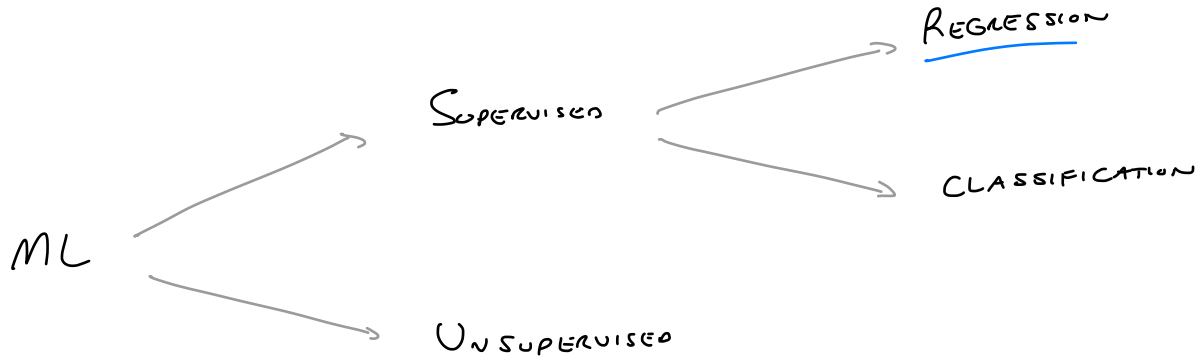
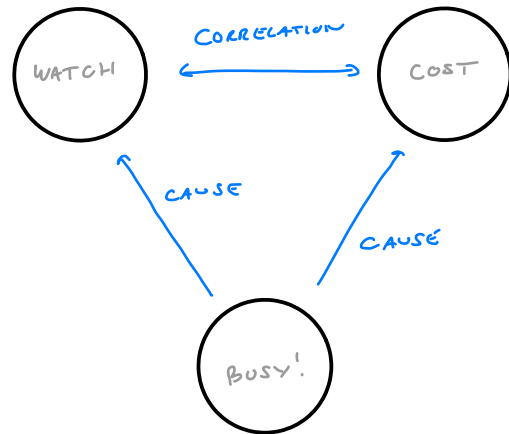
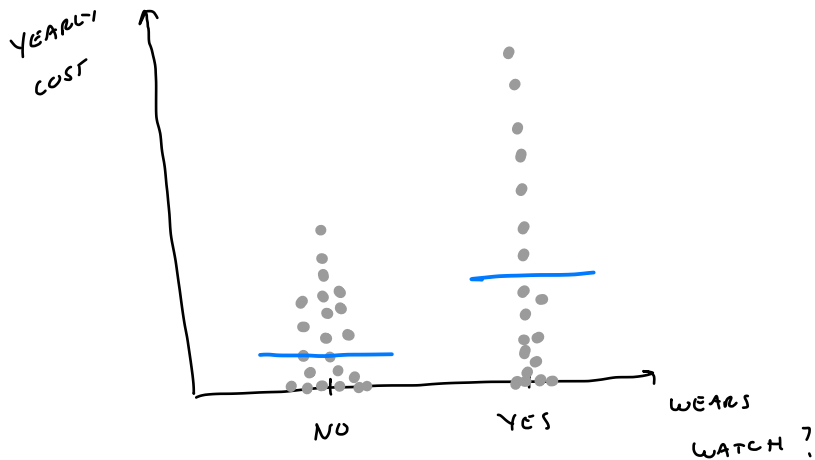
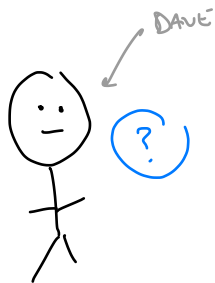
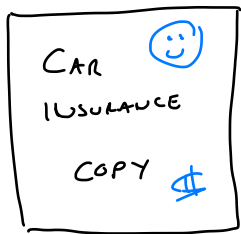


LINEAR REGRESSION

FOR PREDICTION



# EXPLANATION VERSUS PREDICTION



CORRELATION



USEFUL FOR  
PREDICTION

vs

~~CAUSATION~~



NICE TO KNOW.  
HARD TO DO.

