

# Lecture 1 : Linear Regression

Example 1: Predict levels of PSA from various measurements on the prostate

Training Data:  $(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4) \dots, (x_N, y_N)$

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

$x_1, x_2, \dots, x_N$  are measurements on prostate,  
 $x_i \in \mathbb{R}^d$

Example 2: Netflix - (user, movie, rating)

$x \leftarrow (\text{Indecrypt}, \text{Breaking Bad}, \underbrace{\text{Ep 1, Season 1}}_{y_i})$

\$1M Netflix Prize

Goal: Predict  $y$  for a new  $x$ ,  $y \in \mathbb{R}$  (regression)

Example 3: Predict whether an 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  $\rightarrow$  Classification

When  $Y$  is real-valued  $\rightarrow$  Regression

## Regression Problem

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

$$x = \begin{bmatrix} x(1) \\ x(2) \\ \vdots \\ x(d) \end{bmatrix}$$

Prediction (linear):

$$\begin{aligned} y(x) &= w_0 + w_1 x(1) + w_2 x(2) + \dots + w_d x(d) \\ &= w_0 + \bar{w}^T x, \quad \bar{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \\ &= w_0 \cdot 1 + \bar{w}^T x \end{aligned}$$

$$= \mathbf{w}^\top \bar{\mathbf{x}} , \quad \mathbf{x} = \begin{bmatrix} 1 \\ x(1) \\ x(2) \\ \vdots \\ x(d) \end{bmatrix}, \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$$\mathbf{w}, \bar{\mathbf{x}} \in \mathbb{R}^{d+1}$$

$$y(x_i) = \mathbf{w}^\top \bar{\mathbf{x}}_i \approx y_i \rightarrow (\mathbf{w}^\top \bar{\mathbf{x}}_i - y_i)^2 \text{ is small}$$

$$y(x_2) = \mathbf{w}^\top \bar{\mathbf{x}}_2 \approx y_2 \rightarrow (\mathbf{w}^\top \bar{\mathbf{x}}_2 - y_2)^2 \text{ is small}$$

$$y(x_3) = \mathbf{w}^\top \bar{\mathbf{x}}_3 \approx y_3 \rightarrow (\mathbf{w}^\top \bar{\mathbf{x}}_3 - y_3)^2 \text{ is small}$$

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

$$y(x_N) = \mathbf{w}^\top \bar{\mathbf{x}}_N \approx y_N \rightarrow (\mathbf{w}^\top \bar{\mathbf{x}}_N - y_N)^2 \text{ is small}$$

$$F(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^\top \bar{\mathbf{x}}_i - y_i)^2$$

Goal : Minimize  $F(\mathbf{w})$  - Least Squares Regression

$$\mathbf{X} = \underbrace{\begin{bmatrix} \bar{\mathbf{x}}_1 & \bar{\mathbf{x}}_2 & \dots & \bar{\mathbf{x}}_N \end{bmatrix}}_{\mathbb{R}^{d+1 \times N}}, \quad \mathbf{x} \in \mathbb{R}^{(d+1) \times N}$$

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

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1(1) & x_2(1) & \dots & x_N(1) \\ x_1(2) & x_2(2) & \dots & x_N(2) \\ \vdots & \vdots & \ddots & \vdots \\ x_1(d) & x_2(d) & \dots & x_N(d) \end{bmatrix}, \quad \mathbf{x} \in \mathbb{R}^{(d+1) \times N}$$

$$\mathbf{X} \in \mathbb{R}^{d+1 \times N}, \quad \mathbf{y} \in \mathbb{R}^N, \quad \mathbf{w} \in \mathbb{R}^{d+1}$$

$$F(\mathbf{w}) = \frac{1}{2} \| \mathbf{X}^\top \mathbf{w} - \mathbf{y} \|_2^2$$

$$\mathbf{X}^\top \mathbf{w} \in \mathbb{R}^{d+1} \approx \mathbf{y}(\mathbf{x})$$

