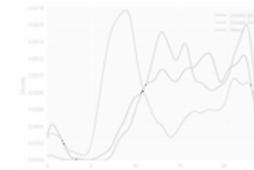
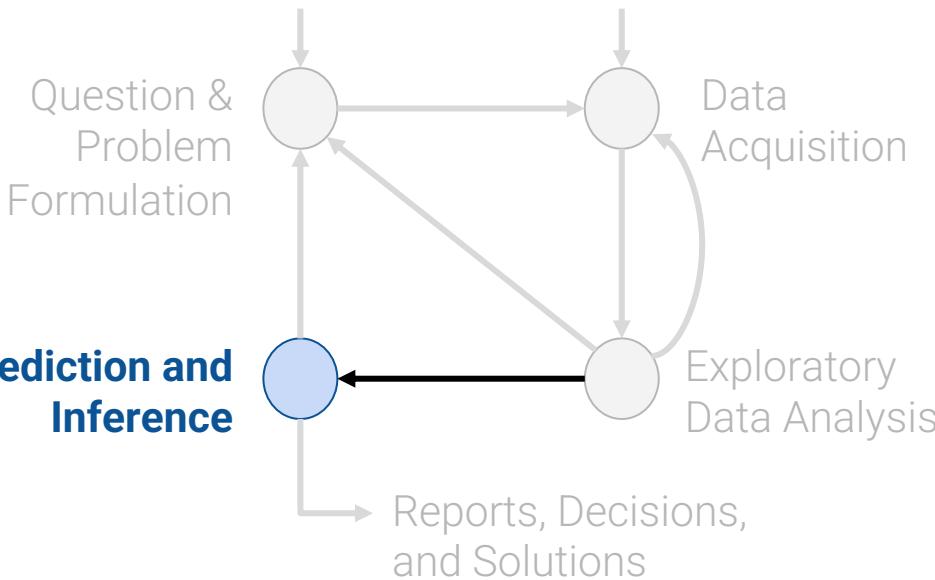


LECTURE 11

# Ordinary Least Squares

Using linear algebra to derive the multiple linear regression model.

# Plan for Next Few Lectures: Modeling



Modeling I:  
Intro to Modeling, Simple  
Linear Regression



Modeling II:  
Different models, loss  
functions, linearization



Modeling III:  
Multiple Linear  
Regression

(today)

# Today's Roadmap

---

## OLS Problem Formulation

- Multiple Linear Regression Model
- Mean Squared Error

Geometric Derivation

Performance: Residuals, Multiple R<sup>2</sup>

OLS Properties

- Residuals
- The Bias/Intercept Term
- Existence of a Unique Solution

# Multiple Linear Regression Model

---

## OLS Problem Formulation

- **Multiple Linear Regression Model**
- Mean Squared Error

Geometric Derivation

Performance: Residuals, Multiple R<sup>2</sup>

OLS Properties

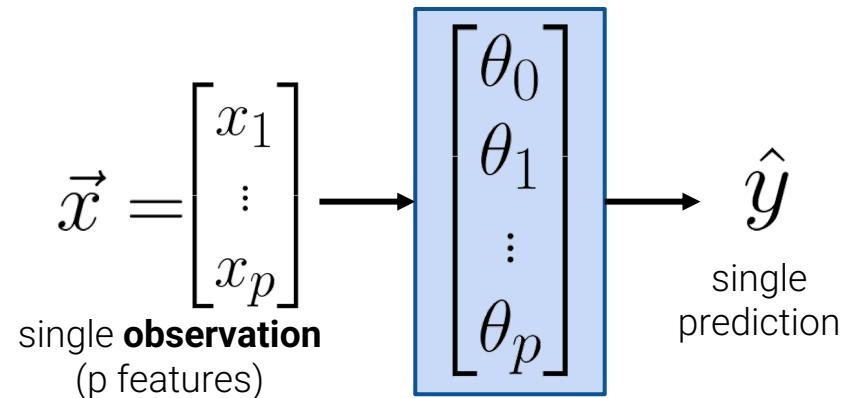
- Residuals
- The Bias/Intercept Term
- Existence of a Unique Solution

## Multiple Linear Regression

Define the **multiple linear regression** model:

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p$$

Predicted  
value of  $y$



# NBA 2018-2019 Dataset

How many points does an athlete score per game?

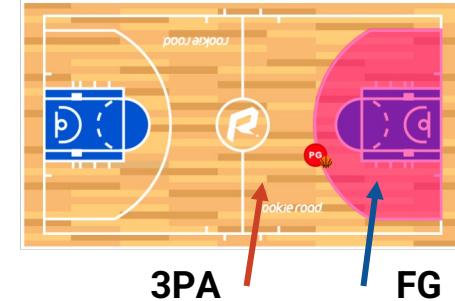
**PTS** (average points/game)

To name a few factors:

- **FG**: average # 2 point field goals
- **AST**: average # of assists
- **3PA**: average # 3 point field goals attempted

	FG	AST	3PA	PTS
1	1.8	0.6	4.1	5.3
2	0.4	0.8	1.5	1.7
3	1.1	1.9	2.2	3.2
4	6.0	1.6	0.0	13.9
5	3.4	2.2	0.2	8.9
6	0.6	0.3	1.2	1.7

Rows correspond to individual players.



**assist**: a pass to a teammate that directly leads to a goal

# Multiple Linear Regression Model

How many points does an athlete score per game?

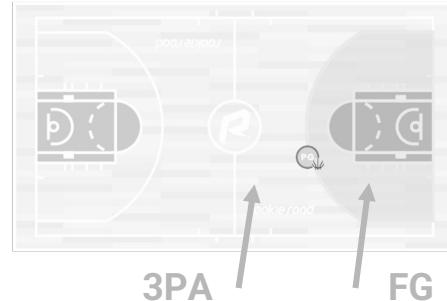
**PTS** (average points/game)

To name a few factors:

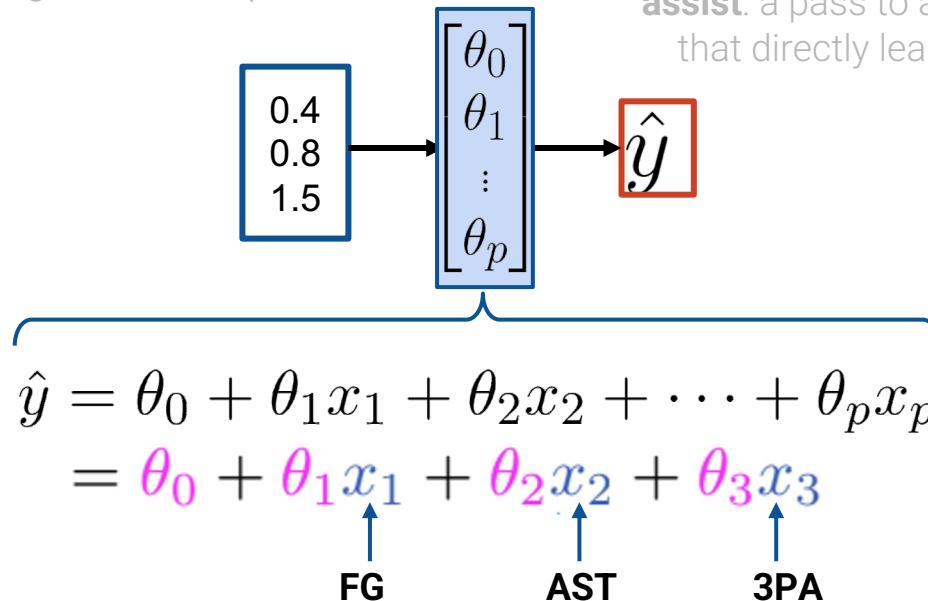
- **FG**: average # 2 point field goals
- **AST**: average # of assists
- **3PA**: average # 3 point field goals attempted

	FG	AST	3PA	PTS
1	1.8	0.6	4.1	5.3
2	0.4	0.8	1.5	1.7
3	1.1	1.9	2.2	3.2
4	6.0	1.6	0.0	13.9
5	3.4	2.2	0.2	8.9
6	0.6	0.3	1.2	1.7

Rows correspond to individual players.



**assist**: a pass to a teammate that directly leads to a goal



# Today's Goal: Ordinary Least Squares

1. Choose a model

2. Choose a loss function

3. Fit the model

4. Evaluate model performance

## Multiple Linear Regression

L2 Loss

## Mean Squared Error (MSE)

Minimize average loss with ~~calculus~~ geometry

Visualize,  
~~Root MSE~~  
Multiple R<sup>2</sup>



In statistics, this model + loss is called **Ordinary Least Squares (OLS)**.

The solution to OLS are the minimizing loss for parameters  $\hat{\theta}$ , also called the **least squares estimate**.

# Today's Goal: Ordinary Least Squares

## 1. Choose a model

Multiple Linear Regression

## 2. Choose a loss function

L2 Loss

Mean Squared Error  
(MSE)

## 3. Fit the model

Minimize average loss with ~~calculus~~ geometry

## 4. Evaluate model performance

Visualize,  
~~Root MSE~~  
Multiple R<sup>2</sup>

For each of our  $n$  data points:

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p$$



$$\hat{\mathbf{Y}} = \mathbf{X}\boldsymbol{\theta}$$

Linear Algebra!!

