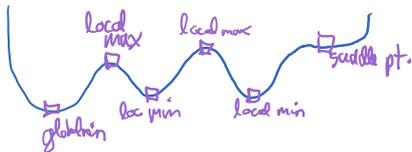


today: • MLE et'ed
• statistical decision theory

optimization comments about MLE

$\min_{\Theta \in \mathcal{G}} f(\theta)$ $\nabla f(\theta^*) = 0$
 "stationary pts."

(if f is differentiable) is a necessary cond. on Θ for θ^* being a local min when θ^* is in the interior of Θ



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

$H \succ 0 \Leftrightarrow u^T H u > 0 \quad \forall u \neq 0 \in \mathbb{R}^d$
 $(S''(f)(\theta^*) \succ 0)$

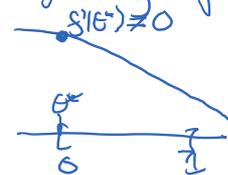
* only local results in general

• but if $\text{Hessian}(f(\theta)) \succ 0 \quad \forall \theta \in \Theta$, $f(\cdot)$ is said to be "convex"
 then $\nabla f(\theta^*) = 0 \Rightarrow \theta^*$ is global min

• otherwise, for smooth $f(\cdot)$, looking at zero gradient pts and boundary pts give you enough information to find global min

⊗ be careful with boundary cases

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



other example:



↑ example where MLE doesn't exist

* 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 nec. unique [e.g. mixture models]



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

example II: multinomial distribution

suppose X_p is discrete R.V. on k choices "multinoulli"

(we could choose $S_{k \times 1} = \{e_1, \dots, e_k\}$)

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

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

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

parameter for discrete RV: $\pi \in \Delta_k$ ($\Theta = \Delta_k$)

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

probability simplex on k choices

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



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

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

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

pmf for X : $p(x|\pi) = \binom{n}{(x_1, \dots, x_k)} \prod_{j=1}^k \pi_j^{x_j}$ $x = (n_1, \dots, n_k)$

multinomial coeff.

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

15h58

multinomial MLE

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

MLE: $\hat{\pi}_{MLE}(x) = \underset{\pi \in \mathbb{R}^k}{\text{argmax}} l(\pi)$
 s.t. $\pi \in \Delta_k$ } constraint



two options:

a) reparameterize problem so that Θ is full dimensional

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

$\{ (\pi_1, \pi_2) : \pi_j \in [0, 1], \sum \pi_j \leq 1 \}$



here $\log \pi_j$ acts as a barrier for. away from $\pi_j = 0$

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



hoping sol'n is in the interior of constraint set (and it usually will)

b) use Lagrange multiplier approach to handle equality constraints on Δ_k

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

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

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

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.

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

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

$$\begin{aligned} l(\pi) = \sum_j n_j \log \pi_j \\ \text{(strictly concave)} \\ \text{set. in } \pi_j \end{aligned}$$

$$\frac{\partial J}{\partial \pi_j} \stackrel{\text{want}}{=} 0$$

$$\frac{\partial J}{\partial \pi_j} - \lambda \stackrel{\text{want}}{=} 0 \Rightarrow \pi_j^* = \frac{n_j}{\lambda}$$

scaling constant

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

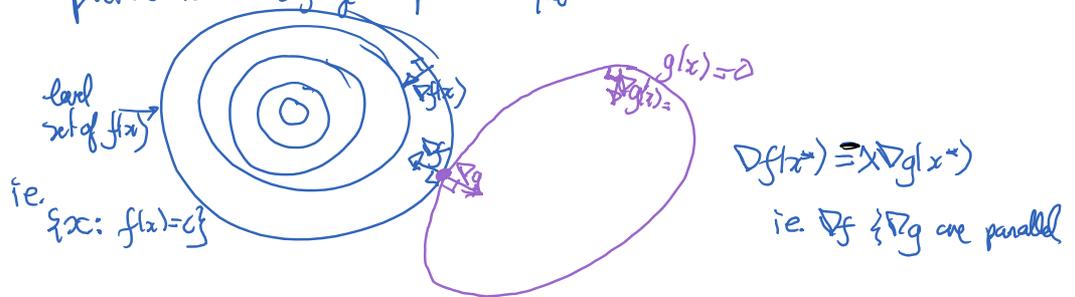
$$\sum_j \frac{n_j}{\lambda} = 1$$

$$\Rightarrow \lambda^* = \sum_j n_j = n$$

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

$$\boxed{\pi_j^* = \frac{n_j}{n}} \quad \text{MLE for multinomial}$$

picture behind Lagrange multiplier technique



Statistical decision theory

A) Bias-variance decomposition for squared loss

estimator: fct. from data (observation) to parameter

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

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

* how do we evaluate these estimators?

$$\text{estimator } \mathcal{S}: \Omega \rightarrow \Theta$$

$$\hat{\theta} = \mathcal{S}(X)$$

most standard tool: frequentist risk of an estimator

$$\boxed{R(\theta, \mathcal{S}) \triangleq \mathbb{E}_X [L(\theta, \mathcal{S}(X))]}$$

average over possible data

(statistical) loss fct.

