

## Linear Regression:

- Classification:  $(x, f(x))$   
    • Half-spaces  $\uparrow \{0, 1\}$   
    • Decision trees

• Real-valued labels  $(x, y)$   $y \in \mathbb{R}$ .

$X$  and  $Y$  two random variables

we want to predict the value/label.

we get to see  $X$ .

- we want to predict  $Y$ ; we don't see  $X$ .  
 $(x, y) \sim \mathcal{D}$ .

- Optimal guess for  $y$  is  $E[Y]$ .

- Measure loss by using square-loss:  $(\text{prediction} - y)^2$

- we observe  $X$ , we want to predict  $Y$ .

optimal prediction  $E[Y|X] = f(x)$

Regression Function.

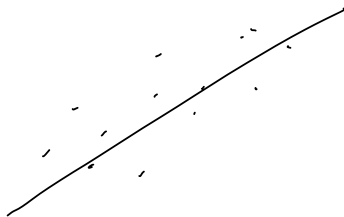
- Obstacle:  $f(x)$  could be unknown,  
or hard to compute.

Linear Regression asks the following question:

\* Given  $X$ , what linear function of  $X$  should we use to predict  $Y$ ?

we want to learn coefficients  $\beta_0$  and  $\beta_1$  such that

$E[(Y - (\beta_0 + \beta_1 X))^2]$  is minimized  
 $(x, y) \sim \mathcal{D}$        $y' = \text{prediction}$



Draw a training set of size  $m$ :

$(x^1, y^1), \dots, (x^m, y^m)$  "simple linear regression"

$$\min_{\beta_0, \beta_1} \underbrace{\frac{1}{m} \sum_{j=1}^m (y^j - (\beta_0 + \beta_1 x^j))^2}_{l = \text{cost function}}.$$

How to find  $\beta_0$  and  $\beta_1$ ?

Take derivative wrt  $\beta_0, \beta_1$ . Set them equal to 0.

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

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

Eliminating -2

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

$$\Rightarrow \beta_0 \text{ in terms of } \beta_1, \bar{y}, \bar{x}: \boxed{\beta_0 = \bar{y} - \beta_1 \bar{x}}$$

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

$$\frac{1}{m} \sum_{j=1}^m (x_j y_j - \beta_0 x_j - \beta_1 (x_j)^2)$$

$\beta_1$  will not involve  $\beta_0$  : substitute  $\beta_0 = \bar{y} - \beta_1 \bar{x}$  :

$$\overline{xy} - \beta_0 \bar{x} - \beta_1 \overline{x^2} = \overline{xy} - (\bar{y} - \beta_1 \bar{x}) \bar{x} - \beta_1 \overline{x^2}$$

$$\overline{xy} - \bar{x} \bar{y} + \beta_1 (\bar{x})^2 - \beta_1 \overline{x^2} = 0$$

$$\beta_1 = \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - (\bar{x})^2} = \boxed{\frac{\text{cov}(x, y)}{\text{var}(x)}} \quad \text{where } \text{cov}(x, y) = E[XY] - E[X]E[Y]$$

This is the simple linear regression. Now let's look at multiple variables/features,  $x_1, \dots, x_n$

$x \in \mathbb{R}^n$ .  $x$  is  $n$  dimensional vector.

$y \in \mathbb{R}$  (scalar).

Fitting a line to  $n$ -dimensional data.

consider  $m$  by  $n$  matrix  $X$ .

$$X = \begin{matrix} m \\ \left[ \begin{array}{c} \\ \\ \end{array} \right] \\ n \end{matrix}$$

- Each row is equal to  $x^i$  drawn from  $D$ ,
- Each column: Each point is in  $\mathbb{R}^n$ .

$\vec{y} \in \mathbb{R}^m \leftarrow$  labels for these  $m$  points ( $m \times 1$ )

