

today: - finish prob.  
 • frequentist vs. Bayesian

binomial distribution:

model  $n$  indep. coin flips

sum of  $n$  indep.  $\text{Bern}(\theta)$  R.V.

let  $X_i \sim \text{Bern}(\theta)$  <sup>iid</sup>  $\rightarrow$  (mutually) independent and identically distributed

implicitly talking  $\rightarrow X_1, X_2, \dots, X_n$  iid

let  $X = \sum_{i=1}^n X_i$  then we have  $X \sim \text{Bin}(n, \theta)$

"binomial with parameter  $n$  &  $\theta$ "

$\Omega_X = \{0, 1, \dots, n\}$

proof:  $p(x; \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$  for  $x \in \Omega_X$

$\binom{n}{x} = \frac{n!}{x!(n-x)!}$

# of ways to choose  $x$  elements out of  $n$  "n choose x"

$p(x_1, \dots, x_n) = \prod_{i=1}^n \text{bern}(x_i; \theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum x_i} (1-\theta)^{n - \sum x_i}$

mean:  $X = \sum_{i=1}^n X_i$

$E[X] = \sum_{i=1}^n E[X_i] = n\theta$

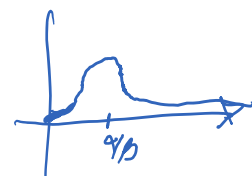
similarly,  $\text{Var}(\sum_{i=1}^n X_i) = \sum_{i=1}^n \text{Var}(X_i) = n\theta(1-\theta)$

other distributions: Poisson  $(\lambda)$   $\Omega_X = \{0, 1, \dots\} = \mathbb{N}$  [count data]  
 mean  $\lambda$  mean variance

Gaussian in 1D  $N(\mu, \sigma^2)$   $\Omega_X = \mathbb{R}$

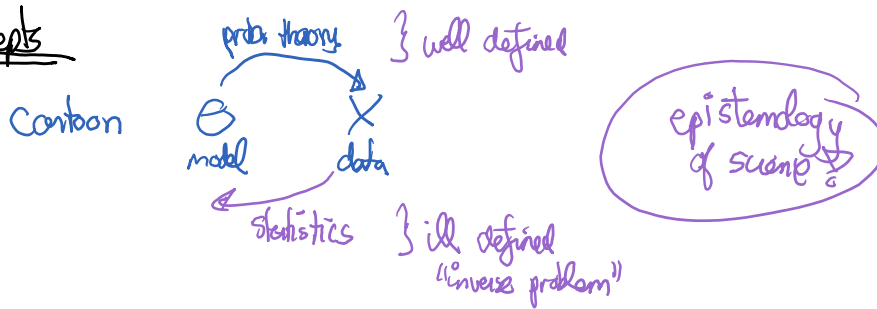
gamma  $\Gamma(\alpha, \beta)$   $\Omega_X = \mathbb{R}^+$

shape  $\alpha$  inverse cka. rate  $\beta$   
 mean  $\frac{\alpha}{\beta}$  variance  $\frac{\alpha}{\beta^2}$



other: Laplace, Cauchy, exponential  $\rightarrow \Gamma(1, \beta)$ , beta  $\rightarrow$  Dirichlet, Dirichlet on 2 demands  $\Omega_X = [0, 1]$

# Statistical Concepts



example: model  $n$  indep. coin flips

prob. theory → <sup>compute</sup> prob.  $k$  heads in a row

statistics: I have observed  $k$  heads and  $n-k$  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  $\infty$  # of iid 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.]

Bayesian approach: (★) 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 answers

Justification for frequentist semantic → see notes from last year ([lecture4 scribbles 2020](#))

coin flips → Bayesian approach

biased coin flip

unknown ⇒ model it as a R.V.

$\theta \sim \text{Bern}(\theta)$

⇒  $x \sim \text{Bern}(\theta)$  (i.i.d.)

biased coin flip

unknown  $\Rightarrow$  model it as a 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 "update" our belief about  $\theta$  using Bayes rule

$$p(\theta = \theta | X=x) = \underbrace{p(x=z | \theta)}_{\substack{\text{observation model /} \\ \text{likelihood}}} \underbrace{p(\theta = \theta)}_{\substack{\text{"prior belief"} \\ p(\theta)}} \underbrace{1}_{\substack{\text{normalization} \\ \text{"marginal likelihood"}}$$

[note:  $p(x|\theta) \rightarrow$  pmf  $p(\theta, x)$  is a "mixed distribution"  
 $p(\theta) \rightarrow$  pdf

example:

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

$$p(\theta|x) \propto \underbrace{p(x|\theta)p(\theta)}_{\substack{\text{"proportional to"} \\ p(x|\theta) \text{ up to a constant}}} \propto \theta^x (1-\theta)^{n-x} \underbrace{\mathbb{1}_{[0,1]}(\theta)}_{p(\theta)}$$

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

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 fct.

