

## Lecture 4: More Linear Regression + KNN

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

Ex. Variable = Hogwarts House = {G, R, H, S}

$$\text{data} = \begin{bmatrix} \dots & H & \dots \\ \dots & R & \dots \\ \dots & G & \dots \\ \dots & H & \dots \\ \dots & S & \dots \\ \vdots & \vdots & \vdots \end{bmatrix} \Rightarrow X = \begin{bmatrix} & H & R & G & \\ 1 & 0 & 0 & 0 & \\ 0 & 1 & 0 & 0 & \\ \dots & 0 & 0 & 1 & \dots \\ 1 & 0 & 0 & 0 & \\ 0 & 0 & 0 & 0 & \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

How do I interpret coefs?

Coef assoc. w/ H dummy is

"holding other vars. constant, what is predicted diff btwn H and S"

Fitting issues:

Recall that I get  $\hat{\beta}$  by setting  $\frac{\partial \text{RSS}}{\partial \beta} = 0$

which yields my normal eqns:

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

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

when can this fail?

$P$  or  $P+1$

$$X^T X \text{ isn't invertible} \Leftrightarrow \text{rank}(X) < \# \text{ cols of } X$$

pf.  $\Leftarrow$  Assume  $\text{rank}(X) < \# \text{ cols}$ .

Then  $\exists v \neq 0$  where  $Xv = 0$ .

Thus  $X^T X v = 0$ .

So  $\exists v \neq 0$  when  $(X^T X)v = 0$

i.e.  $\text{rank}(X^T X) < \# \text{ cols of } X^T X$  i.e. it's not invertible.

$\Rightarrow$  If  $X^T X$  isn't invertible then  $\exists v \neq 0$  where  $X^T X v = 0$ .

for  $a$ ,  $\|a\|^2 = a^T a$

$$\text{thus } 0 = v^T 0 = v^T X^T X v = (Xv)^T (Xv) = \|Xv\|^2$$

$$\text{So } \|Xv\|^2 = 0 \Rightarrow \|Xv\| = 0$$

can only happen if  $Xv = 0$

$\exists v \neq 0$  where  $Xv = 0$ , i.e.  $X$  is rank deficient.

$$\|a\| = 0 \Rightarrow a = 0$$

When does this happen in reality?

- ① Accidentally include a var. twice in my design

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

- ② If one of my vars is a LC of the others.

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

$$\text{Intercept} = M + F$$

- ③ If  $\underbrace{\# \text{ cols of } X}_P > \underbrace{\# \text{ obs.}}_N$

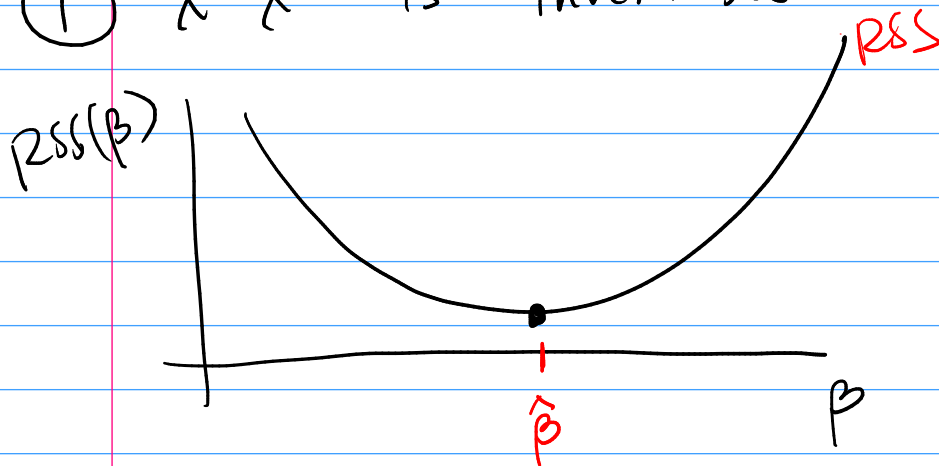
Ex. I measure 50,000 genes for each of 30 patients

