

## Lec 9

Thursday, September 26, 2019 10:53

Recap: density estimation

See observations  $y_1, \dots, y_n$

drawn from distrib w/ density  $f$

want to estimate  $f$

KDE (kernel density estimator)

$$\hat{f}_n^{KDE}(y) = \frac{1}{n} \sum_{i=1}^n K_\lambda(y - y_i)$$

$$K_\lambda(u) = \frac{1}{\lambda} \varphi(u/\lambda)$$

$$\varphi(u) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2}\|u\|^2}$$

KDE  $\rightarrow$  kernel regression

$$\begin{aligned} \mu(x) &= \mathbb{E}[Y|X=x] = \int y f_{Y|X=x}(y) dy \\ &= \int y \frac{f_{Y,X}(y,x)}{f_X(x)} dy \\ &= \frac{\int y f_{Y,X}(y,x) dy}{f_X(x)} \\ &= \frac{\int y f_{Y,X}(y,x) dy}{\int f_{Y,X}(y,x) dy} \end{aligned}$$

All we have to do is estimate the density  $f_{Y,X}(y,x)$

If we use the estimate  $\hat{f}_{Y,X}^{KDE}(y,x)$

we get kernel regression

$$\hat{\mu}_n(x) = \frac{\sum_{i=1}^n K_\lambda(x_i - x) Y_i}{\sum_{i=1}^n K_\lambda(x_i - x)}$$

KDE  $\rightarrow$  Kernel classification

Let  $Y = \begin{cases} +1 \\ 0 \end{cases}$

$$P(Y=1|X=x) \geq .5 \iff \underbrace{E[Y|X=x]}_{\text{apply kernel regression}} \geq .5$$

Kernel classifier declare 1 iff

$$\frac{\sum_{i=1}^n K_\lambda(x_i - x) Y_i}{\sum_{i=1}^n K_\lambda(x_i - x)} \geq .5$$

$$\iff \frac{\sum_{i=1}^n K_\lambda(x_i - x) Y_i}{\sum_{i=1}^n K_\lambda(x_i - x) Y_i + \sum_{i=1}^n K_\lambda(x_i - x) (1 - Y_i)}$$

$\Downarrow$   
declare 1 iff

$$\hat{\pi}_1, \hat{f}_1(x) \geq \hat{\pi}_0, \hat{f}_0(x)$$

$\uparrow$                        $\uparrow$   
 KDE of X              KDE of X  
 among pos ex        among neg ex

Eigen - & singular value decompositions

For a square matrix  $A \in \mathbb{R}^{p \times p}$

$\lambda \in \mathbb{C}, v \in \mathbb{C}^p$  are eigen-value/vector pair of  $A$

$$\text{If } \Lambda V = \lambda V$$

If  $A$  is symmetric ( $A = A^T$ )

then all of its e-vals are real

Any square symmetric matrix can be diagonalized

$$A = \sum_{i=1}^p \lambda_i V_i V_i^T$$

$$\text{s.t. } \lambda_i \in \mathbb{R}$$

$$V_i^T V_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{orw} \end{cases} \quad \leftarrow \text{known as an orthogonal set of vectors}$$

$$A = V \Lambda V^T \quad \Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix}$$

$$V V^T = I_{p \times p} = V^T V$$

known as eigendecomposition of  $A$

- If all of  $A$ 's e-vals are nonneg then  $A$  is called positive semidefinite (PSD) and we can define

$$A^{1/2} = V \Lambda^{1/2} V^T \quad \Lambda^{1/2} = \begin{pmatrix} \lambda_1^{1/2} & & 0 \\ & \ddots & \\ 0 & & \lambda_p^{1/2} \end{pmatrix}$$

$$\begin{aligned} A^{1/2} A^{1/2} &= V \Lambda^{1/2} V^T V \Lambda^{1/2} V^T \\ &= V \Lambda^{1/2} I \Lambda^{1/2} V^T \\ &= V \Lambda V^T = A \end{aligned}$$

- If all of  $A$ 's e-vals are non-zero then  $A$  is called invertible/nonsingular

$$A^{-1} = V \Lambda^{-1} V^T$$

$$A^{-1}A = \cancel{V \Lambda^{-1} V^T V \Lambda V^T} = I$$

- If all of  $A$ 's e-values are positive.  
then  $A$  is both PSD & nonsingular  
and is called positive definite (PD)

What about non-square matrices?

