

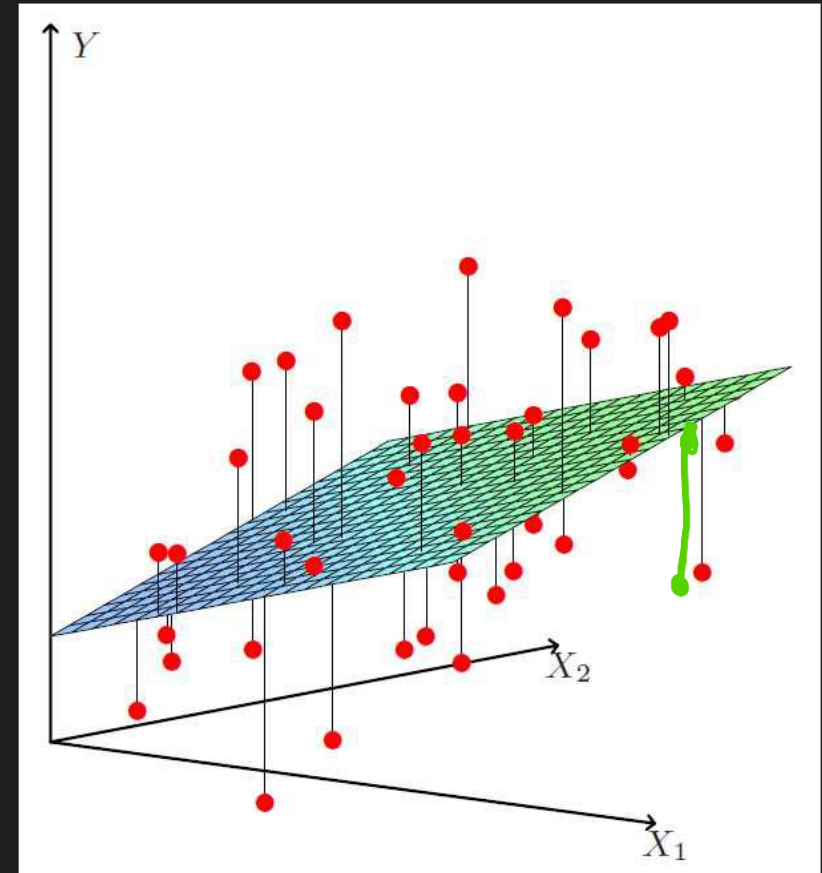
INTRODUCTION TO
MACHINE LEARNING
COMPSCI 4ML3

LECTURE 3

HASSAN ASHTIANI

ORDINARY LEAST SQUARES (D-DIMENSIONS)

- ASSUME $x \in \mathbb{R}^d$, $y \in \mathbb{R}$
- INSTEAD OF A LINE, WE NEED TO FIT A HYPERPLANE!
- WHY ARE THE LINES VERTICAL?
 - ANY DIFFERENT IF WE MINIMIZE THE DISTANCE TO THE HYPERPLANE?



$$\sum (\Delta_i)^2$$

MATRIX FORM OLS

$$\bullet \Delta = \begin{pmatrix} \Delta_1 \\ \dots \\ \Delta_n \end{pmatrix} = \begin{pmatrix} x_1^1 & \dots & x_d^1 \\ \vdots & \ddots & \vdots \\ x_1^n & \dots & x_d^n \end{pmatrix} \begin{pmatrix} w_1 \\ \dots \\ w_d \end{pmatrix} - \begin{pmatrix} y^1 \\ \dots \\ y^n \end{pmatrix} \approx \begin{matrix} \hat{Y} \\ \downarrow \\ Y \end{matrix}_{n \times 1} \approx \begin{pmatrix} y_1^1 - y_1^1 \\ \vdots \\ y_1^n - y_1^n \end{pmatrix}$$

$$\min_{W \in \mathbb{R}^{d \times 1}} \sum_{i=1}^n (\Delta_i)^2 = \min_{W \in \mathbb{R}^{d \times 1}} \langle \Delta, \Delta \rangle = \min_{W \in \mathbb{R}^{d \times 1}} \|\Delta\|_2^2$$

Handwritten annotations: $X_{n \times d}$, $W_{d \times 1}$, $Y_{n \times 1}$

$$\min_{W \in \mathbb{R}^{d \times 1}} \|XW - Y\|_2^2$$

TAKING THE "DERIVATIVE"

- REAL-VALUED FUNCTION OF A VECTOR

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

GRADIENT:

$$\nabla_w f = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}$$

$$f(w) =$$

- VECTOR-VALUED FUNCTION OF A VECTOR

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

JACOBIAN:

$$\nabla_w g = \begin{bmatrix} \frac{\partial g_1}{\partial w_1} & \dots & \frac{\partial g_1}{\partial w_m} \\ \vdots & & \vdots \\ \frac{\partial g_n}{\partial w_1} & \dots & \frac{\partial g_n}{\partial w_m} \end{bmatrix}$$

$n \times m$

unlike gradients,
the columns are
input variables
not the rows.

