

## Lecture 6 - scribbles

Friday, September 22, 2017

13:21

- today:
- properties of estimators
  - gen. vs. disc. methods
  - linear & logistic regression

### properties of estimators: bias-variance

icd setting:  $P = p_0^{\otimes n}$       $D_n = (x_i)_{i=1}^n$   
 $x_i \stackrel{iid}{\sim} p_0$

notation  $\hat{\theta}_n = \delta_n(D_n)$

↳ emphasize dependence on  $n$

study  $R(\theta, \delta_n)$  as a function of  $n$

in particular, would like  $R(\theta, \delta_n) \xrightarrow{n \rightarrow \infty} 0$

"consistency"

for estimation, typical loss: squared loss  $L(\theta, \delta_n(D)) = \|\theta - \delta_n(D)\|_2^2$

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

ie.  $\forall \epsilon > 0, P\{\|\hat{\theta}_n - \theta\| \geq \epsilon\} \xrightarrow{n \rightarrow \infty} 0$  random  
↑ randomness is from  $\hat{\theta}_n = \delta_n(D)$

$$\text{risk } R(\theta, \delta_n) = \mathbb{E}_{D \sim P} [\|\theta - \hat{\theta}_n\|_2^2]$$

$$= \mathbb{E} [\|\theta - \mathbb{E}[\hat{\theta}_n] + \mathbb{E}[\hat{\theta}_n] - \hat{\theta}_n\|_2^2]$$

$$= \mathbb{E} [\|\theta - \mathbb{E}[\hat{\theta}_n]\|_2^2] + \mathbb{E} [\|\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n]\|_2^2]$$

$$= \mathbb{E}[\|\theta - \mathbb{E}[\hat{\theta}_n]\|^2] + \mathbb{E}[\|\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n]\|^2] + 2 \mathbb{E}[\underbrace{\langle \theta - \mathbb{E}[\hat{\theta}_n], \mathbb{E}[\hat{\theta}_n] - \hat{\theta}_n \rangle}_{\text{const.}}]$$

$$2 \langle \theta - \mathbb{E}[\hat{\theta}_n], \mathbb{E}[\mathbb{E}[\hat{\theta}_n] - \hat{\theta}_n] \rangle = 0$$

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

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

bias-variance decomposition tradeoff

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

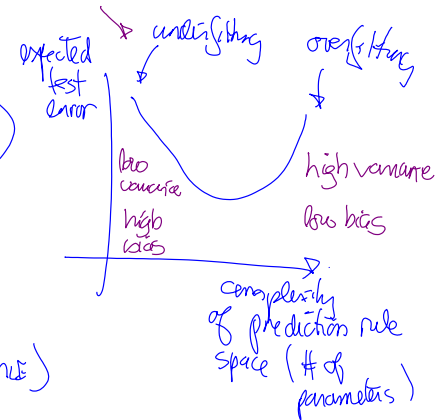
↳ is baised, but lower variance than MLE

and actually, J-S estimator dominates MLE for  $d \geq 3$   
ie.

$$R(\theta, \hat{\theta}_{JS}) \leq R(\theta, S_{MLE})$$

and  $\exists \theta$  st.  $R(\theta, \hat{\theta}_{JS}) < R(\theta, S_{MLE})$

MLE is sometimes inadmissible!



consistency:  $(\hat{\theta}_n \xrightarrow{p} \theta)$

also mentioned

$$R(\theta, \hat{\theta}_n) \xrightarrow{n \rightarrow \infty} 0 \quad \mathbb{E}[\|\hat{\theta}_n - \theta\|^2] \rightarrow 0 \quad \text{"convergence in } \mathcal{L}_2 \text{"}$$

convergence in  $\mathcal{L}_2 \Rightarrow$  convergence in prob.

