

## Lecture 4 - scribbles

Friday, September 14, 2018 13:34

today:

- continue Bayesian approach
- MLE

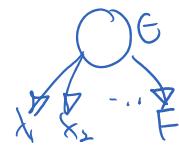
(continuation from last class):

posterior belief  $p(\theta | X=x)$  contains all the info from data that we need to answer new questions  
 observations coin flip outcome

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

as frequentist:  $P(F=1 | \Theta) = \hat{\Theta}$

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



$$p(x, y | z) = p(x|y, z) p(y|z)$$

$$p(x, y | z, w) = p(x|y, z, w) p(y|z, w)$$

$$\begin{aligned} & \int_{\Theta} p(F=1 | X=x, \Theta) p(\Theta | X=x) d\Theta \\ & \quad \text{(always true)} \\ & \quad \text{|| (by our model)} \\ & \quad P(F=1 | \Theta) = \hat{\Theta} \\ & = \int_{\Theta} \Theta p(\Theta | X=x) d\Theta = \mathbb{E}[\Theta | X=x] \\ & \quad \text{"posterior mean" of } \Theta \end{aligned}$$

a meaningful Bayesian estimator of  $\Theta$

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  $\rightarrow$   $\hat{\theta}$ )

$$p(\theta) = \text{unif}(\theta) \\ = \text{Beta}(\theta | \alpha=1, \beta=1)$$

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

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

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

but asymptotically unbiased  
biased i.e.  $\mathbb{E}_X \hat{\theta}(x) \neq \theta$

asympt. unbiased ← compare & contrast with  $\hat{\theta}_{\text{MLE}}(x) = \frac{x}{n}$  ← unbiased

$$\mathbb{E}[\hat{\theta}_{\text{Bayes}}(x)] = \frac{\mathbb{E}X+1}{n+2} = \frac{n\theta+1}{n+2} \xrightarrow{n \rightarrow \infty} \theta$$

$$\text{i.e. } \mathbb{E}_{X|E}[\hat{\theta}_{\text{MLE}}(x)] = \frac{\theta n}{n} \\ = \theta$$

to summarize:

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

consider multiple possible estimators

MLE  
moment matching  
Bayesian Posterior mean

moment matching  
Bayesian posterior mean  
MAP

and then analyze their statistical properties

- biased?
- variance?
- consistent?

14h20

### Maximum Likelihood principle

setup: • given a parametric family  $p(x; \theta)$  for  $\theta \in \Theta$   
           • we want to estimate  $\theta$

$$\hat{\theta}_{ML}(x) \triangleq \underset{\theta \in \Theta}{\operatorname{argmax}} p(x; \theta) \quad \text{ie. } \hat{\theta}_{ML}(x) \text{ maximizes } p(x; \cdot)$$

$\triangleq L(\theta)$   
                         "likelihood function" of  $\theta$

example:  $n$  coin flips       $\mathcal{L}_X = 0 : n$

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

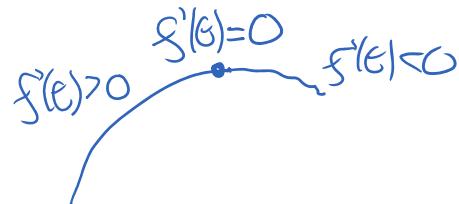
trick: to maximize  $\log L(\theta)$  instead of  $L(\theta)$

$\triangleq \ell(\theta)$   
                         [log likelihood]

Justification:  $\log(\cdot)$  is strictly increasing  
 i.e.  $a < b \Leftrightarrow \log a < \log b$

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

$$\log p(x; \theta) = \underbrace{\log \left( \frac{x}{n} \right)}_{\text{constant}} + x \log \theta + (n-x) \log (1-\theta) = D(\theta)$$



