

LECTURE 13

# Ordinary Least Squares

Using linear algebra to generalize the simple linear regression model.

**Data 100/Data 200, Fall 2021 @ UC Berkeley**

Fernando Pérez and Alvin Wan

(Content by Suraj Rampure, Anil Adhikari, Deborah Nolan, Joseph Gonzalez)

# Recap: Simple Linear Regression

# Simple Linear Regression

In the last lecture, we re-introduced the simple linear regression model from Data 8.

$$\hat{y} = f_{\theta}(x) = \theta_0 + \theta_1 x$$

- Our loss function was squared loss, and so our objective function was mean squared error (MSE).
- To solve for the **optimal parameters** (also known as coefficients or weights), we minimized MSE by hand using calculus.

$$\hat{\theta}_1 = r \frac{\sigma_y}{\sigma_x} \qquad \hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$$

- These are parameter estimates.
- We also looked at  $r$ , the correlation coefficient, and its relation to the optimal coefficients.

# Multiple Regression

We also extended this model to account for multiple **features**.

$$\hat{y} = f_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p = \theta_0 + \sum_{j=1}^p \theta_j x_j$$

Each  $x_j$  is a separate feature.

- We learned about multiple  $R^2$ , an extension of the correlation coefficient  $r$  to multiple features.
- Our loss function yet again was squared loss.
- We didn't minimize MSE by hand – we abstracted away the process of determining the theta values.
- Lastly, we introduced RMSE as a method of comparing model performance.

Today, we will learn how to find the optimal parameters (“thetas”) for multiple regression, for any number of features.

# Agenda

- Use vector dot products to define the multiple regression model.
- Formulate the problem statement using vector norms.
- Use a geometric derivation to solve for the optimal  $\theta$  (which is now a vector).
- Explore properties of residuals.
- Understand when a unique solution exists.

Lots of linear algebra!

- There is a resources post on Piazza if you need a refresher.
- We will try and take things slowly.

Linear algebra formulation

# Dot products

The dot product of two vectors  $\vec{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$   $\vec{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$  is defined as follows:

$$\vec{a} \cdot \vec{b} = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$$

- An alternate way of writing a dot product:  $\vec{a} \cdot \vec{b} = \vec{a}^T \vec{b}$ .
  - This is the form we will use primarily moving forward.
- The dot product between two vectors is a **scalar**, not another vector.
- The dot product is only defined for two vectors of the same length.
- The dot product is a special case of the inner product.

# Multiple regression as a dot product

We previously stated that the multiple regression model was of the form

$$\hat{y} = f_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p = \theta_0 + \sum_{j=1}^p \theta_j x_j$$

This can be restated as a dot product between two vectors.

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}$$

$$\hat{y} = x^T \theta$$

scalar

Even though they don't have arrows on top of them, **x** and **θ** are still **vectors**!



# Design matrix

Our mean squared error involves all observations at once, so it may be valuable to express our model in terms of all observations, instead of just one observation. One step in that process is stacking all of our observations together into a **design matrix**. With  $n$  observations:

$$\mathbb{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & x_{33} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{np} \end{bmatrix}$$

**Rows** correspond to **observations**.

e.g. all features for data point 3

**Columns** correspond to **features**.

e.g. feature 1, for all data points

# Example design matrix

|     | Bias | FG  | AST | PTS  |
|-----|------|-----|-----|------|
| 0   | 1    | 1.8 | 0.6 | 5.3  |
| 1   | 1    | 0.4 | 0.8 | 1.7  |
| 2   | 1    | 1.1 | 1.9 | 3.2  |
| 3   | 1    | 6.0 | 1.6 | 13.9 |
| 4   | 1    | 3.4 | 2.2 | 8.9  |
| ... | ...  | ... | ... | ...  |
| 703 | 1    | 4.0 | 0.8 | 11.5 |
| 704 | 1    | 3.1 | 0.9 | 7.8  |
| 705 | 1    | 3.6 | 1.1 | 8.9  |
| 706 | 1    | 3.4 | 0.8 | 8.5  |
| 707 | 1    | 3.8 | 1.5 | 9.4  |

708 rows × 4 columns

Here,  $n = 708$ , and  $p = 3$ .

- Each of the 4 columns corresponds to a different feature.
- Each of the 708 rows corresponds to a different data point.

# Multiple regression as a matrix multiplication

This allows us to express our linear model on our **entire dataset** (not just one observation) as

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

**Vector** with dimensions  $\mathbf{n} \times 1$ .

**Matrix** with dimensions  $\mathbf{n} \times (\mathbf{p} + 1)$ .

**Vector** with dimensions  $(\mathbf{p} + 1) \times 1$ .

Note: This means that  $\mathbf{Y}$  represents an  $\mathbf{n}$ -length vector containing all of our true  $\mathbf{y}$  values.

# Multiple regression as a matrix multiplication

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & x_{33} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_p \end{bmatrix}$$

Moving forward:

- $\mathbb{X}_{i,:}$  or just simply  $\mathbb{X}_i$  will represent **row i** of our design matrix. **Rows are observations / data points.**
- $\mathbb{X}_{:,j}$  will represent **column j** of our design matrix. **Columns are features.**

# Multiple regression as a matrix multiplication

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & x_{33} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_p \end{bmatrix}$$

For instance, to compute the predicted output for the second observation:

$$\hat{y}_2 = \theta_0 + \theta_1 x_{21} + \theta_2 x_{22} + \theta_3 x_{23} + \dots + \theta_p x_{2p}$$

Or, more compactly:

$$\hat{y}_2 = \mathbb{X}_2^T \theta$$

$\mathbb{X}_2$  is a vector, not a matrix. That is why it is transposed.

