

## Lecture 1 : Linear Regression I

Example : Predicts level of PSA from various measurements on prostate

$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N) \rightarrow \text{Training Data}$

$y_1, y_2, \dots, y_N$  are known,  $y_i \in \mathbb{R}$

$x_1, x_2, \dots, x_N$  are measurements on prostate

$$x_i \in \mathbb{R}^d$$

Example : Netflix - (user, movie, rating) triplets

Jeff      Sorry to  $\overset{\leftarrow}{\$}$  you  
              bother you

\$1MM Netflix

Goal : Predict  $y$  for a  $x$

Example : Predict whether email is spam or not

$X = \text{set of emails} = \{x_1, x_2, \dots, x_N\}, x_i \in \mathbb{R}^d$

$Y = \{\text{spam, normal}\}$

When  $Y$  is categorical - Classification

When  $Y$  is real-valued (continuous) - Regression

## Regression Problem

Given  $(x_i, y_i)$ ,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ ,  $i = 1, 2, 3, \dots, N$

$$x = \begin{bmatrix} x(1) \\ x(2) \\ \vdots \\ x(d) \end{bmatrix} . \text{ Prediction (linear): } w \in \mathbb{R}$$

$$y(x) = w_0 + w_1 x(1) + w_2 x(2) + \dots + w_d x(d).$$

$$= w_0 + \bar{w}^T x, \bar{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$$\text{where } w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}, \bar{x} = \begin{bmatrix} 1 \\ x(1) \\ x(2) \\ \vdots \\ x(d) \end{bmatrix}, w, \bar{x} \in \mathbb{R}^{d+1}$$

$$\bar{x}_1^T w \approx y_1 \rightarrow (\bar{x}_1^T w - y_1)^2 \text{ is small}$$

$$\bar{x}_2^T w \approx y_2 \rightarrow (\bar{x}_2^T w - y_2)^2 \text{ " " }$$

$$\vdots \quad \vdots \quad \vdots$$

$$\bar{x}_N^T w \approx y_N \rightarrow (\bar{x}_N^T w - y_N)^2 \text{ is small}$$

$$F(w) = \sum_{i=1}^N (\bar{x}_i^T w - y_i)^2 . \text{ Minimize } F(w)$$

$$X = d+1 \begin{bmatrix} \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_N \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1(1) & x_2(1) & \dots & x_N(1) \\ x_1(2) & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_1(d) & x_2(d) & \dots & x_N(d) \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, y \in \mathbb{R}^N$$

$$X \in \mathbb{R}^{(d+1) \times N}$$

$$F(w) = \frac{1}{2} \|X^T w - y\|_2^2 \quad - \text{Least Squares Objective}$$

$$= \frac{1}{2} (X^T w - y)^T (X^T w - y)$$

$$= \frac{1}{2} (w^T X - y^T)(X^T w - y)$$

$$= \frac{1}{2} (w^T X X^T w - 2 y^T X^T w + y^T y)$$

$$\nabla_w F(w) = X X^T w - X y$$

$$w^* = \arg \min F(w)$$

$$\nabla F(w^*) = 0$$

$$X X^T w^* - X y = 0$$

$$(X X^T) w^* = X y$$

$w^* = (X X^T)^{-1} X y$  — ok if  $X X^T$  is non-singular

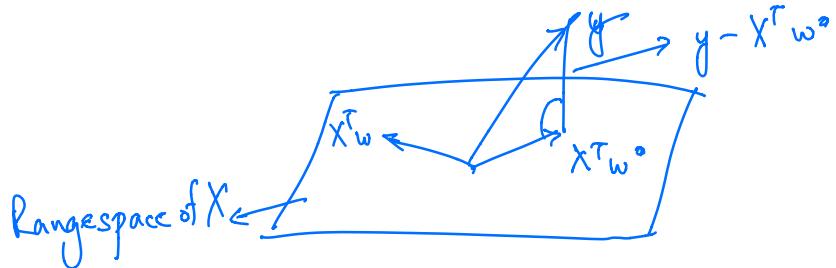
Normal Equations

Geometric View

$$X^T = \begin{bmatrix} 1 & x_1^T \\ 1 & x_2^T \\ \vdots & \vdots \\ 1 & x_N^T \end{bmatrix} = \begin{bmatrix} 1^T \\ x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix}$$