## From one feature to many features

Dataset  
for  
GLR  $y$

$x_1$	$y_1$
$x_2$	$y_2$
$\vdots$	$\vdots$
$x_n$	$y_n$

Dataset for  
Constant Model

$y$
$y_1$
$y_2$
$\vdots$
$y_n$

Dataset for Multiple Linear  
Regression

$x_{:,1}$	$x_{:,2}$	$\dots$	$x_{:,p}$	$y$
$x_{11}$	$x_{12}$	$\dots$	$x_{1p}$	$y_1$
$x_{21}$	$x_{22}$	$\dots$	$x_{2p}$	$y_2$
$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$x_{n1}$	$x_{n2}$	$\dots$	$x_{np}$	$y_n$

	FG	PTS
1	1.8	5.3
2	0.4	1.7
3	1.1	3.2
4	6.0	13.9
5	3.4	8.9
...	...	...

	PTS
1	5.3
2	1.7
3	3.2
4	13.9
5	8.9
...	...

	FG	AST	3PA	PTS
1	1.8	0.6	4.1	5.3
2	0.4	0.8	1.5	1.7
3	1.1	1.9	2.2	3.2
4	6.0	1.6	0.0	13.9
5	3.4	2.2	0.2	8.9
...	...	...	...	...

## From one feature to many features

Dataset for Multiple Linear Regression

$x_{:1}$	$x_{:2}$	$\dots$	$x_{:p}$	$y$
$x_{11}$	$x_{12}$	$\dots$	$x_{1p}$	$y_1$
$x_{21}$	$x_{22}$	$\dots$	$x_{2p}$	$y_2$
$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$x_{n1}$	$x_{n2}$	$\dots$	$x_{np}$	$y_n$

Feature 2  
 $\{x_{12}, x_{22}, \dots, x_{n2}\}$

Observation i  
 $\{x_{i1}, x_{i2}, \dots, x_{ip}, y_i\}$

	FG	AST	3PA	PTS
1	1.8	0.6	4.1	5.3
2	0.4	0.8	1.5	1.7
3	1.1	1.9	2.2	3.2
4	6.0	1.6	0.0	13.9
5	3.4	2.2	0.2	8.9
...	...	...	...	...

Model

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p$$

$$\begin{cases} \hat{y}_1 = \theta_0 + \theta_1 x_{11} + \theta_2 x_{12} + \dots + \theta_p x_{1p} \\ \hat{y}_2 = \theta_0 + \theta_1 x_{21} + \theta_2 x_{22} + \dots + \theta_p x_{2p} \\ \vdots \\ \hat{y}_n = \theta_0 + \theta_1 x_{n1} + \theta_2 x_{n2} + \dots + \theta_p x_{np} \end{cases}$$

The **dot product (or inner product)** is a vector operation that

- can only be carried out on two vectors of the **same length**
- sums up the products of the corresponding entries of the two vectors, and
- returns a single number

$$\vec{u} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad \vec{v} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$\begin{aligned} \boxed{\vec{u} \cdot \vec{v}} &= \vec{u}^\top \vec{v} = \vec{v}^\top \vec{u} \\ &= 1 \cdot 1 + 2 \cdot 1 + 3 \cdot 1 \\ &= 6 \end{aligned}$$

"u dot v"

Sidenote (not in scope): we can interpret dot product geometrically:

- It is the product of three things: the **magnitude** of both vectors, and the **cosine** of the angles between them.  $\vec{u} \cdot \vec{v} = ||\vec{u}|| \cdot ||\vec{v}|| \cdot \cos \theta$
- Another interpretation: [3Blue1Brown](#)

## Vector Notation

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p$$



This part looks a little like a dot product...

$$= \boxed{\theta_0} + \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} = \theta_0 \cdot 1 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p$$

🤔 What about  
this one???

We want to collect  
all the  $\theta_i$ 's into a  
single vector

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p$$

$$= \theta_0 \cdot 1 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p$$

$$= \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} \cdot \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} = \mathbf{x}^\top \boldsymbol{\theta}$$

bias term,  
intercept term

We want to collect  
all the  $\theta_i$ 's into a  
single vector

## Matrix Notation

---

$$\begin{cases} \hat{y}_1 = \theta_0 + \theta_1 x_{11} + \theta_2 x_{12} + \cdots + \theta_p x_{1p} \\ \hat{y}_2 = \theta_0 + \theta_1 x_{21} + \theta_2 x_{22} + \cdots + \theta_p x_{2p} \\ \vdots \\ \hat{y}_n = \theta_0 + \theta_1 x_{n1} + \theta_2 x_{n2} + \cdots + \theta_p x_{np} \end{cases}$$

$$\begin{cases} \hat{y}_1 = \mathbf{x}_1^\top \boldsymbol{\theta} \quad \text{where } \mathbf{x}_1^\top = [1 \quad x_{11} \quad x_{12} \quad \dots \quad x_{1p}] \text{ is datapoint/observation 1} \\ \hat{y}_2 = \mathbf{x}_2^\top \boldsymbol{\theta} \quad \text{where } \mathbf{x}_2^\top = [1 \quad x_{21} \quad x_{22} \quad \dots \quad x_{2p}] \text{ is datapoint/observation 2} \\ \vdots \\ \hat{y}_n = \mathbf{x}_n^\top \boldsymbol{\theta} \quad \text{where } \mathbf{x}_n^\top = [1 \quad x_{n1} \quad x_{n2} \quad \dots \quad x_{np}] \text{ is datapoint/observation n} \end{cases}$$

## Matrix Notation

$$\left\{ \begin{array}{l} \hat{y}_1 = x_1^\top \theta \quad \text{where } x_1^\top = [1 \quad x_{11} \quad x_{12} \quad \dots \quad x_{1p}] \text{ is datapoint/observation 1} \\ \hat{y}_2 = x_2^\top \theta \quad \text{where } x_2^\top = [1 \quad x_{21} \quad x_{22} \quad \dots \quad x_{2p}] \text{ is datapoint/observation 2} \\ \vdots \\ \hat{y}_n = x_n^\top \theta \quad \text{where } x_n^\top = [1 \quad x_{n1} \quad x_{n2} \quad \dots \quad x_{np}] \text{ is datapoint/observation n} \end{array} \right.$$

	FG	AST	3PA	PTS
1	1.8	0.6	4.1	5.3
2	0.4	0.8	1.5	1.7
3	1.1	1.9	2.2	3.2
4	6.0	1.6	0.0	13.9
5	3.4	2.2	0.2	8.9
...	...	...	...	...

For data point/observation 2, we have

$$x_2 = \begin{bmatrix} 1 \\ 0.4 \\ 0.8 \\ 1.5 \end{bmatrix} \quad y_2 = 1.7 \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$\begin{aligned} \hat{y}_2 &= x_2^\top \theta \\ &= \theta_0 + \theta_1 \cdot 0.4 + \theta_2 \cdot 0.8 + \theta_3 \cdot 1.5 \end{aligned}$$

### Dimension check

$$x_2 \in \mathbb{R}^4 \text{ or } \mathbb{R}^{(p+1)}$$

$$\theta \in \mathbb{R}^4 \text{ or } \mathbb{R}^{(p+1)}$$

$$y_2 \in \mathbb{R} \quad \hat{y}_2 \in \mathbb{R}$$

also called scalars

$$\hat{y}_1 = [1 \ x_{11} \ x_{12} \ \dots \ x_{1p}]$$

$$\hat{y}_2 = [1 \ x_{21} \ x_{22} \ \dots \ x_{2p}]$$

$\vdots$                      $\vdots$

$$\hat{y}_n = [1 \ x_{n1} \ x_{n2} \ \dots \ x_{np}]$$

$$\theta = x_1^T \theta$$

$$\theta = x_2^T \theta$$

$\vdots$

$$\theta = x_n^T \theta$$

**n** row vectors, each  
with dimension **(p+1)**

Expand out each datapoint's  
(transposed) input

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \theta$$

**n** row vectors, each  
with dimension **(p+1)**

Vectorize predictions and parameters  
to encapsulate all n equations into a  
single matrix equation.

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \mathbf{X} \theta$$

**Design matrix** with  
dimensions  $n \times (p + 1)$

# The Design Matrix $\mathbb{X}$

We can use linear algebra to represent our predictions of all  $n$  data points at once.

One step in this process is to stack all of our input features together into a **design matrix**:

$$\mathbb{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

What do the **rows** and **columns** of the design matrix represent in terms of the observed data?



Field Goals  
Assists  
3-Point  
Attempts

Bias	FG	AST	3PA	PTS
1	1.8	0.6	4.1	5.3
1	0.4	0.8	1.5	1.7
1	1.1	1.9	2.2	3.2
1	6.0	1.6	0.0	13.9
1	3.4	2.2	0.2	8.9
...	...	...	...	...
1	4.0	0.8	0.0	11.5
1	3.1	0.9	0.0	7.8
1	3.6	1.1	0.0	8.9
1	3.4	0.8	0.0	8.5
1	3.8	1.5	0.0	9.4

Example design matrix  
708 rows x (3+1) cols

# The Design Matrix $\mathbf{X}$

We can use linear algebra to represent our predictions of all  $n$  data points at once.

