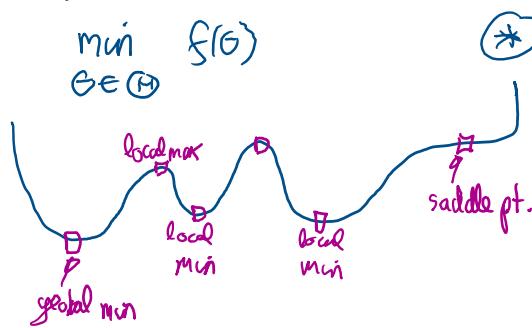


today:

- MLE cf'ed
- Statistical decision theory

optimization comments about MLE



$$\text{(*) } Df(\theta^*) = 0$$

θ^* is a stationary point

If f is differentiable on G ,
is a necessary cond.
for θ^* to be a local min
when θ^* is in the interior of G .

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

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

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

(*) only local results in general

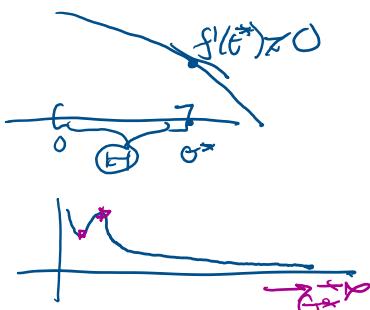
• but if $\text{Hessian}(f(\theta)) \geq 0 \quad \forall \theta \in G$, fct. is said to be "convex"

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

• otherwise, for smooth fct., looking at zero gradient pts & boundary points
give you enough info to find global min.

(*) be careful with boundary cases
i.e. $\theta^* \in \partial G$

other example:



(*) some notes about MLE

• does not always exist $[\theta^* \in \text{bd}(G) \text{ but } G \text{ is open}]$ or when " $\theta^* = \text{ab}$ "

• it is not necessarily unique [i.e. Mixture models]

• is not "admissible" in general [see later]

exists strictly "better" estimator



example II: multinomial distribution

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

(one could choose $\mathcal{Q}_{X_i} = \{1, 2, 3, \dots, k\}$)

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

i.e. $\mathcal{Q}_{X_i} = \{e_1, e_2, \dots, e_k\}$ where $e_j \in \mathbb{R}^k$ "one hot encoding"

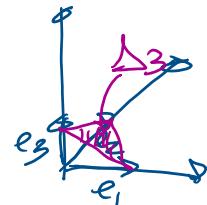
$$e_j = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix} \text{ } j^{\text{th}} \text{ coordinate}$$

parameter for discrete RV $\pi \in \mathbb{R}^k$ ($\mathcal{A} = \Delta_k$)

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

probability simplex on k choices

we will write $X_i \sim \text{Mult}(\pi)$ parametrized by π = Mult(1, π)



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

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

"multinomial distribution"

$$x \rightarrow \in \mathbb{R}^k$$

$$X \in \mathbb{N}^k \quad \mathcal{Q}_X = \left\{ (n_1, n_2, \dots, n_k) : n_j \in \mathbb{N}; \sum_{j=1}^k n_j = n \right\}$$

prnt for X :

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

$$= \frac{p(x|\pi)}{\binom{n}{(n_1, n_2, \dots, n_k)} \prod_{j=1}^k \pi_j^{n_j}}$$

multinomial coeff

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

5h15 Multinomial MLE

$$\text{Log-Likelihood} \quad l(\pi) = \log p(x|\pi) = \log \binom{n}{n_1, n_2, \dots, n_k} + \sum_{j=1}^k n_j \log \pi_j$$

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

constant \rightarrow ignore

$$\text{MLE} \quad \hat{\pi}_{\text{MLE}}(x) = \underset{\substack{\pi \in \mathbb{R}^k \\ \text{st } \pi \in \Delta_k}}{\text{argmax}} l(\pi)$$

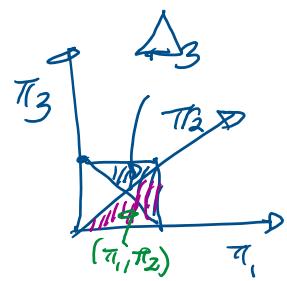
constraint

two options:

a) reparameterize problem so that \mathcal{A} is full dimensional

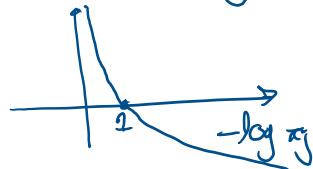
$$\pi_k = 1 - \sum_{j=1}^{k-1} \pi_j$$

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



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

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



can try unconstrained opt. on π_1, \dots, π_{k-1}

$$\text{of } l(\pi_1, \dots, \pi_{k-1}, 1 - \sum_{j=1}^{k-1} \pi_j)$$

hanging soln is in interior of constraint set
(and it is usually case for log-type problem)

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

$$\max f(\pi)$$

$$\text{st. } g(\pi) = 0$$

$$[g(\pi) \triangleq 1 - \sum_{j=1}^k \pi_j]$$

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

Lagrange multiplier

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

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

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

necessary cond. for local opt.