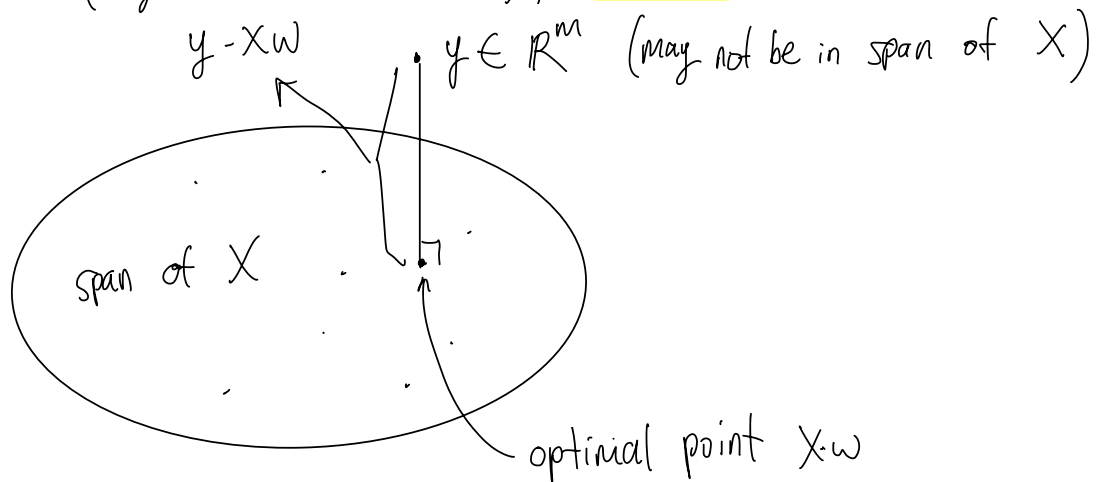
Goal: Find a vector  $\underbrace{w \in \mathbb{R}^n}_{n \times 1} : \min_w \left\| \underbrace{X \cdot w}_{m \times 1} - \vec{y} \right\|_2^2$

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

$$(y - (X_1^1 \omega_1 + \dots + X_n^1 \omega_n))^2$$

$$\min_{\omega} \|X \cdot \vec{\omega} - \vec{y}\|_2^2. \quad \text{How to find } \vec{\omega}?$$

$X \cdot \omega$  is a vector in the span of the columns of  $X$  (e.g. # of features),  $n \times 1$ .



vector  $y - Xw$  is orthogonal to  $X$ :

$$X^T \cdot (y - Xw) = 0. \quad \text{orthogonal} = \text{inner product} = 0.$$

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

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

$$\Rightarrow \boxed{(X^T X)^{-1} X^T y = w} \quad \left\{ \begin{array}{l} \text{Normal} \\ \text{Equations.} \end{array} \right.$$

Issues:

1) what if  $X^T X$  is not invertible?  
"pseudo-inverse".

2) What is the running time for computing  $w$ ?  
crude estimate:  $O(n^3 + mn^2)$ .  
can improve using gradient descent (later),

Maximum Likelihood:

Assume "simple linear regression case".

$X$ ; Assume  $y = \beta_0 + \beta_1 x + \epsilon$   
 $\epsilon$  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  $\beta_0$  and  $\beta_1$ ,  
( $\sigma^2$  is known),

What is the probability that we see  
 $(x^1, y^1), \dots, (x^m, y^m)$ .

Depends on  $\beta_0$  and  $\beta_1$ .

Likelihood Function:

Probability of seeing training set given a choice  
 $\beta_0$  and  $\beta_1$  of our parameters.

$$\prod_{i=1}^m p(y^i | x^i; \beta_0, \beta_1) = \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left[-\frac{(y^i - (\beta_0 + \beta_1 x^i))^2}{2\sigma^2}\right]}_{\text{likelihood of our training set,}} \quad \leftarrow \text{Choose } \beta_0 \text{ and } \beta_1 \text{ that maximize the likelihood.}$$

Gaussian:

Instead of directly maximizing the likelihood,  
we max the log of it.

$$\begin{aligned} \log(L(\beta_0, \beta_1)) &= \log \prod_{i=1}^m p(y^i | x^i; \beta_0, \beta_1) \\ &= \sum_{i=1}^m \log(p(y^i | x^i; \beta_0, \beta_1)) \end{aligned}$$

$$= \underbrace{-\frac{m}{2} \log 2\pi}_{\text{constant}} - \underbrace{m \log \sigma}_{\text{constant}} - \frac{1}{2\sigma^2} \underbrace{\sum_{i=1}^m (y^i - (\beta_0 + \beta_1 x^i))^2}_{\text{least-squares estimate for simple linear regression}}$$

Two interpretations for coefficients in linear regression;

- geometric; coefficients of the line that minimizes squared distance from the line to our labels.
- statistical: coefficients give you the maximum likelihood estimator for a training set generated per the assumption.  
 $y \sim N(\beta_0 + \beta_1 x, \epsilon)$ .

Read ch 9 of ML book.