## Lec 4

Tuesday, September 10, 2019   10:59

# Recap:

Let's pick the best linear model

$$f \in \{ f_\beta : \beta \in \mathbb{R}^p \}$$

$$f_\beta(x) = \beta^T x$$

best in terms of minimal avg squared errors

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2$$

Last time:

$$\hat{R}_n(f_\beta) = \frac{1}{n} \| \underline{Y} - \underline{X}\beta \|_2^2$$

$$\underline{Y} \in \mathbb{R}^n \quad \underline{Y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$\underline{X} \in \mathbb{R}^{n \times p} \quad \underline{X} = \begin{pmatrix} - x_1 - \\ \vdots \\ - x_n - \end{pmatrix}$$

$$\underline{X}\beta = \begin{pmatrix} f_\beta(x_1) \\ \vdots \\ f(x_n) \\ \beta \end{pmatrix} \in \mathbb{R}^n$$

$$\nabla_\beta \hat{R}_n(f_\beta) = \begin{pmatrix} \partial \hat{R}_n(f_\beta) / \partial \beta_1 \\ \vdots \\ \partial \hat{R}_n(f_\beta) / \partial \beta_p \end{pmatrix} = \frac{2}{n} \underline{X}^T (\underline{Y} - \underline{X}\beta)$$

$$\hat{R}_n(f_\beta) \text{ is diff'able \& convex in } \beta$$

So minimizers = critical pts

Want to solve

$$\frac{2}{n} X^T (Y - X \hat{\beta}) = 0 \qquad \text{for } \hat{\beta}$$

$$X^T Y = X^T X \hat{\beta}$$

Get: $\hat{\beta} = \underbrace{(X^T X)^{-1} X^T} Y$

called the (left) pseudoinverse
of $X$

Wanted to solve $Y \approx X\beta$

↳ can't solve exactly
pseudoinverse gives
the closest-possible soln

$\hat{\beta}$ is called OLS
ordinary least squares

Conditional Expectation is the Best
Regression Model

kNN, OLS are some regr models
Which is the BEST?

should minimizes

X,Y r.v.s representing new unseen
random examples

$$R(f) = \mathbb{E}[ \ell(Y, f(x))]$$
$$= \mathbb{E}[(Y - f(x))^2]$$

Preliminary warm up:
Given a r.v. $Y$, which $c \in \mathbb{R}$
minimizes $\mathbb{E}[(Y-c)^2]$ ?      $c = \mathbb{E}[Y]$

$$\frac{\partial}{\partial c} \mathbb{E}[(Y-c)^2] = \mathbb{E}[\frac{\partial}{\partial c}(c-Y)^2]$$

$$= \mathbb{E}[z(c-Y)]$$

$$= zc - z\mathbb{E}[Y] = 0 \Rightarrow c = \mathbb{E}Y$$

In words:

The mean (avg) is the single no. that is simultaneously closest to all vals of a random variable in avg squared dist.

Now consider minimizing

$$R(f) = \mathbb{E}[(Y - f(x))^2]$$

$$= \mathbb{E}[\mathbb{E}[(Y - f(x))^2 | X]]$$

<span style="color:red">expectation over Y drawn from $Y|X$</span>

What should $f(x)$ be?

$$f^*(x) = \mathbb{E}[Y | X = x]$$

the optimal prediction

GLS regression estimating the conditional mean is some sense:
  - we're replacing $\mathbb{E}$ with empirical avgs
  - we're restricting to linear models

recap: Conditional means, modes, probs <u>are</u> the targets of supervised learning.

## Linear Model for Classification

## Log Odds

Focus for now on binary classification

$$G = \{0,1\}$$

Recall, Bayes classifier declare $\hat{Y} = 1$

When $P(Y=1 | X=x) > P(Y=0 | X=x)$

Same as

$$O.R. = \underbrace{\frac{P(Y=1|X=x)}{P(Y=0|X=x)}}_{\substack{\| \\ \frac{P(Y=1|X=x)}{1 - P(Y=1|X=x)}}} > 1 \quad \in [0, \infty]$$

$\log$ odds: $\log(O.R.) = \log\left(\frac{P(Y=1|X=x)}{P(Y=0|X=x)}\right) \in [-\infty, \infty]$

$$= logit(P(Y=1|X=x))$$

for $p \in [0,1]$

$$logit(p) = \log\left(\frac{p}{1-p}\right) = \log\left(\frac{1}{\frac{1}{p}-1}\right) = \log p - \log(1-p)$$

$logit: [0,1] \longrightarrow [-\infty, \infty]$
         $\underbrace{\qquad}_{domain}$   $\underbrace{\qquad}_{co-domain}$

$logit(P(Y=1|X=x))$ is score for declaring $\hat{Y}=1$

that's symmetric & takes vals in $[-\infty, \infty]$

$\rightarrow$ When it's positive $\Rightarrow$ declare $\hat{Y} = 1$

$\rightarrow$ —//— neg $\Rightarrow$ declare $\hat{Y} = 0$

Logistic regression

Posit $logit\, P(Y=1|X=x) = \beta^T x$

$$P(Y=1|X=x) = \ell \cdots^{-1}(\cdots)$$

$$\pi(\cdots \mid X = x) = \text{logit}^{-1}(\beta^T x)$$
$$= \sigma(\beta^T x)$$

Logistic sigmoid fn

$$\sigma(z) = \text{logit}^{-1}(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z}$$

$\sigma$ gives us a way to transform a score into a prob

$$\sigma : \underbrace{[-\infty, \infty]}_{\text{domain}} \to \underbrace{[0, 1]}_{\text{co-domain}}$$

$$\sigma(0) = \frac{1}{2}$$
$$\sigma(\infty) = 1$$
$$\sigma(-\infty) = 0$$
$$\sigma(-z) = 1 - \sigma(z)$$
$$\partial \sigma / \partial z = \frac{1}{1 + e^{-z}} \cdot \left(1 - \frac{1}{1 + e^{-z}}\right)$$
$$= \sigma(z) \cdot (1 - \sigma(z)) = \sigma(z) \cdot \sigma(-z)$$

## Fitting Logistic regression: Maximum Likelihood

Logistic regression provides a generative model for the ~~data~~

— i.e., a model that ~~specifies~~
how (probabilistically) the
~~data~~ was generated

Given $X = x$, logistic regression says

$$Y \sim \text{Bernoulli}(\sigma(\beta_0^T x))$$

for some special $\beta$

What is $\beta_0$, or what $\beta$ makes the data look the most likely under our generative modl.