# Lecture 3: Linear Regression

If $\underset{\sim}{x} = (x^{(1)}, \ldots, x^{(P)})$

and $\beta = (\beta^{(0)}, \ldots, \beta^{(P+1)})$

and $\underline{x} = (1, x^{(1)}, \ldots, x^{(P)})$

LR proposes the model

$$Y = f_\beta(\underset{\sim}{x}) = \underline{x}^T \beta = \beta^{(0)} + \sum_{j=1}^{P} x^{(j)} \beta^{(j)}$$

Looking at least-squares linear regression (OLS)

$$L(y_n, f_\beta(\underset{\sim}{x}_n)) = (y_n - f_\beta(x_n))^2$$

So if we use ERM to get $\hat{f}$

$$\hat{f} = \underset{f \in F}{\text{argmin}} \; \frac{1}{N} \sum_{n=1}^{N} L(y_n, f_\beta(\underset{\sim}{x}_n))$$

$$= \underset{f \in F}{\text{argmin}} \; \frac{1}{N} \sum_{n=1}^{N} (y_n - \underline{x}^T \beta)^2$$

Notice! to determine $f$ — I simply need to determine $\beta$

So equivalently $\hat{f} = f_{\hat{\beta}}$ i.e. $\hat{f}(\underset{\sim}{x}) = \underline{x}^T \hat{\beta}$

where $\hat{\beta}$ minimizes ER:

$$\hat{f} = f_{\hat{\beta}} \quad \text{where} \quad \hat{\beta} = \underset{\beta}{\text{argmin}} \; \frac{1}{N} \sum_{n=1}^{N} (y_n - \underline{x}^T \beta)^2$$

$$= \underset{\beta}{\text{argmin}} \underbrace{\sum_n (y_n - \underline{x}_n^T \beta)^2}_{\text{RSS}(\beta)}$$

So equiv.

$$f(\underline{x}) = \underline{x}^T \hat{\beta} \quad \text{where} \quad \hat{\beta}^{(OLS)} = \underset{\beta}{\text{argmin}} \text{RSS}(\beta)$$

So ERM = Ordinary Least Squares Regression

Practically, How do we find $\hat{\beta}^{(OLS)}$?

Let $X$ be our <u>design matrix</u> so that

$$X = \begin{bmatrix} - \underline{x}_n^T - \\ \vdots \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} & & x_1^{(P)} \\ \vdots & x_2^{(1)} & & \cdots & \\ & \vdots & & & \\ 1 & x_N^{(1)} & x_N^{(2)} & & x_N^{(P)} \end{bmatrix}$$

$N \times (P+1)$

and $y = (y_1, \ldots, y_n)^T \in \mathbb{R}^{N \times 1}$

then I want to minimize $\text{RSS}(\beta)$

$$\text{RSS}(\beta) = \sum_{n=1}^{N} (y_n - \underset{N \times 1}{\underline{x}^T \beta})^2$$

$$= \| \underset{N \times 1}{y} - \underset{N \times (P+1)}{X} \underset{(P+1) \times 1}{\beta} \|^2$$

$N \times 1$

$$\|y - X\beta\|^2$$

$$\rightarrow y - X\beta = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} - x_1^T - \\ \vdots \\ - x_N^T - \end{bmatrix} \beta$$

$$X$$

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} - x_1^T \beta - \\ \vdots \\ - x_N^T \beta - \end{bmatrix}$$

$$N \times 1 \qquad N \times 1$$

$$\|a\| = \sqrt{\sum_i a_i^2}$$

$$\|a\|^2 = \sum_i a_i^2 \quad = \begin{bmatrix} y_1 - x_1^T \beta \\ y_2 - x_2^T \beta \\ \vdots \\ y_N - x_N^T \beta \end{bmatrix}$$

$$\text{So } \|y - X\beta\|^2 = \sum_n (y_n - x_n^T \beta)^2 = RSS(\beta).$$

- - - - - - - - - -

$$\text{So } \hat{\beta} = \underset{\beta}{\arg\min} \underbrace{\|y - X\beta\|^2} \leftarrow$$

$$\underline{RSS : \mathbb{R}^{P+1} \rightarrow \mathbb{R}}$$

How do I find the minimizer of RSS?

