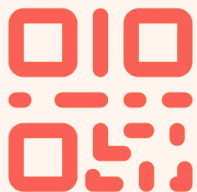slido

Join at slido.com
#2327602

2327602

ⓘ Start presenting to display the joining instructions on this slide.

# Ordinary Least Squares

Using linear algebra to derive the multiple linear regression model.
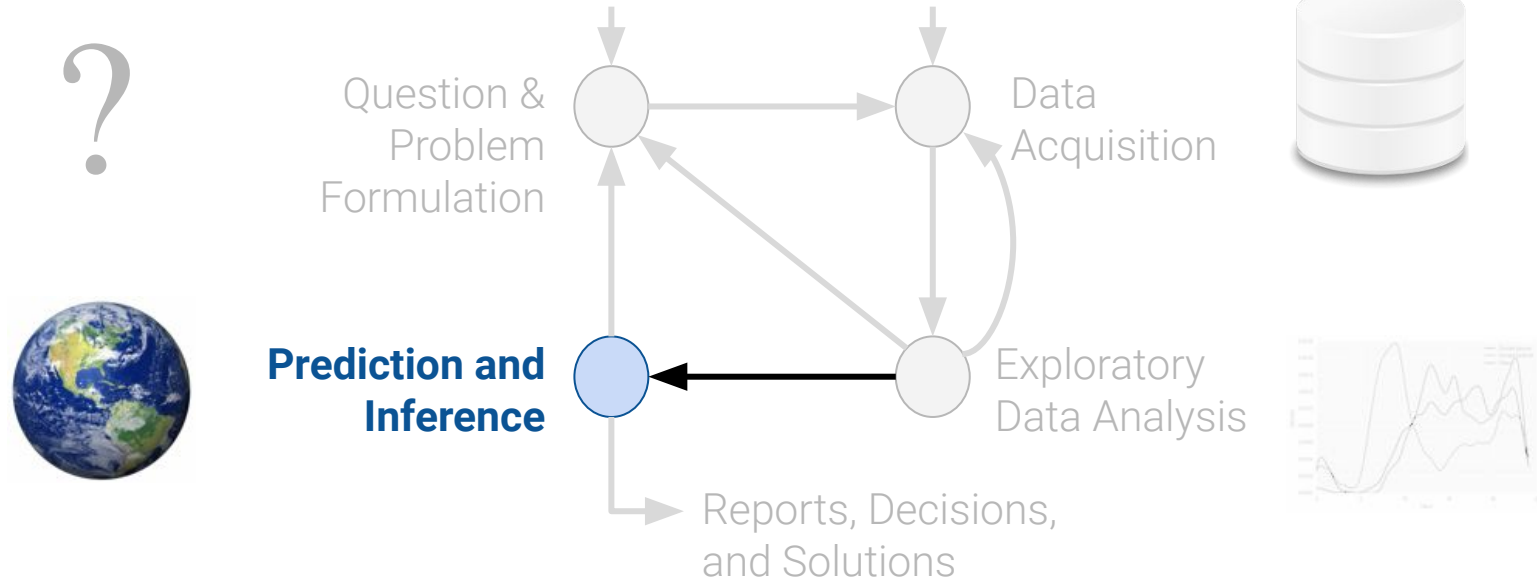
**Data 100/Data 200, Spring 2023 @ UC Berkeley**

Narges Norouzi and Lisa Yan

Content credit: [Acknowledgments](Acknowledgments)

# Plan for Next Few Lectures: Modeling

?

Question & Problem Formulation

Data Acquisition

**Prediction and Inference**

Exploratory Data Analysis

Reports, Decisions, and Solutions

**(today)**

| Modeling I: Intro to Modeling, Simple Linear Regression | Modeling II: Different models, loss functions, linearization | Modeling III: Multiple Linear Regression |

3

- Today's lecture is math heavy.

- If you need more practice, please watch Linear Algebra recordings.

- I will be holding a Regression/Linear Algebra review session **tomorrow at 1 pm**, please do come or watch the recording of it.

