

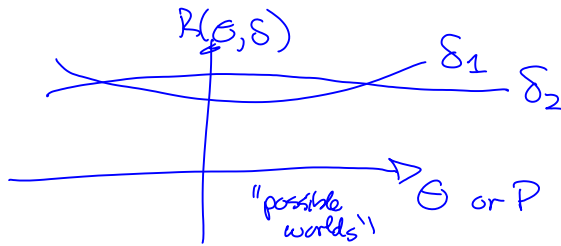
today:

finish decision theory  
 2) prediction models  $\rightarrow$  linear regression  
 logistic regression

(frequentist) statistical decision theory

$L(P, a)$   $\leftarrow$  loss describing task  
 $A$  action space  
 $S: \mathcal{D} \rightarrow A$   
 $a = S(D)$   
 $\uparrow$   
 random observation  $D \sim P$

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



• minimax

'prior' over  $\Theta$  to get weighted summary of performance

$$\int_{\Theta} R(G, S) p(G) d(G)$$

\* Bayesian decision theory:  $\rightarrow$  we condition on data  $D$

instead of averaging

Bayesian posterior risk:  $R_B(a|D) = \int_{\Theta} L(G, a) \underbrace{p(G|D)}_{\text{posterior or } p(G)p(D|G)} dG$

Bayesian optimal action  $S_{\text{Bayesian}}(D) = \arg \min_{a \in A} R_B(a|D)$

example: if  $A = \Theta$  ("estimation")

$$L(G, a) = \|G - a\|_2^2$$

then (exercise)  $S_{\text{Bayesian}}(D) = \mathbb{E}[G|D]$  (posterior mean)

$\rightarrow$  if consider  $\int_{\Theta} p(G) R(G, S_{\text{Bayesian}}) dG$

$\rightarrow$  I use same  $p(G)$  to define it

frequentist  $\mathbb{E}_{D \sim P} R(G, S) \rightarrow \mathbb{E}_{D \sim P} R(G, S)$   
 $\Delta$  minimax nature then this has minimal value

$R(\theta, \delta) \rightarrow$   $\Delta$  regression risk then this has minimal value amongst "all" estimators  $\delta$  (under regularity conditions)

$\mathbb{E}_{\text{emp}} \rightarrow \mathbb{E}_{\text{true}}$  same answer

examples of estimators  $\mathcal{S}: \mathcal{D} \rightarrow \mathcal{H}$

- MLE
- another: MAP, given a prior  $p(\theta)$ , then pick  $\hat{\theta} = \underset{\theta \in \Theta}{\text{argmax}} \underbrace{p(\theta | \mathcal{D})}_{\propto p(\mathcal{D} | \theta) p(\theta)}$
- method of moments:

idea: find injective mapping from  $\mathcal{H}$  to "moments" and surjective on "possible moments"

$\rightarrow \mathbb{E}[X]$   
 $\mathbb{E}[X^2]$   
 etc...

and then invert it from empirical distribution  $\hat{\mathbb{E}}[X]$   
 $\hat{\mathbb{E}}[X^2]$  etc...

example for Gaussian  $X \sim N(\mu, \sigma^2)$

$$\mathbb{E}[X] = \mu$$

$$\mathbb{E}[X^2] = \sigma^2 + \mu^2$$

$$\begin{pmatrix} \mathbb{E}[X] \\ \mathbb{E}[X^2] \end{pmatrix} = f(\mu, \sigma^2)$$

$$\begin{pmatrix} \hat{\mu} \\ \hat{\sigma}^2 \end{pmatrix} = f^{-1}\left(\begin{pmatrix} \hat{\mathbb{E}}[X] \\ \hat{\mathbb{E}}[X^2] \end{pmatrix}\right)$$

- is the context of prediction  $\mathcal{F} = \{f: X \rightarrow Y\}$
- $X \leftarrow$  input  
 $Y \leftarrow$  output space

example of  $\mathcal{S}: \mathcal{D} \rightarrow \mathcal{F}$

is using empirical "risk" minimization (ERM)

$\downarrow$   
 "Vapnik risk"

recall  $L(P, f) = \mathbb{E}_{(x,y) \sim P} [l(y, f(x))] \leftarrow$  "generalization error"

prediction loss

$$\text{ERM} \quad \hat{\mathbb{E}} [l(y, f(x))]$$

$$= \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i)) \text{ for } \{(x_i, y_i)\}_{i=1}^n$$

$$\hookrightarrow \hat{f} = \underset{f \in \mathcal{F}}{\text{argmin}} \hat{\mathbb{E}} [l(y, f(x))]$$

iid setting

suppose that  $P = p \otimes n$  i.e.  $D = \{X_i\}_{i=1}^n$   
 $X_i \stackrel{iid}{\sim} P$

write  $\begin{cases} S_n \\ \hat{\Theta}_n \end{cases}$  to emphasize dependence on  $n$

study  $R(\Theta, S_n)$  as function of  $n$   
 in particular  $R(\Theta, S_n) \xrightarrow{n \rightarrow \infty} 0$   
 "consistency"

for estimation, had  $L(\Theta, S_n(D)) = \|\Theta - S_n(D)\|_2^2$

risk:  $E[\|\Theta - \hat{\Theta}_n\|^2]$

standard statistical consistency:  $\hat{\Theta}_n \xrightarrow{P} \Theta$  "in probability"