MV Calc Problem: get $\underbrace{\text{derivative}}_{\text{gradient}}$ and set equal to zero

Turns Out:

$$\text{gradient of } RSS(\beta) \text{ w.r.t. } \beta = \frac{\partial RSS}{\partial \beta} = -2 \overbrace{(y - X\beta)^T}^{1 \times N} \overbrace{X}^{N \times (P+1)}$$

$$\underline{1 \times (P+1) \text{ row vector}}$$

So if $\dfrac{\partial RSS}{\partial \beta} = -2(y - X\beta)^T X$

then Calc 3 says set equal to zero

$$\frac{\partial RSS}{\partial \beta} = -2(y - X\beta)^T X = 0$$

$$\Rightarrow -2(y^T - \beta^T X^T) X = 0$$

$$\Rightarrow y^T X - \beta^T X^T X = 0$$

$$\Rightarrow \boxed{X^T y = X^T X \beta} \qquad \underline{\text{Normal equations}}$$

If $X^T X$ is invertible the

$\underbrace{\underbrace{p+1 \times N \times p+1}}_{(p+1) \times (p+1)}$

$$(X^T X)^{-1} X^T y = \overbrace{(X^T X)^{-1} X^T X}^{\text{I}} \beta$$

So $\boxed{\overset{\wedge (as)}{\beta} = (X^T X)^{-1} X^T y}$

So $\hat{f}(\underline{x}) = \underline{x}^T \hat{\beta}$ when $\hat{\beta} = (X^T X)^{-1} X^T y$

$$= \underline{x}^T (X^T X)^{-1} X^T y$$

Consider predictions on traing data

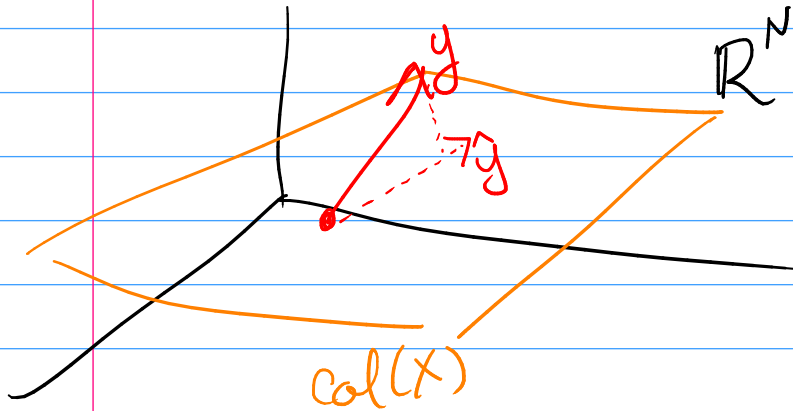$$\hat{y}_n = \hat{f}(\underline{x}_n) = \underline{x}_n^T \underbrace{(X^T X)^{-1} X^T y}_{\hat{\beta}}$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} - & \underline{x}_n^T \hat{\beta} & - \end{bmatrix} = X\hat{\beta} \qquad \boxed{\substack{\text{For traing data} \\ \hat{y} = X\hat{\beta}}}$$

notice then that

$$\hat{y} = X\hat{\beta} = \underbrace{X(X^TX)^{-1}X^T}_{\text{projection mtx onto Col(X)}}y = \text{proj. of } y \text{ onto Col(X)}$$



Col(X)

---

## How flexible is OLS?

Is this regression?

$$y = \beta^{(0)} + \sum_{j=1}^{P} \beta^{(j)} x^{(j)2} \quad ?$$

Yes. This is still linear in $\beta$s.

lets just change the design

$$\underline{x} = (1, x^{(1)2}, x^{(2)2}, \ldots, x^{(P)2})$$

then

$$y = \underline{x}^T \beta$$

this is still a linear operation wrt $\beta$

so as before

$$\hat{\beta} = (X^TX)^{-1}X^Ty \text{ but now I'm using a slightly different } X$$

**Ex.** what about this:

$$Y = f(\underline{x}) = \underline{\beta^{(0)}} + \beta^{(1)} X^{(1)^2} + \beta^{(2)} \log(X^{(2)})$$
$$+ \beta^{(3)} \sin(X^{(1)} X^{(2)}) \; ?$$

Is this <u>linear</u> regression?

<u>Yes.</u> All we need to do is change the design

$$\underline{x} = (1, X^{(1)^2}, \log(X^{(2)}), \sin(X^{(1)} X^{(2)}))$$

So if

$$X = \begin{bmatrix} 1 & & & \\ \vdots & X^{(1)^2} & \lg(X^{(2)}) & \sin(X^{(1)} X^{(2)}) \\ 1 & & & \end{bmatrix}$$

then $\hat{\beta} = (X^T X)^{-1} X^T \underline{y}$

and $\hat{f}(\underline{x}) = \underline{x}^T \hat{\beta}$.

