

Lecture 4 - scribbles

Tuesday, September 13, 2016
14:27

today

- MLE
- statistical decision theory
bias/variance decomposition

Maximum likelihood principle

parametric family $p(x; \theta)$ for $\theta \in \Theta$

how to estimate θ ?

choose $\hat{\theta}_{ML}(x)$ which maximizes $p(x; \theta)$

$$\hat{\theta} \in \underset{\theta \in \Theta}{\text{argmax}} p(x; \theta) \triangleq \underset{\theta \in \Theta}{\text{argmax}} \underbrace{f(\theta)}_{\text{likelihood function (of } \theta)}$$

example: n coin flips

we had $X = \sum_{i=1}^n x_i$ $X \sim \text{Bin}(n, \theta)$
Bernoulli(θ)

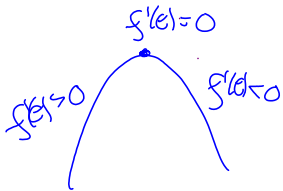
$$p(x; \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \quad [x \in 0:n]$$

trick: $\log(\cdot)$ is strictly increasing fct.

$$\text{i.e. } a < b \Rightarrow \log(a) < \log(b)$$

$$\Rightarrow \underset{\theta \in \Theta}{\text{argmax}} \log(p(x; \theta)) = \underset{\theta \in \Theta}{\text{argmax}} p(x; \theta)$$

log-likelihood: $\log p(x; \theta) = \text{const.} + x \log \theta + (n-x) \log(1-\theta) = \ell(\theta)$
constant (look at context)



look for $\frac{\partial \ell(\theta)}{\partial \theta} = 0$

$$\text{want } \frac{\partial \ell(\theta)}{\partial \theta} = \frac{x}{\theta} - \frac{(n-x)}{1-\theta} = 0$$

$\nabla f(\theta) = 0$
necessary condition
if max is in the interior of Θ

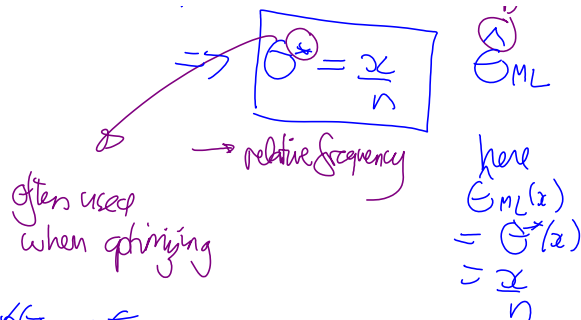
$$\Rightarrow x - \theta x - n\theta + \theta x = 0$$

estimator

$$\Rightarrow \theta = \frac{x}{n}$$

$\hat{\theta}_{ML}$

necessary condition
if max is in the interior of Θ



note: $E[\hat{\theta}_{ML}] = E[\frac{X}{n}] = \frac{E[X]}{n} = \theta$
 \downarrow
 with respect to $X \sim p(x; \theta_0)$ \rightarrow unbiased (here)

some optimization comments:

- $f'(\theta) = 0$ not sufficient for local max \rightarrow also need to check $f''(\theta) < 0$ for max

- only local result in general



but if $f''(\theta) < 0 \forall \theta \in \Theta$ function is "concave"

$\rightarrow f'(\theta) = 0$ is sufficient for global max

- careful with boundary cases $\theta^* \in \text{boundary}(\Theta)$

some notes on MLE

- does not always exist [e.g. $\hat{\theta} \in \text{bd}(\Theta)$ but Θ is open]
- is not necessarily unique
- is not admissible in general (see end of class)

Example 2:

suppose X_i is discrete R.V. on k choices

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

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

i.e. $\Omega_{X_i} = \{e_1, \dots, e_k\}$

where $e_i \in \mathbb{R}^k$ $e_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$ \leftarrow i -th position

parameter $\pi \in \Delta_k \leftarrow$ probability simplex on k choices



$$\Delta_k \triangleq \left\{ \pi \in \mathbb{R}^k : \begin{array}{l} \pi_j \geq 0 \quad \forall j \\ \sum_{j=1}^k \pi_j = 1 \end{array} \right\}$$

(*) Δ_k

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

parameter

[really $\text{Mult}(1, \pi)$]

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

then $X = \sum_{i=1}^n X_i \sim \text{Mult}(n, \pi)$ (analogy of binomial but for k choices)

$$\Omega_X = \left\{ (n_1, \dots, n_k) : \sum_k n_k = n, n_k \in \mathbb{N} \right\}$$

consider MLE for k -vector $\vec{x} = (x_1, \dots, x_n)$

$$p(\vec{x} | \pi)$$

$$p(x_i | \pi) = \text{Mult}(x_i | \pi)$$

$$= \prod_{j=1}^k \pi_j^{x_{ij}}$$

x_{ij} is j^{th} component of vector x_i

[e.g. if $x_i = e_j$, then $p(x_i | \pi) = \pi_j$] "switch statement"

$$p(\vec{x} | \pi) = \prod_{i=1}^n p(x_i | \pi) = \prod_{i=1}^n \left(\prod_{j=1}^k \pi_j^{x_{ij}} \right)$$

by indep.