# Example

Consider the following design matrix and value of  $\theta$ .

| Bias | FG  | AST | PTS  |
|------|-----|-----|------|
| 1    | 1.8 | 0.6 | 5.3  |
| 1    | 0.4 | 0.8 | 1.7  |
| 1    | 1.1 | 1.9 | 3.2  |
| 1    | 6.0 | 1.6 | 13.9 |
| 1    | 3.4 | 2.2 | 8.9  |
| ...  | ... | ... | ...  |
| 1    | 4.0 | 0.8 | 11.5 |

$$\theta = \begin{bmatrix} -5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

The predicted **response** (output) for the second observation:

$$\begin{aligned} \hat{y}_2 &= \begin{bmatrix} 1 & 0.4 & 0.8 & 1.7 \end{bmatrix} \begin{bmatrix} -5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \\ &= 1(-5) + 0.4(2) + 0.8(-1) + 1.7(3) \\ &= 0.1 \end{aligned}$$

# Example

Consider the following design matrix and value of  $\theta$ .

| Bias | FG  | AST | PTS  |
|------|-----|-----|------|
| 1    | 1.8 | 0.6 | 5.3  |
| 1    | 0.4 | 0.8 | 1.7  |
| 1    | 1.1 | 1.9 | 3.2  |
| 1    | 6.0 | 1.6 | 13.9 |
| 1    | 3.4 | 2.2 | 8.9  |
| ...  | ... | ... | ...  |
| 1    | 4.0 | 0.8 | 11.5 |

$$\theta = \begin{bmatrix} -5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

The predicted **response vector**:

$$\hat{\mathbf{Y}} = \begin{array}{c|c|c|c|c} \text{Bias} & \text{FG} & \text{AST} & \text{PTS} & \hat{Y} \\ \hline 1 & 1.8 & 0.6 & 5.3 & 13.9 \\ 1 & 0.4 & 0.8 & 1.7 & 0.1 \\ 1 & 1.1 & 1.9 & 3.2 & 4.9 \\ 1 & 6.0 & 1.6 & 13.9 & 47.1 \\ 1 & 3.4 & 2.2 & 8.9 & 26.3 \\ \dots & \dots & \dots & \dots & \dots \\ 1 & 4.0 & 0.8 & 11.5 & 36.7 \end{array} \begin{bmatrix} -5 \\ 2 \\ -1 \\ 3 \end{bmatrix} =$$

# Summary of notation

When looking at a **single observation**,  
our model is

$$\hat{y} = f_{\theta}(x) = x^T \theta$$

- $x$  is a **vector** of size  $p + 1$ .
- $\hat{y}$  is a **scalar**.
- $\theta$  is a **vector** of size  $p + 1$ .

When looking at **multiple observations**,  
our model is

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

- $\mathbf{X}$  is a **matrix** of size  $n \times (p + 1)$ .
- $\hat{\mathbf{Y}}$  is a **vector** of size  $n$  (i.e.  $\hat{\mathbf{Y}} \in \mathbb{R}^n$ ).
- $\theta$  is a **vector** of size  $p + 1$ .

In many settings, we assume that we have only  $p$  (and not  $p + 1$ ) columns. One of those  $p$  columns may be “1” for each observation.



Problem statement

# Vector norms

- The **norm** of a vector is some measure of that vector's size.
  - More formal definitions exist, but we don't really need them in Data 100.
- The two norms we need to know are the  $L_1$  and  $L_2$  norms (surprise surprise!).
  - $L_2$  will be used today.
  - $L_1$  will appear a few lectures from now.

Consider the vector:

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

$L_2$  vector norm:

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2} = \sqrt{\sum_{i=1}^n x_i^2}$$

$L_1$  vector norm:

$$\|x\|_1 = |x_1| + |x_2| + |x_3| + \dots + |x_n| = \sum_{i=1}^n |x_i|$$

## $L_2$ vector norm

- The  $L_2$  vector norm can be thought of as the “length” of a vector.
  - It is a generalization of the Pythagorean theorem into  $n$  dimensions.

$$||x||_2 = \sqrt{x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2} = \sqrt{\sum_{i=1}^n x_i^2}$$

- The “distance” between two vectors is the  $L_2$  norm of their difference.
  - For instance, if  $a$  and  $b$  are two vectors of the same length, then their distance is

$$||a - b||_2$$

- Note, the square of the  $L_2$  norm of a vector is the sum of the squares of the vector's elements:

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

# Residuals

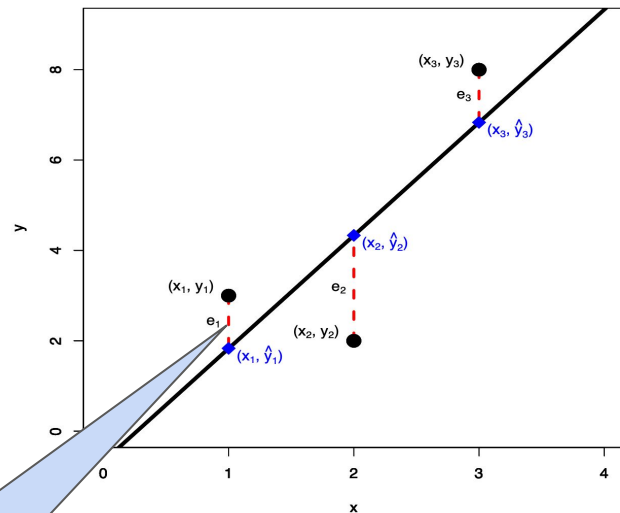
Residuals are defined as being the difference between an actual and predicted value, in the regression context.

- We use the letter  $e$  to denote residuals. The residual  $i$  is

$$e_i = y_i - \hat{y}_i$$

- The MSE of a model is equal to the mean of the squares of its residuals:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n e_i^2$$



# Residual vector

We can stack all  $n$  residuals into a vector, called the residual vector,  $e$ :

$$e = \mathbb{Y} - \hat{\mathbb{Y}} = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix}$$

The residual vector is the “difference” between the two vectors containing our true  $y$  values and predicted  $y$  values.