Look for  $\theta$  s.t.  $\frac{\partial L(\theta)}{\partial \theta} = 0$

$$\text{i.e. } \frac{x}{\theta} - \frac{(n-x)}{1-\theta} = 0$$

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

*often used in optimization*

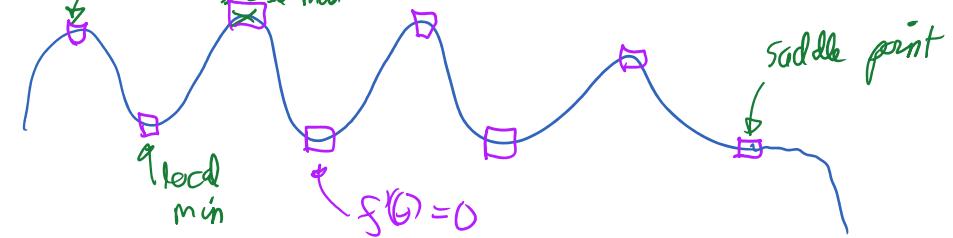
$$\Rightarrow \theta^* = \frac{x}{n}$$

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

### Some optimization comments

in multiple dim.  
 $(\nabla f(\theta)) = 0$ ) •  $f'(\theta) = 0$  is necessary condition for a local max when  $\theta$  is in interior of  $\Theta$   
 → also need to check  $f''(\theta) < 0$  for a local max  
local max      global max

→ also need to check  $f''(E) < 0$  for a local max



dxd matrix



→ only local result in general

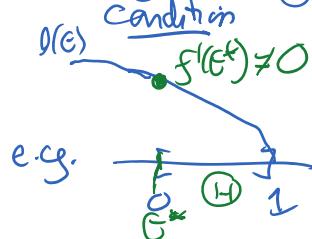
Hessian ( $E$ )  $\leq 0$

i.e. negative  
semi-definite  $\Leftrightarrow u^T A u \leq 0$   
Hence  $A$

but if  $f'(E) \leq 0 \quad \forall E \in \Theta$ , function is said "concave"

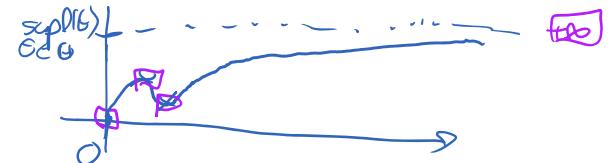
in this case,  $f'(E) = 0$  is sufficient for global max

- Be careful with boundary cases i.e.  $E^* \in \text{boundary}(\Theta)$



another example

$$\Theta = [0, +\infty[$$



### Some notes about MLE

- does not always exist [ $E^* \in \text{bd}(\Theta)$  but  $\Theta$  is open] or when " $E^* = +\infty$ "

e.g.  $\Theta = ]0, 1[$

- is not necessarily unique [i.e. multiple maxima]

- \* is not "admissible" in general [ see next class ]  
 (basically,  $\exists$  strictly "better" estimators)

example 2 : Multinomial distribution

suppose  $X_i$  is discrete R.V. on  $K$  choices "Multinoulli"

(we could choose  $\Omega_{X_i} = \{1, \dots, K\}$ )

but instead, convenient to encode with unit basis in  $\mathbb{R}^K$

i.e.  $\Omega_{X_i} = \{e_1, \dots, e_K\}$  where  $e_j \in \mathbb{R}^K$  "one hot encoding"

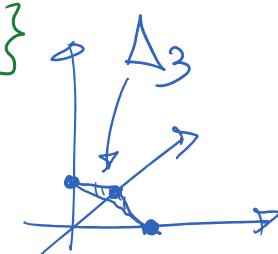
$$e_j = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix} \quad j^{\text{th}} \text{ coordinate}$$

parameters for discrete R.V.:  $\pi \in \Delta_K$

$$(\Theta = \Delta_K)$$

$$\Delta_K \triangleq \left\{ \pi \in \mathbb{R}^K : \pi_j \geq 0 \forall j; \sum_{j=1}^K \pi_j = 1 \right\}$$

probability simplex on  $K$  choices



We will write  $X_i \sim \text{Mult}_q(\pi)$  "multinoulli" parameter

④ consider  $X_i \stackrel{\text{iid.}}{\sim} \text{Mult}(\pi)$

then  $X \triangleq \sum_{i=1}^n X_i \sim \text{Mult}(n, \pi)$

"multinomial" distribution

$$X \in \mathbb{N}^K$$

$$\Omega_X = \left\{ (n_1, \dots, n_k) : \begin{array}{l} n_j \in \mathbb{N} \\ \sum_{j=1}^k n_j = n \end{array} \right\}$$