# Disclaimer

# Today's Roadmap

Lecture 12, Data 100 Spring 2023

**OLS Problem Formulation**

- Multiple Linear Regression Model
- Mean Squared Error

Geometric Derivation

- Lin Alg Review: Orthogonality, Span
- Least Squares Estimate Proof

Performance: Residuals, Multiple $R^2$

OLS Properties

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

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

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

**2.** $\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 x_3 + \theta_3. \log(x_4)$

**3.** $\hat{y} = \theta_0 + \theta_1 x_1 + \log(\theta_2) x_2 + \theta_3 \theta_4$

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

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

Which of the following expressions are linear in theta?

6

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

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

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

**2.** $\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 x_3 + \theta_3 . \log(x_4)$

$$= \begin{bmatrix} 1 & x_1 & x_2 x_3 & \log(x_4) \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

**3.** $\hat{y} = \theta_0 + \theta_1 x_1 + \log(\theta_2) x_2 + \theta_3 \theta_4$

❌

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

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

"**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.
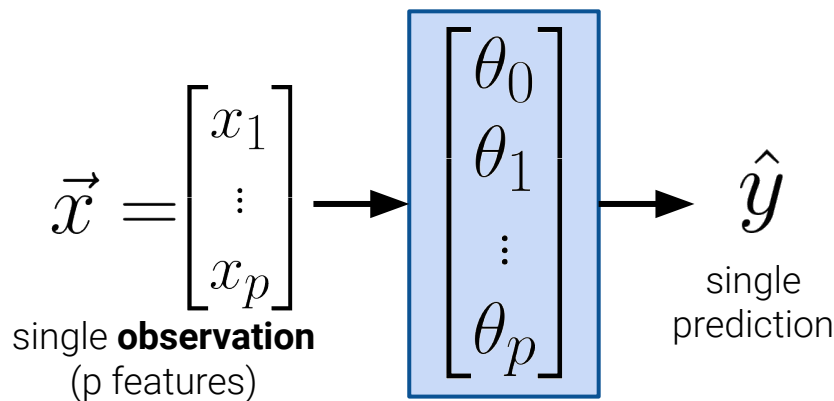
8

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 $y$

This is a linear model because it is a linear combination of parameters $\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$

$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix} \rightarrow \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix} \rightarrow \hat{y}$$

single **observation** (p features)

single prediction
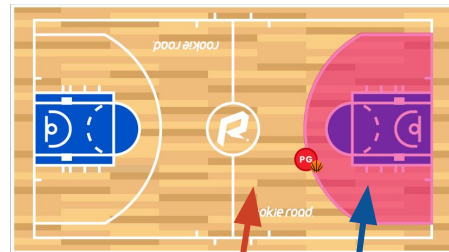
9

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



**3PA**          **FG**

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

|   | 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.

2327602

10

# Multiple Linear Regression Model

How many points does an athlete score per game?
**PTS** (average points/game)
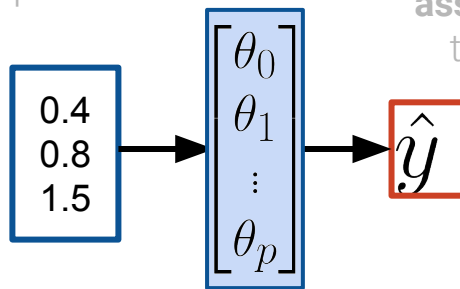
To name a few factors:

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

**3PA**  **FG**

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

| | 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.

$$\begin{bmatrix} 0.4 \\ 0.8 \\ 1.5 \end{bmatrix} \rightarrow \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix} \rightarrow \hat{y}$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p$$
$$= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

**FG**    **AST**    **3PA**

11

# Today's Goal: Ordinary Least Squares

| 1. Choose a model | **Multiple Linear Regression** |
|---|---|

| 2. Choose a loss function | L2 Loss **Mean Squared Error (MSE)** |
|---|---|

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

| 3. Fit the model | Minimize average loss with ~~calculus~~ geometry |
|---|---|

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

| 4. Evaluate model performance | Visualize, ~~Root MSE~~ Multiple $R^2$ |
|---|---|

# Multiple Linear Regression Model

Lecture 12, Data 100 Spring 2023

OLS Problem Formulation

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

Geometric Derivation

- Review: Orthogonality, Span
- Least Squares Estimate Proof

Performance: Residuals, Multiple $R^2$

OLS Properties

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

2327602

**1. Choose a model**

Multiple Linear Regression

For each of our $n$ data points:

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

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

2. Choose a loss function

L2 Loss

Mean Squared Error (MSE)

Linear Algebra!!