One step in this process is to stack all of our input features together into a **design matrix**:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

A **column** corresponds to a **feature**,  
e.g. feature 1 for all  $n$  data points

Special all-ones feature often  
called the **bias/intercept**

A **row** corresponds to one  
**observation**, e.g., all  $(p+1)$   
features for datapoint 3

Field Goals  
Assists  
3-Point  
Attempts

Bias	FG	AST	3PA	PTS
1	1.8	0.6	4.1	5.3
1	0.4	0.8	1.5	1.7
1	1.1	1.9	2.2	3.2
1	6.0	1.6	0.0	13.9
1	3.4	2.2	0.2	8.9
...	...	...	...	...
1	4.0	0.8	0.0	11.5
1	3.1	0.9	0.0	7.8
1	3.6	1.1	0.0	8.9
1	3.4	0.8	0.0	8.5
1	3.8	1.5	0.0	9.4

Example design matrix  
708 rows x (3+1) cols

## The Multiple Linear Regression Model using Matrix Notation

We can express our linear model on our entire dataset as follows:

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$

$$\hat{\mathbf{Y}} = \mathbf{X}\boldsymbol{\theta}$$

Prediction vector  
 $\mathbb{R}^n$

Design matrix  
 $\mathbb{R}^{n \times (p+1)}$

Parameter vector  
 $\mathbb{R}^{(p+1)}$

Note that our  
**true output** is  
also a vector:  
 $\mathbf{Y} \in \mathbb{R}^n$

An expression is “**linear in theta**” if it is a **linear combination** of parameters  $\theta = [\theta_0, \theta_1, \dots, \theta_p]$

1.  $\hat{y} = \theta_0 + \theta_1(2) + \theta_2(4 \cdot 8) + \theta_3(\log 42)$

4. 
$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 5 & 6 & 7 \\ 1 & 8 & 9 & 0 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

2.  $\hat{y} = \theta_0 + \theta_1x_1 + \theta_2x_2x_3 + \theta_3.\log(x_4)$

5. 
$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} \\ 1 & x_{21} & x_{22} & x_{23} \\ 1 & x_{31} & x_{32} & x_{33} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

3.  $\hat{y} = \theta_0 + \theta_1x_1 + \log(\theta_2)x_2 + \theta_3\theta_4$

Which of these expressions are linear in theta?



**Which of the following  
expressions are linear in  
theta?**

- ① Start presenting to display the poll results on this slide.

## Linear in Theta

An expression is “**linear in theta**” if it is a **linear combination** of parameters  $\theta = [\theta_0, \theta_1, \dots, \theta_p]$

1.  $\hat{y} = \theta_0 + \theta_1(2) + \theta_2(4 \cdot 8) + \theta_3(\log 42)$

$$= [1 \ 2 \ 4.8 \ \log(42)] \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

2.  $\hat{y} = \theta_0 + \theta_1x_1 + \theta_2x_2x_3 + \theta_3 \cdot \log(x_4)$

$$= [1 \ x_1 \ x_2x_3 \ \log(x_4)] \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

3.  $\hat{y} = \theta_0 + \theta_1x_1 + \log(\theta_2)x_2 + \theta_3\theta_4$

4.  $\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 5 & 6 & 7 \\ 1 & 8 & 9 & 0 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$

5.  $\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} \\ 1 & x_{21} & x_{22} & x_{23} \\ 1 & x_{31} & x_{32} & x_{33} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$

“**Linear in theta**” means the expression can separate into a matrix product of two terms: a **vector of thetas**, and a matrix/vector not involving thetas.

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 5 & 6 & 7 \\ 1 & 8 & 9 & 0 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix}$$

# Mean Squared Error

---

## OLS Problem Formulation

- Multiple Linear Regression Model
- **Mean Squared Error**

Geometric Derivation

Performance: Residuals, Multiple R<sup>2</sup>

OLS Properties

- Residuals
- The Bias/Intercept Term
- Existence of a Unique Solution

# Today's Goal: Ordinary Least Squares



1. Choose a model

Multiple Linear  
Regression

$$\hat{\mathbb{Y}} = \mathbb{X}\theta$$

## 2. Choose a loss function

L2 Loss

Mean Squared Error  
(MSE)

$$R(\theta) = \frac{1}{n} \|\mathbb{Y} - \mathbb{X}\theta\|_2^2$$

3. Fit the model

Minimize  
average loss  
with ~~calculus~~ geometry

More Linear Algebra!!

4. Evaluate model  
performance

Visualize,  
~~Root MSE~~  
Multiple R<sup>2</sup>

The **norm** of a vector is some measure of that vector's **size/length**.

- The two norms we need to know for Data 100 are the  $L_1$  and  $L_2$  norms (sound familiar?).
- Today, we focus on  $L_2$  norm. We'll define the  $L_1$  norm another day.

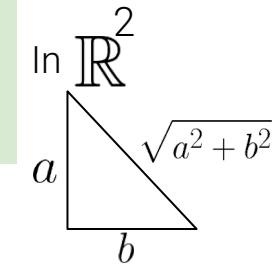
For the n-dimensional vector  $\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ , the **L2 vector norm** is

$$\|\vec{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2} = \sqrt{\sum_{i=1}^n (x_i^2)}$$

## [Linear Algebra] The L2 Norm as a Measure of Length

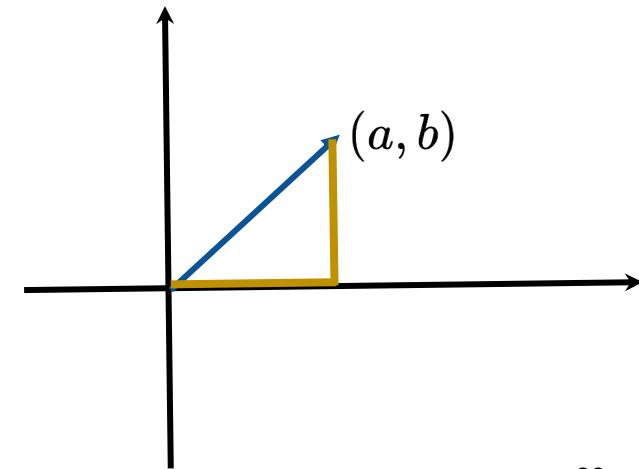
The L2 vector norm is a generalization of the Pythagorean theorem into  $n$  dimensions.

$$\|\vec{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2} = \sqrt{\sum_{i=1}^n (x_i^2)}$$



It can therefore be used as a measure of **length** of a vector

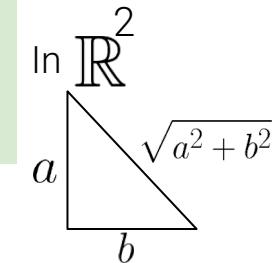
- The vector on the right has length  $\|\vec{v}\|_2 = \sqrt{a^2 + b^2}$



## [Linear Algebra] The L2 Norm as a Measure of Distance

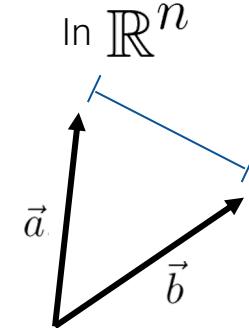
The L2 vector norm is a generalization of the Pythagorean theorem into  $n$  dimensions.

$$\|\vec{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2} = \sqrt{\sum_{i=1}^n (x_i^2)}$$



It can also be used as a measure of **distance** between two vectors.

- For  $n$ -dimensional vectors  $\vec{a}, \vec{b}$ , their distance is  $\|\vec{a} - \vec{b}\|_2$ .



Note: The square of the L2 norm of a vector is the sum of the squares of the vector's elements:

$$(\|\vec{x}\|_2)^2 = \sum_{i=1}^n x_i^2$$

Looks like Mean Squared Error!!

## Mean Squared Error with L2 Norms

---

We can rewrite mean squared error as a squared L2 norm:

$$\begin{aligned} R(\theta) &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \frac{1}{n} \|\mathbb{Y} - \hat{\mathbb{Y}}\|_2^2 \end{aligned}$$

With our linear model  $\hat{\mathbb{Y}} = \mathbb{X}\theta$  :

$$R(\theta) = \frac{1}{n} \|\mathbb{Y} - \mathbb{X}\theta\|_2^2$$

## Ordinary Least Squares

The **least squares estimate**  $\hat{\theta}$  is the parameter that **minimizes** the objective function  $R(\theta)$ :

$$R(\theta) = \frac{1}{n} \|\mathbb{Y} - \mathbb{X}\theta\|_2^2$$

How should we interpret the OLS problem?

- A. Minimize the mean squared error for the linear model  $\hat{\mathbb{Y}} = \mathbb{X}\theta$
- B. Minimize the **distance** between true and predicted values  $\mathbb{Y}$  and  $\hat{\mathbb{Y}}$
- C. Minimize the **length** of the residual vector,  $e = \mathbb{Y} - \hat{\mathbb{Y}} = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$
- D. All of the above
- E. Something else





# How should we interpret the OLS problem?

- ① Start presenting to display the poll results on this slide.

## Ordinary Least Squares

The **least squares estimate**  $\hat{\theta}$  is the parameter that **minimizes** the objective function  $R(\theta)$ :

$$R(\theta) = \frac{1}{n} \|\mathbb{Y} - \mathbb{X}\theta\|_2^2$$

How should we interpret the OLS problem?