# THE REGRESSION TASK

NUMERIC

RESPONSE

FEATURES

GRADE	GPA	MAJOR	YEAR
89	3.23	STAT	3
100	3.51	MATH	3
72	2.61	ECON	4
93	4.00	STAT	4
95	3.62	CS	2

"DATA VIEW"

NEW STUDENT

GPA = 3.42  
MAJOR = MATH  
YEAR = 4  
GRADE = ?

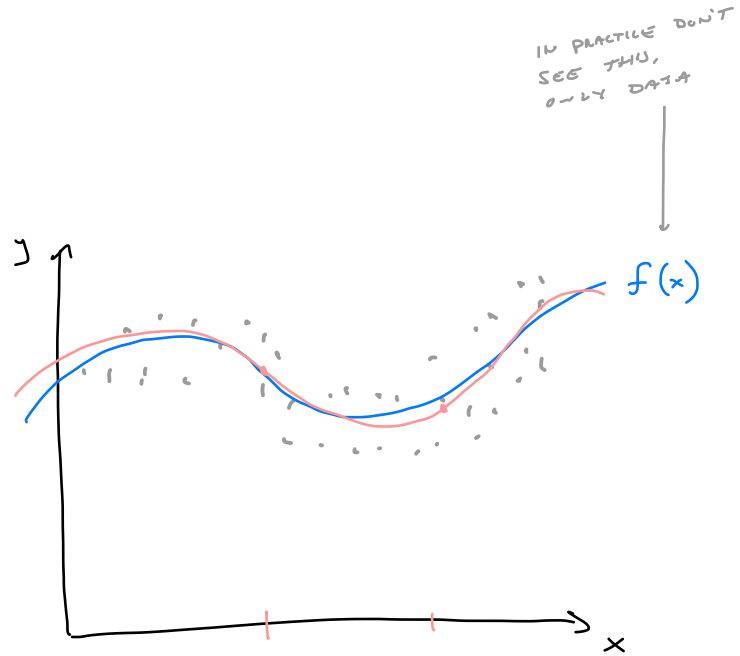
# THE REGRESSION TASK

"MATH VIEW"

$$Y = f(x) + \epsilon$$

↑                      ↑                      ↑  
RESPONSE            SIGNAL            NOISE

WANT TO LEARN THIS FUNCTION



# MAKING GOOD PREDICTIONS

WE WANT  $f(x)$  CLOSE TO  $Y$ .

IN PARTICULAR WE WANT  $(Y - f(x))^2$  TO BE SMALL

THIS HAPPENS WHEN

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

↑  
THE REGRESSION FUNCTION

# LINEAR REGRESSION MODELS

ASSUME

LINEAR COMBINATION OF FEATURES

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

OR

CONDITIONAL MEAN

$$Y | X=x \sim \mathcal{N}(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$$

THAT IS

REGRESSION  
FUNCTION

$$\rightarrow \mathbb{E}[Y | X=x] = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

How TO "FIT" LINEAR MODELS ?

"TRAIN"

USING DATA  $(x_i, y_i)$   
 $i=1, \dots, n$

FIND  $\beta$ s THAT  
MINIMIZE

$$\sum_{i=1}^n \left( y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \right)^2$$

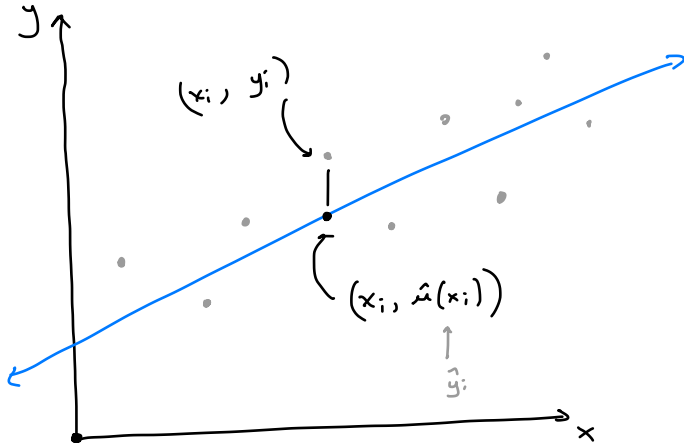
"LEAST SQUARES"