3. Fit the model

Minimize average loss with ~~calculus~~ geometry

4. Evaluate model performance

Visualize, ~~Root MSE~~ Multiple $R^2$

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

$$= \theta_0 + \sum_{j=1}^{p} \theta_j x_j$$

**NBA Data**

| | 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 |

To combine the two terms into one matrix operation, we can assume that there is an additional term $x_0 = 1$ in $\vec{x}$ and hence:

Rows correspond to individual players.

$$= x^T \theta \qquad \vec{x} = \begin{bmatrix} \mathbf{1} \\ 0.4 \\ 0.8 \\ 1.5 \end{bmatrix} \qquad \vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \qquad \vec{x}, \vec{\theta} \in \mathbb{R}^{(p+1)}$$

$$= \begin{bmatrix} \mathbf{1} & 0.4 & 0.8 & 1.5 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} = \boxed{\hat{y}} \in \mathbb{R}$$

Note that:

$$\vec{x}^T \times \vec{\theta} = \vec{x}.\vec{\theta}$$

15

# Matrix Notation

| | 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 |

2327602

To make predictions on all $n$ datapoints in our sample:

$$\hat{y}_1 = x_1^T \theta \qquad \text{where } x_1^T = \begin{bmatrix} 1 & x_{11} & x_{12} \ldots & x_{1p} \end{bmatrix} \quad \text{Datapoint 1}$$

$$\hat{y}_2 = x_2^T \theta \qquad \text{where } x_2^T = \begin{bmatrix} 1 & x_{22} & x_{22} & \ldots & x_{2p} \end{bmatrix} \quad \text{Datapoint 2}$$

$$\vdots \qquad \vdots \qquad\qquad\qquad\qquad \vdots$$

$$\hat{y}_n = x_n^T \theta \qquad \text{where } x_n^T = \begin{bmatrix} 1 & x_{n1} & x_{n2} & \ldots & x_{np} \end{bmatrix} \quad \text{Datapoint n}$$

same
$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$
for all
preds

16

# Matrix Notation

To make predictions on all $n$ datapoints in our sample:

$$\hat{y}_1 = \begin{bmatrix} 1 & x_{11} & x_{12} \dots & x_{1p} \end{bmatrix} \theta \quad = x_1^T \theta$$

$$\hat{y}_2 = \begin{bmatrix} 1 & x_{22} & x_{22} & \dots & x_{2p} \end{bmatrix} \theta \quad = x_2^T \theta$$

$$\vdots \qquad \qquad \vdots \qquad \qquad \vdots$$

$$\hat{y}_n = \begin{bmatrix} 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \theta \quad = x_n^T \theta$$

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

Expand out each data
point's (transposed) input

**Data**

| | 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 |

2327602

same
$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$
for all
preds

17

# Matrix Notation

To make predictions on all $n$ datapoints in our sample:

**Data**

| | 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 |

2327602

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

same $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$ for all preds

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

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

18

# Matrix Notation

To make predictions on all $n$ datapoints in our sample:

| | 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 |

2327602

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

**Design matrix** with
dimensions n x (p + 1)

same
$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$
for all
preds

19

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**:

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

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

2327602

| | | 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

20

# 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**:

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

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

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

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



| 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 |

Field Goals, Assists, 3-Point Attempts

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

21

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

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

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

**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:

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

# Mean Squared Error

Lecture 12, Data 100 Spring 2023

OLS Problem Formulation

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

Geometric Derivation

- Lin Alg Review: Orthogonality, Span
- Least Squares Estimate Proof

Performance: Residuals, Multiple $R^2$

OLS Properties

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

# Today's Goal: Ordinary Least Squares

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

1. Choose a model ✅

Multiple Linear Regression

**2. Choose a loss function**

L2 Loss

Mean Squared Error (MSE)

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

More Linear Algebra!!

