

Lin Repr: (super, regr problem)

$$\text{model form: } \underset{\sim}{Y} = f(\underset{\sim}{x}) = \underset{\sim}{x}^T \underset{\sim}{\beta}$$

$$\text{learn: } \hat{f}(\underset{\sim}{x}) = \underset{\sim}{x}^T \hat{\underset{\sim}{\beta}}$$

How do I "learn"  $\hat{\underset{\sim}{\beta}}$ ?

Want:  $Y \approx \hat{f}(x) = x^T \hat{\beta}$

Need some measure of "goodness"

Ordinary Least Squares (OLS) regression

Going to use a squared-error loss  $L$  to measure goodness of fit

On training data

$$N, \quad \tau, \quad 1^2$$

$$L(\beta) = \sum_{n=1}^N (y_n - \tilde{x}_n^T \beta)^2$$

↑  
target to predict
↑  
prediction

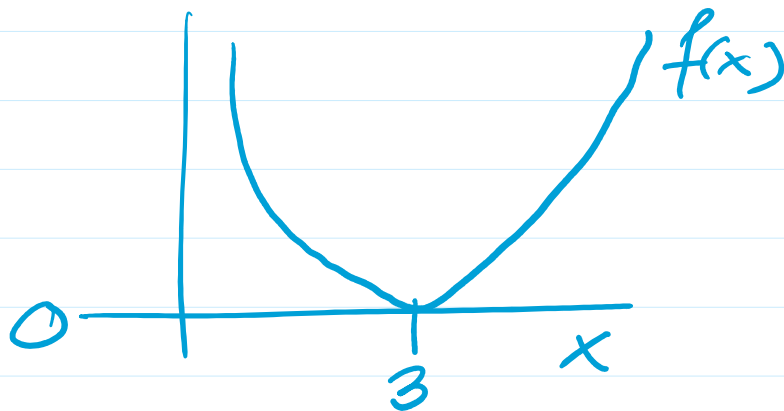
Goal: choose  $\hat{\beta}$  to minimize  $L$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} L(\beta)$$

argmin = value that minimizes

$$\min_x f(x) = 0$$

$$\operatorname{argmin}_x f(x) = 3$$



Simple enough to get closed form soln.

Let our design mtrx be

$$X = \begin{bmatrix} \text{---} & \underline{x}_n^T & \text{---} \end{bmatrix} \in \mathbb{R}^{N \times P}$$

and  $y = (y_1, \dots, y_N)^T \in \mathbb{R}^N$

then  $L(\beta) = \|y - X\beta\|^2$

$$= \sum_{n=1}^N (y_n - \underline{x}_n^T \beta)^2$$

So then  $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \|y - X\beta\|^2$

$L(\beta)$ : fn of  $\beta$

At this point, Calc 3 problem.

Take deriv wrt  $\beta$  and set equal to zero,  
solve.

Solve.

Can show:

$$\text{gradient} = \frac{\partial L}{\partial \beta} = -2(y - X\beta)^T X$$

$1 \times P$

Set equal to zero. (Solve for  $\beta$ )

$$\cancel{-2}(y - X\beta)^T X = 0$$

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

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

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

If  $X^T X$  is invertible, then multiply  
each side by  $(X^T X)^{-1}$  to get

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

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

So, overall for lin Regr.

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

②  $\hat{f}(x_{\text{new}}) = x_{\text{new}}^T \hat{\beta}$

Consider predictions on training data

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_N \end{bmatrix} = \begin{bmatrix} -x_1^T \hat{\beta} - \\ \vdots \\ -x_N^T \hat{\beta} - \end{bmatrix} = \begin{bmatrix} -x_1^T - \\ \vdots \\ -x_N^T - \end{bmatrix} \hat{\beta} = X \hat{\beta}$$

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

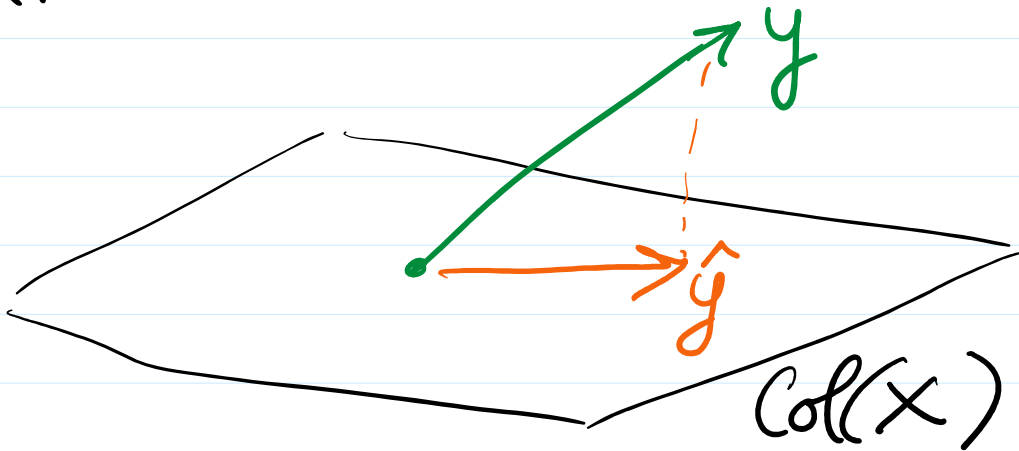
$P_X = \text{proj. onto Col}(X)$

on train.