$\forall \epsilon > 0, P\{\|\Theta - \hat{\Theta}_n\| > \epsilon\} \xrightarrow{n \rightarrow \infty} 0$

$$\begin{aligned} \xrightarrow{D \sim P} E[\|\Theta - \hat{\Theta}_n\|^2] &= E[\|\Theta - E[\hat{\Theta}_n] + E[\hat{\Theta}_n] - \hat{\Theta}_n\|^2] \\ &= E[\|\Theta - E[\hat{\Theta}_n]\|^2] + E[\|\hat{\Theta}_n - E[\hat{\Theta}_n]\|^2] \\ &\quad + 2E[\langle \Theta - E[\hat{\Theta}_n], E[\hat{\Theta}_n] - \hat{\Theta}_n \rangle] \\ &\quad \underbrace{\langle \Theta - E[\hat{\Theta}_n], E[\hat{\Theta}_n] - \hat{\Theta}_n \rangle}_{\text{constants}} \rightarrow E[\hat{\Theta}_n] - E[\hat{\Theta}_n] = 0 \\ &= \underbrace{\|\underbrace{E[\hat{\Theta}_n]}_{\triangleq \text{bias}} - \Theta\|^2}_{\text{bias}} + \underbrace{E[\|\hat{\Theta}_n - E[\hat{\Theta}_n]\|^2]}_{\text{variance}(\hat{\Theta}_n)} \end{aligned}$$

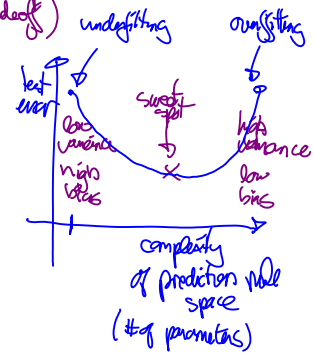
risk for squared loss =  $\|\text{bias}(\hat{\Theta}_n)\|^2 + \text{variance}(\hat{\Theta}_n)$

bias-variance decomposition (tradeoff)

James-Stein estimator: estimating mean of  $N(\mu, \sigma^2 I)$

$\hookrightarrow$  is biased, but lower variance than MLE

and so JS dominates MLE for  $d \geq 3$



note:  $R(\Theta, S) \xrightarrow{n \rightarrow \infty} 0$   
 for squared loss

$E[\|\Theta - \hat{\Theta}_n\|^2] \rightarrow 0$  "convergence in  $L_2$ "

convergence in  $L_2 \Rightarrow$  convergence in prob

so  $R(\hat{\theta}_n) \rightarrow 0 \Rightarrow \hat{\theta}_n \xrightarrow{P} \theta$  i.e. is consistent  
sq. squared loss

$\left. \begin{array}{l} \text{bias} \rightarrow 0 \\ \text{variance} \rightarrow 0 \end{array} \right\} \Rightarrow \text{consistency}$

⊗ asymptotic properties of MLE: under regularity conditions on  $(\mathcal{X}) \{p(x; \theta)\}$

a)  $\hat{\theta}_n \xrightarrow{P} \theta$

b) CLT:  $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \underline{I}(\theta)^{-1/2})$   
information matrix

c) asymptotically optimal  
 i.e. has minimal asymptotic variance (Cramer-Rao lower bound)  
 among all "reasonable" estimators  
↳ consistent

d) invariance:

if  $\Gamma$  reparameterize  $(\mathcal{X})$   
 i.e.  $f: \mathcal{X} \rightarrow \mathcal{X}$   
bijection

then  $\hat{f}(\theta) = f(\hat{\theta})$

if you don't use bijection, can generalize the MLE

with "profile likelihood"  $\mathcal{X} \rightarrow \mathcal{A}$

profile likelihood  $L(\eta) \triangleq \max_{\theta: \eta=g(\theta)} p(\text{data}; \theta)$

define  $\hat{\eta}_{MLE} = \underset{\eta \in g(\mathcal{X})}{\text{argmax}} L(\eta)$

then we have  $\hat{\eta}_{MLE} = g(\hat{\theta}_{MLE})$  } "plug-in"

example:  $\hat{\sigma^2} = (\hat{\sigma})^2$   
 $\hat{\sin \sigma^2} = \sin \hat{\sigma}^2$

$N(\mu, \sigma^2)$