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

$$\begin{aligned} \mathbb{E}_X [\|\theta - \hat{\theta}\|_2^2] &= \mathbb{E} [\|\theta - \underbrace{\mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta} - \mathbb{E}[\hat{\theta}]]}_{\hat{\theta}}\|_2^2] \\ &= \mathbb{E} [\|\theta - \mathbb{E}[\hat{\theta}]\|_2^2] + \mathbb{E} [\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_2^2] \\ &\quad + 2 \mathbb{E} [\langle \theta - \mathbb{E}[\hat{\theta}], \mathbb{E}[\hat{\theta}] - \hat{\theta} \rangle] \\ &\quad \quad \quad \downarrow \text{by linearity} \\ &\quad \quad \quad \downarrow \text{constant} \\ &\quad \quad \quad 2 \langle \theta - \mathbb{E}[\hat{\theta}], \mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}] \rangle \\ R(\theta, \delta) = \mathbb{E}_X [\|\theta - \hat{\theta}\|_2^2] &= \underbrace{\|\theta - \mathbb{E}[\hat{\theta}]\|_2^2}_{\text{bias}^2} + \underbrace{\mathbb{E} [\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|_2^2]}_{\text{variance}} \end{aligned}$$

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

(freq.) risk for squared loss = bias² + variance

bias-variance decomposition "tradeoff"

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

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

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

Statistical decision theory - formal setup

- a random observation $D \sim P$ ✓ unknown distribution which models the world / phenomenon (often P_θ)
- action space \mathcal{A}
- loss (statistical) $L(P, a) =$ statistical loss of doing action $a \in \mathcal{A}$ when the world is P } describe the goal / task
- ↳ often write $L(\theta, a)$ if we have a parametric model of world i.e. P is $P_\theta / \text{pref } P_\theta$ for some $\theta \in \Theta$
- $\delta: \mathcal{D} \rightarrow \mathcal{A}$ "decision rule"
 \downarrow
 Ω_D

examples: a) parameter estimation:

$\mathcal{A} = \Theta$ for parametric family P_θ

δ is a parameter estimator from data

typical loss $L(\theta, a) = \|\theta - a\|_2^2$

$\mathcal{D} \stackrel{\text{typically}}{=} (X_1, \dots, X_n)$
 [usually $X_i \stackrel{\text{i.i.d.}}{\sim} P_\theta$?
 unknown]

"squared loss"

but other losses are $k_L(p_G \| p_A)$

b) $\mathcal{A} = \{0, 1\}$; this is hypothesis testing

δ describes a statistical test

loss \rightarrow usually 0-1 loss $L(\theta, a) = \mathbb{1}\{\theta \neq a\}$

c) prediction in ML (function estimation) learn a prediction fct. in supervised learning

here $\mathcal{D} = (x_i, y_i)_{i=1}^n$ $x_i \in \mathcal{X}$ (input space) $y_i \in \mathcal{Y}$ (output space) $\mathcal{Y} = \{0, 1\} \rightarrow$ classifiers $\mathcal{Y} = \mathbb{R} \rightarrow$ regression

p_G gives joint (x, y)

$\mathcal{D} \sim \mathcal{P}$ where $\mathcal{P} = p_G \otimes p_G \otimes \dots \otimes p_G$

$\mathcal{F} = \mathcal{Y}^{\mathcal{X}}$ (set of fct.'s from \mathcal{X} to \mathcal{Y}) n times

"generalization error"
"classification error"

in ML

$L(p_G, f) \triangleq \mathbb{E}_{p_G} [l(y, f(x))]$
 $(x, y) \sim p_G$ prediction loss

in ML is often called the "risk"

eg classification $l(y, f(x)) = \mathbb{1}\{y \neq f(x)\}$

0-1 error

Simon calls it "Vapnik risk" to distinguish it from frequentist risk

frequentist risk $\mathbb{E}_{\mathcal{P}} [L(p_G, \delta(\mathcal{D}))]$

* decision rule

$\delta = \delta(\mathcal{D})$
prediction fct.
classifier
etc.

"learning algorithm"