

Lecture 5 - scribbles

Tuesday, September 18, 2018 14:40

today: finish multinomial MLE
 • Statistical decision theory + properties of estimator

multinomial (cited from last class)

$$X_i \sim \text{Mult}(\pi)$$

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

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

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

$$p((x_1, \dots, x_n) | \pi)$$

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

$$= \prod_{j=1}^k \pi_j^{\sum_{i=1}^n x_{ij}} = \prod_{j=1}^k \pi_j^{n_j = (x)_j}$$

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

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

← jth component of vector x_i

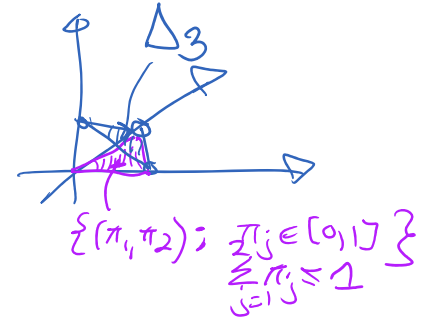
$$= \prod_{j=1}^k \pi_j^{n_j}$$

log-likelihood:
$$l(\pi) = \log p(x|\pi) = \underbrace{\log \binom{n}{(n_1, \dots, n_k)}}_{\text{const. with respect to } \pi} + \sum_{j=1}^k n_j \log \pi_j$$

(can ignore for MLE)

we want: $\max_{\pi \in \mathbb{R}^k} l(\pi)$
 s.t. $\pi \in \Delta_k$ } constraint

two options: a) would reparameterize $\pi_k \triangleq 1 - \sum_{j=1}^{k-1} \pi_j$
 with $\pi_1, \dots, \pi_{k-1} \in [0, 1]$
 with constraints: $\sum_{j=1}^{k-1} \pi_j \leq 1$



and then do unconstrained optimization on π_1, \dots, π_{k-1}
hoping solution is in interior of constraint set

b) use Lagrange multiplier approach to handle equality constraint on Δ_k necessary but not sufficient conditions

$$\begin{aligned} \max_{\pi} f(\pi) \\ \text{s.t. } g(\pi) = 0 \end{aligned} \quad \begin{aligned} \sum_{j=1}^k \pi_j = 1 \\ \downarrow \\ 1 - \sum_{j=1}^{k-1} \pi_j = 0 \\ \triangleq g(\pi) \end{aligned}$$

look for stationary points (0-gradient) of

$$J(\pi, \lambda) \triangleq f(\pi) + \lambda \underbrace{g(\pi)}_{1 - \sum_{j=1}^k \pi_j}$$

↑
Lagrange multiplier

ie. want $\nabla_{\pi} J(\pi, \lambda) = 0$
 and

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

$0(\pi) = \sum_{j=1}^k \pi_j, 0 \text{ on } \pi_i \quad \partial J = 0 \Rightarrow \pi_i - \lambda = 0$

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

concave in π

$$\frac{\partial^2 l}{\partial \pi_j^2} = \delta_{ij} \left(-\frac{1}{\pi_j^2} \right)$$

$$\frac{\partial J}{\partial \pi_j} = 0 \Rightarrow \frac{n_j}{\pi_j} - \lambda = 0$$

$$\Rightarrow \pi_j^* = \frac{n_j}{\lambda} \text{] scaling constant}$$

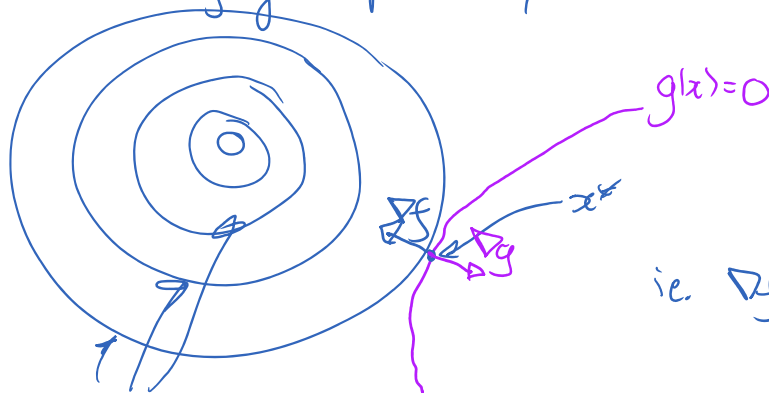
$$\text{want } \sum_j \pi_j^* = 1 \Rightarrow \frac{1}{\lambda} \sum_j n_j = 1 \Rightarrow \lambda^* = \sum_{j=1}^k n_j = n$$

$$\Rightarrow \pi_j^* = \frac{n_j}{n}$$

MLE for a multinomial

side note: $\pi_j^* = \frac{n_j}{n} \geq 0$

picture behind Lagrange multiplication technique:



level sets of $f(x)$ [ie. $\{x : f(x) = c\}$]

$$a^T x \leq b$$

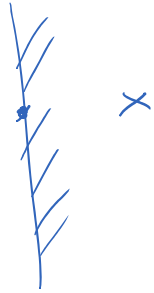
$$J(x, \lambda) = f(x) + \lambda g(x)$$

$$\nabla_x J(x, \lambda) = 0$$

$$\Rightarrow \nabla f(x^*) = -\lambda \nabla g(x^*)$$

$$\text{ie. } \nabla f(x^*) = \lambda \nabla g(x^*)$$

ie. ∇f & ∇g are parallel



Bias-variance decomposition for squared loss

estimator: function from data to parameters

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

MAP
"maximum a posteriori"

