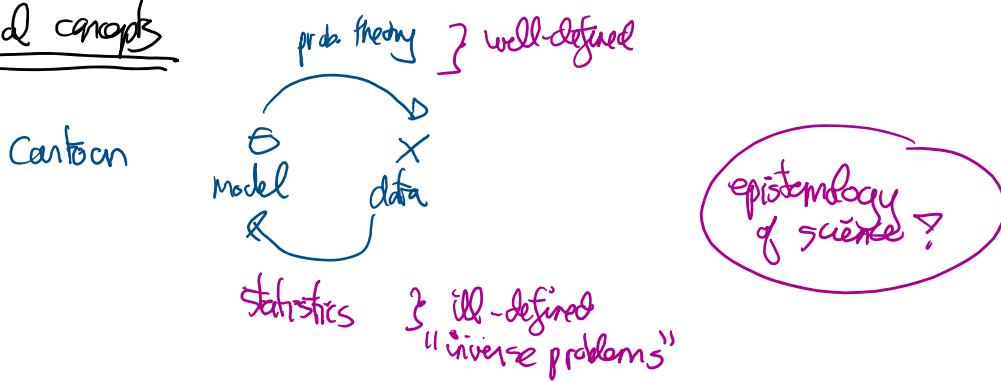


today: statistical frequentist vs. Bayesian

## Statistical concepts



Example : Model  $n$  independent coin flips

prob. theory  $\rightarrow$  prob.  $k$  heads in a 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  $X=x$

I could repeat  $N$  # of i.i.d experiments

b) Bayesian (subjective) semantic

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

Views 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:  $\textcircled{S}$  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

as a "prior" on probabilistic models

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

justification for frequentist semantics:

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

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

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

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

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

by L.L.N.  
Law of Large  
numbers

$$\frac{1}{n} \sum_{i=1}^n B_i \xrightarrow{\text{a.s.}} \mathbb{E}[B_i] = \epsilon$$

↑ limiting frequency

by C.L.T.

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n B_i - \epsilon \right) \xrightarrow{d} N(0, \epsilon(1-\epsilon))$$

$\sim \text{Bin}(n, \epsilon)$

### Coin-flips - Bayesian approach

biased coin-flips

unknown  $\Rightarrow$  model it as R.V.

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

$$\Omega_H = [0, 1]$$

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

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

$$p(H=\epsilon | X=x) = \frac{p(X=x | H=\epsilon) p(H=\epsilon)}{p(X=x)} \rightarrow \text{prior belief}$$

$\downarrow$  posterior belief       $\downarrow$  observation model / likelihood       $\leftarrow$  "marginal likelihood" (normalization)

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

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

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

<sup>↑</sup>  
"proportional to"

$$p(\theta|x) \propto \theta^x (1-\theta)^{n-x} \underbrace{1_{[0,1]}(\theta)}$$

up to a constant in  $\theta$

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

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

$$\underbrace{B(a,b)}_{\text{beta fn.}} \triangleq \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

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

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

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

- uniform dist:  $B(\theta|1,1)$   
density

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

exercise to the reader: If use  $B(\theta_0, \beta_0)$  as prior

posterior will be  $B(\theta+x_0, n-x+\beta_0)$

13h34

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

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

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

$$\begin{aligned} \text{as a Bayesian } P(F=1 | X=x) &= \int_{\theta} P(F=1, \Theta=\theta | X=x) d\theta \\ &= \int_{\theta} P(F=1 | \Theta=\theta, X=x) P(\theta | X=x) d\theta \\ &\quad \text{product rule} \quad \text{posterior} \\ &\quad \text{by own model} \\ &= \int \theta p(\theta | X=x) d\theta = \mathbb{E}[\theta | X=x] \end{aligned}$$

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

*"posterior mean" of  $\theta$*

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

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

notation:  $\hat{\theta}$ : observation  $\rightarrow \Omega_{\Theta}$

An coin example:  $p(\theta|x) = \text{Beta}(\theta|\alpha=x+\alpha_0, \beta=n-x+\beta_0)$

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

$$\text{thus } \hat{\theta}_{\text{Bayes}}(x) = \mathbb{E}[\theta|x] = \frac{x+\alpha_0}{n+\alpha_0+\beta_0}$$

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

$$\mathbb{E}\left[\frac{X+\alpha_0}{n+\alpha_0+\beta_0}\right] = \frac{\mathbb{E}X + \alpha_0}{n+\alpha_0+\beta_0} \neq \theta \text{ unless } \alpha_0 = \beta_0 = 0 \quad (\text{not allowed})$$

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

compare & contrast  $\hat{\theta}_{\text{MLE}}(x) = \bar{x}$  (unbiased)

To summarize:

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

consider multiple estimators

MLE  
 moment matching  
 Bayesian posterior mean  
 MAP  
 Regularized MLE  
 ...

and then analyze the statistical properties of estimator?

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

Maximum Likelihood principle

## Maximum Likelihood principle

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

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

$$\hat{\theta}_{MLE}(x) \triangleq \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n p(x_i; \theta)$$

"Likelihood fct." of  $\theta$

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

MLE example I: binomial

$n$  coin flips  $\sum x_i = 0:n$

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

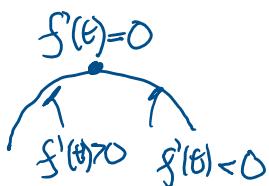
trick: to maximize  $\log L(\theta)$  instead of  $L(\theta)$   
 $\triangleq l(\theta)$  Log-Likelihood

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

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

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

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



Look for  $\theta$  s.t.  $\frac{dl}{d\theta} = 0$

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

Used as  
solution in  
optimization

$$\cancel{x(1-\theta)} = \theta(n-x) \Rightarrow \cancel{x - \theta x} = n\theta - \theta x$$

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

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