

Lecture 7 - scribbles

Tuesday, September 22, 2020 14:29

- today:
- finish estimators
 - linear regression

other examples of estimators

4) in the context of prediction $\mathcal{F} = \{f: X \rightarrow Y\}$ X input space
 Y output

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

is using empirical "risk" minimization (ERM)

↳
 Vapnik risk i.e. generalization error

$$L(P, f) = \mathbb{E}_{(X,Y) \sim P} [l(Y, f(X))]$$

replace with $\hat{\mathbb{E}} [l(Y, f(X))] = \frac{1}{n} \sum_{i=1}^n l(y^{(i)}, f(x^{(i)}))$

$$\hat{S}_{ERM} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{\mathbb{E}} [l(Y, f(X))]$$

↳ hypothesis class

James-Stein estimator:

estimator to estimate the mean of $N(\vec{\mu}, \sigma^2 I)$ ← d independent Gaussian variables
 $x_i \overset{\text{indep}}{\sim} N(\mu_i, \sigma^2)$

S_{JS} is baised, but ^{much} lower variance than MLE

recall bias-variance decomposition for squared loss

$$R(\theta, \hat{\theta}) = \mathbb{E} \|\theta - \hat{\theta}\|_2^2 = \underbrace{\mathbb{E} \|\hat{\theta} - \theta\|_2^2}_{\text{bias}} + \underbrace{\mathbb{E} [(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2]}_{\text{variance}}$$

S_{JS} actually strictly dominates S_{MLE} for $d \geq 3$
 ↑ dimension of $\vec{\mu}$

i.e. $R(\theta, S_{JS}) \leq R(\theta, S_{MLE}) \forall \theta$
 and $\exists \theta$ st. $R(\theta, S_{JS}) < R(\theta, S_{MLE})$

→ MLE is inadmissible in this case

(can interpret the S_{JS} as an "empirical" Bayesian method)

(asymptotic)
properties of MLE:

under suitable regularity conditions on $\Theta \ni p(x; \theta)$ $\hat{\Theta}_n = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log p(x_i; \theta)$

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

$$D_n \sim p_{\theta}^{\otimes n}$$

b) CLT
 (central limit theorem)

$$\sqrt{n}(\hat{\Theta}_n - \theta) \xrightarrow{\text{dist.}} N(0, I(\theta)^{-1/2})$$

information matrix

$$\begin{aligned} E\|\hat{\Theta}_n - \theta\|^2 &\xrightarrow{L} 0 \\ \hat{\Theta}_n &\xrightarrow{L} \theta \\ \hat{\Theta}_n &\xrightarrow{P} \theta \end{aligned}$$

c) asymptotically optimal

(Cramer-Rao lower bound)

ie. it has minimal asymptotic variance among all "reasonable" estimators

d) invariance: MLE is preserved under reparameterization

Suppose have a bijection $f: \Theta \rightarrow \Theta'$

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

Example: $\widehat{(\sigma^2)} = (\hat{\sigma})^2$

$$\widehat{\sin \sigma^2} = \sin \hat{\sigma}^2$$

* If not a bijection, can generalize MLE with "profile likelihood"

suppose $g: \Theta \rightarrow \Lambda$

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

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

then we have

$$\hat{\eta}_{MLE} = g(\hat{\theta}_{MLE})$$

"plug-in estimator"

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

eg. $g(\mu) = \mu^2$

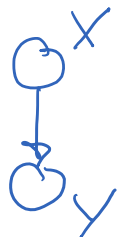
15h34

Prediction

went to Dima's prediction (i.e. $\mu, \sigma^2 \rightarrow \mu, \sigma^2$)

want to learn prediction fct. $h: X \rightarrow Y$
 $x \in \mathbb{R}^d$

$Y = \{0,1\} \rightarrow$ binary class
 $Y = \{0,1, \dots, k-1\} \rightarrow$ multiclass
 $Y = \mathbb{R} \rightarrow$ regression



"prediction model"
 $p(x,y) = p(y|x) p(x)$ (model over X)
 $= p(x|y) p(y)$ ("class conditional", prior over class)

"generative perspective" (in context of classification) \rightarrow Model $p(x)$ as well