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

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  $B(\theta | x+1, n-x+1)$

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

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

17h37

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

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"

as a Bayesian  $P(F=1 | X=x) = \int p(F=1, \theta | X=x) d\theta$

$$= \int_{\theta} \underbrace{p(F=1 | \theta, X=x)}_{\text{by an model}} \underbrace{p(\theta | X=x)}_{\text{posterior}} d\theta$$

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

"posterior mean of  $\theta$ "

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

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

relation:  $\hat{\theta} : \text{observation} \rightarrow \Theta$

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  $\mathbb{E}_X[\hat{\theta}(X)] \neq \theta$

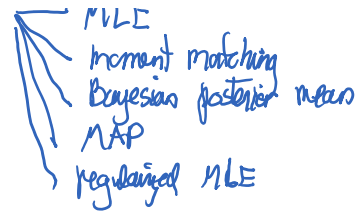
but asymptotically unbiased  $\xrightarrow{n \rightarrow \infty} \mathbb{E}[\frac{X+1}{n+2}] = \frac{\mathbb{E}[X]+1}{n+2} = \frac{n\theta+1}{n+2} \rightarrow \theta$

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

to summarize:

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

Various multiple estimators



and then analyze the statistical properties:   
 • biased?   
 • variance?   
 • consistent?

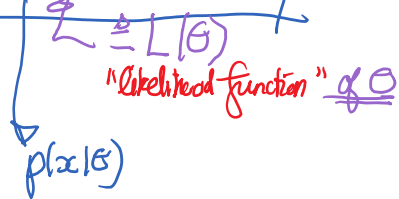
Maximum likelihood principle

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

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

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

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



MLE example I: binomial:

n coin flips

$$\Omega_x = 0:n$$

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

$$p(x; \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

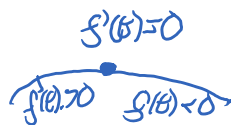
trick: to maximize  $\log L(\theta)$  instead of  $L(\theta)$   
 $\triangleq \ell(\theta)$  log likelihood

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

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

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

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



look for  $\theta$  s.t.  $\frac{\partial \ell}{\partial \theta} = 0$

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

$$x(1-\theta) - \theta(n-x) = 0$$

$$\theta = x \quad \left( \text{use often as } \frac{x}{n} \right)$$

$$x(1-x) - 0(1-x) = 0$$

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

use often as soln in optimization

hence  $\hat{\theta}_{ML}(x) = \frac{x}{n}$  i.e. relative frequency

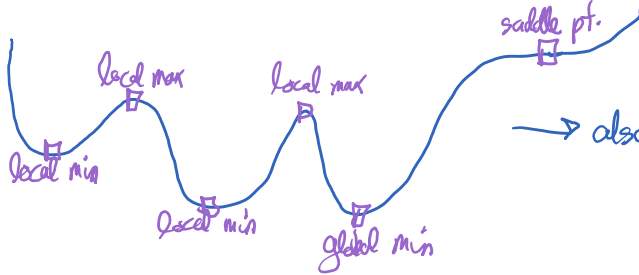
some optimization comments:

$$\min_{\theta \in \Theta} f(\theta)$$

$$(*) \nabla f(\hat{\theta}) = 0$$

"stationary pts."

(if  $f$  is diff) is a necessary condition for  $\hat{\theta}$  being a local min when  $\hat{\theta}$  is in the interior of  $\Theta$



$\hat{\theta}$  being a local min when  $\hat{\theta}$  is in the interior of  $\Theta$

→ also need to check  $\text{Hessian}f(\hat{\theta}) \succ 0$  for a local min

$$\text{scalar}(S^{-1}(\hat{\theta}) \succ 0)$$

$$H \succ 0$$

$$\Leftrightarrow u^T H u > 0 \quad \forall u \neq 0 \in \mathbb{R}^d$$

(\*) only local result in general

• but if  $\text{Hessian}f(\hat{\theta}) \succeq 0 \quad \forall \theta \in \Theta$ ,  $f(\theta)$  is said "convex"

and in this case,  $\nabla f(\hat{\theta}) = 0 \Rightarrow$  sufficient for  $\hat{\theta}$  to be a global min

• otherwise, for smooth  $f(\theta)$ , look at zero gradient pts and boundary pt.

give you enough information to find global optima