3. Fit the model

Minimize average loss with ~~calculus~~ geometry

4. Evaluate model performance

Visualize,
~~Root MSE~~
Multiple $R^2$

24

2327602

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

- 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

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

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

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

In $\mathbb{R}^2$



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

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

In $\mathbb{R}^n$



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

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

Looks like Mean Squared Error!!

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

$$R(\theta) = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y_i})^2$$

$$= \frac{1}{n} ||\mathbb{Y} - \hat{\mathbb{Y}}||_2^2$$

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

2327602

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, $e = \mathbb{Y} - \hat{\mathbb{Y}} = \begin{bmatrix} y_1 - \hat{y_1} \\ y_2 - \hat{y_2} \\ \vdots \\ y_n - \hat{y_n} \end{bmatrix}$

Important
for today

**D.** All of the above

**E.** Something else

30

# Geometric Derivation

Lecture 12, Data 100 Spring 2023

2327602

2327602

| | | |
|---|---|---|
| 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 |

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**.

| | |
|---|---|
| 4. Evaluate model performance | Visualize, ~~Root MSE~~ Multiple $R^2$ |

32

The set of all possible linear combinations of the columns of $\mathbb{X}$ is called the **span** of the columns of $\mathbb{X}$ (denoted $span(\mathbb{X})$), 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 $n$, $span(\mathbb{X})$ is a subspace of $\mathbb{R}^n$



$\mathbb{X}_{:,1}$   $\mathbb{X}_{:,2}$

Subspace of $\mathbb{R}^n$ spanned by $\mathbb{X}$

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

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

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

We can also think of $\hat{\mathbb{Y}}$ as a **linear combination of feature vectors**, scaled by **parameters**.

$$n\begin{bmatrix} | \\ \hat{\mathbb{Y}} \\ | \\ \scriptstyle 1 \end{bmatrix} = n\begin{bmatrix} | & | \\ \mathbb{X}_{:,1} & \mathbb{X}_{:,2} \\ | & | \\ & \scriptstyle p+1 \end{bmatrix} \begin{bmatrix} | \\ \theta \\ | \\ \scriptstyle 1 \end{bmatrix}^{p+1} = \theta_1 \mathbb{X}_{:,1} + \theta_2 \mathbb{X}_{:,2} \; + \; \ldots$$
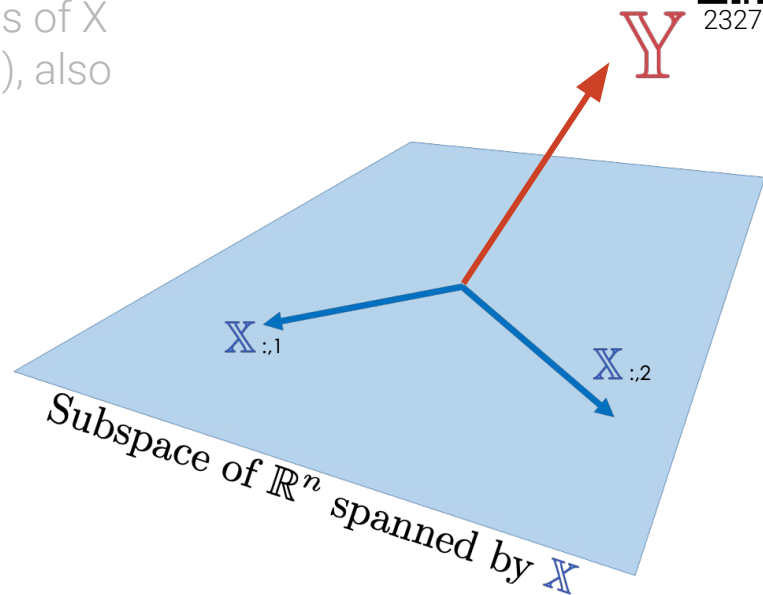
# A Linear Combination of Columns

The set of all possible linear combinations of the columns of X is called the **span** of the columns of X (denoted $span(\mathbb{X})$), 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 $n$, $span(\mathbb{X})$ 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 $\hat{\mathbb{Y}} \in span(\mathbb{X})$.