"conditional perspective" \rightarrow only model $p(y|x)$
 ("more discriminative" \rightarrow traditionally called "discriminative")

generative	conditional	"fully discriminative"
model $p_G(x,y)$ MLE	model $p_G(y x)$ max. conditional likelihood	model $h_G: X \rightarrow Y$ (not nec. derived from $p(y x)$) reg. ERM; etc... $\sum_{i=1}^n \ell(y^{(i)}, h_G(x^{(i)}))$
more assumptions \Rightarrow less robust for prediction		less assumptions more robust

$\hookrightarrow \hat{h}(x) \triangleq \underset{\tilde{y} \in Y}{\text{argmin}} \sum_y p_G(y|x) \ell(y, \tilde{y})$

if $\ell(y, \tilde{y}) = \mathbb{1}\{y \neq \tilde{y}\}$ (or loss) then $\hat{h}(x) = \underset{\tilde{y} \in Y}{\text{argmax}} p_G(\tilde{y}|x)$

Linear regression: derive/motivate with conditional approach to regression ($X \in \mathbb{R}$)

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

$$p(y_i | x_i; w) = \mathcal{N}(y_i | w^T x_i, \sigma^2)$$

\uparrow parameter $x_i \in \mathbb{R}^d$
 $w \in \mathbb{R}^d$

$$\mathcal{N}(\mu, \sigma^2)$$

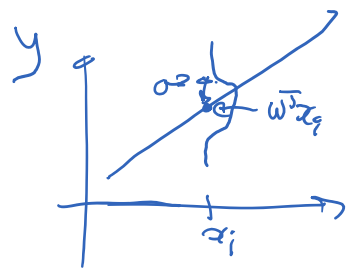
$$\mathcal{N}(y | \mu, \sigma^2)$$

equivalently; $Y_i = w^T X_i + \epsilon_i$ where $\epsilon_i | X_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$

[aside: we use "offset" notation for x

ie. $x = \begin{pmatrix} \tilde{x} \\ 1 \end{pmatrix}$ $\tilde{x} \in \mathbb{R}^{d-1}$
 "constant feature" \downarrow "bias/offset"

thus $\langle w, x \rangle = \langle w_{1:d-1}, \tilde{x} \rangle + w_d$



- dataset $(x_i, y_i)_{i=1}^n$ $X_i \sim$ whatever
 $Y_i | X_i \stackrel{ind\&ap}{\sim} \mathcal{N}(w^T X_i, \sigma^2)$

conditional likelihood $p(y_{1:n} | x_{1:n}) \stackrel{ind\&ap}{=} \prod_{i=1}^n p(y_i | x_i)$

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

\nwarrow not a concave fcn of σ^2

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

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

obj $\rightarrow -\infty$ as $\sigma \rightarrow 0$ or $\sigma \rightarrow +\infty$ so conclude that this is correct global max for w fixed

"design matrix" $X \triangleq$ $\begin{pmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{pmatrix}$ $\begin{matrix} \text{matrix (??)} \\ n \times d \\ \text{matrix} \end{matrix}$ $\text{vector } y \triangleq \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

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

$$\sum_{i=1}^n (y_i - w^T x_i)^2 = \|y - Xw\|_2^2$$

$$\begin{pmatrix} \vdots \\ x_n^T w \end{pmatrix} \quad \sum_{i=1}^n (y_i - w^T x_i)^2 = \|y - Xw\|_2^2$$

can rewrite $-\log p(y_{\text{train}} | X) = \frac{\|y - Xw\|_2^2}{2\sigma^2} + \text{const.}$

design matrix

MCL $\rightarrow \min_w \|y - Xw\|_2^2 \Leftrightarrow$ projecting y on the column space of design matrix X



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

\downarrow
jth column of X

$$\hat{w}_{MLE} = \underset{w \in \mathbb{R}^d}{\text{argmin}} \|y - Xw\|_2^2 \quad \text{"least square"}$$

"Bayesian Methods for Hackers" -- cute book on practical Bayesian approach with probabilistic programming: <https://camdavidsonpilon.github.io/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers/>

(thanks to Dora Jambor for the reference!)

- note about σ^2 being a global max

(aside: showing that the σ^2 above is the **global max** is subtle because the objective is not concave in σ^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 $-\infty$ 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 (μ, σ^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
- 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](#))