CALL THESE  $\hat{\beta}$ s

GIVES 
$$\hat{\mu}(x) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$$

AS ESTIMATE OF 
$$\mu(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$



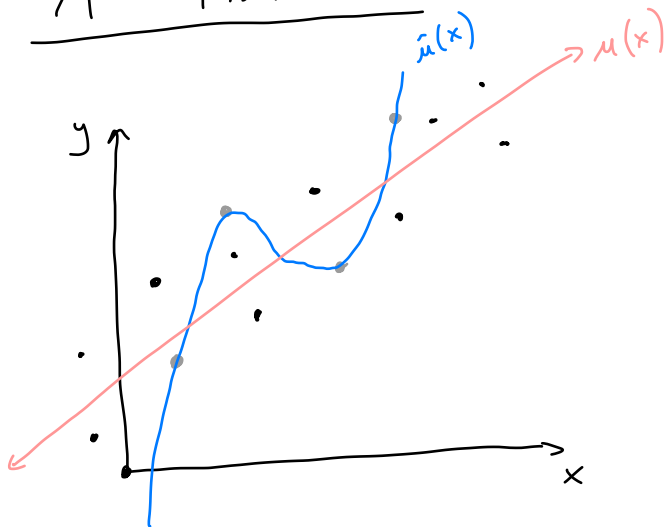


$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}(x_i))^2}$$

RMSE

ROOT MEAN SQUARED ERROR

A PROBLEM



Assume

$$\mu(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

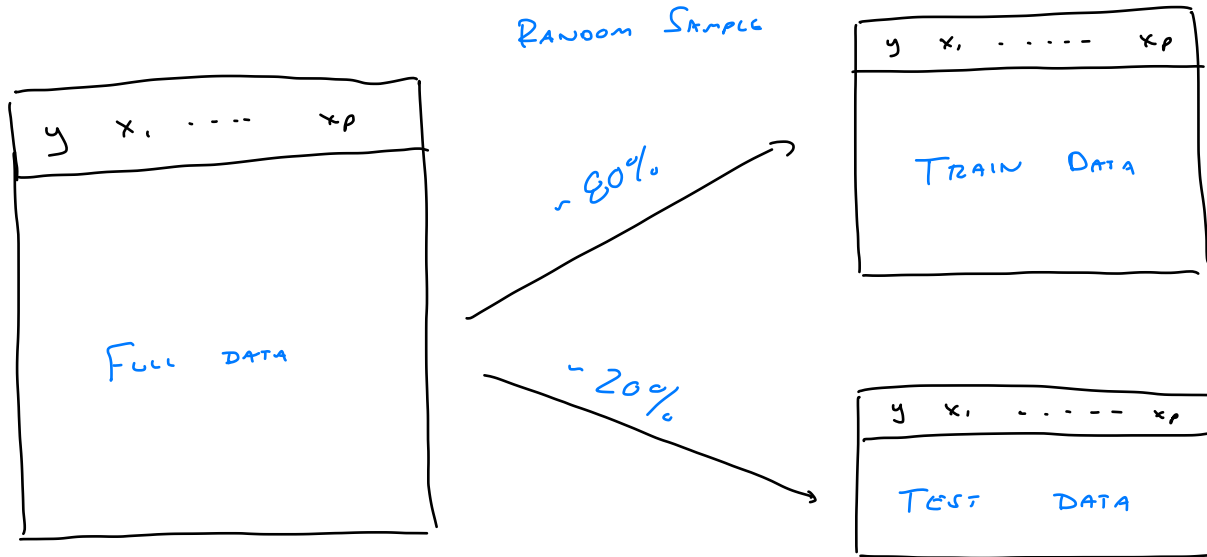
~~PERFECT!~~ ?