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

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



$\mathbb{Y}$

$\mathbb{X}_{:,1}$    $\mathbb{X}_{:,2}$

Subspace of $\mathbb{R}^n$ spanned by $\mathbb{X}$

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

$n$ ... $1$

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

$\mathbb{Y}$

$\mathbb{Y}$
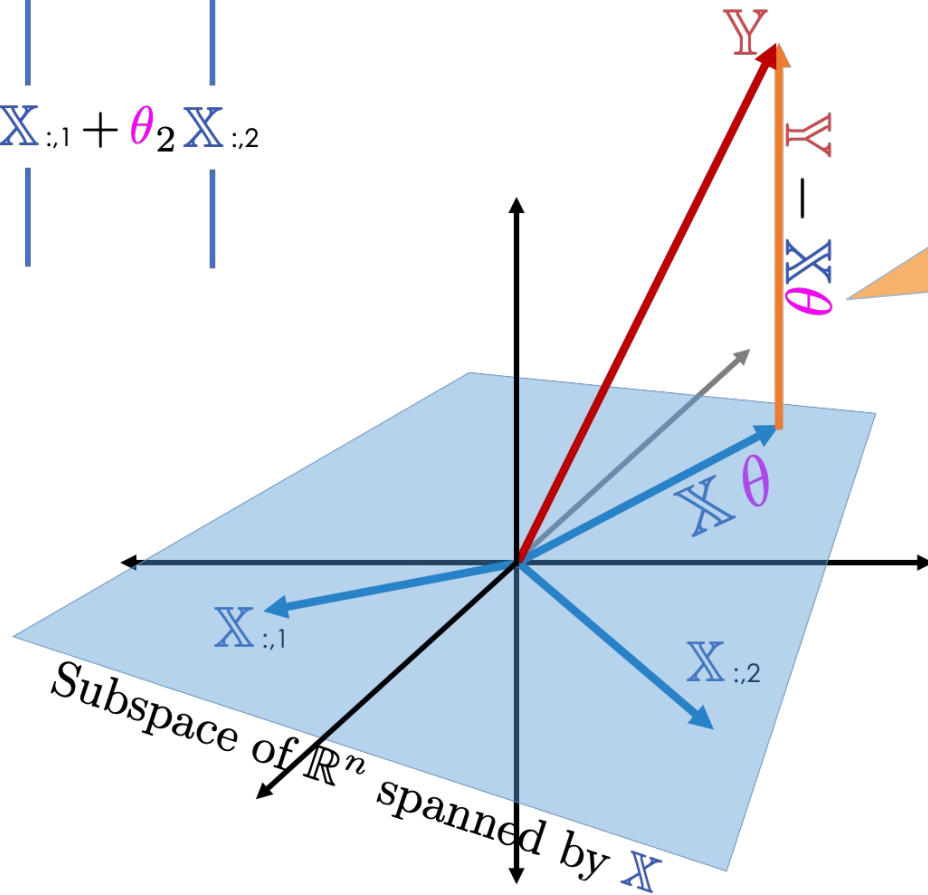
$-$

$\mathbb{X}\theta$

$\mathbb{X}\theta$

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.

$\mathbb{X}_{:,1}$

$\mathbb{X}_{:,2}$

Subspace of $\mathbb{R}^n$ spanned by $\mathbb{X}$

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

36

37

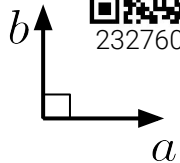$$\begin{bmatrix} | \\ \hat{\mathbb{Y}} \\ | \end{bmatrix} = \theta_1 \mathbb{X}_{:,1} + \theta_2 \mathbb{X}_{:,2}$$

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 $span(\mathbb{X})$.

Subspace of $\mathbb{R}^n$ spanned by $\mathbb{X}$

We will not prove this property of orthogonal projection: see Khan Academy.

