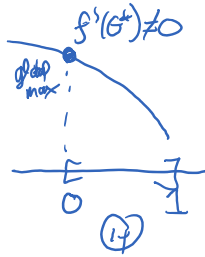


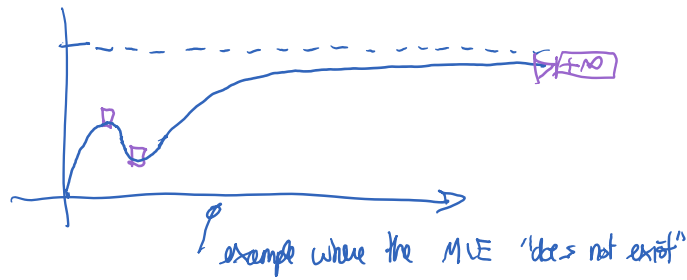
today: • MLE of θ
• decision theory

⊗ be careful with boundary cases

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



other example



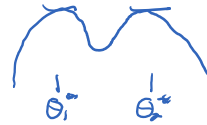
⊗ Some notes about MLE

• does not always exist [$\theta^* \in \text{bd}(\Theta)$ but Θ is open] or when " $\theta^* < \theta_0$ "

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

• is not rec. unique [i.e. multiple maxes]

e.g. mixture models



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

∃ strictly "better" estimators

example II: multinomial distribution

suppose X_i is a discrete R.V. on K choices "multinoulli"

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

but instead, convenient to encode with unit basis vectors in \mathbb{R}^K "one-hot encoding"

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

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

$$\begin{pmatrix} X \\ \vdots \\ X \end{pmatrix} = \sum_{i=1}^n X_i$$

multinomial

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

$$\Delta_K \triangleq \left\{ \pi \in \mathbb{R}^K : \pi_j \geq 0 \text{ } \forall j \text{ } \sum_{j=1}^K \pi_j = 1 \right\}$$

probability simplex on K choices

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



$\sum_{i=1}^K x_i = 0$

parameter

↖

pmf: $p(X_i = e_j) = \pi_j = \prod_{l=1}^k \pi_l^{X_{i,l}}$
 "X_i = j"

⊗ Consider X_i iid Mult(π)

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

"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:

$p(x|\pi) = \text{coeff.} \cdot \prod_{j=1}^k \pi_j^{x_j}$

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

$\binom{n}{x_1, x_2, \dots, x_k} = \binom{n}{n_1, n_2, \dots, n_k} \triangleq \frac{n!}{n_1! n_2! \dots n_k!}$

"multinomial coeff."

multinomial MLE:

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 \rightarrow ignore for MLE

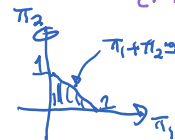
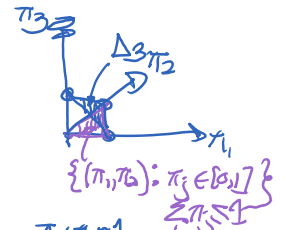
MLE: $\hat{\pi}_{ML}(x) = \underset{\pi \in \mathbb{R}^k}{\text{argmax}} l(\pi)$
 s.t. $\pi \in \Delta_k$ } constraint

two options:

a) reparameterize problem so that Θ is full dimensional

$\pi_1, \dots, \pi_{k-1} \in [0, 1]$

$\pi_k \triangleq 1 - \sum_{j=1}^{k-1} \pi_j$ with constraint $\sum_{j=1}^{k-1} \pi_j \leq 1$



here magic is that $\log \pi_j$ act as a barrier $\text{Exp. away from } \pi_j = 0$

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

hope

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



[note: here $l(x)$ is a concave fct.]

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

max $f(\pi)$
 s.t. $g(\pi) = 0$

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

$\sum_j \pi_j = 1 \Rightarrow g(\pi) \triangleq 1 - \sum_j \pi_j$

Lagrange multiplier

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

necessary for local opt.

