

Lecture 8 - linear & logistic regression

Friday, September 25, 2020 13:23

today: finish linear regression
logistic regression

back to least squares

$$\hat{w}_{MLE} = \arg \min_{w \in \mathbb{R}^d} \|y - Xw\|^2$$

algebra: want $\nabla_w \rightarrow 0$

$$\frac{\partial}{\partial w} [(y - Xw)^T (y - Xw)] = 0$$

$$\frac{\partial}{\partial w} [\|y\|^2 - 2y^T Xw + w^T X^T X w] = 0$$

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

$$\Rightarrow (X^T X) w^* = X^T y$$

vector

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

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

"normal equation"

a) if $X^T X$ is invertible, then have unique sol'n

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

$$X \in \mathbb{R}^{n \times d} \Rightarrow \text{rank}(X) \leq \min\{n, d\}$$

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

$X^T X$ is invertible

$$\Rightarrow n \geq d$$

prediction on training set

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

projection matrix on column space of X

(recall geometric perspective)

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

b) what if $X^T X$ is not invertible?

→ there is no unique sol'n

any \hat{w} s.t. $(X^T X) \hat{w} = X^T y$ is a MLE estimate

could choose $\hat{w} = \arg \min_w \|w\|_2$ s.t. $(X^T X) w = X^T y$

Moore-Penrose pseudo-inverse (see Wikipedia)

$$X^+ = (X^T X)^+ X^T$$

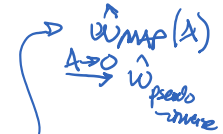
$$X = U \Sigma V^T \quad X^T = V \Sigma^T U^T$$

$n \times d$ $n \times d$
 $\begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_d \\ & & & 0 \end{pmatrix}$

when X is full rank

Problem: pseudo-inverse is not numerically stable

instead it is better to regularize to get similar effect



regularization (can be motivated from MAP pt. of view) dxd identity matrix

suppose we put a prior $p(w) = N(w|0, \sigma^2 I)$

$\lambda \leftarrow$ "precision" parameter

log posterior: $\log p(w|\text{data}) = \log p(y_{\text{train}}|X, w) + \log p(w) + \text{cs.}$

$$= -\frac{1}{2\sigma^2} \|y - Xw\|_2^2 + \text{cs.} - \frac{\lambda}{2\sigma^2} \|w\|_2^2 + \text{cs.}$$

MAP here

$$\hat{w}_{\text{MAP}} = \underset{w}{\text{argmin}} \frac{1}{2} \|y - Xw\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

"ridge regression"

same as "regularized" ERM empirical risk minimization

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

this objective is strongly convex

in w
 \Rightarrow a unique soln

$f(\cdot)$ is λ -strongly convex
 $\Leftrightarrow f(\cdot) - \frac{\lambda}{2} \|w\|_2^2$ is convex in (\cdot)

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

always invertible for $\lambda > 0$

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

ridge regression

no problem for $d > n$

14h39

good practice: to either standardize features i.e. make each feature zero mean & unit empirical variance
or
normalize \rightarrow make x_i unit norm $\|x_i\|_2 = 1$
or
scale features to $[0, 1]$ or $[-1, 1]$

Logistic regression

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

generative model motivation:

suppose only assumption is there exists a pdf (densities) in \mathbb{R}^d for the class conditions:

$$p(x|Y=1) \text{ \& } p(x|Y=0)$$

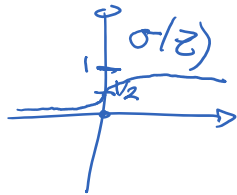
$$P(Y=1|X=x) = \frac{p(Y=1, X=x)}{p(Y=1, X=x) + p(Y=0, X=x)} \Bigg\} p(x=x)$$

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

where $f(x) \triangleq \log \frac{p(x=x|Y=1)}{p(x=x|Y=0)} + \log \frac{p(Y=1)}{p(Y=0)}$
 "log odds" class conditional ratio prior pdf ratio

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

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



some properties of $\sigma(z)$:

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

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

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

$$p(x|\eta) \triangleq \underbrace{h(x)}_{\text{"canonical parameter"}} \exp(\underbrace{\eta^T T(x)}_{\text{"sufficient statistics"}} - A(\eta))$$

scalar η^T : log partition fun.
normalized

these specify the "flat" exponential family that we are considering

Gaussian: $\log p(x|\mu, \sigma^2) = -\frac{1}{2} \log 2\pi\sigma^2 - \frac{(x-\mu)^2}{2\sigma^2}$

$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$-\left[\frac{x^2}{2\sigma^2} - x\frac{\mu}{\sigma^2} + \frac{\mu^2}{2\sigma^2}\right]$$

let $T(x) = \begin{bmatrix} -\frac{x^2}{2} \\ x \end{bmatrix}$

$$\eta(\mu, \sigma^2) = \begin{bmatrix} 1/\sigma^2 \\ \mu/\sigma^2 \end{bmatrix}$$

$$A(\eta) = \frac{1}{2} \log(2\pi\sigma^2) + \frac{\mu^2}{2\sigma^2}$$

$$p(x|Y=1) = p(x|m_1)$$

$$p(x|Y=0) = p(x|m_0)$$

log odds $\xi(x) = \log \frac{\frac{p(x|m_1)}{p(x|m_0)}}{\frac{p(x|Y=1)}{p(x|Y=0)}} + \log \frac{\frac{\pi}{1-\pi}}{\frac{p(Y=1)}{p(Y=0)}}$

$$= (m_1 - m_0)^T T(x) + A(m_0) - A(m_1) + \log \frac{\pi}{1-\pi}$$

$$\triangleq w^T \phi(x)$$

where $w = \begin{pmatrix} m_1 - m_0 \\ A(m_0) - A(m_1) + \log \frac{\pi}{1-\pi} \end{pmatrix}$ $\phi(x) = \begin{pmatrix} T(x) \\ 1 \end{pmatrix}$

get logistic regression model

$$P_w(Y=1|X=x) = \sigma(w^T \phi(x))$$

"feature map"

exercise to reader:

try argument above with $p(x|y) = N(x|\mu_y, \Sigma_y)$

if $\Sigma_0 = \Sigma_1$, then $\phi(x) = \begin{pmatrix} x \\ 1 \end{pmatrix}$

otherwise $\phi(x) = \begin{pmatrix} x x^T \\ x \\ 1 \end{pmatrix}$
optional