on train,  
reg. is just  
proj.  $y$  onto  
 $\text{Col}(X)$ .

$$= P_X y$$

$\text{Col}(X)$



Another view  $\text{rank}(X) = p$

$$X = U D V^T$$

$N \times P$

$$D = \begin{bmatrix} D_* \\ 0 \end{bmatrix}$$

$N \times P$

$$X^T X = V D^T U^T U D V^T$$

$$= V D^T D V^T$$

$$= V D_*^2 V^T$$

$$D^T D = [D_* \ 0] \begin{bmatrix} D_* \\ 0 \end{bmatrix}$$

$$= D_*^2$$

$$(X^T X)^{-1} = (VD_*^2 V^T)^{-1} = V D_*^{-2} V^T$$

$$= D_*^{-1}$$

$$X(X^T X)^{-1} X^T = U \cancel{D V^T} V \cancel{D_*^{-2}} V^T V D^T U^T$$

$$= U D D_*^{-2} D^T U^T$$

$$\rightarrow \begin{bmatrix} D_* \\ 0 \end{bmatrix} D_*^{-2} \begin{bmatrix} D_* & | & 0 \end{bmatrix}$$

$$= \begin{bmatrix} I_p & | & 0 \\ \hline 0 & | & 0 \end{bmatrix}$$

$$P_X = U_{1:p} U_{1:p}^T$$

$$\hat{y} = P_X y = U_{1:p} U_{1:p}^T y$$

- Proj. of  $y$  onto  $u_j$

$$y = P_X y = U_{1:p} U_{1:p}^T y$$

$$= \sum_{j=1}^p u_j u_j^T y$$

Proj. onto  $u_j$   
 j<sup>th</sup> col of  $U$

Another way of getting  $\hat{y}$  is

- ① let  $X = UDV^T$
  - ② proj.  $y$  onto each of first  $p$  cols of  $U$
  - ③ Sum them up
- 

Linear regression is quite flexible.

e.g.  $\hat{f}(x) = \hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j x_j^2$

this is still lin. regr.



Still linear in  $\hat{\beta}$ .

$$X = \begin{bmatrix} 1 & X_1 & X_2 & \dots \\ 1 & X_1^2 & X_2^2 & \dots \\ 1 & X_1^3 & X_2^3 & \dots \end{bmatrix}$$

e.g.  $\hat{f}(x) = \beta_1 X_1^3 + \beta_2 \sin(X_2) + \beta_3 X_1 X_2 + \dots$

$$X = \begin{bmatrix} X_1^3 & \sin(X_2) & X_1 X_2 \\ 1 & 1 & 1 \end{bmatrix}$$

---

What about categorical inputs/features?

e.g. race, gender, color, ...

Can I build a model like

$$\hat{f}(x) = \beta_0 + \beta_1(\text{color})$$

# One-hot encoding

$$\begin{bmatrix} R \\ Y \\ B \\ R \\ B \end{bmatrix} \rightarrow \begin{array}{c} R \quad Y \quad B \\ \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \end{array}$$

Issue: What if I include a intercept?

$$X = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ \vdots & 0 & 0 \\ \vdots & \vdots & \vdots \\ 1 & \vdots & \vdots \end{bmatrix}$$

← one of my  
cols is a LC  
of others

⇒  
 $X^T X$  not  
invertible

So what people do typically is throw  
out one col:

$$X = \begin{bmatrix} 1 & 0 & 0 \\ \vdots & 1 & 0 \\ \vdots & 0 & 1 \\ \vdots & \vdots & \vdots \\ 1 & \vdots & \vdots \end{bmatrix}$$

---

More general issue is that to get  $\hat{\beta}$  I solved normal eqns

$$\boxed{X^T X \beta = X^T y}$$

and sometimes  $X^T X$  not invertible.