A. Minimize the mean squared error for the linear model  $\hat{\mathbb{Y}} = \mathbb{X}\theta$

B. Minimize the **distance**  
between true and predicted values  $\mathbb{Y}$  and  $\hat{\mathbb{Y}}$

C. Minimize the **length** of the residual vector,  $e = \mathbb{Y} - \hat{\mathbb{Y}} =$

$$\left[ \begin{array}{c} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{array} \right]$$

}  
Important  
for today

D. All of the above

E. Something else

# Interlude

---

LEAST SQUARES  
REGRESSION



MOST SQUARES  
REGRESSION



# Geometric Derivation

---

## OLS Problem Formulation

- Multiple Linear Regression Model
- Mean Squared Error

## Geometric Derivation

Performance: Residuals, Multiple R<sup>2</sup>

## OLS Properties

- Residuals
- The Bias/Intercept Term
- Existence of a Unique Solution

# Today's Goal: Ordinary Least Squares



1. Choose a model

Multiple Linear  
Regression



2. Choose a loss  
function

L2 Loss  
Mean Squared Error  
(MSE)

## 3. Fit the model

Minimize  
average loss  
with ~~calculus~~ geometry

4. Evaluate model  
performance

Visualize,  
~~Root MSE~~  
Multiple R<sup>2</sup>

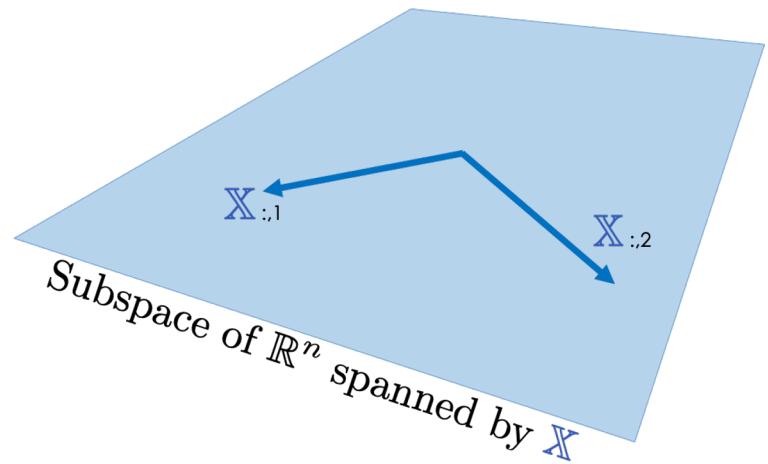
$$\hat{\mathbb{Y}} = \mathbf{X}\theta$$

$$R(\theta) = \frac{1}{n} \|\mathbb{Y} - \mathbf{X}\theta\|_2^2$$

The calculus derivation requires matrix calculus (out of scope, but here's a [link](#) if you're interested). Instead, we will derive  $\hat{\theta}$  using a **geometric argument**.

The set of all possible linear combinations of the columns of  $\mathbb{X}$  is called the **span** of the columns of  $\mathbb{X}$  (denoted  $\text{span}(\mathbb{X})$ ), also called the **column space**.

- Intuitively, this is all of the vectors you can "reach" using the columns of  $\mathbb{X}$ .
- If each column of  $\mathbb{X}$  has length  $n$ ,  $\text{span}(\mathbb{X})$  is a subspace of  $\mathbb{R}^n$ .



**Approach 1:** So far, we've thought of our model as horizontally stacked predictions per datapoint:

$$\begin{matrix} \text{n} \\ \hat{\mathbf{Y}} \\ \text{---} \\ \text{1} \end{matrix} = \begin{bmatrix} \text{x}_1^T \\ \text{x}_2^T \\ \vdots \\ \text{x}_n^T \end{bmatrix} \begin{matrix} \text{p+1} \\ \theta \\ \text{---} \\ \text{1} \end{matrix}$$

**Approach 1:** So far, we've thought of our model as horizontally stacked predictions per datapoint:

$$\begin{bmatrix} \hat{Y}_1 \\ \vdots \\ \hat{Y}_n \end{bmatrix} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$

$$\hat{y}_1 = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} \cdot \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} = x^\top \theta$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$

**Approach 1:** So far, we've thought of our model as horizontally stacked predictions per datapoint:

$$n \begin{bmatrix} \hat{Y} \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{bmatrix} \theta \\ \vdots \\ 1 \end{bmatrix}^{p+1}$$

**Approach 2:** However, it is helpful sometimes to think of matrix-vector multiplication as performed by columns. We can also think of  $\hat{Y}$  as a **linear combination of feature vectors**, scaled by **parameters**.

$$n \begin{bmatrix} \hat{Y} \\ \vdots \\ 1 \end{bmatrix} = n \begin{bmatrix} | & | \\ \mathbb{X}_{:,1} & \mathbb{X}_{:,2} \\ | & | \end{bmatrix} \begin{bmatrix} \theta \\ \vdots \\ 1 \end{bmatrix}^{p+1} = \theta_1 \mathbb{X}_{:,1} + \theta_2 \mathbb{X}_{:,2}$$

## Prediction is a Linear Combination of Columns

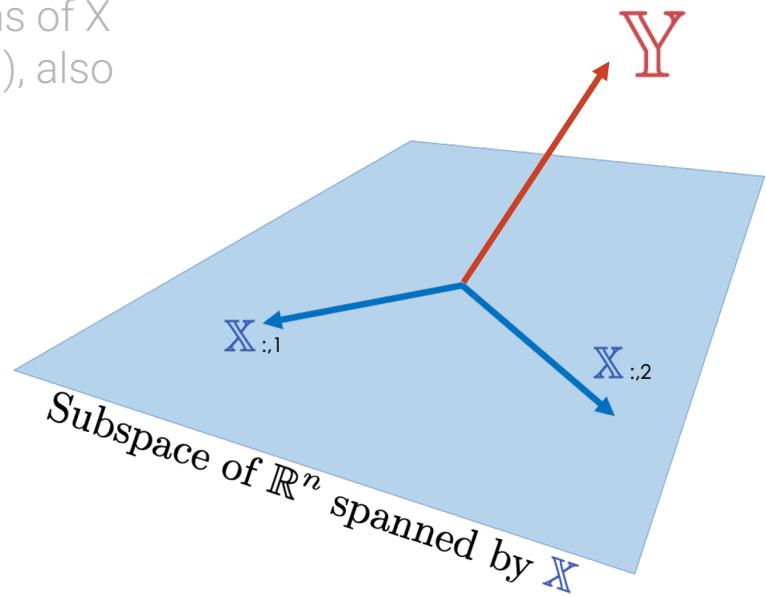
The set of all possible linear combinations of the columns of  $X$  is called the **span** of the columns of  $X$  (denoted  $\text{span}(\mathbb{X})$ ), also called the **column space**.

- Intuitively, this is all of the vectors you can “reach” using the columns of  $X$ .
- If each column of  $X$  has length  $n$ ,  $\text{span}(\mathbb{X})$  is a subspace of  $\mathbb{R}^n$ .

Our prediction  $\hat{\mathbb{Y}} = \mathbb{X}\theta$  is a **linear combination** of the columns of  $\mathbb{X}$ . Therefore  $\hat{\mathbb{Y}} \in \text{span}(\mathbb{X})$ .

Interpret: Our linear prediction  $\hat{\mathbb{Y}}$  will be in  $\text{span}(\mathbb{X})$ , even if the true values  $\mathbb{Y}$  might not be.

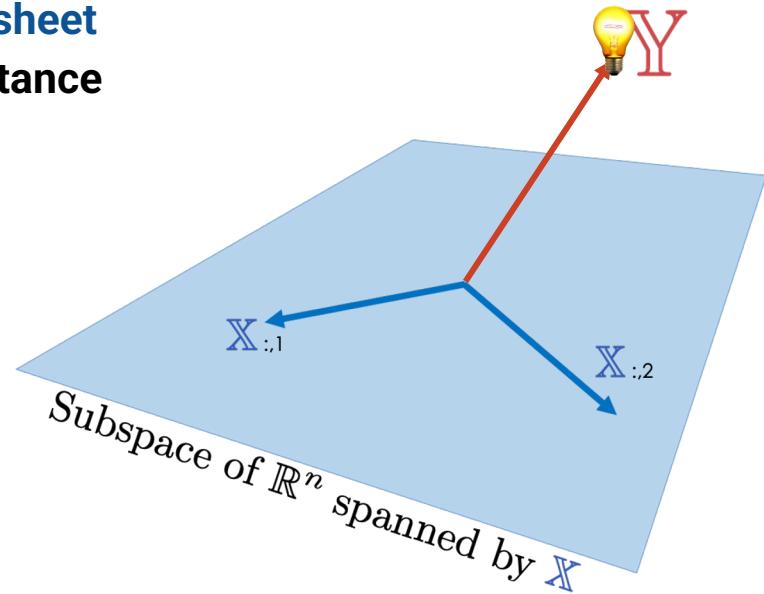
Goal: Find the vector  $\hat{\mathbb{Y}}$  in  $\text{span}(\mathbb{X})$  that is **closest** to  $\mathbb{Y}$



## A thought experiment

If you're a human being who can only stand on the **blue sheet of paper**, and you need to get as close as possible in **distance** to the **light bulb** located at the tip of the **red arrow**.

Where do you stand on the blue sheet?

