theta) = l(\theta)$$

$$f'(\theta) = 0$$

constant w.r.t.  $\theta$

$$f'(0) > 0 \quad f'(1) < 0$$

look for  $\theta$  s.t.  $\frac{d l}{d \theta} = 0$

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

$$x(1-\theta) = \theta(n-x)$$

$$x - x\theta = n\theta - x\theta$$

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

$\boxed{\theta = x/n}$   
Used often  
as solution  
in optimization

hence

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