$$= \prod_{j=1}^k \pi_j^{\left(\sum_{i=1}^n x_{ij} \right)}$$

note $n_j(\vec{x})$

log-likelihood: $l(\pi) = \log p(\vec{x} | \pi) = \sum_{j=1}^k n_j \log \pi_j$

we want $\max_{\pi \in \Delta_k} l(\pi)$ constraints

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

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



full 2-dimensional set

here instead demonstrate simple Lagrange multiplier approach for equality constraint

look for stationary points (∇ gradient) of

$$J(\pi, \lambda) \triangleq l(\pi) + \lambda \left(1 - \sum_{j=1}^k \pi_j \right)$$

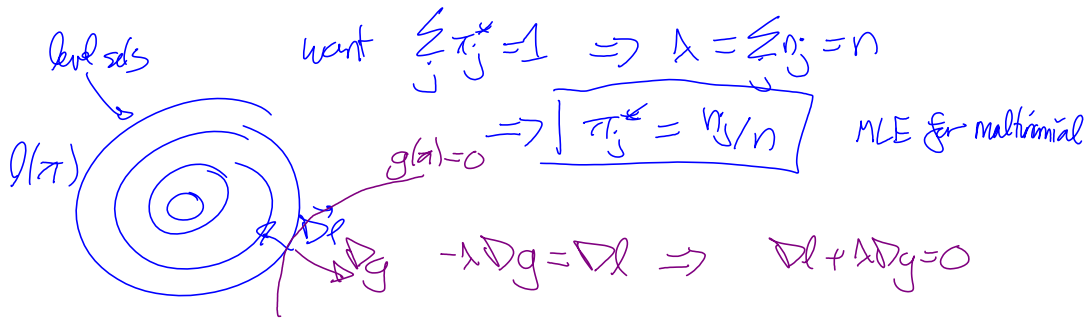
Lagrange

equality constraint $a(\pi) = 0$

want $\nabla_{\pi} J(\pi, \lambda) = 0$
 and $\nabla_{\lambda} J(\pi, \lambda) = 0 \rightarrow$ equivalent to ask $g(\pi) = 0$

$$\partial_{\pi_j} J = 0 \Rightarrow \frac{n_j}{\pi_j} - \lambda = 0 \quad \forall j$$

$$\Rightarrow \pi_j^* = \frac{n_j}{\lambda} \quad \text{Scaling constant}$$



Statistical decision theory (frequentist)

general setup:

- $D \sim P$ unknown distribution which model the "world"
random data
- A action space
- $L(P, a)$ loss of doing action a when "true world" is P } describes goal/task

[if have parametric model; often write $L(\theta, a)$ where $P = P_{\theta}$]

- $\mathcal{S}: \mathcal{D} \rightarrow A$ decision rule

Examples: a) $A = \mathbb{R}$ \mathcal{S} is then parameter estimator from data

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

b) $A = \{0, 1\}$ hypothesis test

c) say $D = ((x_i, y_i)_{i=1}^n)$ $x_i \in X, y_i \in Y$

if P_{θ} gives joint on (x_i, y_i) ; then $P = P_{\theta}^{\otimes n}$

$A = Y^X$ (set of fct's from $X \rightarrow Y$)

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

in M.L., $L(P, f) = \mathbb{E}_P [l(Y, f(X))]$

↑ prediction loss
 ↑ $(x, y) \sim P$
 "generalization error"

decision rule: $f = \mathcal{S}(D)$
 ↑ prediction fct.
 ↑ training dataset
 "learning algorithm"

How to compare procedures? (given \mathcal{S})

(frequentist) risk $R(P, \mathcal{S}) \triangleq \mathbb{E}_P [L(P, \mathcal{S}(D))]$

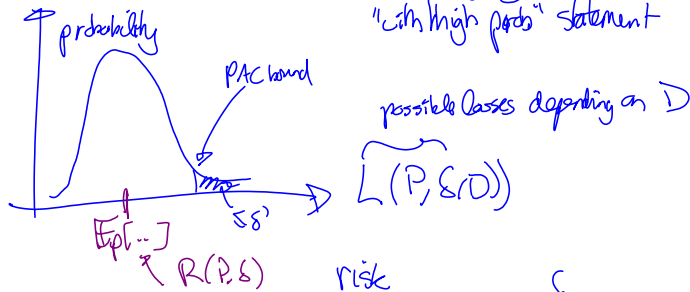
↑ $D \sim P$
 ↑ observations is random

"how well does it do in average"?

vs. PAC analysis (learning theory) → look at tail bound small

$P\{L(P, \mathcal{S}(D)) \geq \text{stuff}\} \leq \delta$
 "with high prob" statement

cartoon:



⊗ now $R(P, \mathcal{S})$

depends P/θ unknown

