

today :

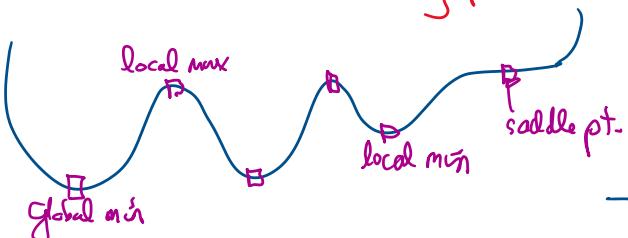
- MLE cf. ex
- statistical decision theory

optimization comments about MLE

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

$$\textcircled{S} \quad \nabla f(\theta^*) = 0$$

"stationary pts."



(f is differentiable)

is a necessary cond.:

for θ^* to be a local min
when θ^* is in interior of Θ

→ also check that $\text{Hessian}(f)(\theta^*) > 0$
for a local min

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

$$(f''(\theta^*) > 0)$$

\textcircled{S} only local results in general

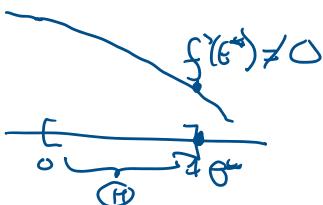
- but if $\text{Hessian}(f(\theta)) > 0 \quad \forall \theta \in \Theta$, fct. is said to be "convex"

then $\nabla f(\theta^*) = 0 \Rightarrow \theta^*$ is a global min

- otherwise, for smooth fct., looking at all zero gradient pts. and boundary pts
give you enough information to find global min $\textcircled{S} \textcircled{S}$

\textcircled{S} be careful with boundary case

i.e. $\theta^* \in \text{boundary}(\Theta)$ e.g.



other example:



example where MLE does not exist

\checkmark some notes about MLE

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

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

- It is not rec. unique [e.g. mixture models] 

- is not "admissible" in general [see later.]
 \exists strictly "better" estimators

Example II: multinomial distribution

Suppose X_i is a discrete R.V. on k choices "multinomial"

(we could choose $\Omega_{X_i} = \{1, 2, \dots, k\}$)

but instead, convenient to encode the k possibilities using unit basis in \mathbb{R}^k

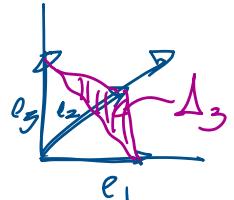
i.e. $\Omega_{X_i} = \{e_1, \dots, e_k\}$ where $e_j \in \mathbb{R}^k$ "one hot encoding"

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

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

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

probability simplex on k choices



We will write $X_i \sim \text{Mult}(\pi)$ $\underset{\text{parameters}}{P} = \text{Mult}(1, \pi)$

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

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

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

pmf for X :

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

↳ multinomial coeff

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

$$(x)_j \triangleq n_j$$

[vs. x_i of X_i]

$$\binom{n}{n_1, \dots, n_k} \triangleq \frac{n!}{n_1! n_2! \dots n_k!}$$

multinomial MLE

$$\text{log-likelihood } l(\pi) = \log p(x|\pi) = \log \binom{n}{n_1, \dots, n_k} + \sum_{j=1}^k n_j \log \pi_j$$

constant during MLE
→ ignore MLE

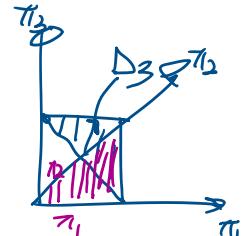
$$\text{MLE: } \hat{\pi}_{\text{MLE}}(x) = \underset{\pi \in \Delta^K}{\operatorname{arg\,max}} l(\pi)$$

s.t. $\pi \in \Delta_K$ } Constraint

two options:

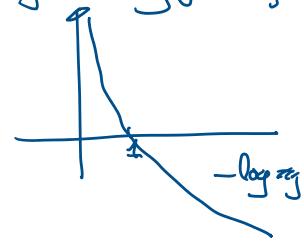
a) reparameterize problem so that Δ is full dimensional

$$\begin{aligned} \pi_K &\triangleq 1 - \sum_{j=1}^{K-1} \pi_j \\ \Rightarrow \pi_1, \dots, \pi_{K-1} &\in [0,1] \text{ with constraint } \sum_{j=1}^{K-1} \pi_j \leq 1 \end{aligned}$$



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

can try unconstrained optimization on π_1, \dots, π_{K-1}
of $l(\pi_1, \dots, \pi_{K-1})$



hoping sol'n is in the interior of constraint set
(and it usually will for log-type problems)

b) use Lagrange multiplier approach to handle equality constraint on Δ_K

[and still ignoring $\pi_j \in [0,1]$]

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

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

Lagrange
multiplier

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

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

necessary cond. for local opt.

$$\begin{aligned} \nabla_\lambda J(\pi, \lambda) = 0 \\ \Leftrightarrow g(\pi) = 0 \end{aligned}$$