## Mean squared error, again

We are choosing our loss to be squared loss. This means, the average loss across our dataset is mean squared error.

$$R(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbb{X}_i^T \theta)^2$$

We can write this in terms of the norm of the residual vector!

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

This is the residual vector!

# Optimization procedure

As we did in the last lecture, we note that the value of  $\theta$  that minimizes  $\frac{1}{n} ||\mathbb{Y} - \mathbb{X}\theta||_2^2$  is the same value that minimizes

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

**Therefore, our goal is to find the value of  $\theta$  that minimizes the squared  $L_2$  norm of the residual vector.** In other words, we want the “distance” between  $\mathbb{Y}$  and  $\hat{\mathbb{Y}}$  to be minimized.

There are two ways we can determine the optima  $\hat{\theta}$  here.

- Using calculus, like we’ve done earlier.
  - Requires a lot of matrix calculus. Out of scope, but here’s a [link](#) if you’re interested.
- Using a **geometric argument**. This is what we’ll do here.

# Geometric derivation



A linear combination of columns

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

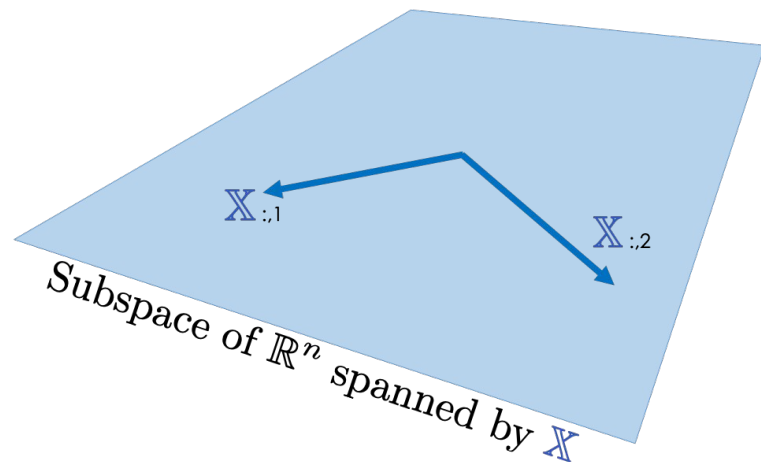
$$\begin{matrix} n \\ \left[ \begin{array}{c} | \\ \hat{\mathbf{Y}} \\ | \end{array} \right] \\ 1 \end{matrix} = \begin{matrix} n \\ \left[ \begin{array}{cc} | & | \\ \mathbf{X}_{:,1} & \mathbf{X}_{:,2} \\ | & | \end{array} \right] \\ p+1 \end{matrix} \begin{matrix} \left[ \begin{array}{c} | \\ \boldsymbol{\theta} \\ | \end{array} \right] \\ p+1 \\ 1 \end{matrix} = \theta_1 \begin{matrix} | \\ \mathbf{X}_{:,1} \\ | \end{matrix} + \theta_2 \begin{matrix} | \\ \mathbf{X}_{:,2} \\ | \end{matrix}$$

The linear model represents  $\hat{\mathbf{Y}}$  as a linear combination of the columns of  $\mathbf{X}$ .

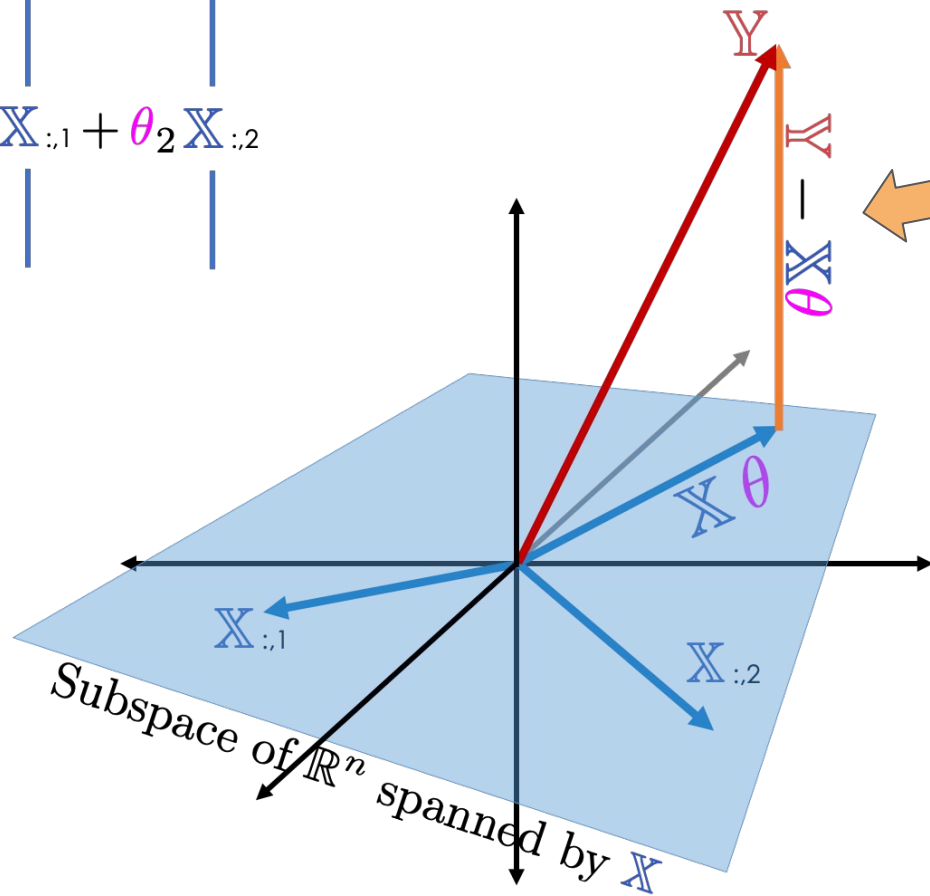
# Span

**Our prediction is a linear combination of the columns of  $X$ .**

- The set of all possible linear combinations of the columns of  $X$  is called the **span** of the columns of  $X$  (denoted  $\text{span}(X)$ ).
  - Also called the **column space**.
- Intuitively, this is all of the vectors you can “reach” using the columns of  $X$ .
- Since each column of  $X$  has length  $n$ ,  $\text{span}(X)$  is a subspace of  $\mathbb{R}^n$ .
- Our goal is to find the vector in  $\text{span}(X)$  that is closest to  $Y$ .