Least Squares Objective

Goal: Find  $\mathbf{w}$  such that  $F(\mathbf{w})$  is minimized

$$\|z\|_2^2 = z^T z \quad a^T b = b^T a$$

$$\begin{aligned}
 F(\omega) &= \frac{1}{2} (x^T \omega - y)^T (x^T \omega - y) \\
 &= \frac{1}{2} (\omega^T x - y^T)(x^T \omega - y) \\
 &= \frac{1}{2} (\omega^T x x^T \omega - \underline{y^T x^T \omega} - \underline{\omega^T x y} + y^T y) \\
 &= \frac{1}{2} (\omega^T x x^T \omega - 2 y^T x^T \omega + y^T y)
 \end{aligned}$$

$$\nabla_{\omega} F(\omega) = \frac{1}{2} (2 x x^T \omega - 2 x y) = x x^T \omega - x y$$

$$\omega^* = \arg \min_{\omega} F(\omega)$$

$$\nabla_{\omega} F(\omega^*) = 0$$

$$x x^T \omega^* = x y$$

d+1 linear equations

d+1 unknowns

$$\omega^* = (x x^T)^{-1} x y$$

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

$$x x^T \in \mathbb{R}^{(d+1) \times (d+1)}$$

$$\omega \in \mathbb{R}^{d+1}$$

- ok if  $x x^T$  is non-singular

Normal Equations

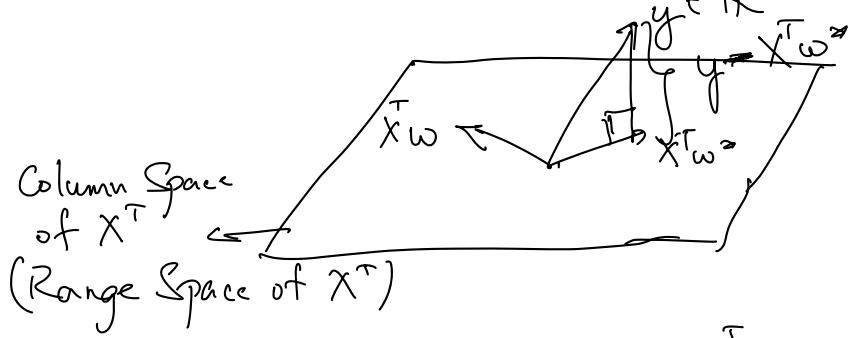
Geometric View

$$x^T = \begin{bmatrix} 1 & x_1^T \\ 1 & x_2^T \\ \vdots & \vdots \\ 1 & x_N^T \end{bmatrix} = \begin{bmatrix} 1 \\ x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} \quad x^T \in \mathbb{R}^{N \times (d+1)}$$

$$x^T \omega \approx y$$

$$\begin{aligned}
 x^T \omega &\leftarrow \text{Linear combination of columns of } x^T \\
 &= \omega_0 \times \text{col1} + \omega_1 \times \text{col2} + \dots + \omega_d \times \text{col}(d+1)
 \end{aligned}$$

$$\in \mathbb{R}^N$$



Column Space  
of  $X^T$   
(Range Space of  $X^T$ )

$$X^T w \perp y - X^T w^* \quad \forall w \quad (\text{to minimize } F(w))$$

$w^* = \text{argmin } F(w)$

$$a \perp b \quad a^\top b = 0$$

$$(X^T w)^\top (y - X^T w^*) = 0 \quad \forall w$$

$$w^\top X(y - X^T w^*) = 0 \quad \forall w$$

$$w^\top (Xy - \underbrace{XX^T w^*}_{0}) = 0 \quad \forall w$$

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

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

$d+1$  equations in  $d+1$  unknowns ( $w_0, w_1, \dots, w_d$ )

Solve  $XX^T w^* = Xy$  - linear system of equations

How?

① Do Gaussian Elimination + solve  $\underline{XX^T w^*} = \underline{Xy}$

- Form the matrix  $\underline{XX^T} = \underline{A}$  -  $O(d^2N)$  operations

- Form the right hand side  $Xy = b$  -  $O(dN)$

$$Aw^* = b$$

$A = L U$  decomposition (Gaussian Elimination)

$\curvearrowleft A = XX^T$  - positive semi-definite matrix

$\curvearrowright A = LL^T$  (Cholesky Decomposition) -  $L$  is lower triangular

$\curvearrowright O(d^3)$  operations

$L^T$  is upper triangular

$$LL^T w^* = b$$

Solve  $L_2 = b$        $\boxed{L} = \boxed{I}$       Forward Substitution  
 $L^T w^* = z$        $\boxed{L^T} = \boxed{I}$       Backward Substitution  
 $O(d^2)$   
 $O(d^2)$

Solve normal equations by Gaussian Elimination  
in  $O(d^2N + d^3)$  operations

Inverse ~~of~~ of  $A (x x^T)$  exists only if  $A$  is non-singular  
Condition Number of  $A$  is large  $\equiv A$  is close to singularity  
Solving normal equations can yield large error when  
 $A$  is poorly conditioned

Other methods which are better when  $A$  is close to  
singular:

① Use QR decomposition of  $X$

XXXXX ② Use SVD of  $X$ .

↳ Singular value decomposition

↳ Most accurate but more expensive to compute.