Generically we still have an OLS method in any basis expansion so long as the bases don't depend on $\beta$s:

$$\underline{x}^{(j)} = \varphi_j(\underline{x}) \quad \text{for arb. } \varphi_j \text{ doesn't depend on } \beta s$$

then

$$X = \begin{bmatrix} \varphi_1(\underline{x}) & \varphi_2(\underline{x}) & \cdots \\ & & \end{bmatrix} \quad \text{and} \quad \hat{\beta} = (X^T X)^{-1} X^T \underline{y}$$

What about categorical vars? (factors in R)

e.g. race, color, gender etc.

How do I do something like

$$Y = \beta^{(0)} + \beta^{(1)} \overbrace{\underline{Gender}}^{X^{(1)}} ?$$

Can do this using dummy variable encoding.

$$\underline{x} = (1, 0) \quad \text{if female}$$

$$X^{(1)} = \begin{cases} 0 & \text{fem} \\ \underline{1} & \text{male} \end{cases}$$

$$\underline{x} = (1, \underline{1}) \quad \text{if male}$$

then my design mtx X will look somethy like

$$X = \begin{bmatrix} \underline{1} & 0 \\ \vdots & \vdots \\ & 0 \\ \underline{1} & \vdots \end{bmatrix} \quad \leftarrow \text{females} \\ \leftarrow \text{males}$$

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

How do I interpret $\hat{\beta}$?

Ex. $(\hat{\beta}^{(0)}, \hat{\beta}^{(1)})$ interpret?

If gender = F then $Y \approx \hat{\beta}^{(0)} + \hat{\beta}^{(1)} \cdot 0 = \hat{\beta}^{(0)}$

gender = M then $Y \approx \hat{\beta}^{(0)} + \hat{\beta}^{(1)} \cdot 1$

So $\hat{\beta}^{(0)}$ as the typ. val. for $Y$ (w/ other vars fixed) for a female

$\hat{\beta}^{(1)} = \underbrace{\hat{\beta}^{(0)} + \hat{\beta}^{(1)}}_{\text{typ val for M}} - \underbrace{\hat{\beta}^{(0)}}_{\text{typ val for F}} = $ contrast $=$ diff. btwn typ. vals for M and F

Generically I can encode a $K$ level factor using $K-1$ dummy vars.

data $=$

$\Rightarrow X = $


How do I interpret $\hat{\beta}$?

Holding other vars. constant what is the diff btwn ___ and S     in typ. vals for Y

---

## Fitting Issues

Recall that $\dfrac{\partial RSS}{\partial \beta} = 0$ yielded "normal equations"

$$X^T X \beta = X^T y.$$

IF $X^T X$ is invertible then $\hat{\beta} = (X^T X)^{-1} X^T y.$

When can this fail? when $X^TX$ isn't invertible.

$$\boxed{X^TX \text{ isn't invertible} \iff \text{rank}(X) < \text{\# cols } X}$$

$\Leftarrow$ rank$(X) <$ \# cols $X$ then $\exists v \neq 0$ s.t. $Xv = 0$.

Consequently $X^TXv = X^T 0 = 0$

So $\exists v \neq 0$ s.t. $(X^TX)v = 0$ i.e. $X^TX$ is

rank defficient so its not invertible.

$\Rightarrow$ If $X^TX$ isn't invertible then $\exists v \neq 0$ s.t.

$$X^TXv = 0 \qquad\qquad \|u\| = \sqrt{u^Tu}$$

hence
$$0 = v^T 0 = v^T X^T X v = (Xv)^T(Xv) = \|Xv\|^2$$

So $\underline{\exists v \neq 0}$ when $\|Xv\| = 0$ hence $\underline{Xv = 0}$

$$(\|a\| = 0 \iff a = 0)$$

So $X$ is rank defficient ie. rank$(X) <$ \# cols.

when does this happen in reality?

① If I accidentically incl. a var. twice in design

$$X = \begin{bmatrix} 1 & 1 & 1 \\ \vdots & 2 & 2 \\ \vdots & 3 & 3 \\ 1 & 4 & 4 \end{bmatrix}$$

lin. dep. $\Rightarrow X$ rank def.

this is like saying

$$\hat{Y} = \hat{\beta}^{(0)} + \hat{\beta}^{(1)} X^{(1)} + \hat{\beta}^{(2)} X^{(1)}$$

e.g. $\quad = 1 + 5 X^{(1)} + 7 X^{(1)}$

this has exactly same preds as

$$= 1 + \underbrace{3 X^{(1)} + 9 X^{(1)}}_{12 X^{(1)}}$$

or any model where sum of $\hat{\beta}^{(1)} + \hat{\beta}^{(2)} = 12$

②  If # cols of $X > N$

if $X$ is $N \times P+1$

and $P+1 > N$

then $\text{rank}(X) \leq \min\{\#\text{rows}, \#\text{cols}\} = N < P+1$