\* bias  $\rightarrow 0$   
and  
variance  $\rightarrow 0$  }  $\Rightarrow$  consistency

note: in hwk 1, for simplicity, by consistency, I mean that  $R(\theta, \hat{\theta}_n) \rightarrow 0$  (so no need to worry about the subtleties of convergence of probability to show that an estimator is not consistent)

variance  $\leftrightarrow$  )

## properties (asymptotic) of MLE

under regularity conditions on  $\Theta$  &  $p(x; \theta)$

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

b) CLT:  $\sqrt{n}(\hat{\Theta}_n - \theta) \xrightarrow{d} N(0, \underbrace{I(\theta)^{-1}}_{\text{information matrix}})$

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

d) invariance: MLE is preserved under reparameterization

suppose have bijection  $f: \Theta \rightarrow \Theta'$

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

\* if not a bijection, can generalize the MLE  
with "profile likelihood"

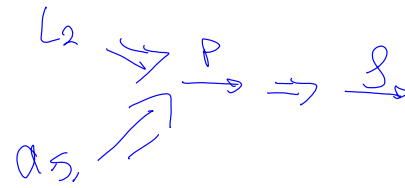
suppose  $g: \Theta \rightarrow \Lambda$  profile likelihood  $L(n) \triangleq \max_{\theta: g(\theta) = \eta} p(\text{data}; \theta)$

define  $\hat{\eta}_{MLE} \triangleq \arg \max_{\eta \in g(\Theta)} L(n)$

then we have  $\hat{\eta}_{MLE} = g(\hat{\theta}_{MLE})$  invariance

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

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

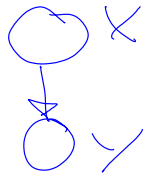


prediction:

want learn a prediction function  $h: X \rightarrow Y$   
 $x \in \mathbb{R}^d$

$Y = \{0,1\} \rightarrow$  binary classification  
 $\{0,1,\dots,k\} \rightarrow$  multiclass "

$\mathbb{R} \rightarrow$  regression



$$p(x,y) = \underbrace{p(y|x)}_{\text{"prediction model"}} \underbrace{p(x)}_{\text{model over } X}$$

$$= \underbrace{p(x|y)}_{\text{"class-conditional"}} \underbrace{p(y)}_{\text{prior over classes}}$$

generative perspective  $\rightarrow$  model  $p(x)$  as well

conditional perspective  $\rightarrow$  only models  $p(y|x)$

(traditionally called "discriminative")  
 more disc  $\rightarrow$

gen	conditional	"fully discrimin"
model $p(y x)$	model $p(y x)$	model $\tilde{h}_\theta: X \rightarrow \mathcal{Y}$ (not nec. $p(y x)$ ) use $l(y, \tilde{y})$ for estimation
more assumptions $\Rightarrow$ less robust for prediction		more robust

Linear regression : conditional approach to regression ( $Y \in \mathbb{R}$ )

$$p(y|x; w) = N(y | \underbrace{\langle w, x \rangle}_{w^T x}, \sigma^2)$$

↑  
parameter

$w \in \mathbb{R}^d$   
 $x \in \mathbb{R}^d$

equivalently :  $Y = w^T X + \epsilon$  where  $\epsilon \stackrel{\text{indep.}}{\sim} N(0, \sigma^2)$

[ aside : we'll use "offset" notation for  $x$  i.e.  $x = \begin{pmatrix} \tilde{x} \\ 1 \end{pmatrix}$   $\tilde{x} \in \mathbb{R}^{d-1}$   
"constant features" ]

thus  $\langle w, x \rangle = \langle w_{1:(d-1)}, \tilde{x} \rangle + \underbrace{w_d}_{\text{bias/offset}}$

\* dataset  $(x_i, y_i)_{i=1}^n$

$x_i \sim$  whatever

$y_i | x_i \stackrel{\text{iid}}{\sim} N(w^T x_i, \sigma^2)$

$$N(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right)$$

conditional likelihood  $p(y_1, \dots, y_n | x_1, \dots, x_n) = \prod_{i=1}^n p(y_i | x_i)$

log-likelihood  $\log p(y_1, \dots, y_n | x_1, \dots, x_n) = \sum_{i=1}^n \log p(y_i | x_i)$

$$\log \pi \quad \log p(y_{1:n} | x_{1:n}) = \sum_{i=1}^n \left[ -\frac{(y_i - w^T x_i)^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right]$$

$$\frac{\partial}{\partial \sigma^2} (\quad) = 0 \quad \sum_{i=1}^n \left( -\frac{(y_i - w^T x_i)^2}{2\sigma^2} - \frac{1}{2} \frac{1}{\sigma^2} \right) = 0$$