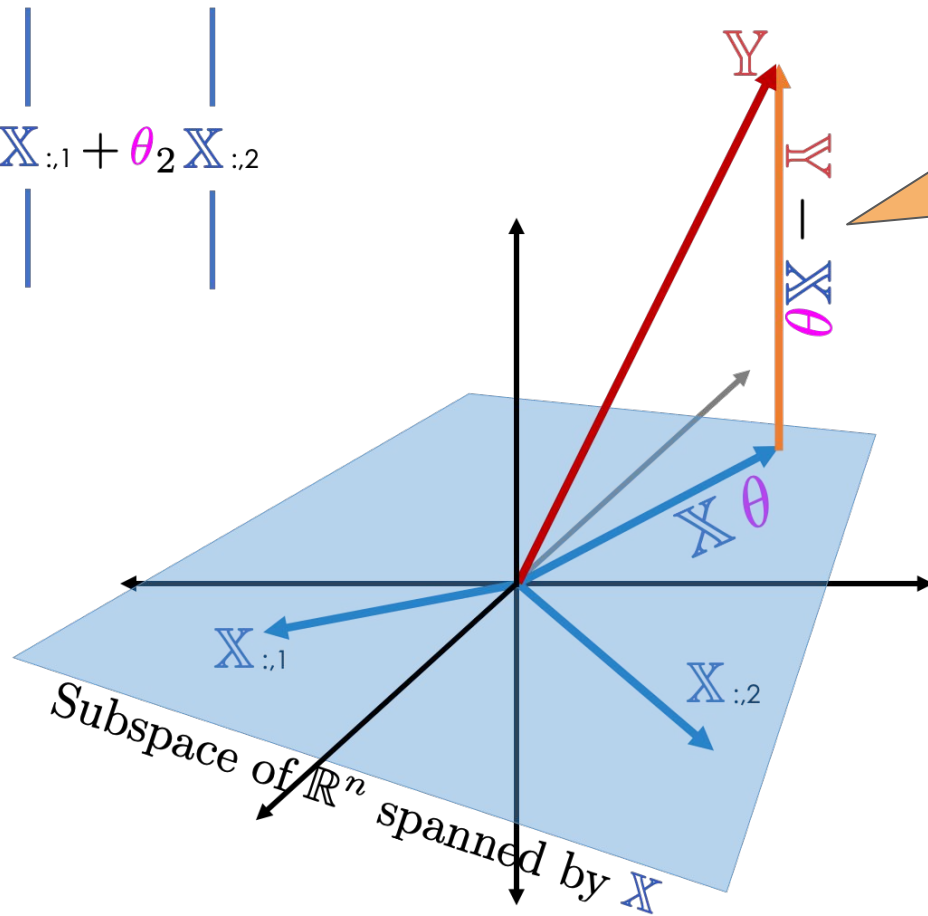
$$\begin{matrix} n \\ \left[ \begin{array}{c} | \\ \hat{Y} \\ | \end{array} \right] \\ 1 \end{matrix} = \theta_1 \begin{matrix} | \\ X_{:,1} \\ | \end{matrix} + \theta_2 \begin{matrix} | \\ X_{:,2} \\ | \end{matrix}$$



Recall, this is the residual vector,  $e = Y - \hat{Y}$ .

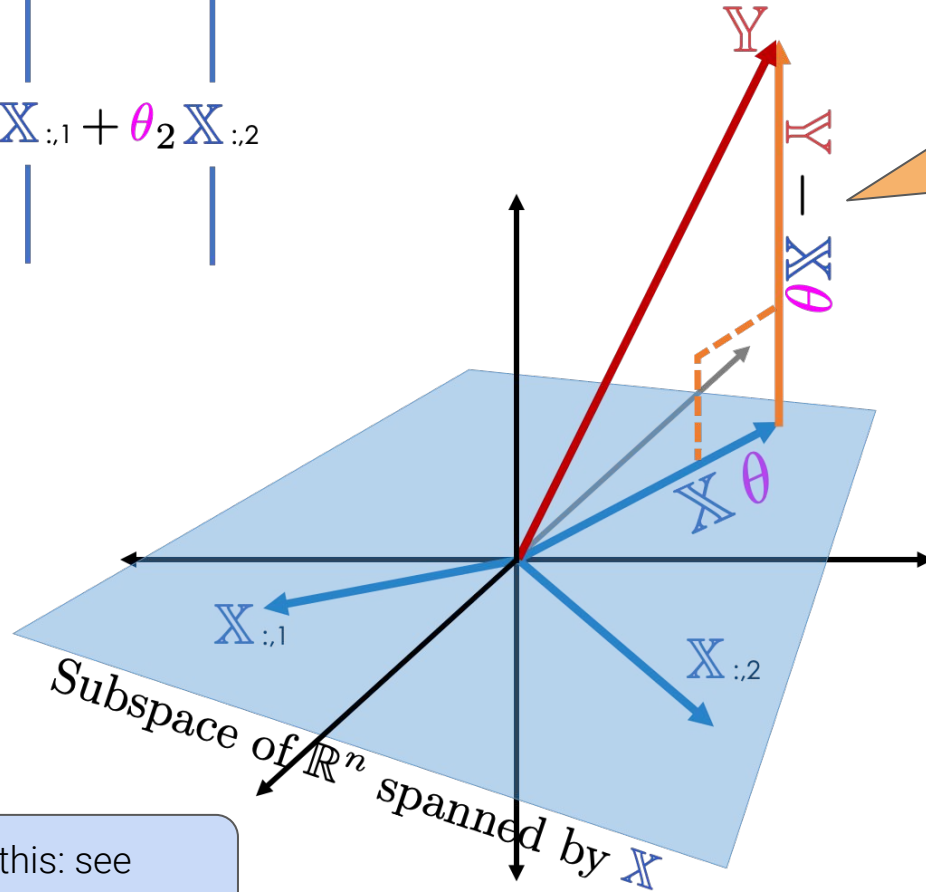
Our goal is to minimize the  $L_2$  norm of the residual vector, i.e. we want our predictions to be “as close” to our true  $y$  values as possible.