(check "banded Hessian" to get local or max)

(see [Wikipedia](#))

$$l(\pi) = \sum_j n_j \log \pi_j$$

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

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

scaling constant

(strictly concave sf. in π_j)

$$\text{want } g(\pi^*) = 0 \text{ i.e. } \sum_j \pi_j^* = 1$$

$$\text{i.e. } \sum_j n_j = 1$$

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

note: $\pi_j^* \in [0, 1]$

$$\boxed{\pi_j^* = \frac{n_j}{\eta}}$$

MLE for multinomial

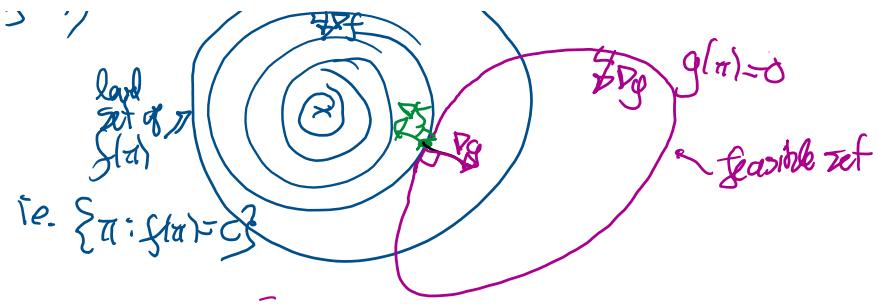
$$\nabla_\pi L(\pi, \lambda) = 0 \Rightarrow \nabla_\pi f(\pi) + \lambda \nabla_\pi g(\pi) = 0$$

(max $J(\pi)$)



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

$$\nabla g(\pi) = 0$$



statistical decision theory

A) Bias-variance decomposition for squared loss

estimator: fd. from data (observation) to parameter

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

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

likelihood prior
marginal
 $\log p(x|\theta) + \log p(\theta)$
regularization

how do we evaluate these estimators? estimator $\delta: \mathcal{X} \rightarrow \Theta$
 $\hat{\theta} = \delta(x)$

most standard tool: frequentist risk of an estimator

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

average over possible data / training set "statistical loss fn."

$$\text{squared loss, } L(\theta, \delta) \triangleq \|\theta - \delta\|_2^2$$

$$\begin{aligned} \mathbb{E}_x [\|\theta - \hat{\delta}\|_2^2] &= \mathbb{E} [\|\theta - \mathbb{E}\hat{\delta} + \mathbb{E}\hat{\delta} - \hat{\delta}\|_2^2] \quad \hat{\theta} = \delta(x) \\ &= \mathbb{E} [\|\theta - \mathbb{E}\hat{\delta}\|_2^2] + \mathbb{E} [\|\hat{\delta} - \mathbb{E}\hat{\delta}\|_2^2] \quad \|a+b\|^2 = \|a\|^2 + \|b\|^2 + 2\langle a, b \rangle \\ &\quad + 2\mathbb{E} [\langle \theta - \mathbb{E}\hat{\delta}, \mathbb{E}\hat{\delta} - \hat{\delta} \rangle] \\ &\quad \underbrace{\quad}_{\text{constant}} \\ &= 2 \langle \theta - \mathbb{E}\hat{\delta}, \mathbb{E}(\hat{\delta} - \hat{\delta}) \rangle \end{aligned}$$

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

$$\text{bias} \triangleq \|\theta - \mathbb{E}(\hat{\delta})\|_2$$

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

"bias-variance."

$$\text{bias} = \| \theta - E(\theta') \|_2$$

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

"bias-variance tradeoff"

④ consistency: informally "do right thing as $n \rightarrow \infty$ " where n is the training set size

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

$$X \sim (x_i)_{i=1}^n$$

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

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

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

i.e.
 $\hat{\theta}_n$ is consistent