

Lecture 6 - scribbles

Tuesday, September 20, 2016
14:27

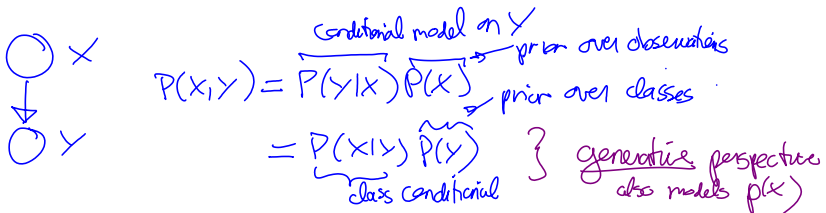
today: gen. vs. discriminative methods
linear & logistic regression

prediction

want to learn a prediction function $f: X \rightarrow Y$

\uparrow
 $\in \mathbb{R}^d$

$\{0, 1\} \rightarrow$ binary classification
 $\{0, 1, \dots, k\} \rightarrow$ multiclass
 $\mathbb{R} \rightarrow$ regression



conditional perspective: only models $P(y|x)$
 (traditionally called "discriminative")

gen	conditional	disc. continuum for prediction "fully disc"
model $p(x,y)$	model $p(y x)$	model $f: X \rightarrow Y$ (not nec. $p(y x)$)
more assumptions \Rightarrow less robust for prediction		more robust

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

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

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

$y \in \mathbb{R}$

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

[aside: we'll use "offset" notation for x $x = \begin{pmatrix} \tilde{x} \\ 1 \end{pmatrix}$ $\tilde{x} \in \mathbb{R}^{d-1}$]
 ϵ constant feature bias/offset

thus $w^T x = w_{(d-1)}^T \tilde{x} + w_d$

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

x_i : N , whatever
 $y_i | x_i \sim N(w^T x_i, \sigma^2)$

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

conditional likelihood $P(y_1, \dots, y_n | x_1, \dots, x_n)$
 $= \prod p(y_i | x_i)$

log-conditional likelihood:

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

$$\frac{\partial}{\partial \sigma^2} (\) = 0$$

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

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

design matrix $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \in \mathbb{R}^{n \times d}$

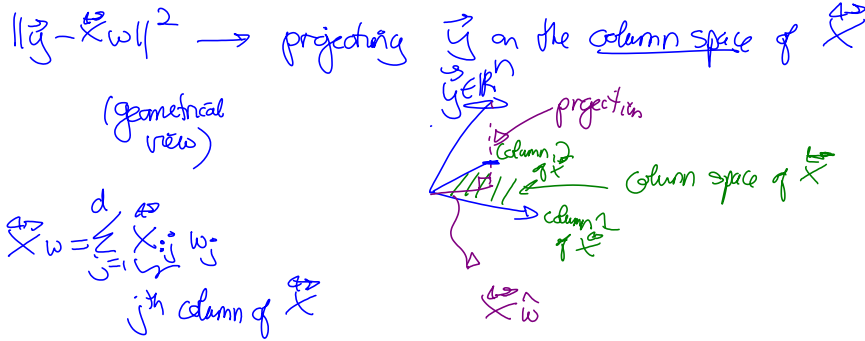
$$\vec{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

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

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

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

minimizing with respect to w



algebra:

$$\frac{\partial}{\partial w} (y - Xw)^T (y - Xw)$$

$$\nabla_w (w^T A w) = (A + A^T) w$$

$$\frac{\partial}{\partial w} [\|y\|^2 - 2y^T Xw + w^T X^T X w]$$

$$= 0 - 2X^T y + 2X^T X w = 0$$

$$\Rightarrow (X^T X) w = X^T y \quad \text{"normal equations!"}$$

if $X^T X$ is invertible, unique solution is

$$\hat{w} = (X^T X)^{-1} X^T y$$

$X^T X \rightarrow d \times d$ matrix
dim $n \times d$

maximum conditional likelihood

$$\text{rank}(X) \leq \min\{n, d\}$$

if $n < d$, then X is not full rank and so $X^T X$ is not invertible

⊗ what if $X^T X$ is not invertible?