MATRIX/VECTOR CALCULUS

- $u, v \in \mathbb{R}^n$

- $g(u) = u^T v$

$$g: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\boxed{u_1 v_1} + u_2 v_2 + \dots$$

- $\nabla u(g) =$
gradient

$$\begin{bmatrix} \frac{\partial g}{\partial u_1} \\ \vdots \\ \frac{\partial g}{\partial u_n} \end{bmatrix}$$

=

$$\begin{bmatrix} \frac{\partial \sum u_i v_i}{\partial u_1} \\ \vdots \\ \frac{\partial \sum u_i v_i}{\partial u_n} \end{bmatrix}$$

=

$$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$$

= v

MATRIX/VECTOR CALCULUS

- $A \in \mathbb{R}^{m \times n}, u \in \mathbb{R}^n$
- $g(u) = \underline{Au}$

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

$$A = \begin{bmatrix} a_{1,1} & \dots & \dots \\ \vdots & & \vdots \\ a_{m,1} & \dots & a_{m,n} \end{bmatrix} \vec{u}$$

- $\nabla u(g) = A_{m \times n}$

Jacobian

$$\frac{\partial g_i}{\partial u_j} = \frac{\partial \left(\sum_{k=1}^n a_{i,k} u_k \right)}{\partial u_j} = a_{i,j}$$

MATRIX/VECTOR CALCULUS

- $A \in \mathbb{R}^{n \times n}, u \in \mathbb{R}^n$

- $g(u) = \underline{u^T A u}$ $g: \mathbb{R}^n \rightarrow \mathbb{R}$

- $\nabla u(g) = \underbrace{u^T}_{1 \times n} \underbrace{(A + A^T)}_{n \times n}$

Jacobian

show this as
an exercise,

SOLVING OLS

$f(W) =$ $\|XW - Y\|_2^2$. WHAT IS ∇f ?

$$\frac{\partial \|XW - Y\|_2^2}{\partial W} = \frac{\partial [(XW - Y)^T (XW - Y)]}{\partial W}$$

$$\frac{\partial [(W^T X^T - Y^T)(XW - Y)]}{\partial W} = \frac{\partial W^T X^T X W}{\partial W} + \frac{\partial Y^T X W}{\partial W}$$

$$= \frac{\partial W^T X^T X W}{\partial W} + 0$$

↻ ↻

↗ 0

$$\frac{\partial Y^T X W}{\partial W}$$

$Y^T X W$
Scalar

~~$\frac{\partial Y^T X W}{\partial W}$~~

$$\Rightarrow w^T (x^T x + x^T x) + 0 - 2 Y^T x = 0$$

$$\Rightarrow 2 x^T x w = 2 x^T Y$$

* if $x^T x$ is invertible then

$$\boxed{\bar{w} = (x^T x)^{-1} x^T Y}$$

SOLVING OLS

$$W^{LS} = (X^T X)^{-1} X^T Y$$

- DEGENERATE CASE WHEN $X^T X$ IS NOT INVERTIBLE?

BIAS/INTERCEPT TERM

- WE ARE MISSING THE BIAS TERM (w_0)