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

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

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

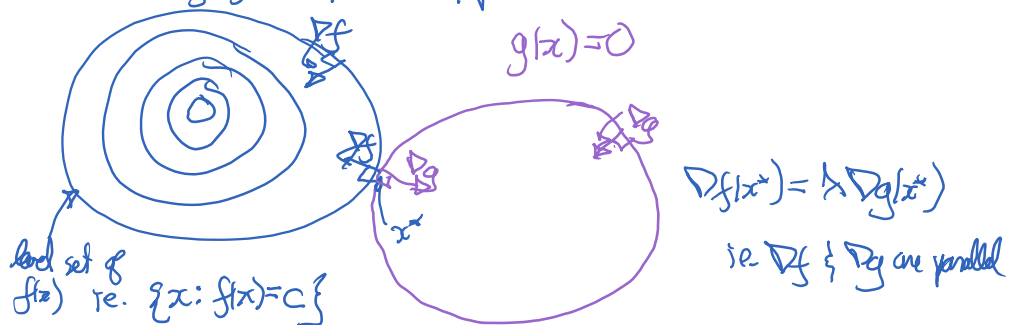
(strictly concave fct. of π)

want $\frac{\partial J}{\partial \pi_j} = 0 \Rightarrow \frac{n_j}{\pi_j} - \lambda = 0 \Rightarrow \pi_j^* = \frac{n_j}{\lambda}$ (see Wikipedia)
scaling constant

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

notice: $\pi_j^* = \frac{n_j}{N} \in [0, 1]$ $\pi_j^* = \frac{n_j}{N}$ MLE for multinomial $\Rightarrow \lambda^* = \sum_j n_j = N$

picture behind Lagrange multiplier technique



16h10

statistical decision theory

A) Bias-variance decomposition for squared loss

estimator: function from data (observation) to parameters

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

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

⊛ how do we evaluate these estimators?

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

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

most standard tool: frequentist risk of estimator

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

average over possible datasets

(statistical) loss fct.

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

$$\mathbb{E}_X [\|\theta - \delta\|_2^2] = \mathbb{E} [\underbrace{\|\theta - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \hat{\theta}\|_2^2}_{\text{variance}}]$$

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

constant

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

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

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

risk for squared loss = bias² + variance

bias-variance decomposition
"tradeoff"

* consistency: informally "do right thing as $n \rightarrow \infty$ "

where n is training set size
 $X \rightarrow (X_i)_{i=1}^n$

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

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

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

formal setup for statistical decision theory

- a random observation $D \sim \mathcal{P}$ unknown distribution which models the world / phenomena (often related to θ)
- action space \mathcal{A}
- loss $L(\mathcal{P}, a) =$ statistical loss of doing action $a \in \mathcal{A}$ when the world is \mathcal{P} describe the goal / task
- $\delta: \mathcal{D} \rightarrow \mathcal{A}$ "decision rule" often write $L(\theta, a)$ if we have a parametric model of world i.e. \mathcal{P} has pmf/pdf p_θ for some $\theta \in \Theta$

examples: a) parameter estimation

$\mathcal{A} = \Theta$ for a parametric family \mathcal{P}_Θ

δ is a parameter estimator from data

typically $D = (X_1, \dots, X_n)$

[usually, $X_i \stackrel{iid}{\sim} p_\theta$]
unknown

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

but other losses are used e.g. $KL(p_\theta || p_a)$

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

\mathcal{S} describes a statistical test

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

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

have $D = (x_i, y_i)_{i=1}^n$ $x_i \in X$ (input space) $y_i \in \mathcal{Y}$ (output space)

$\mathcal{Y} = \{0, 1\}$ \rightarrow ^{binary} classification
 $\mathcal{Y} = \mathbb{R}$ \rightarrow regression

\mathcal{P}_S gives joint (X, Y)

$D \sim P$ where $P = \underbrace{\mathcal{P}_S \otimes \mathcal{P}_S \dots \otimes \mathcal{P}_S}_{n \text{ times}} = \mathcal{P}_S^{\otimes n}$

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

in machine learning

$$L(\mathcal{P}_S, f) \triangleq \mathbb{E}_{(X, Y) \sim \mathcal{P}_S} [l(Y, f(X))]$$

"generalization error"
 "classification error"

prediction loss

e.g. classification

$$l(Y, f(X)) = \mathbb{1}\{Y \neq f(X)\}$$

0-1 error

* decision rule $\hat{f} = \mathcal{S}(D)$

"learning algorithm"

prediction fct.
 classification
 etc...

in ML is often called "risk"

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

frequentist risk $\mathbb{E}_D [L(\mathcal{P}_S, \mathcal{S}(D))]$