(see [note below](#) about  $\sigma^2$  being true global max)

$$\Rightarrow \hat{\sigma}_{MSE}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2$$

"design matrix"

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \in \mathbb{R}^{n \times d} \quad \vec{y} \triangleq \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

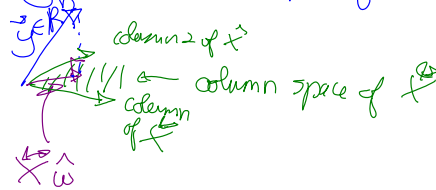
$$Xw = \begin{pmatrix} x_1^T w \\ \vdots \\ x_n^T w \end{pmatrix} \in \mathbb{R}^{n \times 1}$$

$$\|\vec{y} - Xw\|_2^2 = \sum_i (y_i - w^T x_i)^2$$

can rewrite:

$$-\log p(\vec{y} | X) = \frac{\|\vec{y} - Xw\|_2^2}{2\sigma^2} + \text{function}(\sigma^2)$$

minimization  $\|\vec{y} - Xw\|_2^2 \rightarrow$  projecting  $\vec{y}$  on the column space of  $X$   
(geometric view)



$$Xw = \sum_{j=1}^d X_{:j} w_j$$

$j^{\text{th}}$  column of  $X$

- **note about sigma^2 being a global max**

(**aside:** showing that the sigma^2 above is the **global max** is subtle because the objective is not concave in sigma^2. I give more info here for your curiosity, but it is not required for the assignment.)

- Formally, to find a global max of a \*differentiable objective\*, you need to check all **stationary points** (zero gradient points), **as well as the values at the boundary of the domain.**

Thus here, you would need to show that the objective cannot take higher value anywhere at the boundary of the domain (which is the case here (exercise!), as the objective goes to -infinity at the boundary), so you are done (this is the only possible global optimum -- a maximum here, as it should be, given that there are no other stationary points and all values are lower at the boundary, but one could also explicitly check the Hessian to see that it is strictly negative definite at the stationary point, i.e. it looks like a local maximum).

Note that we will see later in the class that the Gaussian is in the exponential family, with a log-concave likelihood in the right ("natural") parameterization, and thus using the invariance principle of the MLE, we could also easily deduce the MLE in the "moment" parameterization which is the usual (mu, sigma^2) one, without having to worry about local optima...

- for a cute counter-example illustrating that a differentiable function could have only one stationary point which is a local min but \*not a global min\* (and thus why one need to look at the values at the boundary), see:

- [https://en.wikipedia.org/wiki/Maxima\\_and\\_minima#Functions\\_of\\_more\\_than\\_one\\_variable](https://en.wikipedia.org/wiki/Maxima_and_minima#Functions_of_more_than_one_variable)

- i.e.

$$f(x, y) = x^2 + y^2(1 - x)^3, \quad x, y \in \mathbb{R},$$

shows. Its only critical point is at (0,0), which is a local minimum with  $f(0,0) = 0$ . However, it cannot be a global one, because  $f(2,3) = -5$ .

(see picture of function [here](#))

(and note that the "[Mountain pass theorem](#)" which basically says that if you have a strict local optimum with another point somewhere with the same value, then there must be a saddle point somewhere (a "mountain pass") i.e. another stationary point, **does not hold for this counter-example** as one of the required regularity condition, the "Palais-Smale compactness condition" fails. Here, the saddle point (which should intuitively exist) "happens at infinity", which is why it only has one stationary point despite (0,0) not being a global minimum)

- the moral of the story: intuitions for multivariate optimization are often misleading! (this counter-example would not work in 1d because of [Rolle's theorem](#))