(pseudo-inverse would choose the minimum norm $\|w\|$ solution)
 $\hookrightarrow \hat{w} = X^+ y$ amongst $\operatorname{argmin}_w \|y - Xw\|$

problem: pseudo-inverse is not numerically stable

instead better to regularize

MAP estimator viewpoint or the regularized ERM

$$\underbrace{\sum_{i=1}^n \ell(y_i, f_w(x_i))}_{\text{empirical error}} + \underbrace{\lambda \frac{\|w\|_2^2}{2}}_{\text{regularizer}}$$

add prior over w $p(w) = N(w | 0, \frac{1}{\lambda} I)$
 identity in dimension d
 λ precision
 $N(w | \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2}(w-\mu)^T \Sigma^{-1}(w-\mu))$

$$\log p(w | y, x) = \log p(y_{1:n} | x_{1:n}; w) + \log p(w) + \text{cst.}$$

as before $+ \text{cst.} - \frac{\lambda}{2} \|w\|_2^2$

"ridge regression"
 [with $\|w\|_2^2$ penalty]

$$\nabla_w = 0 \Rightarrow (X^T X + \lambda I) w = X^T y$$

$-\log p(w | y, x)$

is strongly convex in w
 \Rightarrow unique global minimum

$$\hat{w}_{\text{MAP}} = (X^T X + \lambda I)^{-1} X^T y$$

ridge regression always invertible

one comment: good practice to standardize or normalize features $(x_1, \dots, x_d) \rightarrow$ make them empirical zero mean & unit std \rightarrow put on $[0, 1]$ scale or something

Logistic regression

binary classification $Y \in \{0, 1\}$ $X \in \mathbb{R}^d$

no assumptions apart that $p(x | Y=1)$ & $p(x | Y=0)$ are densities

$$P(Y=1 | X=x) = \frac{P(Y=1, X=x)}{P(Y=1, X=x) + P(Y=0, X=x)}$$

$$= \frac{1}{1 + \frac{P(Y=0, X=x)}{P(Y=1, X=x)}} = \frac{1}{1 + \exp(-f(x))}$$

where $f(x) \triangleq \log \frac{P(X=x | Y=1)}{P(X=x | Y=0)} + \log \frac{P(Y=1)}{P(Y=0)}$
 class conditional ratio \quad prior odds ratio

in general

$$P(Y=1 | X=x) = \sigma(f(x))$$

"log odds" $\sigma(\cdot) = \frac{1}{1 + \exp(-\cdot)}$

$$\sigma(z) \equiv \frac{1}{1 + \exp(-z)}$$


$$\sigma(-z) = 1 - \sigma(z) \quad [\sigma(z) + \sigma(-z) = 1]$$

$$\frac{d\sigma}{dz}(z) = \sigma(z)(1 - \sigma(z)) = \sigma(z)\sigma(-z)$$

for example: consider class conditional in the exponential family

$$p(x | \eta) = h(x) \exp(\underbrace{\eta^T T(x)}_{\text{linear part in } \eta} - A(\eta))$$

canonical parameters

linear part in η

sufficient statistics

$$\begin{aligned} \xi(x) &= \log \frac{p(x | \eta_1)}{p(x | \eta_0)} + \log \frac{p(y=1)}{p(y=0)} \\ &= (\eta_1 - \eta_0)^T T(x) + A(\eta_0) - A(\eta_1) + \log \left(\frac{\pi}{1-\pi} \right) \\ &= w^T \phi(x) \end{aligned}$$

$$\text{where } w = \begin{pmatrix} \eta_1 - \eta_0 \\ A(\eta_0) - A(\eta_1) + \log \frac{\pi}{1-\pi} \end{pmatrix} \quad \phi(x) = \begin{pmatrix} T(x) \\ 1 \end{pmatrix}$$

$$\text{(logistic regression model } p_w(y=1 | x=z) = \sigma(w^T \phi(x)))$$