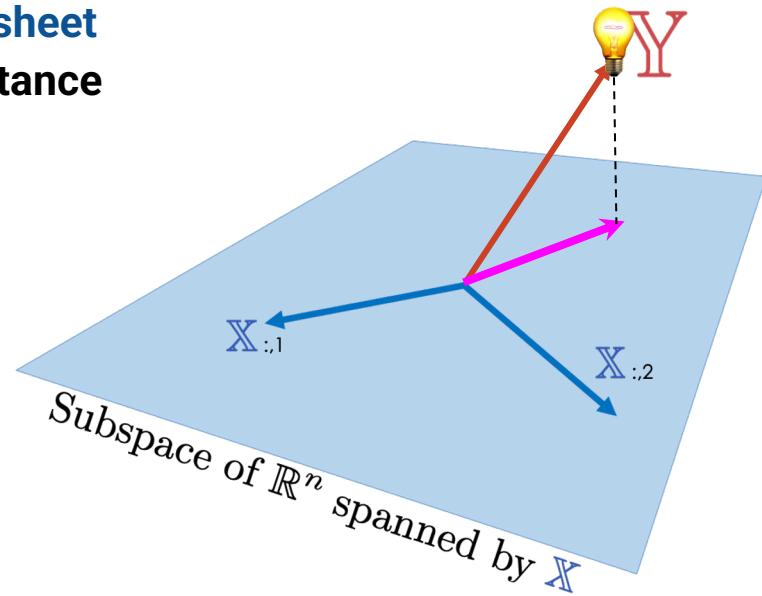


## A thought experiment

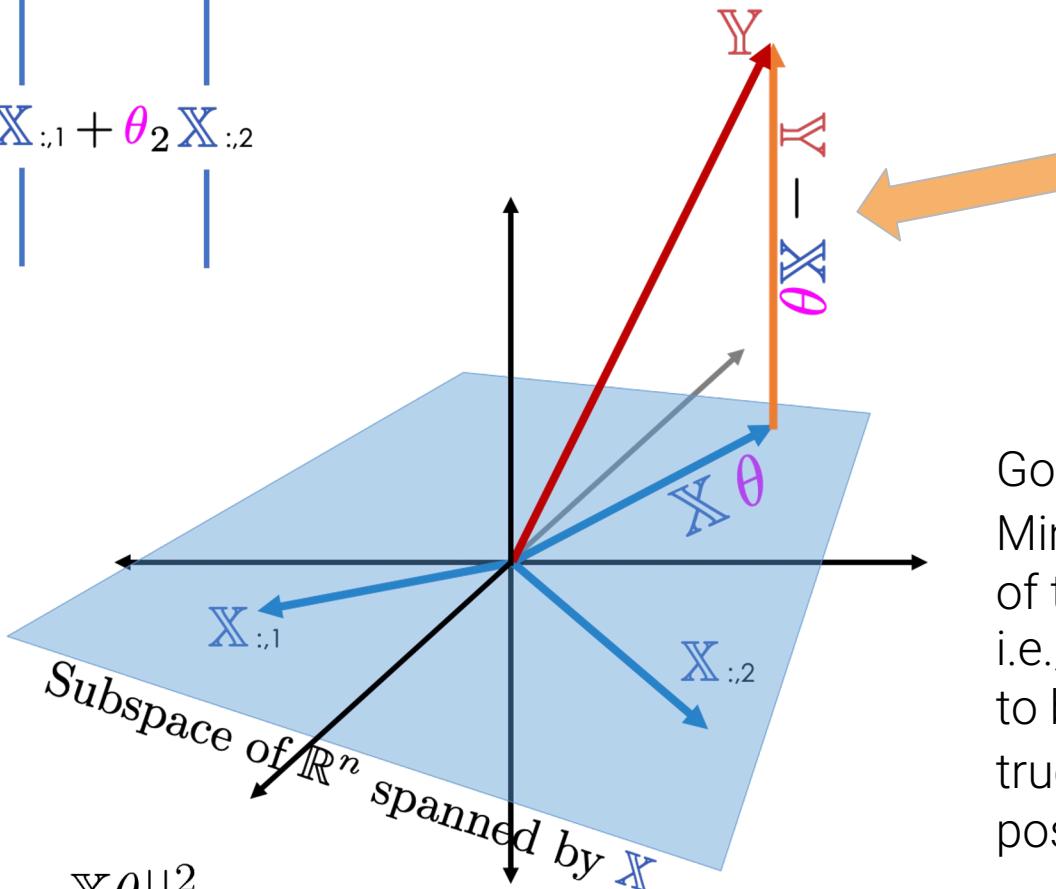
If you're a human being who can only stand on the **blue sheet of paper**, and you need to get as close as possible in **distance** to the **light bulb** located at the tip of the **red arrow**.

Where do you stand on the blue sheet?

**Right below the lightbulb - that's the closest you can get because you can't travel vertically!**



$$\begin{bmatrix} n \\ \hat{\mathbb{Y}} \\ 1 \end{bmatrix} = \theta_1 \mathbb{X}_{:,1} + \theta_2 \mathbb{X}_{:,2}$$



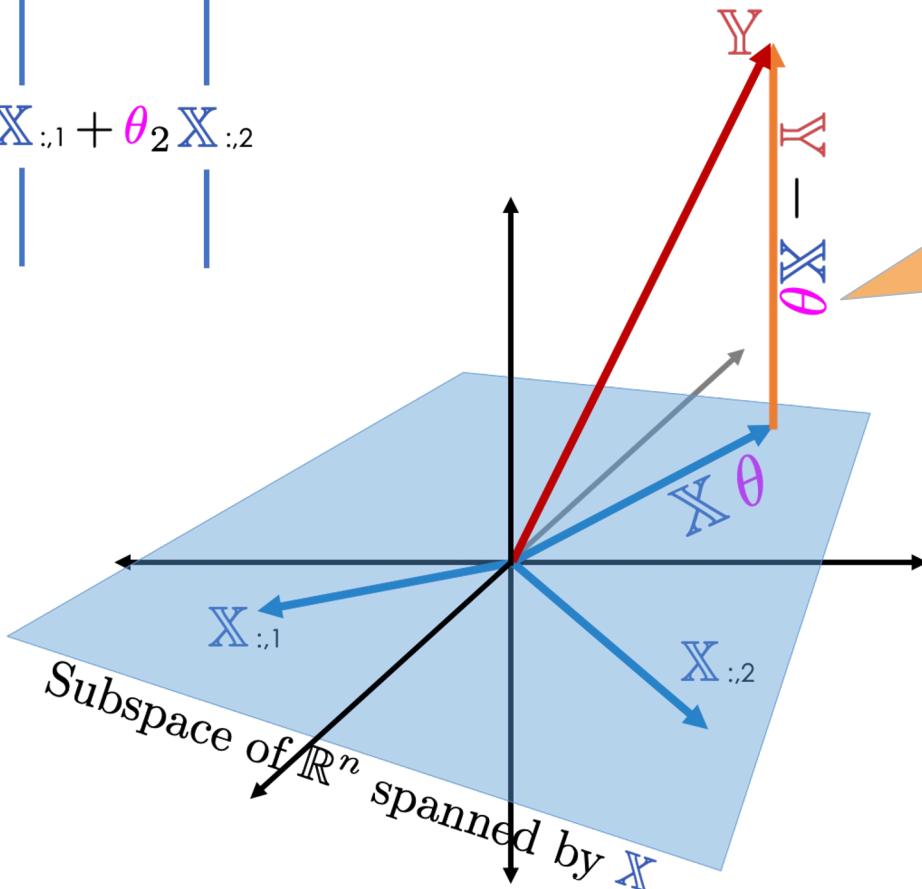
$$R(\theta) = \frac{1}{n} \|\mathbb{Y} - \mathbb{X}\theta\|_2^2$$

This is the residual vector,  
 $e = \mathbb{Y} - \hat{\mathbb{Y}}$ .

Goal:

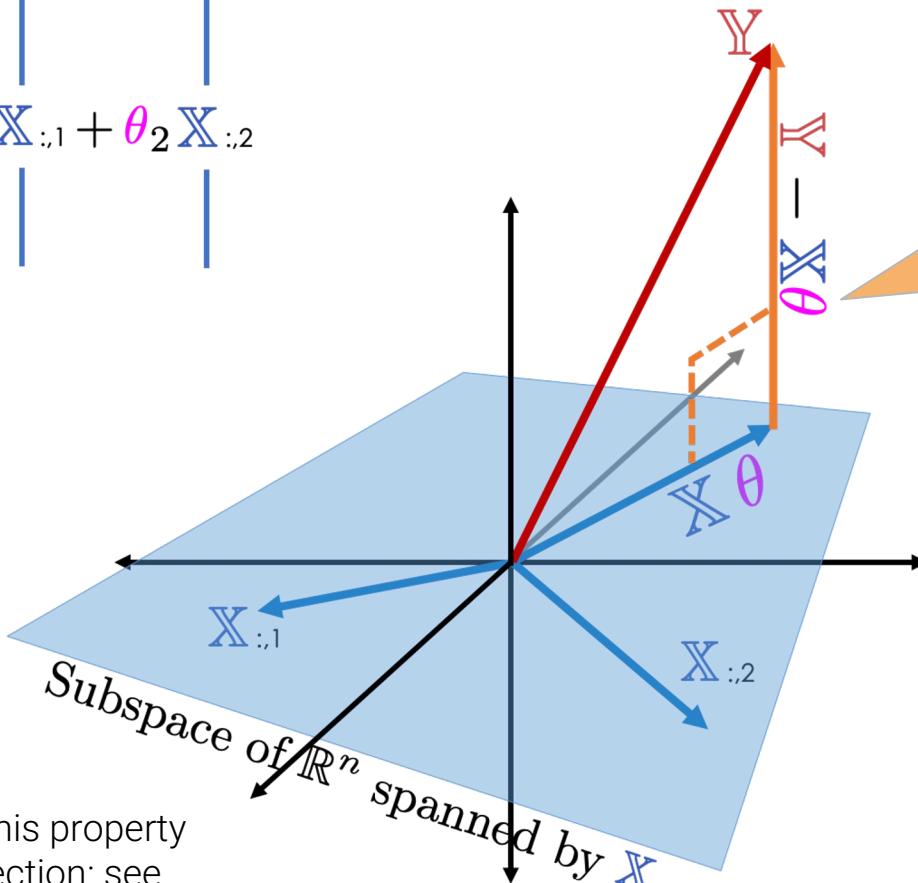
Minimize the  $L_2$  norm of the residual vector.  
 i.e., get the predictions  $\hat{\mathbb{Y}}$  to be “as close” to our true  $\mathbb{Y}$  values as possible.

$$\begin{bmatrix} n \\ \hat{Y} \\ 1 \end{bmatrix} = \theta_1 \mathbf{X}_{:,1} + \theta_2 \mathbf{X}_{:,2}$$



How do we minimize this distance – the norm of the residual vector (squared)?

$$\begin{bmatrix} n \\ \vdots \\ \hat{\mathbb{Y}} \\ \vdots \\ 1 \end{bmatrix} = \theta_1 \mathbb{X}_{:,1} + \theta_2 \mathbb{X}_{:,2}$$