$$\begin{matrix} n \\ \left[ \begin{array}{c} | \\ \hat{Y} \\ | \end{array} \right] \\ 1 \end{matrix} = \theta_1 \begin{matrix} | \\ X_{:,1} \\ | \end{matrix} + \theta_2 \begin{matrix} | \\ X_{:,2} \\ | \end{matrix}$$



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

$$\begin{bmatrix} \vdots \\ \hat{Y} \\ \vdots \end{bmatrix}_n = \theta_1 \begin{bmatrix} \vdots \\ X_{:,1} \\ \vdots \end{bmatrix} + \theta_2 \begin{bmatrix} \vdots \\ X_{:,2} \\ \vdots \end{bmatrix}$$

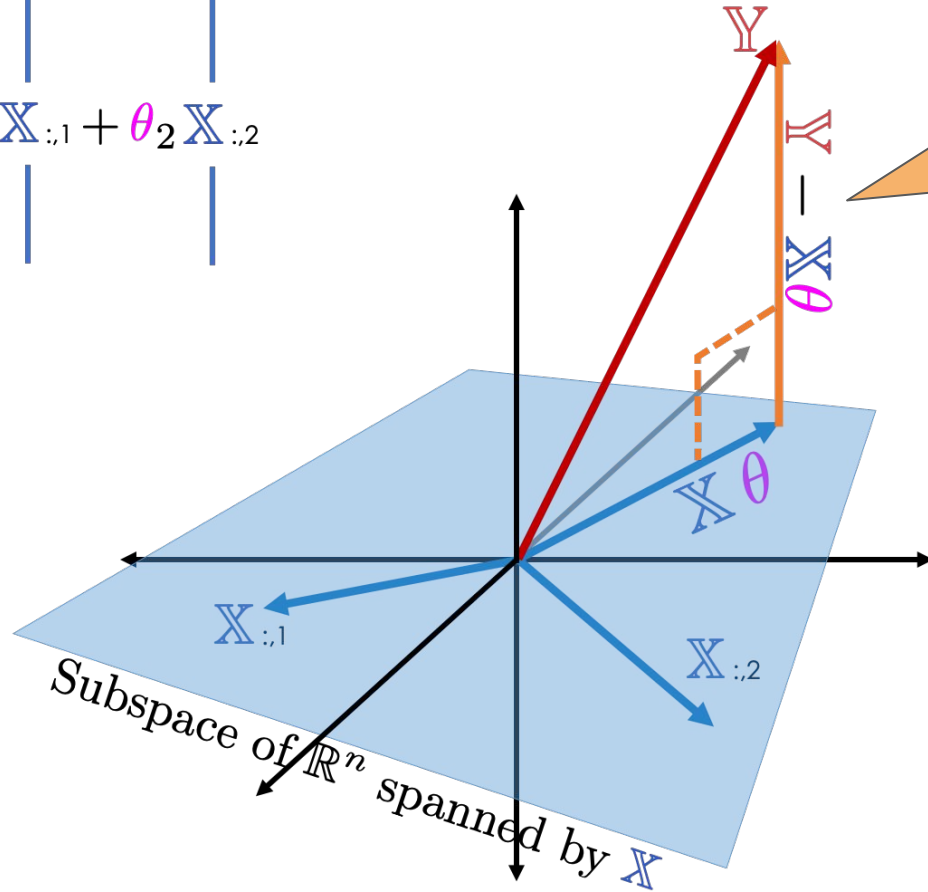


How do we minimize this distance – the norm of the residual vector (squared)?

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

We will not prove this: see [Khan Academy](https://www.khanacademy.org/linear-algebra/a/orthogonal-projection/v/orthogonal-projection).

$$\begin{bmatrix} \vdots \\ \hat{Y} \\ \vdots \end{bmatrix}_n = \theta_1 \begin{bmatrix} \vdots \\ X_{:,1} \\ \vdots \end{bmatrix} + \theta_2 \begin{bmatrix} \vdots \\ X_{:,2} \\ \vdots \end{bmatrix}$$



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

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

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

# Orthogonality

We say two vectors are **orthogonal** if and only if their dot product is 0.

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

$$a^T b = 0 \iff a, b \text{ are orthogonal}$$

Suppose a vector is orthogonal to the span of the columns of a matrix.

- This is true if and only if it is orthogonal to each column individually. ([proof](#))
- Let  $M \in \mathbb{R}^{n \times d}$ ,  $v \in \mathbb{R}^{n \times 1}$ . Suppose  $v$  is orthogonal to the span of the columns of  $M$ .