38

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

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 $span(\mathbb{X})$.

We will not prove this property of orthogonal projection: see Khan Academy.

Subspace of $\mathbb{R}^n$ spanned by $\mathbb{X}$

$\mathbb{Y}$

$\mathbb{Y} - \mathbb{X}\theta$

$\mathbb{X}\theta$

$\mathbb{X}_{:,1}$
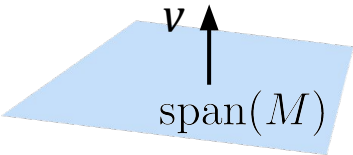
$\mathbb{X}_{:,2}$

2327602

39

**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 $\mathrm{span}(M)$, 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 $v \in \mathbb{R}^{n \times 1}$, $M \in \mathbb{R}^{n \times d}$ where $M = \begin{bmatrix} | & | & & | \\ m_1 & m_2 & \dots & m_d \\ | & | & & | \end{bmatrix}$:

$$m_1^T v = 0$$
$$m_2^T v = 0$$
$$\vdots$$
$$m_d^T v = 0$$

$$\begin{bmatrix} m_1^T v \\ m_2^T v \\ \vdots \\ m_d^T v \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

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

$v$ is orthogonal to each column of $M$, $m_j \in \mathbb{R}^{n \times 1}$

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

40

2327602

$b$

$a$

$v$

$\mathrm{span}(M)$

2327602

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 $\longrightarrow$ $M^T v = \vec{0}$ $\longleftarrow$ Residual vector

Equivalently, this is the $\hat{\theta}$ such that the residual vector $\mathbb{Y} - \mathbb{X}\hat{\theta}$ is orthogonal to $span(\mathbb{X})$.

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

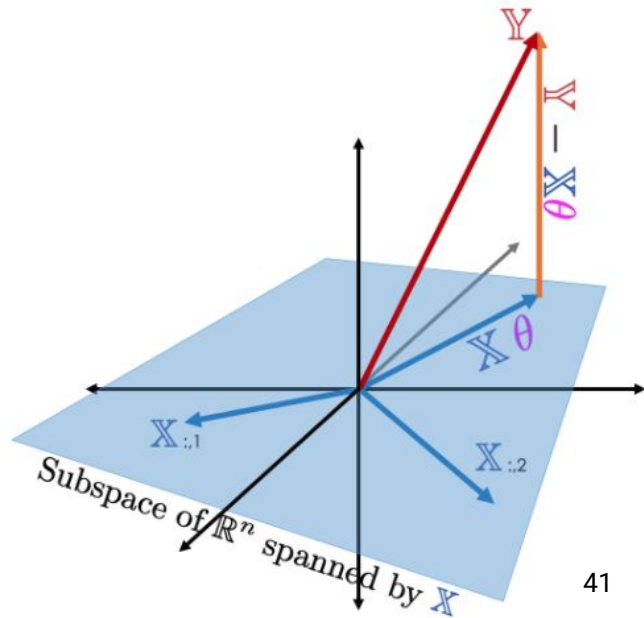$$\mathbb{X}^T\left(\mathbb{Y} - \mathbb{X}\hat{\theta}\right) = 0$$

Rearrange terms

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

The **normal equation**

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

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

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



Subspace of $\mathbb{R}^n$ spanned by $\mathbb{X}$

41

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

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

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

🎉 🎉 🎉

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

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

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

# Least Squares Estimate

| 1. Choose a model | Multiple Linear Regression | $$\hat{\mathbb{Y}} = \mathbb{X}\theta$$ |

| 2. Choose a loss function | L2 Loss<br><br>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,<br>~~Root MSE~~<br>Multiple $R^2$ | |