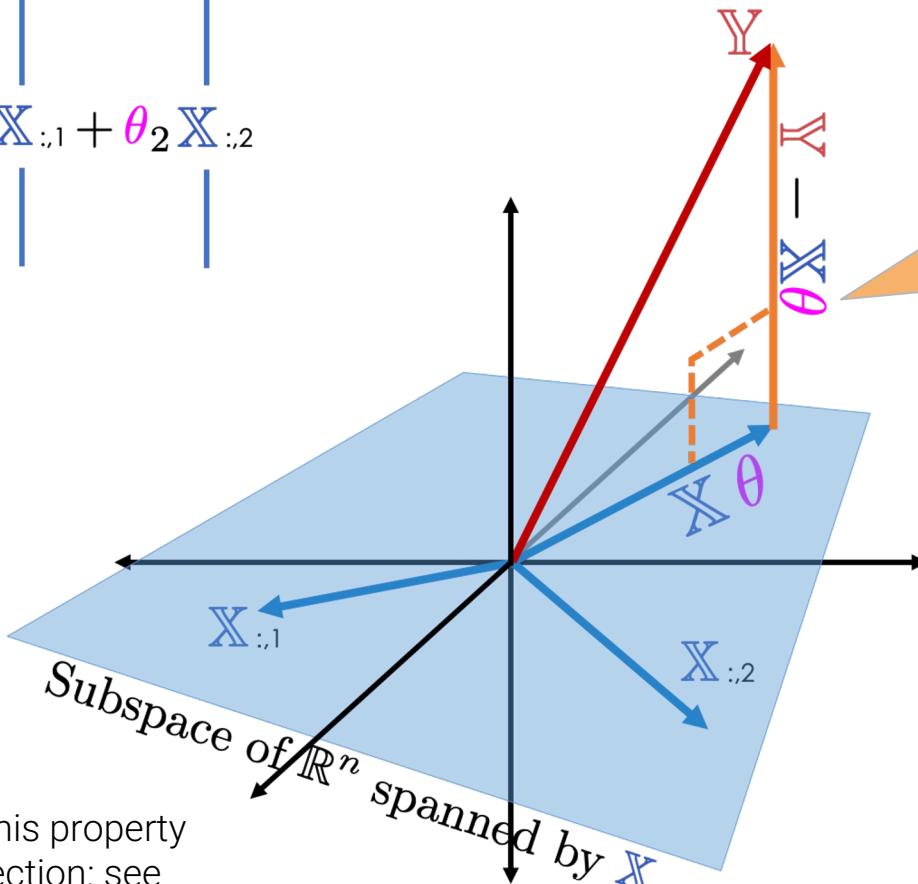


We will not prove this property of orthogonal projection: see [Khan Academy](#).

How do we minimize this distance – the norm of the residual vector (squared)?

The vector in  $span(\mathbb{X})$  that is closest to  $\mathbb{Y}$  is the **orthogonal projection** of  $\mathbb{Y}$  onto  $span(\mathbb{X})$ .

$$\begin{bmatrix} n \\ \vdots \\ \hat{\mathbb{Y}} \\ \vdots \\ 1 \end{bmatrix} = \theta_1 \mathbb{X}_{:,1} + \theta_2 \mathbb{X}_{:,2}$$



We will not prove this property of orthogonal projection: see [Khan Academy](#).

How do we minimize this distance – the norm of the residual vector (squared)?

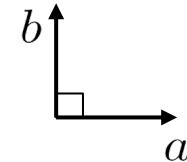
The vector in  $span(\mathbb{X})$  that is closest to  $\mathbb{Y}$  is the **orthogonal projection** of  $\mathbb{Y}$  onto  $span(\mathbb{X})$ .

Thus, we should choose the  $\theta$  that makes the residual vector **orthogonal** to  $span(\mathbb{X})$ .

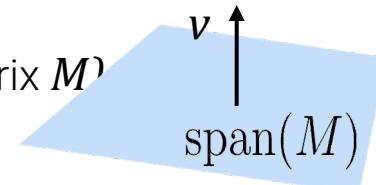
## [Linear Algebra] Orthogonality

1. Vector  $a$  and Vector  $b$  are **orthogonal** if and only if their dot product is 0:  $a^T b = 0$

This is a generalization of the notion of two vectors in 2D being perpendicular.



2. A vector  $v$  is **orthogonal** to  $\text{span}(M)$ , (the span of the columns of a matrix  $M$ ) if and only if  $v$  is orthogonal to **each column** in  $M$ .



Let's express 2 in matrix notation. Let  $v \in \mathbb{R}^{n \times 1}$   $M \in \mathbb{R}^{n \times d}$

$$\begin{aligned} M_{:1}^T v &= 0 \\ M_{:2}^T v &= 0 \\ &\vdots \\ M_{:d}^T v &= 0 \end{aligned}$$



$$\begin{bmatrix} M_{:1}^T v \\ M_{:2}^T v \\ \vdots \\ M_{:d}^T v \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$



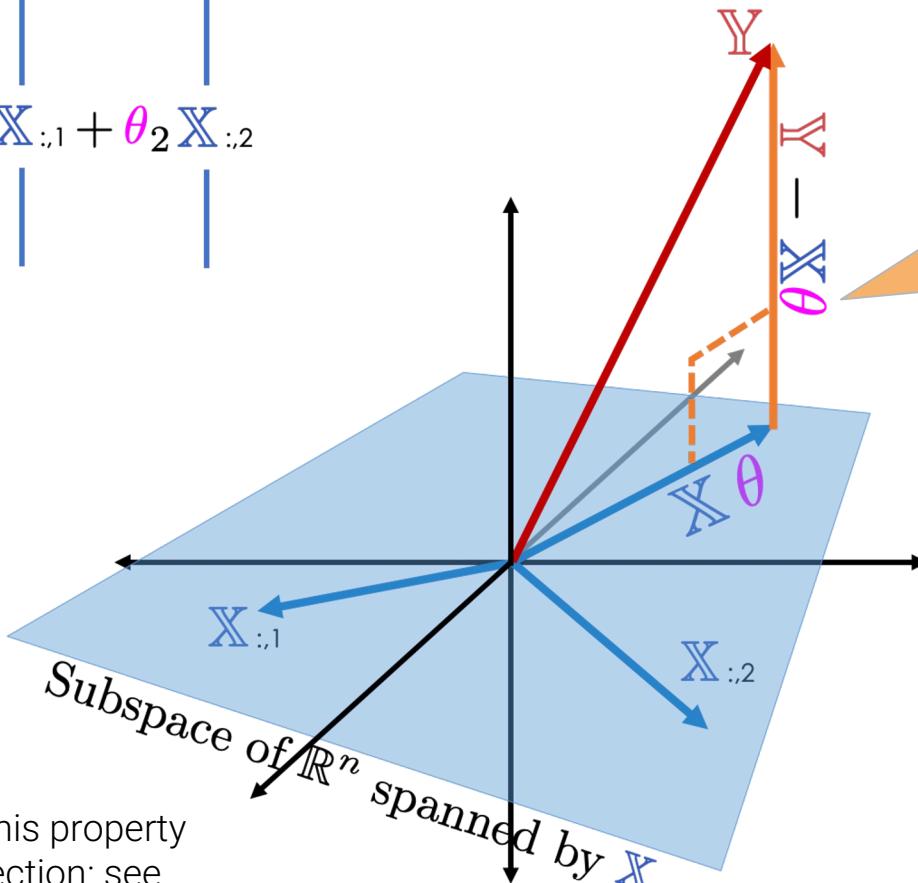
$$\text{where } M = \begin{bmatrix} | & | & | \\ M_{:1} & M_{:2} & \vdots & M_{:d} \\ | & | & & | \end{bmatrix}$$

$$\underbrace{M^T}_{M^T \in \mathbb{R}^{d \times n}} v = \underbrace{\vec{0}}_{\text{zero vector}}$$

$v$  is orthogonal to each column of  $M, M_{:j} \in \mathbb{R}^n$

**zero vector** ( $d$ -length vector full of 0s).

$$\begin{bmatrix} n \\ \vdots \\ \hat{\mathbb{Y}} \\ \vdots \\ 1 \end{bmatrix} = \theta_1 \mathbb{X}_{:,1} + \theta_2 \mathbb{X}_{:,2}$$

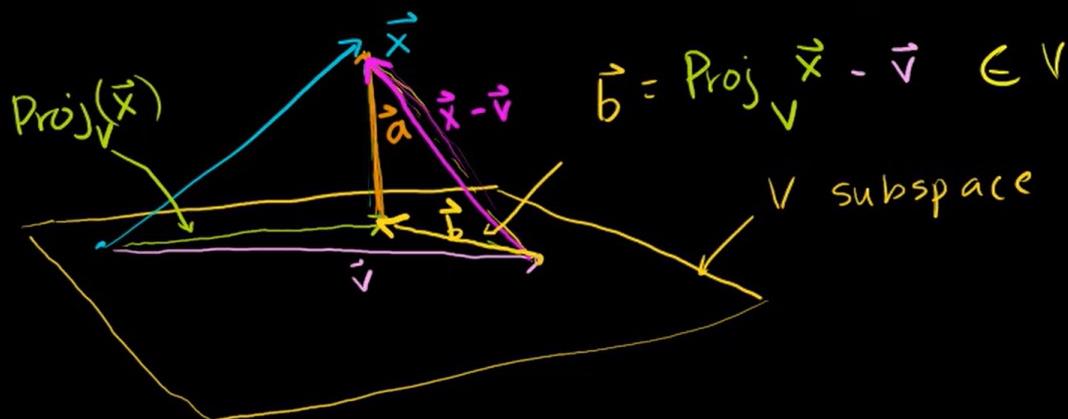


We will not prove this property of orthogonal projection: see [Khan Academy](#).