(check "bordered Hessian" to
get local min or max)

$$l(\pi) = \sum_j n_j \log \pi_j \quad \frac{\partial l}{\partial \pi_j} = 0 \quad \frac{n_j}{\pi_j} \xrightarrow{\text{want}} \lambda = 0 \Rightarrow \pi_j^* = \frac{n_j}{n} \quad \boxed{\pi_j^* = \frac{n_j}{n}} \quad \begin{matrix} \text{scaling} \\ \text{constant} \end{matrix}$$

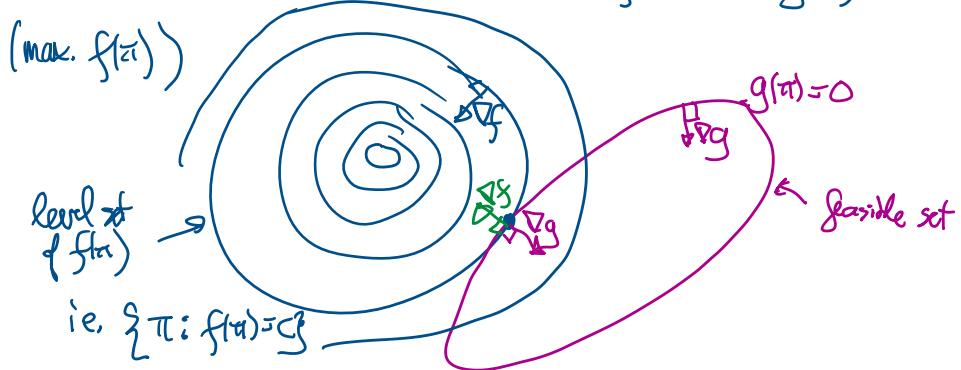
(strictly concave
fn. in π_j)

want $g(\pi^*) = 0$ i.e. $\sum_j \pi_j^* = 1$
 $\sum_j \frac{n_j}{n} = 1$
 $\Rightarrow \lambda^* = \sum_j n_j = n$

notice $\pi_j^* = \frac{n_j}{n} \in [0, 1]$ $\pi_j^* = \frac{n_j}{n}$ $n \in \mathbb{N}$
for multinomial

$$\nabla_{\pi} L(\pi, \lambda) = 0 \Rightarrow \nabla_{\pi} f(\pi) + \lambda \nabla g(\pi) = 0$$

$$\Rightarrow \nabla_{\pi} f(\pi) = -\lambda \nabla g(\pi)$$



Statistical decision theory

A) Bias-variance decomposition for squared loss

estimator: fn. from data (observation) to parameter

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

likelihood prior

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

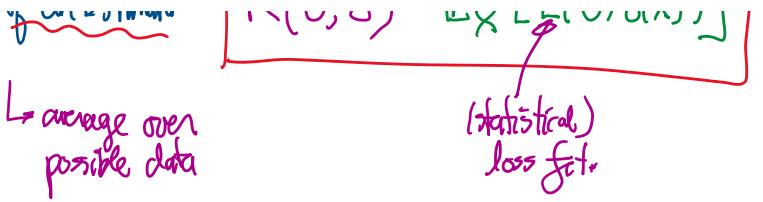
$\log p(x|\theta) + \log p(\theta)$

* how do we evaluate these estimators? estimator $\delta: \Omega_x \rightarrow \Theta$

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

most standard tool: frequentist risk of an estimator

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



$$\text{Squared Loss : } L(\theta, \hat{\theta}) \triangleq \|\theta - \hat{\theta}\|_2^2 \quad \hat{\theta} = \delta(X)$$

$$\begin{aligned} \mathbb{E}_X [\|\theta - \hat{\theta}\|_2^2] &= \mathbb{E} \left[\|\theta - \mathbb{E}\hat{\theta} + \mathbb{E}\hat{\theta} - \hat{\theta}\|_2^2 \right] && \|\alpha\|^2 = \langle \alpha, \alpha \rangle \\ &= \mathbb{E} [\|\theta - \mathbb{E}\hat{\theta}\|^2] + \mathbb{E} [\|\hat{\theta} - \mathbb{E}\hat{\theta}\|^2] \\ &\quad + 2 \mathbb{E} [\langle \theta - \mathbb{E}\hat{\theta}, \mathbb{E}\hat{\theta} - \hat{\theta} \rangle] \\ &= 2 \langle \theta - \mathbb{E}\hat{\theta}, \mathbb{E}(\mathbb{E}\hat{\theta} - \hat{\theta}) \rangle \end{aligned}$$

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

$(\text{freq. risk for squared loss}) = \text{bias}^2 + \text{variance}$

bias-variance "tradeoff"

* consistency: informally "do right thing as $n \rightarrow \infty$ " where n is training set size

$$\begin{aligned} X &\rightarrow (x_i)_{i=1}^n \\ \hat{\theta}_n &(\text{data of size } n) \end{aligned}$$

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

$$\begin{aligned} \text{and} \\ \text{variance}(\hat{\theta}_n) &\rightarrow 0 \quad (\hat{\theta}_n \xrightarrow{P} \theta) \end{aligned}$$