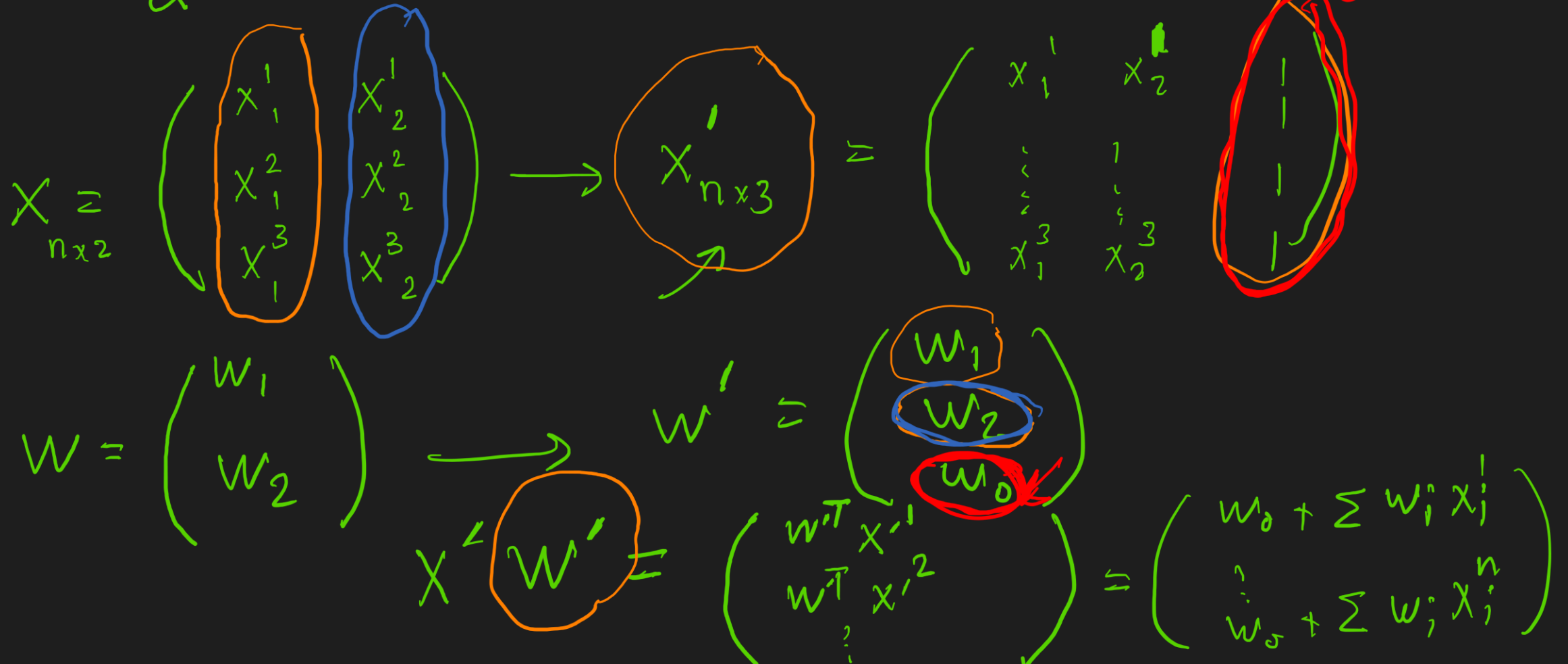
$$\text{MIN}_{w_0, w_1, \dots, w_d \in \mathbb{R}} \sum_{i=1}^n (w_1 x_1^i + \dots + w_d x_d^i + \underline{w_0} - y^i)^2$$

- MATRIX FORM WITH THE BIAS TERM?

$$\text{MIN}_{W \in \mathbb{R}^{d \times 1}, w_0 \in \mathbb{R}} \|XW + \begin{pmatrix} w_0 \\ w_0 \\ \dots \\ w_0 \end{pmatrix} - Y\|_2^2$$

EXAMPLE

$$\hat{y} = \omega_0 + \omega_1 x_1 + \omega_2 x_2$$



BIAS/INTERCEPT TERM

- ADD A NEW AUXILIARY DIMENSION TO THE DATA

- $X'_{n \times (d+1)} = \begin{pmatrix} x_1^1 & \dots & x_d^1 & 1 \\ \vdots & \ddots & \vdots & 1 \\ x_1^n & \dots & x_d^n & 1 \end{pmatrix}, W'_{(d+1) \times 1} = \begin{pmatrix} w_1 \\ \dots \\ w_d \\ w_0 \end{pmatrix}$

- SOLVE OLS: $\min_{W' \in \mathbb{R}^{(d+1) \times 1}} \|X'W' - Y\|_2^2$

- w_0 WILL BE THE BIAS TERM!

test point
 $x = (x_1, x_2)$

$d=2$, you

know w'
 $\hat{y} = ?$ $x' = (x_1, x_2, 1)$
 $\hat{y} = w'^T x'$

SOME EXAMPLES

- OLS NOTEBOOK