$$X^T w \approx y$$

$X^T w \leftarrow$  Linear Combination of columns of  $X^T$



$$X^T w \perp y - X^T w^* + w$$

$$(X^T w)^T (y - X^T w^*) = 0$$

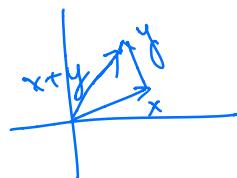
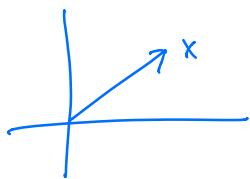
$$\Rightarrow w^T (Xy - X X^T w^*) = 0$$

$$\Rightarrow \boxed{X X^T w^* = Xy}$$

Normal Equations (again) but derived from  
a geometric viewpoint)

Linear Algebra Background.

$$x \in \mathbb{R}^d, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}, \quad x^T = [x_1, x_2, \dots, x_d]$$



Vector Norms  $\|\cdot\| : \mathbb{R}^d \rightarrow \mathbb{R}_+$

1.  $\|x\| \geq 0$  &  $\|x\| = 0 \Leftrightarrow x = 0$  - positivity
2.  $\|\alpha x\| = |\alpha| \|x\|$  - Homogeneity
3.  $\|x+y\| \leq \|x\| + \|y\|$  - Triangle inequality

$$\|x\|_p = \left( |x_1|^p + |x_2|^p + \dots + |x_d|^p \right)^{\frac{1}{p}} \text{ - } L_p \text{ norm}$$

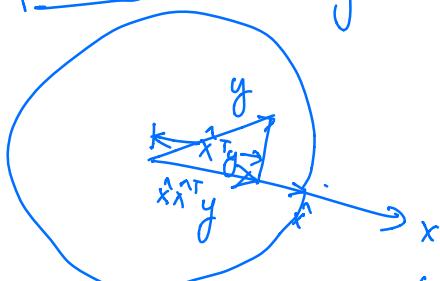
$$\|x\|_1 = |x_1| + |x_2| + \dots + |x_d| \text{ - } L_1 \text{ norm}$$

$$\|x\|_2 = \sqrt{|x_1|^2 + |x_2|^2 + \dots + |x_d|^2} \text{ - } L_2 \text{ norm}$$

$$\|x\|_\infty = \max_{1 \leq i \leq d} |x_i| \text{ - } L_\infty \text{ norm}$$

$$\hat{x} = \frac{x}{\|x\|_2} \text{ - unit vector in direction of } x$$

$\begin{bmatrix} \hat{x} & \hat{x}^\top \end{bmatrix} \rightarrow$   $d \times 1$  matrix  
 $\rightarrow$  orthogonal projector onto the vector  $x$  (or  $\hat{x}$ )



$$(\hat{x} \hat{x}^\top) y = \underbrace{\hat{x}}_{\substack{\text{unit vector}}} (\underbrace{\hat{x}^\top y}_{\substack{\text{scalar}}})$$

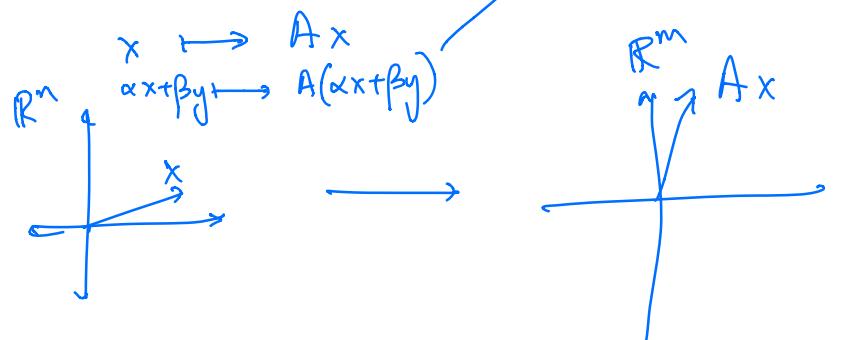
$$\hat{x}^\top y = \left( \frac{x}{\|x\|} \right)^\top y = \frac{x^\top y}{\|x\|}$$

$$-1 \leq \cos \theta = \frac{x^\top y}{\|x\| \cdot \|y\|} \leq 1 \Rightarrow |x^\top y| \leq \|x\| \cdot \|y\| \quad (\text{Cauchy-Schwarz Inequality})$$

