

today: finish MLE for multinomial  
 statistical decision theory + properties of estimator

multinomial MLE:

$$X \sim \text{Mult}(n, \pi) \quad \pi \in \Delta_K$$

$$\text{log-likelihood: } \ell(\pi) = \log p(x|\pi) = \log \binom{n}{n_1, \dots, n_K} + \sum_{j=1}^K n_j \log \pi_j$$

$x = (n_1, \dots, n_K)$ 
const. with respect to  $\pi$  → ignore for MLE

$$\text{MLE: } \hat{\pi} = \underset{\pi \in \mathbb{R}^K}{\text{argmax}} \ell(\pi)$$

s.t.  $\pi \in \Delta_K$  } constraint

two options:

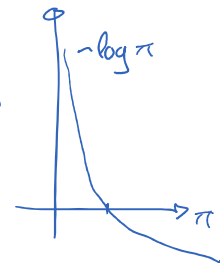
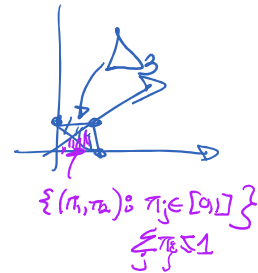
a) reparameterize problem to be full dimensional

$$\pi_K \triangleq 1 - \sum_{j=1}^{K-1} \pi_j$$

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

here magic is that  $\log \pi_j$  acts as barrier fct. away from  $\pi_j = 0$

can try unconstrained opt. on  $\pi_1, \dots, \pi_{K-1}$   
hoping solution is in interior of constraint set



b) use Lagrange multiplier approach to handle equality constraint on  $\Delta_K$   
 [and still not worry about  $\pi_j \in [0, 1]$ ]

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

method: look at stationary pts. (0-gradient) of

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

Lagrange multiplier

↓ necessary but not sufficient conditions for  $(\pi^*, \lambda^*)$  to be optimal

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

$$\ell(\pi) = \sum_j n_j \log \pi_j \quad \frac{\partial J}{\partial \pi_j} = 0 \Rightarrow \frac{n_j}{\pi_j} - \lambda = 0 \Rightarrow \pi_j^* = \frac{n_j}{n}$$

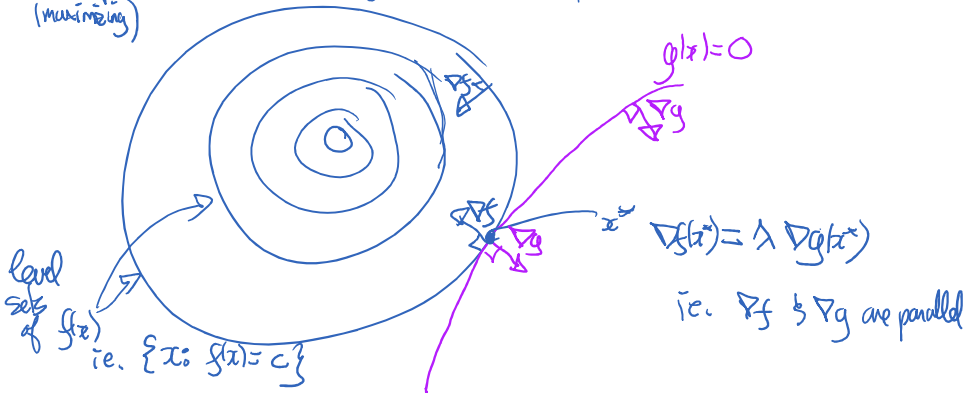
scaling constant

want  $g(x^*)=0$  i.e.  $\sum_j \pi_j^* = 1$   
 $\Rightarrow \sum_j \pi_j = 1 \Rightarrow x^* = \sum_j \pi_j = \eta$

$\Rightarrow \pi_j^* = \frac{n_j}{n}$  MLE for a multinomial

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

picture behind Lagrange multiplier technique:  
 (maximizing)



statistical decision theory

A) bias-variance decomposition for squared loss

estimator: function from data (observation) to parameters

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

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

\* how do we evaluate these estimators?

