

LINEAR REGRESSION

- CLASSIFICATION

$(x, f(x))$
 \uparrow
 $\{0, 1\}$

- HALFSPACES
- DECISION TREES

- REAL-VALUED LABELS (x, y) $y \in \mathbb{R}$

X, Y two RANDOM VARIABLES

WE WANT TO PREDICT THE VALUE / LABEL
WE GET TO SEE X.

- WE WANT TO PREDICT Y; WE DON'T SEE X

(XY) optimal GUESS FOR Y IS $E[Y]$

- MEASURE OUR LOSS USING SQUARE-LOSS: $(\text{PREDICTION} - Y)^2$

- WE OBSERVE X WE WANT TO PREDICT Y

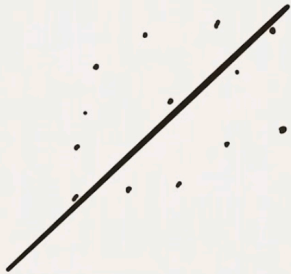
OPTIMAL PREDICTION $E[Y|X] = f(X)$

REGRESSION FUNCTION

OBSTACLE: $f(X)$ COULD BE UNKNOWN OR VERY HARD TO COMPUTE!

LINEAR REGRESSION ASKS THE FOLLOWING QUESTION:

GIVEN X WHAT LINEAR FUNCTION OF X
SHOULD WE USE TO PREDICT Y ?



• WE WANT TO LEARN
COEFFICIENTS β_0 AND β_1 ,

TO MINIMIZE

$$\mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[(Y - (\beta_0 + \beta_1 X))^2 \right]$$

DRAW A TRAINING SET OF SIZE m

$$(x^1, y^1), \dots, (x^m, y^m)$$

"SIMPLE LINEAR
REGRESSION"

$$\min_{\beta_0, \beta_1} \frac{1}{m} \sum_{j=1}^m \left(y^j - (\beta_0 + \beta_1 x^j) \right)^2$$

• TAKE DERIVATIVE w.r.t. β_0, β_1 SET THEM EQUAL TO 0.

$$\frac{\partial \mathcal{L}}{\partial \beta_0} = \sum_{j=1}^m (y^j - \beta_0 - \beta_1 x^j)(-2) = 0$$

$$\frac{\partial \mathcal{L}}{\partial \beta_1} = \sum_{j=1}^m (y^j - \beta_0 - \beta_1 x^j)(-2x^j) = 0$$

$$\frac{1}{n} \sum_{j=1}^n (y^j - \beta_0 - \beta_1 x^j) = 0$$

$$\frac{1}{n} \sum_{j=1}^n (y^j - \beta_0 - \beta_1 x^j) (x^j) = 0$$

SOLVE FOR

β_0

β_0 IN TERMS OF $\beta_1, \bar{y}, \bar{x}$;

↑
avg y
value

↑
avg x
value

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

β_1 WILL NOT INVOLVE β_0

$$\beta_1 = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - (\bar{x})^2}$$

$$\overline{xy} - \beta_0 \bar{x} - \beta_1 \overline{x^2}$$

$$\overline{xy} - \bar{x} \bar{y} + \beta_1 (\bar{x})^2 - \beta_1 \overline{x^2} \quad \beta_1 = \frac{\text{COV}(X, Y)}{\text{VAR}(X)}$$

$$\text{COV}(X, Y) = E[XY] - E[X] \cdot E[Y]$$

x

REGRESSION WITH MULTIPLE VARIABLES.

0 5 2

$X \in \mathbb{R}^n$ $y \in \mathbb{R}$ FITTING A LINE TO n -dimensional DATA.

$X \leftarrow$ matrix X IS GOING TO BE AN $m \times n$ MATRIX

- \cdot m rows \leftarrow EACH ROW IS EQUAL TO X^i DRAWN FROM \mathcal{D}
- \cdot n columns \leftarrow EACH POINT IS IN \mathbb{R}^n

$y \in \mathbb{R}^m$ \leftarrow LABELS FOR THESE m POINTS.