Then:

$$M = \begin{bmatrix} | & | & & | \\ m_1 & m_2 & \dots & m_d \\ | & | & & | \end{bmatrix} \quad \begin{matrix} m_1^T v = 0 \\ m_2^T v = 0 \\ \vdots \\ m_d^T v = 0 \end{matrix} \quad \Rightarrow \quad M^T v = \vec{0}$$

Let  $\mathbf{M} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{v} \in \mathbb{R}^{n \times 1}$ .

Suppose  $\mathbf{v}$  is orthogonal to the span of the columns of  $\mathbf{M}$ .

$$\mathbf{M} = \begin{bmatrix} | & | & & | \\ m_1 & m_2 & \dots & m_d \\ | & | & & | \end{bmatrix}$$

$$m_1^T \mathbf{v} = 0$$

$$m_2^T \mathbf{v} = 0$$

$$\vdots$$

$$m_d^T \mathbf{v} = 0$$



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

$\mathbf{v}$  is orthogonal to each column of  $\mathbf{M}$  separately.

(Note, each column of  $\mathbf{M}$  has length  $n$ , and  $\mathbf{v}$  also has length  $n$ ).

This product encapsulates all  $d$  of the equations on the left into a single equation. The quantity on the right is the **zero vector** ( $d$ -length vector full of 0s).



# Residuals are orthogonal to the span of $\mathbf{X}$

We want the  $\theta$  such that the residual vector is orthogonal to  $\text{span}(\mathbf{X})$ .

By the definition of orthogonality:  $\mathbf{X}^T (\mathbf{Y} - \mathbf{X}\hat{\theta}) = 0$

*residual vector*

Still the zero vector!

Rearranging:

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

The **normal equation**:

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

Assuming  $\mathbf{X}^T \mathbf{X}$  is full rank:

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

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

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

It is the **least squares estimate** for  $\theta$ .

# Summary

1. **Choose a model.** We chose the multiple linear regression model, formulated using a matrix.

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

2. **Choose a loss function.** We chose squared loss, and hence our average loss was

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

3. **Minimize average loss to determine optimal model parameters.** Done!

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

# Residuals

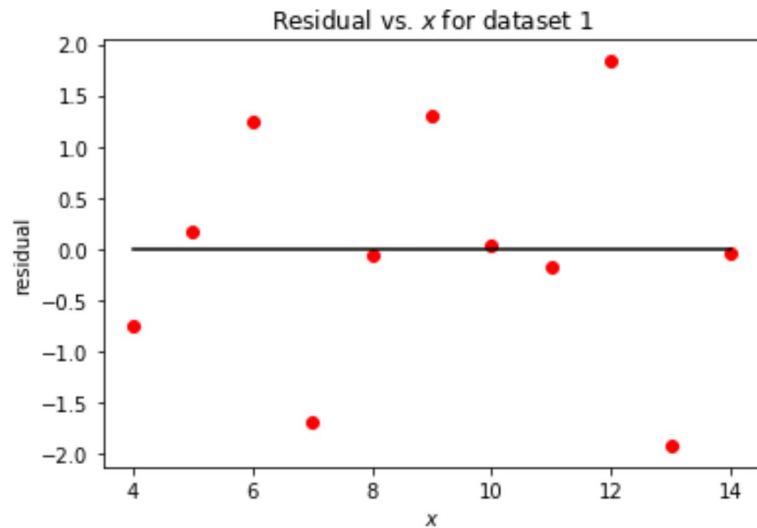
# Residual plots

Residual plots can tell us about the quality of our model.

- In the **simple linear regression** case, with only one independent variable, we typically plot residuals vs.  $x$ .
- More generally, a residual plot is of **residuals vs. fitted values**.

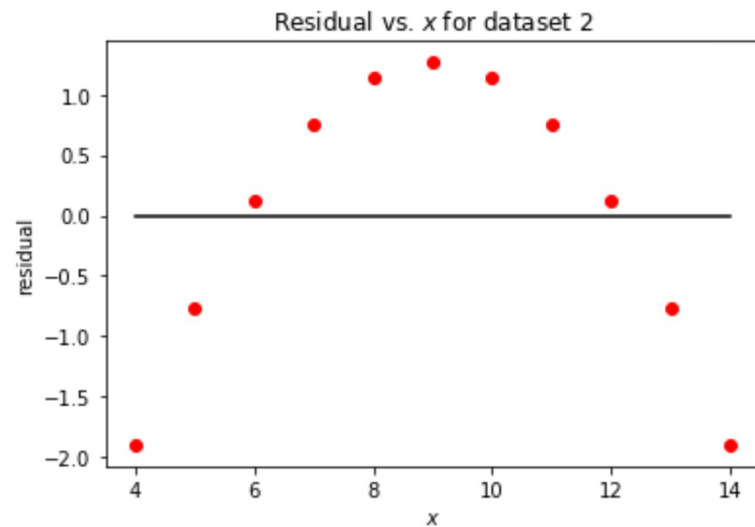
Properties:

- A good residual plot **has no pattern**. This means that our model represents the relationship in the data well.
  - If you see a curve, it is a sign that transformations or additional variables could help.
- A good residual plot also has a **similar vertical spread** throughout the entire plot.
  - If this is not the case, the accuracy of the predictions is not reliable.