- ④ This can "approx." happen if some vars are almost a LC of each other

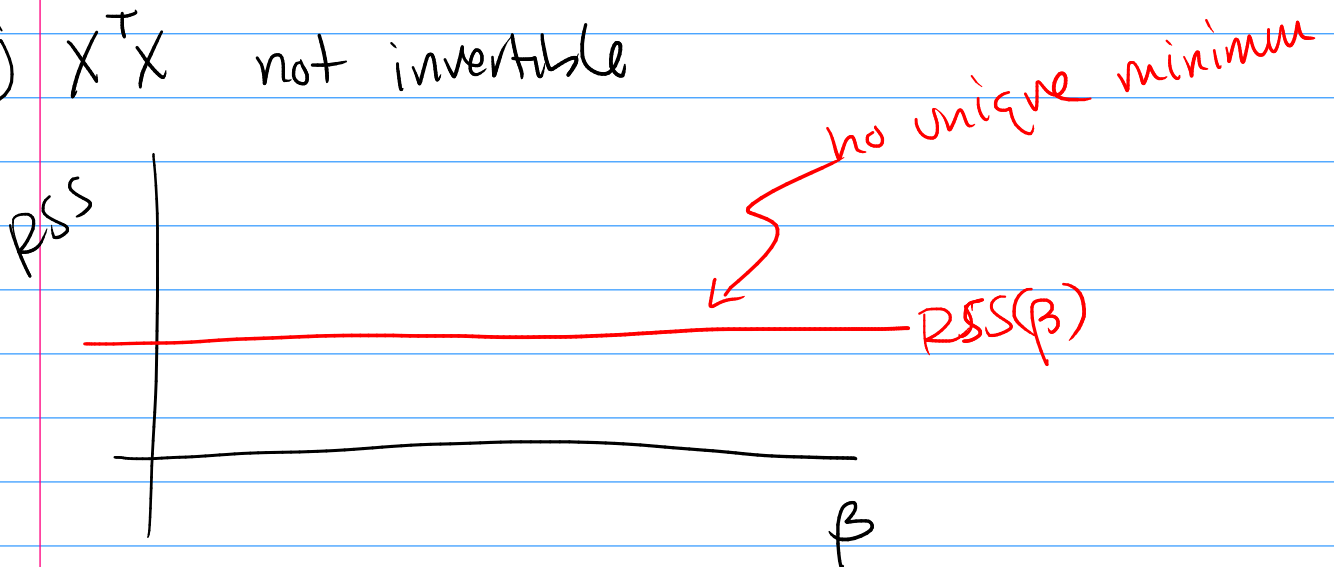
Ex. two vars. are highly correlated

What  $RSS(\beta)$  looks like in various situations

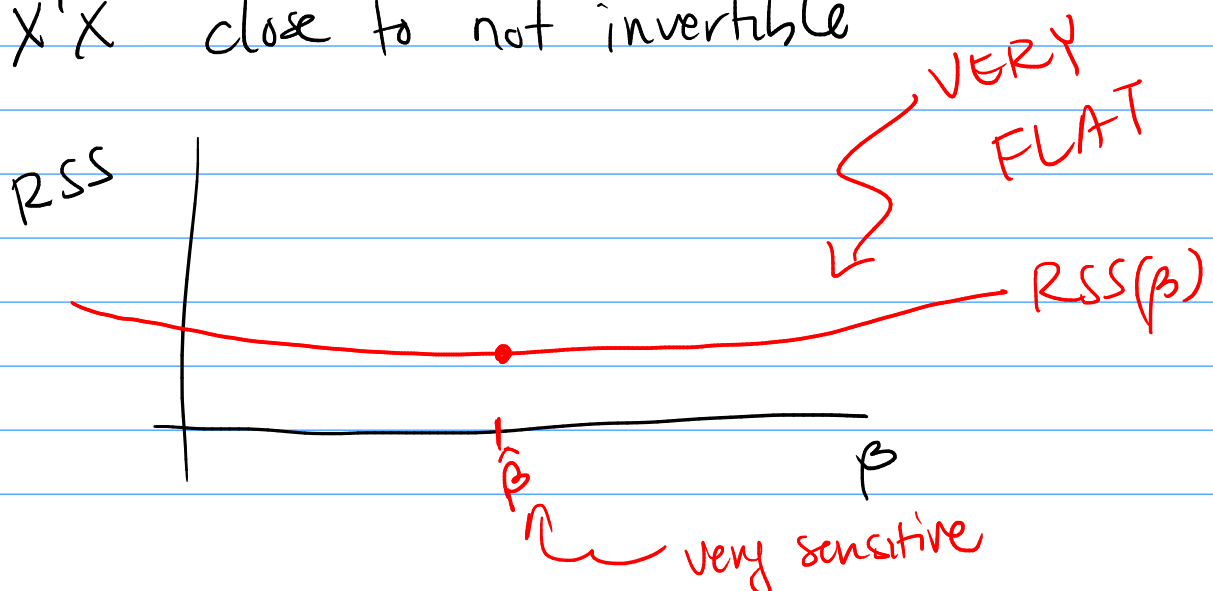
①  $X^T X$  is invertible



②  $X^T X$  not invertible



③  $X^T X$  close to not invertible



# KNN Regression (K nearest neighbors regression)

For LR we have a really strong global assumption about the form of  $f$

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

e.g. in 1-D



Predictions follow linear fn in entire space

Training data affects fit very far away

Benefit: strong global assumption makes  $\hat{f}$  practical to find (opt.  $\beta$  over a  $p$ -dim'l space)

KNN makes a weaker local assumption about the form of  $\hat{f}$

where  $\hat{f}(\underline{x})$  only depend on nearby training pts

In particular  $\hat{f}(\underline{x}) = \frac{1}{K} \sum_{n \in N_K(\underline{x})} y_n =$  avg.  $y_n$  over  $K$  nearest training pts to  $\underline{x}$

$N_K(x)$  = K neighborhood of  $x$   
= indices of K  
nearest training pts  
to  $x$