# Performance

Lecture 12, Data 100 Spring 2023

# 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 $R^2$

46

2327602

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

**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:

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

**Demo**

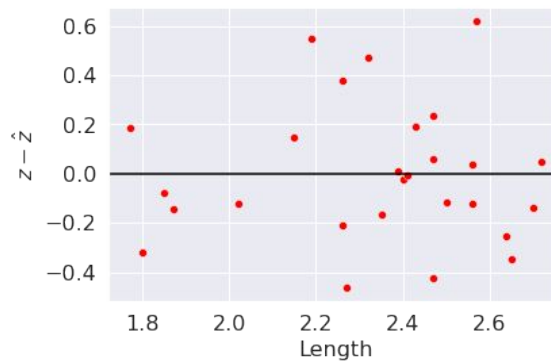$$R(\theta) = \frac{1}{n}||\mathbb{Y} - \mathbb{X}\theta||_2^2$$

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

47

**Simple linear regression**
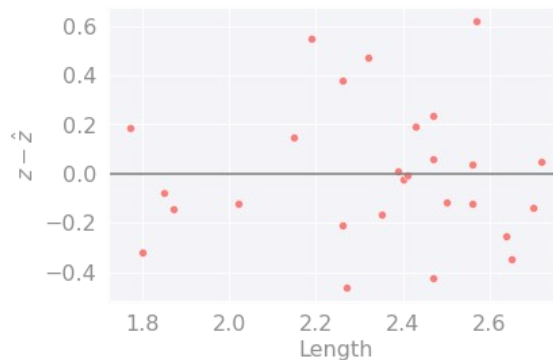
Plot residuals vs
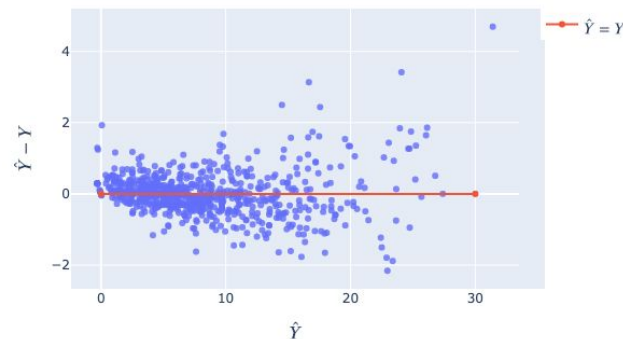   the single feature $x$.



**Compare**

# [Visualization] Residual Plots

## Compare

See notebook

### Simple linear regression

Plot residuals vs
the single feature $x$.



### Multiple linear regression

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



Same interpretation as before  (Data 8 [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.

2327602

2327602

**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$^2$**, 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**

# [Metrics] Multiple R^2

We define the **multiple R²** 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**.

$R^2$ 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. $\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_p x_p$),

$R^2 = [r(y, \hat{y})]^2$ is equal to the square of correlation between $y, \hat{y}$.

- For SLR, $R^2 = r^2$, the correlation between $x, y$.
- The proof of these last two properties is beyond this course. 51

2327602

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

R² = 0.457

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

R² = 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 $R^2$**, 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, R² 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 are a few weeks away from understanding why.

52

# OLS Properties

Lecture 12, Data 100 Spring 2023

OLS Problem Formulation

- Multiple Linear Regression Model
- Mean Squared Error

Geometric Derivation

- Lin Alg Review: Orthogonality, Span
- Least Squares Estimate Proof

Performance: Residuals, Multiple $R^2$

**OLS Properties**

**We will cover in the review session on Friday**

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

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

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

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

For all linear models:

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

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

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

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

You will prove both properties in homework.

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

54

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

$$\sum_{i=1}^{n} e_i = 0 \quad \text{(previous slide)}$$

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

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

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

56

2327602

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

In most settings,
**# observations        # features**
**n       >>       p**

2327602

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**:

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



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 estimate $\hat{\theta}$.

- $\hat{\theta}$ 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) x (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).

2327602

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 >> p so this is less of an issue.

    p + 1 features

    n data points

    $\mathbb{X}$

2.  If we 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,
        ■  Perimeter = 2 * Width + 2 * Height  →  $\mathbb{X}$ is not full rank.
    ○  Important with one-hot encoding (to discuss in later).

60

# Does a Unique Solution Always Exist?

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

61

# Ordinary Least Squares

Content credit: [Acknowledgments](Acknowledgments)