How do we minimize this distance – the norm of the residual vector (squared)?

The vector in  $\text{span}(\mathbb{X})$  that is closest to  $\mathbb{Y}$  is the **orthogonal projection** of  $\mathbb{Y}$  onto  $\text{span}(\mathbb{X})$ .

Thus, we should choose the  $\theta$  that makes the residual vector **orthogonal** to  $\text{span}(\mathbb{X})$ .



We will not prove this property of orthogonal projection: see [Khan Academy](#).

$$\|\vec{x} - \underbrace{\text{Proj}_V \vec{x}}_{\vec{a}}\| \leq \|\vec{x} - \vec{v}\|$$

$$\begin{aligned} \|\vec{x} - \vec{v}\|^2 &= \|\vec{b} + \vec{a}\|^2 = (\vec{b} + \vec{a}) \cdot (\vec{b} + \vec{a}) = \\ &= \vec{b} \cdot \vec{b} + 2 \vec{a} \cdot \vec{b} + \vec{a} \cdot \vec{a} \end{aligned}$$

$$\|\vec{x} + \vec{a}\|^2 = \|\vec{b}\|^2 + \|\vec{a}\|^2$$

## Ordinary Least Squares Proof

The **least squares estimate**  $\hat{\theta}$  is the parameter  $\theta$  that minimizes the objective function  $R(\theta)$ :

$$R(\theta) = \frac{1}{n} \|\mathbb{Y} - \mathbb{X}\theta\|_2^2$$

Design matrix  $M^T v = \vec{0}$  Residual vector

Equivalently, this is the  $\hat{\theta}$  such that the residual vector  $\mathbb{Y} - \mathbb{X}\hat{\theta}$

Definition of orthogonality  
of  $\mathbb{Y} - \mathbb{X}\hat{\theta}$  to  $\text{span}(\mathbb{X})$   
(0 is the  $\vec{0}$  vector)

Rearrange terms

$$\mathbb{X}^T (\mathbb{Y} - \mathbb{X}\hat{\theta}) = 0$$

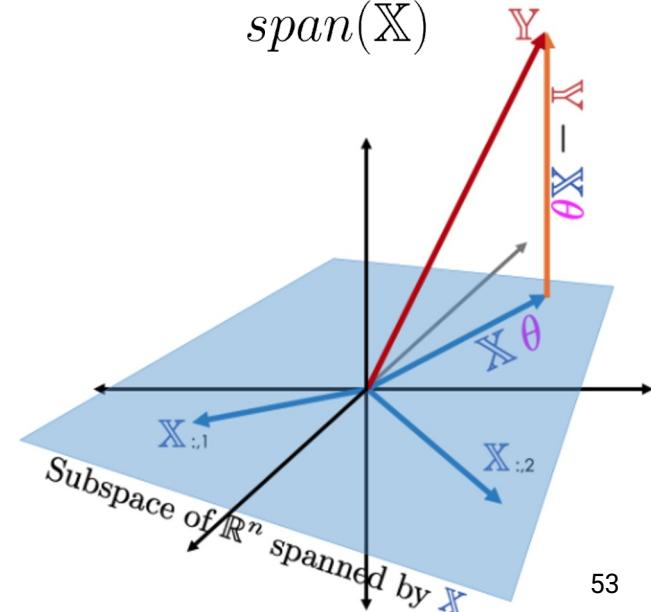
The **normal equation**

$$\mathbb{X}^T \mathbb{X} \hat{\theta} = \mathbb{X}^T \mathbb{Y}$$

If  $\mathbb{X}^T \mathbb{X}$

is invertible

$$\hat{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}$$



$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

This result is so important that it deserves its own slide.

It is the **least squares estimate** and the solution to the normal equation  $\mathbf{X}^T \mathbf{X} \hat{\theta} = \mathbf{X}^T \mathbf{Y}$



$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

---

This result is so important that it deserves its own slide.

It is the **least squares estimate** and the solution to the normal equation  $\mathbf{X}^T \mathbf{X} \hat{\theta} = \mathbf{X}^T \mathbf{Y}$

# Least Squares Estimate

1. Choose a model

Multiple Linear  
Regression

$$\hat{\mathbb{Y}} = \mathbb{X}\theta$$

2. Choose a loss  
function

L2 Loss

Mean Squared Error  
(MSE)

$$R(\theta) = \frac{1}{n} \|\mathbb{Y} - \mathbb{X}\theta\|_2^2$$

3. Fit the model



Minimize  
average loss  
with ~~calculus~~ geometry

$$\hat{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}$$

4. Evaluate model  
performance

Visualize,  
~~Root MSE~~  
Multiple R<sup>2</sup>

# Performance

---

## OLS Problem Formulation

- Multiple Linear Regression Model
- Mean Squared Error

## Geometric Derivation

## **Performance: Residuals, Multiple R<sup>2</sup>**

## OLS Properties

- Residuals
- The Bias/Intercept Term
- Existence of a Unique Solution

## Least Squares Estimate



1. Choose a model

Multiple Linear  
Regression



2. Choose a loss  
function

L2 Loss  
Mean Squared Error  
(MSE)

$$\hat{Y} = X\theta$$

$$R(\theta) = \frac{1}{n} ||Y - X\theta||_2^2$$



3. Fit the model

Minimize  
average loss  
with ~~calculus~~ geometry

$$\hat{\theta} = (X^T X)^{-1} X^T Y$$

**4. Evaluate model  
performance**

Visualize,  
~~Root MSE~~  
Multiple R<sup>2</sup>

## Multiple Linear Regression

$$\hat{Y} = X\theta$$

Prediction  
vector

$$\mathbb{R}^n$$

Design matrix

$$\mathbb{R}^{n \times (p+1)}$$

Parameter  
vector

$$\mathbb{R}^{(p+1)}$$

Note that our  
**true output** is  
also a vector:

$$Y \in \mathbb{R}^n$$

$$R(\theta) = \frac{1}{n} ||Y - X\theta||_2^2$$

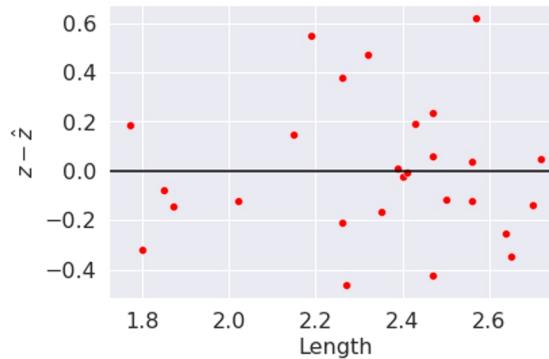
## Demo

$$\hat{\theta} = (X^T X)^{-1} X^T Y$$

## [Visualization] Residual Plots

### Simple linear regression

Plot residuals vs  
the single feature  $x$ .

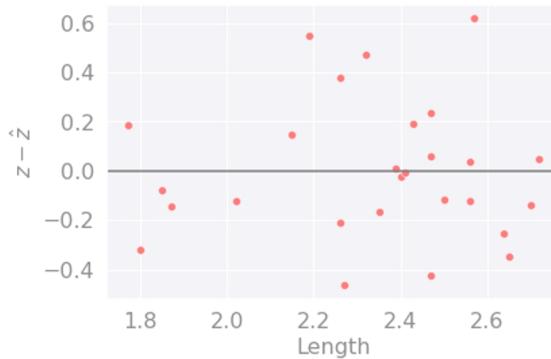


## Compare

## [Visualization] Residual Plots

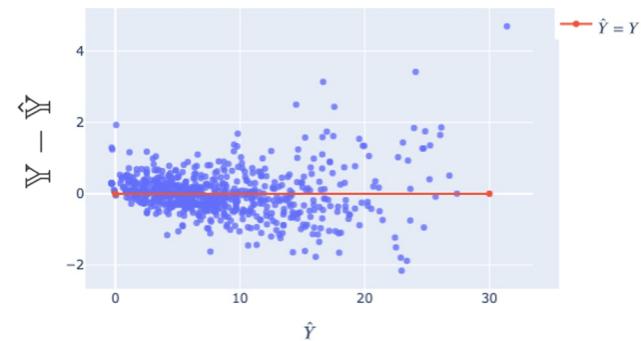
### Simple linear regression

Plot residuals vs  
the single feature  $x$ .



### Multiple linear regression

Plot residuals vs  
**fitted (predicted) values**  $\hat{y}$



## Compare

See notebook

Same interpretation as before ([textbook](#)):

- A good residual plot shows no pattern.
- A good residual plot also has a similar vertical spread throughout the entire plot. Else (heteroscedasticity), the accuracy of the predictions is not reliable.

## [Metrics] Multiple R<sup>2</sup>

### Simple linear regression

Error

RMSE

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Linearity

Correlation coefficient,  $r$

$$r = \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma_x} \right) \left( \frac{y_i - \bar{y}}{\sigma_y} \right)$$

### Multiple linear regression

Error

RMSE

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Linearity

**Multiple R<sup>2</sup>**, also called the **coefficient of determination**

$$R^2 = \frac{\text{variance of fitted values}}{\text{variance of } y} = \frac{\sigma_{\hat{y}}^2}{\sigma_y^2}$$

## Compare

We define the **multiple R<sup>2</sup>** value as the **proportion of variance** or our **fitted values** (predictions)  $\hat{y}$  to our true values  $y$ .

$$R^2 = \frac{\text{variance of fitted values}}{\text{variance of } y} = \frac{\sigma_{\hat{y}}^2}{\sigma_y^2}$$

Also called the **correlation of determination**.

R<sup>2</sup> ranges from 0 to 1 and is effectively  
“the proportion of variance that the **model explains**.”