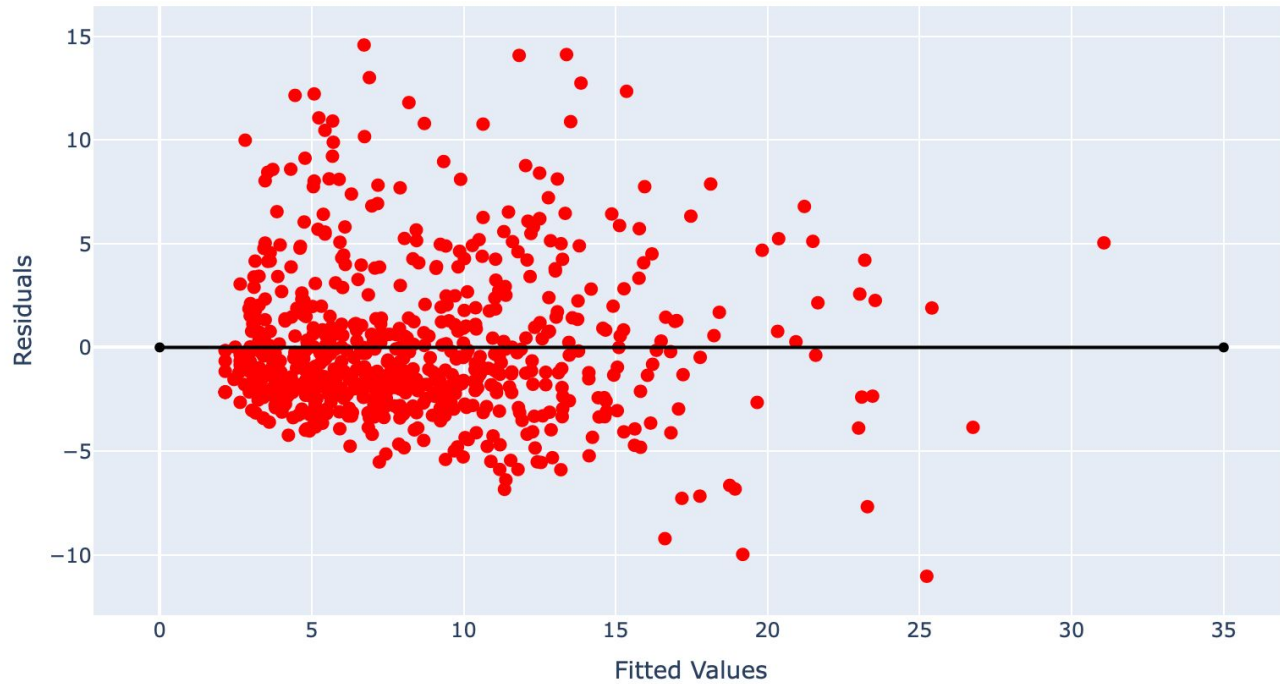


No pattern, even spread.

*Good.*



Clear quadratic relationship in the residuals.



No clear relationship, but uneven spread.

# Residuals are orthogonal to the span of $\mathbb{X}$

When using the optimal parameter vector, our residuals are orthogonal to  $\text{span}(\mathbb{X})$ .

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

- Since our predicted response  $\hat{\mathbb{Y}}$  is in  $\text{span}(\mathbb{X})$ , it is orthogonal to the residuals.
  - This is true, no matter what the features of model are.

$$\hat{\mathbb{Y}}^T \mathbf{e} = 0$$

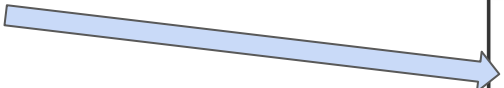
- When our model contains an intercept term, things become slightly more interesting!

We denote this column with  $\mathbb{1}$

Since  $\mathbb{X}^T \mathbf{e} = \mathbf{0}$ , it is also true that

$$\mathbb{1}^T \mathbf{e} = 0$$

**Why?**


$$\mathbb{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & x_{33} & \dots & x_{3p} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{np} \end{bmatrix}$$



# Properties when our model has an intercept term

When our linear model has an intercept term (i.e. when our design matrix has a column of all 1s), the following properties hold true:

- The sum of the residuals is 0.
  - The mean of the residuals is also 0!
  - This is why the positive and negative residuals cancel out in any residual plot where the (linear) model contains an intercept term, even if the model is terrible.
- The average true  $y$  value is equal to the average predicted  $y$  value.
  - This follows from the property above.

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

**You will prove them in a homework.**

Existence of a unique solution

# Does a solution always exist?

- For all models so far, our goal has been to determine the value of  $\theta$  that minimizes some average loss.
- The minimum value of both mean squared error and mean absolute error is 0.

$$MSE(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \qquad MAE(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

- This means, there is always **at least one** model parameter that minimizes average loss.

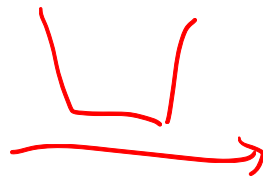
# Does a **unique** solution always exist?

- **Constant** model with **squared** loss:
  - Any set of values has a unique mean.
  - Thus, in this case, a unique solution always exists.
- **Simple** linear model with **squared** loss:
  - Any set of non-constant\* values has a unique mean, SD, and correlation coefficient.
- **Constant** model with **absolute** loss:
  - This is unique when there is an odd number of  $y$  values.
  - But, from Lecture 11, when there is an even number of  $y$  values, there are infinitely many solutions!
    - Recall, in such a case, any value of  $\theta$  between the “middle two” values minimized MAE.

$$\hat{\theta} = \mathbf{mean}(y)$$

$$\hat{\theta}_1 = r \frac{\sigma_y}{\sigma_x} \quad \hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$$

$$\hat{\theta} = \mathbf{median}(y)$$



# Understanding the solution matrices

Typically, **n** is much larger than **p**.

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

The diagram illustrates the dimensions of the matrices in the least squares solution formula  $\hat{\theta} = (X^T X)^{-1} X^T Y$ . The matrix  $X$  is shown as a blue rectangle with dimensions  $p+1$  (rows) by  $n$  (columns). The matrix  $X^T$  is shown as a blue rectangle with dimensions  $n$  (rows) by  $p+1$  (columns). The matrix  $Y$  is shown as a red vertical rectangle with dimensions  $n$  (rows) by  $1$  (column). The inverse matrix  $(X^T X)^{-1}$  is shown as a blue rectangle with dimensions  $p+1$  (rows) by  $p+1$  (columns). The resulting vector  $\hat{\theta}$  is shown as a pink vertical rectangle with dimensions  $p+1$  (rows) by  $1$  (column). Arrows indicate the multiplication of  $X^T$  and  $X$  to form  $(X^T X)^{-1}$ , and the multiplication of  $(X^T X)^{-1}$  and  $Y$  to form  $\hat{\theta}$ .