estimator  $S: \Omega \rightarrow \Theta$   $\theta = S(x)$

most standard tool: frequentist risk of an estimator

$\mathbb{E}_X [L(\theta, S(x))] = R(\theta, S)$   
 (statistical) loss function

average over data

squared loss:  $L(\theta, \theta') \triangleq \|\theta - \theta'\|^2$

$\text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2$

$\theta = S(x)$

$\mathbb{E}(S(x) - \mathbb{E}S(x))^2$

$\mathbb{E}_X [\|\theta - \hat{\theta}\|^2] = \mathbb{E} [\|\theta - \underbrace{\mathbb{E}[\hat{\theta}]}_0 + \mathbb{E}[\hat{\theta}] - \hat{\theta}\|^2]$

$\|a+b\|^2 \leq \|a\|^2 + \|b\|^2 + 2\langle a, b \rangle$   
 $= \|a\|^2 + \|b\|^2 + 2\langle a, b \rangle$

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

$+ 2 \mathbb{E} [\langle \underbrace{\theta - \mathbb{E}[\hat{\theta}]}_{\text{constant}}, \mathbb{E}[\hat{\theta}] - \hat{\theta} \rangle]$   
 $2 \langle \theta - \mathbb{E}[\hat{\theta}], \mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}] \rangle$

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

$\text{bias} \triangleq \|\theta - \mathbb{E}_X[\hat{\theta}]\|_2$

risk for squared loss = bias<sup>2</sup> + variance

bias-variance decomposition  
"tradeoff"

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

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

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

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

15h42

### Statistical decision theory

formal setup:

• a random observation  $D \sim P$

unknown distribution within models the world (often  $P_\theta$ )

• action space  $\mathcal{A}$

• loss  $L(P, a) = \text{loss of doing action } a \in \mathcal{A} \text{ when "the world" is } P$  } describe the goal/task

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

•  $\delta: \mathcal{D} \rightarrow \mathcal{A}$  "decision rule"  
 $\uparrow$   
 $\mathcal{D}$

examples: a) parameter estimation:

$\mathcal{A} = \Theta$  for parametric family  $P_\theta$

$\delta$  is parameter estimator from data  $D = (X_1, \dots, X_n)$  typically

typical loss  $L(\theta, a) = \|\theta - a\|^2$  [usually,  $X_i \stackrel{iid}{\sim} P_\theta$ ]  
 "squared loss" O unknown

but other losses are used,  $KL(P_\theta || P_a)$

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

$\delta$  describes a statistical test

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

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

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

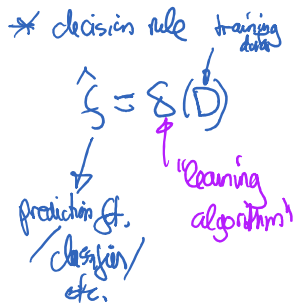
$P_\theta$  gives joint on  $(X, Y)$

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

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

in machine learning  $L(P_\theta, S) \triangleq \mathbb{E}_{P_\theta} [l(Y, f(X))]$

"generalization error"  
"classification error"



e.g. classification  $l(Y, f(X)) = \mathbb{1}\{Y \neq f(X)\}$  0-1 error  
 in ML,  $\rightarrow$  is often called the "risk"

Sumin calls it "Vapnik risk" to distinguish it from statistical frequentist risk

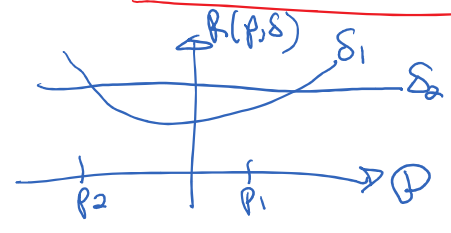
frequentist risk  $\mathbb{E}_D [L(P_\theta, S(D))]$

comparing procedures?

$S_1$  vs.  $S_2$

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

"risk profiles"



$S_i: \mathcal{D} \rightarrow \mathcal{F}$

incompatible functions?  
 $R(\cdot, S_1)$  vs.  $R(\cdot, S_2)$

\* transform to scalar:

• "minimax" analysis:  $\max_{P \in \mathcal{P}} R(P, S)$  "worst case"

• weighted average  $\int_{\mathcal{P}} R(\theta, S) \pi(\theta) d\theta$