$$x \perp y \Leftrightarrow x^T y = 0$$

$A \in \mathbb{R}^{m \times n}$

$$A: \mathbb{R}^n \rightarrow \mathbb{R}^m$$



Matrices Norms,  $A \in \mathbb{R}^{m \times n}$ ,  $\|\cdot\|: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}_+$

$$1. \|A\| \geq 0 \quad \& \quad \|A\| = 0 \Leftrightarrow A = 0$$

$$2. \|\alpha A\| = |\alpha| \cdot \|A\|$$

$$3. \|A+B\| \leq \|A\| + \|B\|$$

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |A_{ij}|^2} \rightarrow \text{Frobenius Norm}$$

$$\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \max_{\|x\|_2 \leq 1} \|Ax\|_2 = \sigma_1$$

↓  
maximum singular value of A

$$\|A\|_1 = \max_{x \neq 0} \frac{\|Ax\|_1}{\|x\|_1} = \max_{\|x\|_1 \leq 1} \|Ax\|_1$$



$$\|A\|_{\infty} = \max_{\|x\|_{\infty} \leq 1} \|Ax\|_{\infty}$$

$$\|A\|_{p,q} = \max_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_q}$$

$$\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sigma_1 \text{ (max sing. value)}$$

$$\min_{x \neq 0} \frac{\|Ax\|}{\|x\|_2} = \sigma_n \text{ minimum singular value of } A$$

$\kappa(A)$  - condition number of  $A$

$$\kappa(A) = \|A\| \|A^{-1}\|$$

$$\kappa_2(A) = \frac{\sigma_1}{\sigma_n} = (\|A\|_2 \|A^{-1}\|_2)$$

Eigenvalues / Eigenvectors

$$A\bar{x} = \lambda \bar{x}, \quad \bar{x} \neq 0, \quad A \in \mathbb{R}^{n \times n}$$

↓  
eigenvalue

eigenvector ( $\|\bar{x}\|_2 = 1$ )

$$A\bar{x} - \lambda \bar{x} = 0$$

$$(A - \lambda I)\bar{x} = 0 \Rightarrow A - \lambda I \text{ is singular}$$

$$\Rightarrow \det(A - \lambda I) = 0$$

$\det \begin{pmatrix} a_{11} - \lambda & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda & \dots & a_{2n} \\ \vdots & & \ddots & \\ a_{n1} & \dots & a_{n(n-1)} - \lambda & \end{pmatrix}$  is polynomial of degree  $n$  in  $\lambda$

If  $A$  is  $n \times n$ , it has  $n$  eigenvalues

If  $A = A^T$  (symmetric)  $\Rightarrow$  all eigenvalues are real and it has a full set ( $n$ ) of mutually orthogonal eigenvectors

$$A = A^T$$

$$Av_1 = v_1 \lambda_1$$

$$Av_2 = v_2 \lambda_2$$

:

$$\underbrace{Av_n = v_n \lambda_n}_{\lambda_i \in \mathbb{R}}$$

$$\lambda_i \in \mathbb{R}$$

$$\boxed{v_i^T v_j = 0, i \neq j}$$

$$\boxed{v_i^T v_i = 1}$$

$$\Downarrow V = [v_1 \ v_2 \ \dots \ v_n]$$

$$V^T V = I = V V^T$$

$$A[v_1 \ \dots \ v_n] = [v_1 \ \dots \ v_n] \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix}$$

$$AV = V\Lambda$$

$$, \quad \Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix}$$

$$\boxed{A = V\Lambda V^T}$$

Eigenvalue Decomposition

$A$  is positive semi-definite if  $x^T A x \geq 0 \forall x \neq 0$  ( $A \succeq 0$ )

$A$  is positive definite if  $x^T A x > 0 \forall x \neq 0$  ( $A \succ 0$ )

$$A = A^T, \quad x^T A x = (x^T V) D (V^T x)$$

Let  $z = V^T x$ , then  $x^T A x = z^T D z = \sum_{i=1}^n \lambda_i z_i^2$

$$x^T A x \geq 0 \Leftrightarrow \lambda_i \geq 0$$

$$A = Q D Q^T, \quad x^T A x = x^T Q D Q^T x = z^T z \geq 0 \quad (z = Q^T x)$$