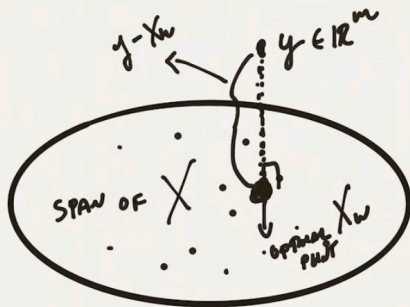
GOAL: FIND A VECTOR $w \in \mathbb{R}^n$ $\|X \cdot w - y\|_2^2$

$$X^1 = X_1^1, \dots, X_n^1 \quad y^1$$

$$\left(y - (X_1^1 w_1 + \dots + X_n^1 w_n) \right)^2$$

$$\min_w \|Xw - y\|_2^2$$

Xw IS A VECTOR IN THE SPAN OF THE COLUMNS OF X .



VECTOR $y - Xw$ IS ORTHOGONAL TO X

$$X^T \cdot (y - Xw) = 0$$

$$X^T y - X^T X w = 0 \quad \text{1st ISSUE: WHAT IF}$$

$$X^T y = X^T X w$$

$(X^T X)$ IS NOT INVERTIBLE?

$$(X^T X)^{-1} X^T y = w \quad \text{"PSEUDO-INVERSE"}$$

WHAT IS THE RUNNING TIME FOR COMPUTING w ?

NORMAL EQUATIONS

$$O(n^3 + m \cdot n^2)$$

MAXIMUM LIKELIHOOD

ASSUMPTION "SIMPLE LINEAR REGRESSION CASE"

$$X; \text{ ASSUME } Y = \beta_0 + \beta_1 X + \epsilon$$

RANDOM NOISE
VARIABLE
 $\epsilon \sim N(0, \sigma^2)$

DRAWN X^1, \dots, X^m AND y^1, \dots, y^m

WE WANT TO UNDERSTAND: FOR A FIXED CHOICE OF
 β_0 AND β_1 (σ^2 IS KNOWN)

WHAT IS THE PROBABILITY THAT WE SEE $(x^1, y^1) \dots (x^m, y^m)$

LET'S WRITE DOWN LIKELIHOOD FUNCTION

PROBABILITY OF SEEING TRAINING SET GIVEN A
CHOICE β_0 AND β_1 OF OUR PARAMETERS

$$\prod_{i=1}^m p(y^i | x^i; \beta_0, \beta_1) =$$

$$\prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma^2} \cdot e^{-\frac{(y^i - (\beta_0 + \beta_1 x^i))^2}{2\sigma^2}}$$

LIKELIHOOD
OF
OUR TRAINING
SET.

CHOOSE β_0 AND β_1 THAT
MAXIMIZES THIS LIKELIHOOD

$$L(b_0, b_1)$$

INSTEAD OF DIRECTLY MAXIMIZING LIKELIHOOD
WE WILL MAXIMIZE LOG-LIKELIHOOD

$$\text{LOG } (L(\beta_0, \beta_1)) = \log \prod_{i=1}^m p(y^i | x^i; \beta_0, \beta_1)$$

$$= \sum_{i=1}^m \log (p(y^i | x^i; \beta_0, \beta_1))$$

$$= \underbrace{-\frac{m}{2} \log 2\pi - m \log \sigma}_{\text{LEAST-SQUARES ESTIMATE}} - \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^m (y^i - (\beta_0 + \beta_1 x^i))^2}_{\text{FOR SIMPLE LINEAR REGRESSION}}$$

LEAST-SQUARES ESTIMATE
FOR SIMPLE LINEAR REGRESSION

TWO INTERPRETATIONS FOR COEFFICIENTS IN LINEAR REGRESSION:

- GEOMETRIC ; COEFFICIENTS OF THE LINE THAT MINIMIZES SQUARED DISTANCE FROM LINE TO OUR LABELS
- STATISTICAL: COEFFICIENTS GIVE YOU THE MAXIMUM LIKELIHOOD ESTIMATOR FOR A TRAINING SET GENERATED $y \sim N(\beta_0 + \beta_1 x, \sigma^2)$