---

Thrm:

$$X^T X \text{ is invertible} \iff \text{rank}(X) = p$$

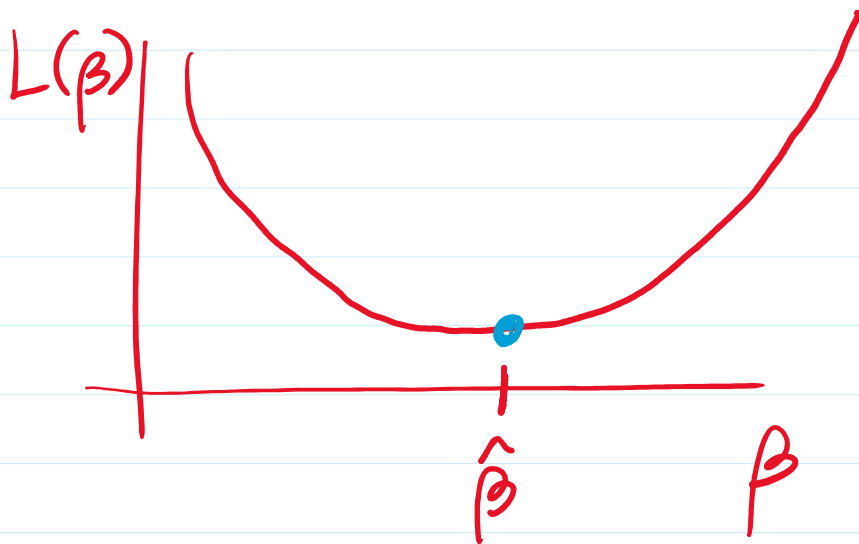
= # cols  
of  $X$

---

When does this happen in reality?

- ① Accidentally include a var twice
- ② One var is a LC of others  
(e.g. one-hot encoding)
- ③ If # cols ( $P$ )  $>$  # rows ( $N$ )
- ④ This can be problematic if  $X^T X$  is "almost" not invertible  
e.g. one var is approx. to a LC of others.

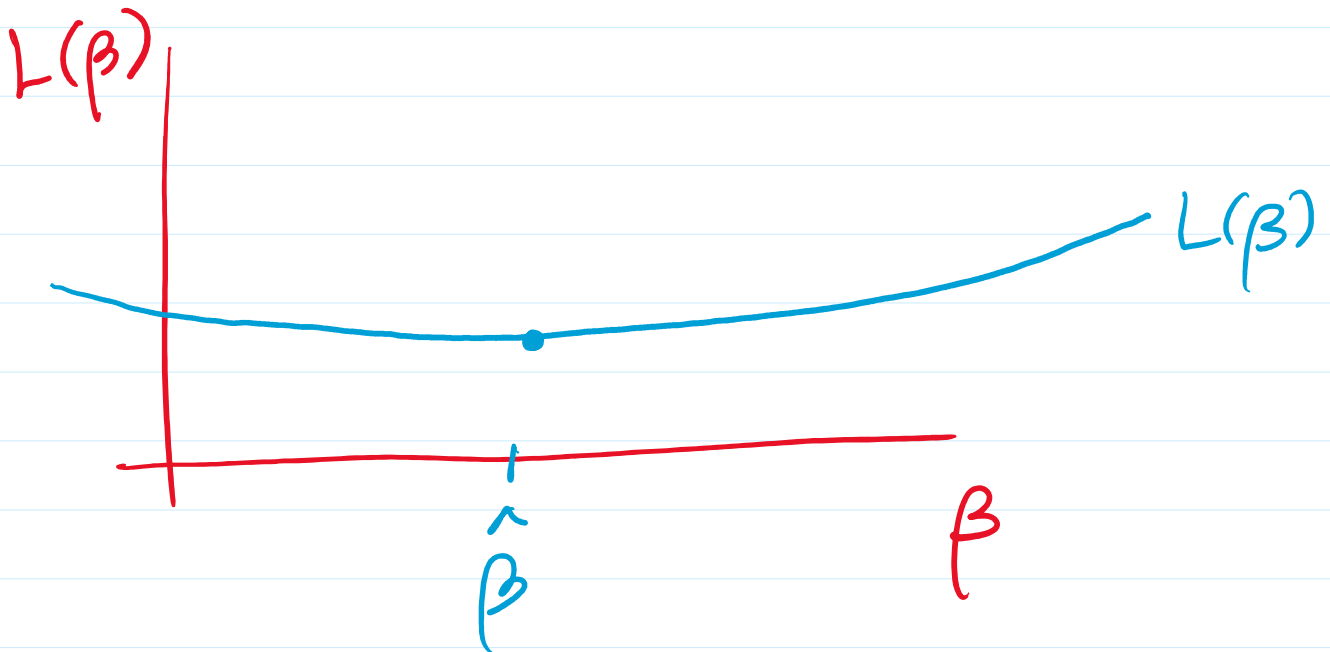
①  $X^T X$  is invertible



②  $X^T X$  not invertible



③  $X^T X$  almost not invertible



$\hat{\beta}$  very sensitive