# Understanding the solution matrices

Typically, **n** is much larger than **p**.

$$\hat{\theta} = \left( \begin{matrix} \text{p+1} \\ \text{p+1} \end{matrix} \right) \left( \begin{matrix} \text{p+1} \\ \text{p+1} \end{matrix} \right)^{-1} \begin{matrix} \text{1} \\ \text{p+1} \end{matrix} \begin{matrix} \text{1} \\ \text{p+1} \end{matrix}$$

The diagram illustrates the dimensions of the matrices and vectors in the least squares solution formula. The matrix  $(X^T X)^{-1}$  is shown as a square with dimensions  $(p+1) \times (p+1)$ . The matrix  $X^T Y$  is shown as a column vector with dimensions  $1 \times (p+1)$ . Arrows indicate the dimensions of the matrices and vectors.

# Understanding the solution matrices

The **Normal Equation**:

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

$$\begin{pmatrix} \overset{p+1}{\square} \\ \text{p+1} \end{pmatrix} \hat{\theta} = \begin{pmatrix} \overset{1}{\square} \\ \text{p+1} \end{pmatrix} \mathbf{b}$$

Our optimal parameter vector can be thought of as the solution to a set of  $p + 1$  equations, with  $p + 1$  unknowns.

# Does a **unique** solution always exist?

Let's consider our optimal  $\theta$  for the multiple linear regression model:

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

- As mentioned previously, **at least one solution always exists**.
  - Intuitively, we can always draw a line of best fit for a given set of data. There may be multiple lines that are “equally good”.
- When does a unique solution for  $\hat{\theta}$  exist?
  - When  $\mathbb{X}^T \mathbb{X}$  is invertible. If it is **not invertible**, a unique solution does **not exist**.
    - In such a case, there will be infinitely many values of theta that minimize average squared loss.
    - If there are infinitely many “optimal” choices of coefficients, it's unclear which to use.
    - We want a **unique solution**.



# Invertibility of $\mathbf{X}^T\mathbf{X}$

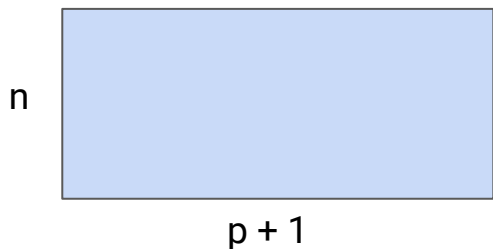
When is  $\mathbf{X}^T\mathbf{X}$  invertible?

- $\mathbf{X}^T\mathbf{X}$  is invertible if and only if it is full rank.
  - The shape of  $\mathbf{X}^T\mathbf{X}$  is  $(p + 1) \times (p + 1)$ . Invertibility is only defined for square matrices.
  - The rank of a matrix is the number of linearly independent columns (or rows) it contains.
- $\mathbf{X}^T\mathbf{X}$  This is one of several conditions of the “invertible matrix theorem.”
- $\mathbf{X}^T\mathbf{X}$  and  $\mathbf{X}$  have the same rank.
  - The proof is beyond the scope of this class.
- The maximum possible rank of  $\mathbf{X}^T\mathbf{X}$  is  $p + 1$ .
- Thus,  $\mathbf{X}^T\mathbf{X}$  is invertible if and only if  $\mathbf{X}$  has rank  $p + 1$  (full column rank).
  - **That is, a unique solution for the least squares estimate exists if and only if all columns of  $\mathbf{X}$  are linearly independent.**

# Invertibility of $\mathbb{X}^T \mathbb{X}$

When does our design matrix  $\mathbb{X}$  **not** have full column rank?

- When some features in our design matrix are linear combinations of other features.
  - If “Width”, “Height”, and “Perimeter” are all columns,  $\mathbb{X}$  will not have full rank, since  $\text{Perimeter} = 2 * \text{Width} + 2 * \text{Height}$  (linear combination).
  - When we discuss one-hot encoding, this is something to be aware of.
- When our design matrix has more columns than rows (i.e. it is “fat”).
  - In the normal setting,  $n > p + 1$  (we typically have more observations than features).



Since the row rank and column rank of a matrix are equal,  $\text{rank}(\mathbb{X}) \leq \min(n, p+1)$ .  
If  $n < p + 1$ , then the column rank cannot possibly be  $p + 1$ , meaning it cannot be full column rank.

Summary, what's next

# Summary

- We defined the multiple linear regression model in terms of matrices.

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

- We used a geometric argument to derive the optimal parameter vector  $\hat{\boldsymbol{\theta}}$ , that minimizes average squared loss. This value is called the **least squares estimate**.

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

- We discussed residuals and their properties.
- We explored when a unique solution for  $\hat{\boldsymbol{\theta}}$  exists, and when one does not.

# Moving forward

Think of this lecture as the third in a trilogy of lectures that establish the foundations of regression. From here, we will cover:

- Feature engineering.
  - The process of extracting and creating more sophisticated features from our data.
- The Bias-Variance tradeoff.
  - Discussing this idea of a “true model,” and where the errors in our predictions come from.
- Regularization and Cross-Validation.
  - Making our models more “general”, by updating our objective function.
- Gradient descent.
  - Minimizing average loss using numerical methods.
- Logistic regression.
  - What if we want to predict categories (1 and 0), instead of numbers?