

Lecture 5 - scribbles

Tuesday, September 15, 2020 14:20

- today:
- MLE
 - statistical decision theory

MLE example I: binomial

n coin flips

$$\Omega_X = 0:n$$

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

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

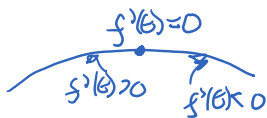
trick: to maximize $\log L(\theta)$ instead $L(\theta)$
 $\stackrel{\Delta}{=} \ell(\theta)$ log likelihood

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

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

$$\Rightarrow \arg \max_{\theta \in \Theta} \log p(x; \theta) = \arg \max_{\theta \in \Theta} p(x; \theta)$$

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



look for θ st. $\frac{\partial \ell}{\partial \theta} = 0$

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

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

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

used often
as solution
in optimization

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

i.e. relative frequency

some optimization comments

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

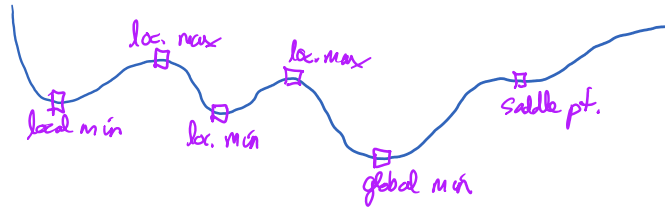
$$\nabla f(\theta^*) = 0$$

"stationary pts."

(if f is diff.) is a necessary condition for

θ^* being a local min when θ^* is in the interior of Θ

→ also need to check $\text{Hessian}(f)(\sigma^*) \succ 0$ for a local min ($f''(\sigma^*) > 0$)



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

⊗ only local result in general

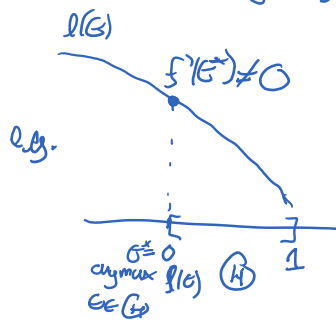
• but if $\text{Hessian}(f(\sigma)) \succeq 0 \quad \forall \sigma \in \Theta$, $f(\cdot)$ is said "convex"

and in this $\nabla f(\sigma) = 0 \Rightarrow$ sufficient for σ^* to be global min

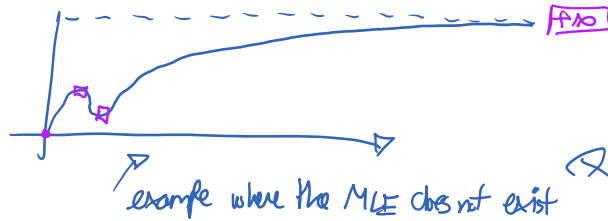
• otherwise, for smooth $f(\cdot)$, looking at zero gradient pts and boundary pts

give you enough information to find global optima

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



another example:



* Some notes about MLE

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

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

• is not nec. unique [i.e. multiple maximax]
e.g. mixture models



• is not "admissible" in general [see later]
⇒ strictly "better" estimators

example II: 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 "one hot encoding"

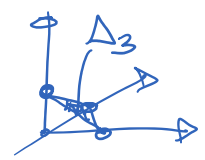
i.e. $\Omega_{X_i} = \{e_1, \dots, e_k\}$ where $e_j \in \mathbb{R}^k$

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

parameter for discrete R.V. : $\pi \in \Delta_k$ ($\Delta_k = \Delta_k$)

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

probability simplex on k choices



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

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

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

↑ "multinomial distribution"

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

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

pmf for X :

$$p(x|\pi) = \binom{n}{(x)_1, \dots, (x)_k} \prod_{j=1}^k \pi_j^{(x)_j} \quad x = (n_1, \dots, n_k)$$

↑ multinomial coeff.

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

15h37

multinomial MLE :

$$\log\text{-likelihood} : l(\pi) = \log p(x|\pi) = \underbrace{\log \binom{n}{n_1, \dots, n_k}}_{\text{constant} \rightarrow \text{ignore for MLE}} + \sum_{j=1}^k n_j \log \pi_j$$

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

$$\text{MLE} : \hat{\pi}_{MLE}(x) = \text{argmax}_{\pi} l(\pi)$$

$$\pi \in \mathbb{R}^k$$

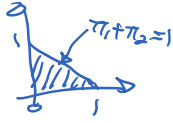
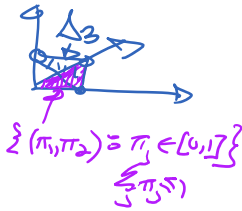
$$\text{s.t. } \pi \in \Delta_k \} \text{ constraint}$$

two options:

a) reparameterize problem so that Θ is full dimensional

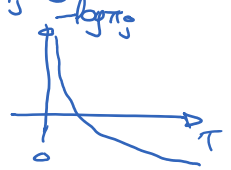
$$\pi_k = 1 - \sum_{j=1}^{k-1} \pi_j$$

$$\rightarrow \pi_1, \dots, \pi_{k-1} \in [0, 1] \text{ with constraint } \sum_{j=1}^{k-1} \pi_j \leq 1$$



here magic is that $\log \pi_j$ acts as barrier fct. away from $\pi_j = 0$

can try unconstrained opt. on π_1, \dots, π_{k-1} of $l(\pi_1, \dots, \pi_{k-1})$
hoping sol'n is on the exterior of constraint set (and it usually will)



b) use Lagrange multiplier approach to handle equality constraints on Δ_k
 [and still ignore $\pi_j \in [0, 1]$]

