

Lecture 4 - scribbles

Friday, September 13, 2019 14:03

today:
 • Bayesian approach
 • MLE

coin flips - Bayesian approach:

biased coin flips

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

↑ unknown \Rightarrow model it as a R.V.

\Rightarrow need a $p(\theta)$ "prior distribution"

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

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

then we can "update" our belief about θ using Bayes rule

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

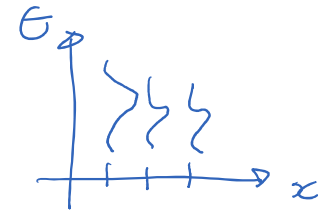
posterior belief prior belief observation model / likelihood normalization "marginal likelihood"

$$P(\{\theta = \theta\})$$

$p(\theta = \theta)$ I mean pdf $p(\theta)$ of R.V. θ

[note: $p(x|\theta) \rightarrow$ is a pmf $p(x, \theta)$ is a "mixed" dist. θ

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



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

$$p(\theta|x) \propto \underbrace{\theta^x (1-\theta)^{n-x}}_{p(x|\theta) \text{ up to scaling}} \underbrace{\mathbb{1}_{[0,1]}(\theta)}_{p(\theta)} \binom{n}{x}$$

↑
"proportional to"

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

normalization constant
 so that $\int_{\theta} p(\theta|x) d\theta = 1$

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

Beta function

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

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

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

↑
parameters

- uniform distribution: $B(\theta | 1, 1)$
- posterior here was $B(\theta | x+1, n-x+1)$

exercise to the reader: if use $B(\alpha_0, \beta_0)$ as prior
posterior will be $B(x+\alpha_0, n-x+\beta_0)$

⇒ posterior belief $p(\theta | X=x)$ contains all the info from data x that we need to answer new queries

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

as a frequentist: $P(F=1 | \text{data}) = \hat{\theta}$

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

$\stackrel{\text{(product rule)}}{=} \int_{\Theta} \underbrace{p(F=1 | X=x, \theta)}_{\text{by our model}} \underbrace{p(\theta | X=x)}_{\text{posterior}} d\theta$

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

"posterior mean" of θ

* a meaningful "Bayesian" estimator of θ

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

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

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

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

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

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

but asymptotically unbiased

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

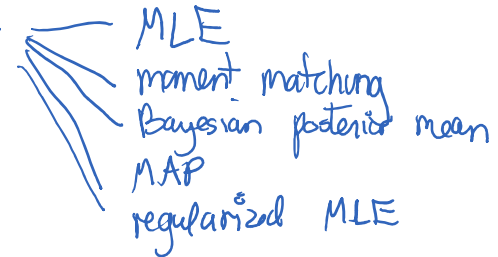
compare { contrast with $\hat{\theta}_{\text{MLE}}(x) = \frac{x}{n}$

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

to summarize

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

consider multiple possible estimators



and then analyze their statistical properties

- biased?
- variance?
- consistent?

15n35

Maximum likelihood principle

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

we want to estimate/learn θ

$$\hat{\theta}_{ML}(x) \triangleq \underset{\theta \in \Theta}{\text{arg max}} p(x; \theta)$$

$(\triangleq L(\theta))$

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

"likelihood function" of θ

example: n coin flips
 $X \sim \text{Bin}(n, \theta)$

$$\rightarrow X \in 0:n$$

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

$\forall x = 0:n$

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

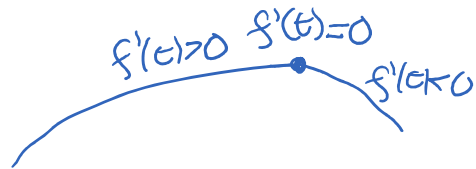
trick: to maximize $\log L(\theta)$ instead of $L(\theta)$
 $\triangleq \ell(\theta)$ log-likelihood

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

$$\text{i.e. } a < b \Leftrightarrow \log a < \log b$$

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

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



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

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

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

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

used often
as solution
in optimization

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

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

Some optimization comments

$\min_{\theta \in \Theta} f(\theta)$ \circledast $\nabla f(\theta) = 0$ is a necessary condition for a local min when θ is in the interior of Θ (when f is differentiable on Θ)

"stationary pts."

\rightarrow also need to check $\text{Hessian}(f) > 0$ for a local min ($f''(\theta) > 0$)

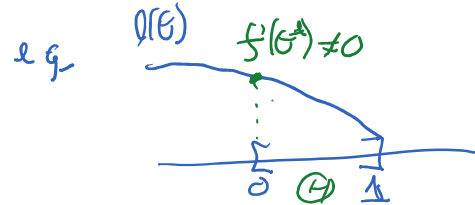


\circledast only local result in general

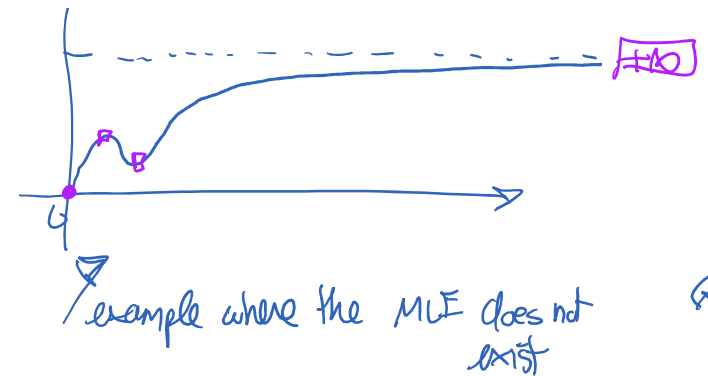
• but if $\text{Hessian}(f(\theta)) \geq 0 \forall \theta \in \Theta$, $f(\cdot)$ is said "convex" and in this case $\nabla f(\theta) = 0 \Rightarrow$ sufficient for global min

- otherwise, for a smooth fct., looking at zero gradient & boundary gives you enough information

⊗ be careful with boundary cases i.e. $\theta^* \in \text{boundary}(\Theta)$



Another example



⊗ Some notes about MLE

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

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

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

- is not "admissible" in general [see later]
 ∃ strictly "better" estimators

15h36

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

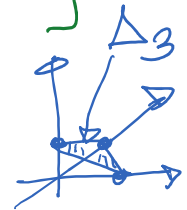
ie. $\Omega_{X_i} = \{e_1, \dots, e_k\}$ where $e_j \in \mathbb{R}^k$ "one hot encoding"

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

parameter for discrete R.V.: $\pi \in \Delta_k$ ($\mathbb{H} = \Delta_k$)

$$\Delta_k \triangleq \left\{ \pi \in \mathbb{R}^k : \pi_j \geq 0 \forall j, \sum_{j=1}^k \pi_j = 1 \right\}$$

probability simplex on k choices



we will write $X_i \sim \text{Mult}(\pi)$ "multinoulli"

parameter

consider iid
 $\otimes X_i \sim \text{Mult}(\pi)$

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

"multinomial distribution"

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

pmf for X :

$$p(x|\pi) = \binom{n}{x_1, \dots, x_k} \prod_{j=1}^k \pi_j^{x_j}$$

$$x = (n_1, \dots, n_k)$$

↳ multinomial coef.

$$\binom{n}{n_1, \dots, n_k} \stackrel{!}{=} \frac{n!}{n_1! \dots n_k!}$$