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

* analyzing property of estimator:

frequentist risk of an estimator

$$\begin{array}{l} \text{estimator } \delta: \Omega \rightarrow \Theta \quad \hat{\Theta} = \delta(x) \\ \boxed{\mathbb{E}_x [L(\Theta, \delta(x))] = R(\Theta, \delta)} \\ \text{on average over data} \\ \text{loss function} \end{array}$$

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

$$\begin{array}{c} \hat{\Theta} \\ \Theta \quad (\hat{\Theta}(x)) \end{array}$$

$$\begin{aligned}
\mathbb{E}_x [\|\theta - \hat{\theta}\|^2] &= \mathbb{E}_x [\|\underbrace{\theta - \mathbb{E}[\hat{\theta}]}_0 + \mathbb{E}[\hat{\theta}] - \hat{\theta}\|^2] \\
&= \mathbb{E} [\|\theta - \mathbb{E}[\hat{\theta}]\|^2] + \mathbb{E} [\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|^2] \\
&\quad + 2 \mathbb{E} [\underbrace{\langle \theta - \mathbb{E}[\hat{\theta}] \rangle}_{\text{constant}}, \mathbb{E}[\hat{\theta}] - \hat{\theta}] \\
&= 2 \langle \theta - \mathbb{E}[\hat{\theta}], \underbrace{\mathbb{E}[\mathbb{E}[\hat{\theta}] - \hat{\theta}]}_{\mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}]} \rangle
\end{aligned}$$

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

$$\text{bias} \triangleq \|\theta - \mathbb{E}_x [\hat{\theta}(x)]\|_2$$

risk for squared loss = bias² + variance

bias-variance decomposition
"tradeoff"

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

here: training set of size n
 $\hat{\theta}_n$ is $\hat{\theta}$ (data of size n)

assignment: $\text{bias}(\hat{\theta}_n) \xrightarrow{n \rightarrow \infty} 0 \Rightarrow \hat{\theta}_n$ is consistent $[\hat{\theta}_n \xrightarrow{P} \theta]$
 $\underbrace{\hspace{1cm}}_{\text{variance}(\hat{\theta}_n) \xrightarrow{n \rightarrow \infty} 0}$

(frequentist) statistical decision theory

, unknown distribution which models the "world"

- formal setup :
- a random observation $D \sim P$ (perhaps P_0)
 - action space A
 - loss $L(P, a) = \text{loss of doing action } a \in A \text{ when "the world" is } P$ } describes the goal/task

unknown distribution which models the "world"

(if have a parametric model of world, often write $L(\theta, a)$ where θ is s.t. $P = P_\theta$)

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

examples: a) estimation of parameters:

$A = \Theta$ for a parametric family P_θ

S is then parameter estimator from data

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

[more specifically, often $\mathcal{D} = (X_1, \dots, X_n)$ where $X_i \stackrel{i.i.d.}{\sim} P_\theta$ (θ is unknown)]

$$S(\mathcal{D}) = \hat{\theta} \quad L(\theta, S(\mathcal{D})) = \|\theta - \hat{\theta}\|_2^2$$

- b) $A = \{0, 1\}$; this is hypothesis testing
 here S describes a statistical test

c) prediction in machine learning: learn a prediction function in supervised learning

here $D = ((x_i, y_i))_{i=1}^n$

$x_i \in X$ (input space)
 $y_i \in Y$ (output space)

$Y = \{0, 1\} \rightarrow$ classification
 $Y = \mathbb{R} \rightarrow$ regression
 ...

$P \in$ gives joint on (X, Y)

$D \sim P$ where $P = P \otimes \dots \otimes P$ (n times)

$A = Y^X$ (set of functions from X to Y)

prediction loss

eg. classification: $l(y, f(x)) = 1_{\{y \neq f(x)\}}$
 0 = error

in machine learning

$L(P, f) \triangleq \mathbb{E}_{P \otimes P} [l(Y, f(X))]$ "generalization error"

* decision rule $\hat{f} = \mathcal{S}(D)$
 (training data D)
 "learning algorithm"
 prediction fct. / classifier / etc.

$(X, Y) \sim P \otimes P$
 in M.L., often called "risk"
 Simon calls it "Vapnik risk" to distinguish it from statistical frequentist risk