$$\sum_{j=1}^k \pi_j = 1 \Rightarrow$$

$$\begin{aligned} \max f(\pi) \\ \text{s.t. } g(\pi) = 0 \\ 1 - \sum_{j=1}^k \pi_j = 0 \\ \triangleq g(\pi) \end{aligned}$$

$$J(\pi, \lambda) \triangleq f(\pi) + \lambda g(\pi)$$

Lagrange multiplier

method: look at stationary pts, (0-gradient) of $J(\pi, \lambda)$

$$\text{i.e. } \nabla_{\pi} J(\pi, \lambda) = 0$$

necessary for local opt. (check "bordered Hessian") to get local min or max

$$\nabla_{\lambda} J(\pi, \lambda) = 0 \Rightarrow g(\pi) = 0$$

(see Wikipedia)

$$l(\pi) = \sum_{j=1}^k n_j \log \pi_j$$

(strictly concave fct. in π_j)

$$\frac{\partial J}{\partial \pi_j} = 0$$

$$\frac{n_j}{\pi_j} - \lambda = 0 \Rightarrow \pi_j^* = \frac{n_j}{\lambda} \text{ scaling constant}$$

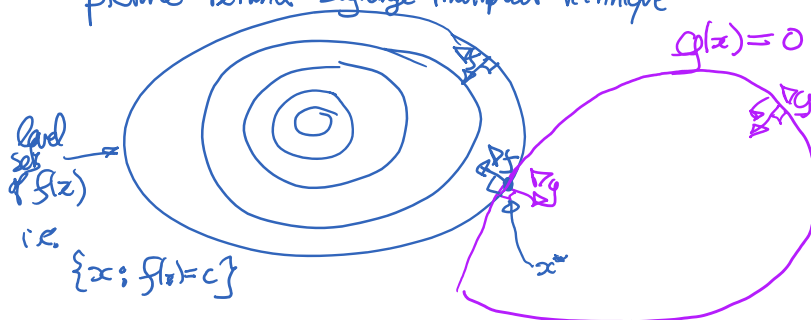
$$\text{want } g(\pi^*) = 0 \text{ i.e. } \sum_{j=1}^k \pi_j^* = 1 \Rightarrow \sum_{j=1}^k \frac{n_j}{\lambda} = 1$$

$$\Rightarrow \lambda^* = \sum_{j=1}^k n_j = n$$

$$\text{notice: } \pi_j^* = \frac{n_j}{n} \in [0, 1]$$

$$\pi_j^* = \frac{n_j}{n} \text{ MLE of multinomial}$$

Picture behind Lagrange multiplier technique



$$\nabla f(x^*) = \lambda \nabla g(x^*)$$

i.e. ∇f & ∇g are parallel

Statistical decision theory

A) Bias-variance decomposition for squared loss

estimator: function from data (observation) to parameters

$$\text{MLE} : \hat{\theta}_{\text{MLE}}(z) = \underset{\theta \in \Theta}{\text{argmax}} p(z|\theta)$$

$$\text{MAP} : \hat{\theta}_{\text{MAP}}(z) = \underset{\theta \in \Theta}{\text{argmax}} p(\theta|z) = \underset{\theta \in \Theta}{\text{argmax}} p(z|\theta) \cdot p(\theta)$$

weibched prior kann

* how do we evaluate those estimators?

estimator $\delta : \Omega \rightarrow \Theta$

$$\hat{\theta} = \delta(x)$$

most standard tool: frequentist risk of an estimator

$$R(\theta, \delta) \triangleq \mathbb{E}_x [L(\theta, \delta(x))]$$

(statistical) loss fct.

Average over possible data

squared loss: $L(\theta, \hat{\theta}) \triangleq \|\theta - \hat{\theta}\|_2^2$ $\hat{\theta} = \delta(x)$

$$\mathbb{E}_x [\|\theta - \hat{\theta}\|_2^2] = \mathbb{E} [\|\theta - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \hat{\theta}\|_2^2]$$

$$= \mathbb{E} [\|\theta - \mathbb{E}[\hat{\theta}]\|_2^2] + \mathbb{E} [\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_2^2]$$

$$+ 2 \mathbb{E} [\langle \theta - \mathbb{E}[\hat{\theta}], \mathbb{E}[\hat{\theta}] - \hat{\theta} \rangle]$$

constant

$$2 \langle \theta - \mathbb{E}[\hat{\theta}], \mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}] \rangle$$

$$R(\theta, \delta) = \mathbb{E}_x [\|\theta - \hat{\theta}\|_2^2] = \underbrace{\|\theta - \mathbb{E}[\hat{\theta}]\|_2^2}_{\text{bias}^2} + \underbrace{\mathbb{E} [\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_2^2]}_{\text{variance}}$$

bias²

variance

$$\text{bias} \triangleq \|\theta - \mathbb{E}[\hat{\theta}]\|$$

$$\boxed{\text{risk for squared loss} = \text{bias}^2 + \text{variance}}$$

bias-variance decomposition "tradeoff"

* consistency: informally "do right thing as $n \rightarrow \infty$ " where n is training set size
 $X \rightsquigarrow (X_i)_{i=1}^n$

$$\hat{\theta}_n \quad \hat{\theta}(\text{data of size } n)$$

assignment: if $\text{bias}(\hat{\theta}_n) \xrightarrow{n \rightarrow \infty} 0$

and $\text{variance}(\hat{\theta}_n) \rightarrow 0$

$$\Rightarrow R(\theta, \hat{\theta}_n) \xrightarrow{n \rightarrow \infty} 0 \Rightarrow \hat{\theta}_n \text{ is consistent}$$

$$(\hat{\theta}_n \xrightarrow{P} \theta)$$

variance $(\epsilon_n) \rightarrow 0$

$$(\hat{\theta}_n \xrightarrow{P} \theta)$$