OVERFITTING

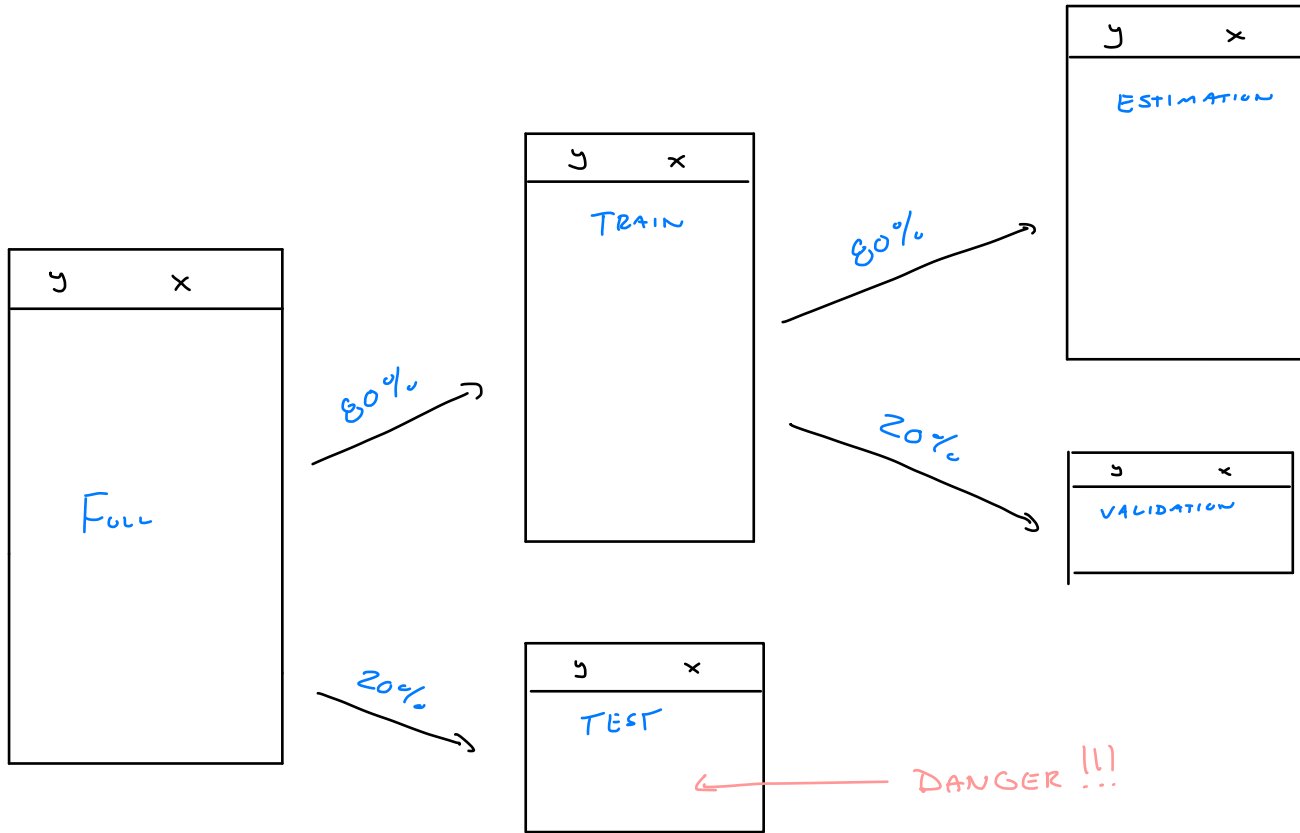
A

SOLUTION

USE "UNSEEN" DATA TO EVALUATE MODELS

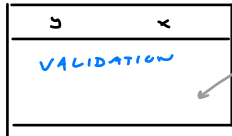


# MORE COMPLICATED SOLUTION





FIT "CANDIDATE"  
MODELS HERE



EVALUATE "CANDIDATE"  
MODELS HERE

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{u}(x_i))^2}$$

VALIDATION RMSE

FROM VALIDATION DATA

FIT TO ESTIMATION DATA

MODEL SELECTION →

ERROR ESTIMATION →

- PICK MODEL WITH SMALLEST VALIDATION RMSE
- RE-FIT TO TRAIN DATA
- REPORT METRIC WITH THIS MODEL ON TEST DATA

# FRAMEWORK

- "LOOK AT THE DATA"
- SPLIT DATA      TST-TRN      EST-VAL
- PICK CANDIDATE MODELS
- FIT MODELS      TO ESTIMATION DATA
- EVALUATE AND SELECT      ON VALIDATION DATA
- FIT CHOSEN MODEL      TO TRAIN DATA
- ESTIMATE ERROR OF CHOSEN MODEL      ON TEST DATA
- USE!