Answer: singular value decomp (SVD)

For a rectangular matrix  $A \in \mathbb{R}^{p \times q}$

$\sigma \in \mathbb{R}$ ,  $v \in \mathbb{R}^p$ ,  $u \in \mathbb{R}^q$  are a  
sing val & sing left- & right-sing vector triple

$$\text{if } A^T v = \sigma u \quad Au = \sigma v$$

Any matrix can SVD'ed:

$$A = \sum_{i=1}^{\min(p,q)} \sigma_i u_i v_i^T \quad \text{s.t.} \quad \begin{aligned} u_i^T u_j &= \begin{cases} 1 & i=j \\ 0 & \text{or } u \end{cases} \\ v_i^T v_j &= \begin{cases} 1 & i=j \\ 0 & \text{or } v \end{cases} \end{aligned}$$

$$= U \Sigma V^T$$

$$U \in \mathbb{R}^{p \times \min(p,q)} \quad V \in \mathbb{R}^{q \times \min(p,q)}$$

$$\Sigma \in \mathbb{R}^{\min(p,q) \times \min(p,q)}$$

$$U^T U = I$$

$$\Sigma = \begin{pmatrix} \cdot & & & \\ & \cdot & & \\ & & \cdot & \\ 0 & & & \sigma_{\min(p,q)} \end{pmatrix} \quad V^T V = I$$

Note:  $A^T A = V \Sigma U^T U \Sigma V^T = V \Sigma^2 V^T$   
 $AA^T = U \Sigma V^T V \Sigma U^T = U \Sigma^2 U^T$

## Principal Component Analysis (PCA)

Have many hi-dim observation  $\Sigma \in \mathbb{R}^{n \times p}$

Want to represent these in fewer dims  $q < p$   
 in the most faithful way.

I.e. want to find:

- vectors  $A_1, \dots, A_q \in \mathbb{R}^p$  i.e.  $A = \begin{pmatrix} | & & | \\ A_1 & \dots & A_q \\ | & & | \end{pmatrix} \in \mathbb{R}^{p \times q}$
- loadings  $z_i \in \mathbb{R}^q$  for each  $x_i$   
 i.e.  $Z \in \mathbb{R}^{n \times q}$

such that the approximation

$$\hat{x}_i = \sum_{j=1}^q z_{ij} A_j \in \text{span}(\{A_1, \dots, A_q\})$$

is as close as possible to  $x_i$  (simultaneously across  $i=1, \dots, n$ )

We may as well let  $A_1, \dots, A_q$  be orthonormal

$$A^T A = I$$

$$\min_{Z \in \mathbb{R}^{n \times q}} \sum_{i=1}^n \|x_i - \hat{x}_i\|_2^2 = \sum_{i=1}^n \|x_i - A z_i\|_2^2$$

$$= \|\Sigma - Z A^T\|_F^2$$

$A \in \mathbb{R}^{n \times q}$   
 $A^T A = I$

diff by  $z_i$

$\|M\|_F^2 = \sum_{ij} M_{ij}^2$   
 $= \|\text{vec}(M)\|_2^2$

get  
 $= A^T (x_i - A z_i) = 0$

$\sum_{i=1}^n \|x_i - A A^T x_i\|_2^2 \Rightarrow z_i = A^T x_i$

min  
 $A^T A = I$   
 $A \in \mathbb{R}^{n \times q}$   
 $\|X - \sum A A^T\|_F^2$

Soln:  $A = V_q \leftarrow$  first  $q$  cols of right-sing vecs  
 $V$  from the SVD  $A = U \Sigma V^T$

PCA for noncentered data:

PCA: encoders  $z = e(x) = A^T x$   
 decoders  $\hat{x} = d(z) = A z$

PCA find  $e, d$  s.t.  $x \approx d(e(x))$