

Lecture 7 - scribbles

Tuesday, September 26, 2017

14:35

today: • finish linear regression
• logistic regression

Linear regression:

$$\hat{w}_{MLE} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \| \vec{y} - Xw \|^2$$

algebra: want $\vec{\nabla}_w = 0$ ↙ "vector" convention

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

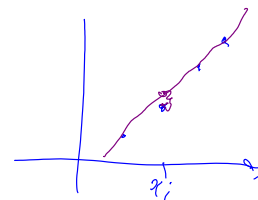
$$\frac{\partial}{\partial w} (\|y\|^2 - 2y^T Xw + w^T X^T X w) \stackrel{\text{const } b}{=} 0$$

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

$$\boxed{(X^T X)w = X^T y}$$

"normal equation"

$$\begin{aligned} \nabla_w (w^T A w) &\stackrel{\text{constant}}{\downarrow} \\ &= (A + A^T)w \end{aligned}$$



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

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

prediction on training set $\hat{y} = X \hat{w} = X (X^T X)^{-1} X^T y$

projection operator on column space of X

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

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

if $n < d$ (i.e. high dimension) or low data regime then $X^T X$ is not invertible

⊛ what if $X^T X$ is not invertible?

⊛ what if $X^T X$ is not invertible?

any \hat{w} s.t. $(X^T X) \hat{w} = X^T y$ is a MLE
 could choose $\hat{w} = \arg \min_{w: X^T X w = X^T y} \|w\|$ ← Moore-Penrose pseudo-inverse

problem: pseudo-inverse is not numerically stable

instead it is better to regularize to get similar effect

regularization (from MAP point view):

suppose we put prior $p(w) = N(w | 0, \frac{\lambda}{2} I)$
 ← "precision" parameter (dxd identity)
 ← Multivariate $N(\vec{\mu}, \Sigma)$ (dxd)

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

$$= - \underbrace{\frac{1}{2} \|\vec{y} - Xw\|^2}_{\text{conditional likelihood}} + \underbrace{\text{fct.}(n)}_{\text{fct.}(n)} - \frac{\lambda}{2} \|w\|^2 + \text{cst.}$$

MAP here $\hat{w}_{\text{MAP}} = \arg \min_w \frac{1}{2} \|\vec{y} - Xw\|^2 + \frac{\lambda}{2} \|w\|^2$ ← "ridge regression"

same as "regularized ERM"

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, f_w(x_i)) + \frac{\lambda}{2} \|w\|^2$$

← squared loss empirical error
 ← regularization

is "strongly convex" w.r.t. w

$f(\cdot)$ is strongly convex
 $\Leftrightarrow f(\cdot) - \frac{\lambda}{2} \|\cdot\|^2$ is convex for some $\lambda > 0$

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

← always invertible

⇓
 unique solution

$$\nabla_w = 0 \Rightarrow \underbrace{(X^T X + \lambda I)}_{\substack{\text{always invertible} \\ (\text{for } \lambda > 0)}} w = X^T y$$

\Downarrow
unique solution

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

ridge regression

one comment:

good practice to either standardize features i.e. make each feature zero mean and unit empirical variance

or
normalize features

make x_i unit norm ($\|x_i\| = 1$)

or
scale features to $[0, 1]$ or $[-1, 1]$

Logistic regression:

setup: binary classification $y = \{0, 1\}$, $X \in \mathbb{R}^d$

motivation: suppose only assumption is \exists pdf (densities) in \mathbb{R}^d

$p(x|y=1)$ and $p(x|y=0)$

$\nwarrow \nearrow$ "class conditionals"

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

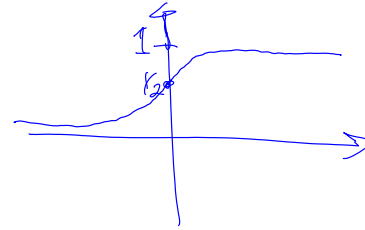
$$= \frac{1}{1 + \frac{P(Y=0, X=x)}{P(Y=1, X=x)}} = \frac{1}{1 + \exp(-f(x))}$$

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

"log odds"

and thus in general $P(Y=1 | X=z) = \sigma(f(z))$

where $\sigma(z) \triangleq \frac{1}{1 + \exp(-z)}$
 "sigmoid function"



some properties of $\sigma(z)$:

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

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

* to motivate linear logistic regression, consider class conditions in the exponential family

$$p(x|\eta) \triangleq h(x) \exp(\eta^T T(x)) \sim A(\eta)$$

"canonical parameter"
"sufficient statistics"
scalar function = log partition function

these specify the family
linear in η

log odds $f(x) = \log \frac{p(x|Y=1)}{p(x|Y=0)} + \log \frac{p(Y=1)}{p(Y=0)}$

$\frac{p(x|\eta_1)}{p(x|\eta_0)}$ $\frac{\pi}{1-\pi}$
 $\eta_0 \rightarrow$ parameter for class 0

$$= (\eta_1 - \eta_0)^T T(x) + A(\eta_0) - A(\eta_1) + \log \frac{\pi}{1-\pi}$$

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

def logistic regression model

"feature map"

get logistic regression model

$$p_w(y=1|x=z) = \sigma(w^T \phi(x))$$

↑ "feature map"

logic regression model: $\gamma = \{0, 1\}$

$$p(y=1|x) = \sigma(w^T x)$$

$$p(y=0|x) = 1 - \sigma(w^T x) = \sigma(-w^T x)$$

$$\left[\text{if } \gamma = \{\pm 1\} \text{ encode } p(y|x) = \sigma(y w^T x) \right]$$

$Y|X=z$ is a Bernoulli($\sigma(w^T x)$)

$$p(y|x) = \sigma(w^T x)^y \sigma(-w^T x)^{1-y}$$

given $(x_i, y_i)_{i=1}^n$, maximum conditional log-likelihood:

$$l(w) = \sum_{i=1}^n \log p(y_i|x_i; w) = \sum_{i=1}^n [y_i \log \sigma(w^T x_i) + (1-y_i) \log \sigma(-w^T x_i)]$$

$$\sigma'(z) = \sigma(z)\sigma(-z)$$

$$\nabla_w \sigma(w^T x_i) = x_i [\sigma(w^T x_i) \sigma(-w^T x_i)] \quad \text{let } v_i \triangleq w^T x_i$$

$$\nabla l(w) = \sum_{i=1}^n x_i \left[y_i \frac{\sigma'(v_i)}{\sigma(v_i)} \sigma(-v_i) - (1-y_i) \frac{\sigma'(-v_i)}{\sigma(-v_i)} \sigma(v_i) \right]$$

$$\left[y_i (\underbrace{\sigma(-v_i) + \sigma(v_i)}_1) - \sigma(v_i) \right]$$

$$\nabla l(w) = \sum_{i=1}^n x_i [y_i - \sigma(w^T x_i)]$$

need to use
numerical optimization

solution for $\nabla l(w) = 0 \Rightarrow$ need to solve a transcendental eq.

solving for $\nabla l(w) = 0 \Rightarrow$ need to solve a transcendental eq.

because $\frac{1}{1 + \exp(-\sigma x_i)} = 0$
(need to solve for this)

(contrast to least square

linear in w

$$\nabla l(w) = \sum_i x_i [y_i - w x_i]$$