## Compare

For OLS with an intercept term (e.g.  $\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p$ )

$$R^2 = [r(y, \hat{y})]^2$$

$R^2 = r^2$  is equal to the square of correlation between  $y$  and  $\hat{y}$ .

- For SLR,  $R^2$  is the correlation between  $x_1$  and  $y$ .

## [Metrics] Multiple R<sup>2</sup>

predicted PTS =  $3.98 + 2.4 \cdot \text{AST}$

R<sup>2</sup> = 0.457

predicted PTS =  $2.163 + 1.64 \cdot \text{AST} + 1.26 \cdot \text{3PA}$

R<sup>2</sup> = 0.609

## Compare

### Simple linear regression

Error

RMSE

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Linearity

Correlation coefficient,  $r$

$$r = \frac{1}{n} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{\sigma_x} \right) \left( \frac{y_i - \bar{y}}{\sigma_y} \right)$$

### Multiple linear regression

Error

RMSE

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Linearity

**Multiple R<sup>2</sup>**, also called the **coefficient of determination**

$$R^2 = \frac{\text{variance of fitted values}}{\text{variance of } y} = \frac{\sigma_{\hat{y}}^2}{\sigma_y^2}$$

As we add more features, our fitted values tend to become closer and closer to our actual  $y$  values. Thus, R<sup>2</sup> increases.

- The SLR **model** (AST only) explains 45.7% of the variance in the true  $y$ .
- The AST & 3PA **model** explains 60.9%.

Adding more features doesn't always mean our model is better, though! We will see why after the midterm.

# OLS Properties

---

OLS Problem Formulation

- Multiple Linear Regression Model
- Mean Squared Error

Geometric Derivation

Performance: Residuals, Multiple R<sup>2</sup>

## **OLS Properties**

- Residuals
- The Bias/Intercept Term
- Existence of a Unique Solution

## Residual Properties

When using the optimal parameter vector, our residuals  $e = \mathbb{Y} - \mathbb{X}\hat{\theta}$  are orthogonal to  $\text{span}(\mathbb{X})$ .

$$\mathbb{X}^T e = 0$$

Proof: First line of our OLS estimate proof ([slide](#)).

For all linear models:

Since our predicted response  $\hat{\mathbb{Y}}$  is in  $\text{span}(\mathbb{X})$  by definition,  $\hat{\mathbb{Y}}^T e = 0$ , and hence it is orthogonal to the residuals.

For all linear models with an **intercept term**,  $\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p$ , the **sum of residuals is zero**.

$$\sum_{i=1}^n e_i = 0$$

$$\mathbb{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

(Proof hint)  $\mathbb{1}^T e = 0$

You will prove both properties in homework.

## Properties When Our Model Has an Intercept Term

---

For all linear models with an **intercept term**, the **sum of residuals is zero**.

- This is the real reason why we don't directly use residuals as loss.

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) = \frac{1}{n} \sum_{i=1}^n e_i = 0 \quad \sum_{i=1}^n e_i = 0 \quad (\text{previous slide})$$

- This is also why positive and negative residuals will cancel out in any residual plot where the (linear) model contains an intercept term, even if the model is terrible.

It follows from the property above that for linear models with intercepts, the average predicted  $\hat{y}$  value is equal to the average true  $y$  value.

$$\bar{y} = \bar{\hat{y}}$$

These properties are true when there is an intercept term, and not necessarily when there isn't.

## Does a Unique Solution Always Exist?

	Model	Estimate	Unique?
Constant Model + MSE	$\hat{y} = \theta_0$	$\hat{\theta}_0 = \text{mean}(y) = \bar{y}$	<b>Yes.</b> Any set of values has a unique mean.
Constant Model + MAE	$\hat{y} = \theta_0$	$\hat{\theta}_0 = \text{median}(y)$	<b>Yes</b> , if odd. <b>No</b> , if even. Return average of middle 2 values.
Simple Linear Regression + MSE	$\hat{y} = \theta_0 + \theta_1 x$	$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$ $\hat{\theta}_1 = r \frac{\sigma_y}{\sigma_x}$	<b>Yes.</b> Any set of non-constant* values has a unique mean, SD, and correlation coefficient.
<b>Ordinary Least Squares</b> (Linear Model + MSE)	$\hat{\mathbb{Y}} = \mathbb{X}\theta$	$\hat{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}$	???

## Understanding The Solution Matrices

$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

In most settings,  
**# observations**  $\gg$  **# features**

The diagram illustrates the dimensions of the matrices involved in the normal equation solution. It shows two main parts: the left side of the equation and the right side.

**Left Side:**  $(\mathbf{X}^T \mathbf{X})^{-1}$

- A large blue bracket groups two matrices:
  - A blue square matrix of size  $n \times p+1$ .
  - A blue vertical vector of size  $n \times 1$  with entries 1, 1, 1, 1, ... followed by a green ellipsis.
- A curly brace below this group indicates the overall size is  $p+1 \times p+1$ .
- An arrow points from the top of this group to the term  $(\mathbf{X}^T \mathbf{X})^{-1}$ .

**Right Side:**  $\mathbf{X}^T \mathbf{Y}$

- A blue horizontal vector of size  $n \times p+1$  with entries 1, 1, 1, 1, ... followed by a green ellipsis.
- An arrow points from the top of this vector to the term  $\mathbf{X}^T \mathbf{Y}$ .
- A red arrow points from the right side of the vector to a red vertical vector of size  $n \times 1$  labeled '1'.
- A curly brace below the vector indicates the overall size is  $p+1 \times 1$ .

## Understanding The Solution Matrices

In practice, instead of directly inverting matrices, we can use more efficient numerical solvers to directly solve a system of linear equations.

The **Normal Equation**:

$$\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\theta}} = \mathbf{X}^T \mathbf{Y}$$

$$\left( \begin{array}{c|c} & p+1 \\ \hline p+1 & A \end{array} \right) \hat{\boldsymbol{\theta}} = \begin{matrix} 1 \\ p+1 \end{matrix} \mathbf{b}$$

Note that at least one solution always exists:

Intuitively, we can always draw a line of best fit for a given set of data, but there may be multiple lines that are “equally good”. (Formal proof is beyond this course.)

## Uniqueness of a Solution: Proof

---

### Claim

The Least Squares estimate  $\hat{\theta}$  is **unique** if and only if  $\mathbb{X}$  is **full column rank**.

### Proof

- The solution to the normal equation  $\mathbb{X}^T \mathbb{X} \hat{\theta} = \mathbb{X}^T \mathbb{Y}$  is the least square  $\hat{\theta}$  estimate .
- $\hat{\theta}$  has a **unique** solution if and only if the square matrix  $\mathbb{X}^T \mathbb{X}$  is **invertible**, which happens if and only if  $\mathbb{X}^T \mathbb{X}$  is full rank.
  - The **rank** of a square matrix is the max **# of linearly independent columns** it contains.
  - $\mathbb{X}^T \mathbb{X}$  has shape  $(p+1) \times (p+1)$ , and therefore has max rank  $p+1$ .
- $\mathbb{X}^T \mathbb{X}$  and  $\mathbb{X}$  **have the same rank** (proof out of scope).
- Therefore  $\mathbb{X}^T \mathbb{X}$  has rank  $p+1$  if and only if  $\mathbb{X}$  has rank  $p+1$  (full column rank).

## Uniqueness of a Solution: Interpretation

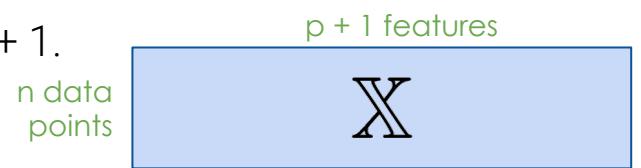
Claim:

The Least Squares estimate  $\hat{\theta}$  is **unique** if and only if  $\mathbb{X}$  is **full column rank**.

When would we **not** have unique estimates?

1. If our design matrix  $\mathbb{X}$  is “**wide**”:

- (property of rank) If  $n < p$ , rank of  $\mathbb{X} = \min(n, p + 1) < p + 1$ .
- In other words, if we have way more features than observations, then  $\hat{\theta}$  is not unique.
- Typically we have  $n \gg p$  so this is less of an issue.



2. If our design matrix  $\mathbb{X}$  has features that are **linear combinations of other features**.

- By definition, rank of  $\mathbb{X}$  is number of linearly independent columns in  $\mathbb{X}$ .
- Example: If “Width”, “Height”, and “Perimeter” are all columns,
  - $\text{Perimeter} = 2 * \text{Width} + 2 * \text{Height} \rightarrow \mathbb{X}$  is not full rank.
- Important with one-hot encoding (to discuss in later).

## Does a Unique Solution Always Exist?

	Model	Estimate	Unique?
Constant Model + MSE	$\hat{y} = \theta_0$	$\hat{\theta}_0 = \text{mean}(y) = \bar{y}$	<b>Yes.</b> Any set of values has a unique mean.
Constant Model + MAE	$\hat{y} = \theta_0$	$\hat{\theta}_0 = \text{median}(y)$	<b>Yes</b> , if odd. <b>No</b> , if even. Return average of middle 2 values.
Simple Linear Regression + MSE	$\hat{y} = \theta_0 + \theta_1 x$	$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$ $\hat{\theta}_1 = r \frac{\sigma_y}{\sigma_x}$	<b>Yes.</b> Any set of non-constant* values has a unique mean, SD, and correlation coefficient.
<b>Ordinary Least Squares</b> (Linear Model + MSE)	$\hat{\mathbb{Y}} = \mathbb{X}\theta$	$\hat{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}$	<b>Yes</b> , if $\mathbb{X}$ is full col rank (all cols lin independent, #datapoints>> #features)

Lecture 11

# Ordinary Least Squares