

Lecture 12

# Multiple Linear Regression

DSC 40A, Fall 2025



# Recap: Regression and linear algebra

## Regression and linear algebra (Solution 1)

- Define the **design matrix**  $\mathbf{X} \in \mathbb{R}^{n \times 2}$ , **observation vector**  $\vec{y} \in \mathbb{R}^n$ , and parameter vector  $\vec{w} \in \mathbb{R}^2$  as:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \begin{array}{l} \nearrow \text{intercept} \\ \searrow \text{slope} \end{array}$$

- How do we make the **hypothesis vector**,  $\vec{h} = \mathbf{X}\vec{w}$ , as close to  $\vec{y}$  as possible? Use the parameter vector  $\vec{w}^*$ :

$$\vec{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$$

- Solution: We chose  $\vec{w}^*$  so that  $\vec{h}^* = \mathbf{X}\vec{w}^*$  is the **projection** of  $\vec{y}$  onto the **span of the columns of the design matrix**,  $\mathbf{X}$  and minimized the length of the projection error  $\|\vec{e}\| = \|\vec{y} - \mathbf{X}\vec{w}\|$ .

## Regression and linear algebra (Solution 2)

- Define the **design matrix**  $\mathbf{X} \in \mathbb{R}^{n \times 2}$ , **observation vector**  $\vec{y} \in \mathbb{R}^n$ , and parameter vector  $\vec{w} \in \mathbb{R}^2$  as:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

- How do we minimize the mean squared error  $R_{\text{sq}}(\vec{w}) = \frac{1}{n} \|\vec{y} - \mathbf{X}\vec{w}\|^2$ ? Using calculus the optimal parameter vector  $\vec{w}^*$  is:

$$\nabla R_{\text{sq}}(\vec{w}) = 0$$

$$\vec{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$$

- Solution: we computed the gradient of  $R_{\text{sq}}(\vec{w})$ , set it to zero and solved for  $\vec{w}$ .

# Multiple linear regression

	$x^{(1)}$	$x^{(2)}$	$\tilde{y}$	commute
	departure_hour	day_of_month	minutes	
0	10.816667	15	68.0	
1	7.750000	16	94.0	
2	8.450000	22	63.0	
3	7.133333	23	100.0	
4	9.150000	30	69.0	
...	...	...	...	...

So far, we've fit **simple** linear regression models, which use only **one** feature ( 'departure\_hour' ) for making predictions.







## The setup

- Suppose we have the following dataset.

row	departure_hour	day_of_month	minutes
1	8.45	22	63.0
2	8.90	28	89.0
3	8.72	18	89.0

commute time

- We can represent each day with a **feature vector**,  $\vec{x}$ :

$$x_1 = \begin{bmatrix} 8.45 \\ 22 \end{bmatrix}$$

$$x_2 = \begin{bmatrix} 8.90 \\ 28 \end{bmatrix}$$

$$y_1 = ?$$

$$y_2 = ?$$

$$H(\vec{x}_i) = y_i$$



## Finding the optimal parameters

- To find the optimal parameter vector,  $\vec{w}^*$ , we can use the **design matrix**  $\mathbf{X} \in \mathbb{R}^{n \times 3}$  and **observation vector**  $\vec{y} \in \mathbb{R}^n$ :

$$\mathbf{X} = \begin{bmatrix} 1 & \text{departure hour}_1 & \text{day}_1 \\ 1 & \text{departure hour}_2 & \text{day}_2 \\ \dots & \dots & \dots \\ 1 & \text{departure hour}_n & \text{day}_n \end{bmatrix} \quad \vec{y} = \begin{bmatrix} \text{commute time}_1 \\ \text{commute time}_2 \\ \vdots \\ \text{commute time}_n \end{bmatrix}$$

- Then, all we need to do is solve the **normal equations**:

$$\mathbf{X}^T \mathbf{X} \vec{w}^* = \mathbf{X}^T \vec{y}$$

$$\vec{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \end{pmatrix}$$

If  $\mathbf{X}^T \mathbf{X}$  is invertible, we know the solution is:

$$\vec{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$$

$3 \times 3$

$x_1$  - 1st data point  $\in \mathbb{R}^2$

## Notation for multiple linear regression

$x^{(1)}$  - 1st feature  $\in \mathbb{R}^n$

- We will need to keep track of multiple features for every individual in our dataset.
  - In practice, we could have hundreds or thousands of features!
- As before, subscripts distinguish between individuals in our dataset. We have  $n$  individuals, also called **training examples**.
- Superscripts distinguish between **features**. We have  $d$  features.

$x^{(1)}, x^{(2)}, \dots, x^{(d)}$

departure hour:  $x^{(1)} \in \mathbb{R}^n$

day of month:  $x^{(2)} \in \mathbb{R}^n$

Think of  $x^{(1)}, x^{(2)}, \dots$  as new variable names, like new letters.

↑  
not exponent!

$$x^2 = x \cdot x$$

$x^{(7)} \rightarrow 7^{\text{th}}$  column  
 $7^{\text{th}}$  feature

$x_{21} \rightarrow 21^{\text{st}}$  data point

## Augmented feature vectors

- The augmented feature vector  $\text{Aug}(\vec{x})$  is the vector obtained by adding a 1 to the front of feature vector  $\vec{x}$ :

$$\vec{x} = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(d)} \end{bmatrix} \quad \text{Aug}(\vec{x}) = \begin{bmatrix} 1 \\ x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(d)} \end{bmatrix} \quad \vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

Annotations in blue:

- $x^{(1)}$  is labeled "departure time"
- $x^{(2)}$  is labeled "day of the month"
- $x^{(d)}$  is labeled "temperature"

- Then, our hypothesis function is:

$$H(\vec{x}) = w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_d x^{(d)}$$

Annotations in green:

- $w_0$  is labeled  $\in \mathbb{R}$
- $w_1 x^{(1)}$  is labeled  $(d+1) \times 1$
- $w_2 x^{(2)}$  is labeled  $(d+1) \times 1$
- $w_d x^{(d)}$  is labeled  $(d+1) \times 1$

## The general problem

- We have  $n$  data points,  $(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_n, y_n)$ , where each  $\vec{x}_i$  is a feature vector of  $d$  features:

$$\vec{x}_i = \begin{bmatrix} x_i^{(1)} \\ x_i^{(2)} \\ \vdots \\ x_i^{(d)} \end{bmatrix} \in \mathbb{R}^d$$

- We want to find a good linear hypothesis function:

$$\begin{aligned} H(\vec{x}) &= w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_d x^{(d)} \\ &= \vec{w} \cdot \text{Aug}(\vec{x}) \end{aligned}$$

How to find  $\vec{w}^* = (w_0, w_1, \dots, w_d)$

SLR model  
 $(x_1, y_1), \dots, (x_n, y_n)$   
 $x_i \in \mathbb{R}$   
scalar

# The general solution

- Define the design matrix  $X \in \mathbb{R}^{n \times (d+1)}$  and observation vector  $\vec{y} \in \mathbb{R}^n$ :

*data point*  $\vec{x}_1$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(d)} \\ 1 & x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(d)} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(d)} \end{bmatrix} = \begin{bmatrix} \text{Aug}(\vec{x}_1)^T \\ \text{Aug}(\vec{x}_2)^T \\ \vdots \\ \text{Aug}(\vec{x}_n)^T \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

- Then, solve the normal equations to find the optimal parameter vector,  $\vec{w}^*$ :

*feature*  $\vec{x}_1$

$$X^T X \vec{w}^* = X^T \vec{y}$$

If  $X^T X$  invertible  
optimal  $\vec{w}^* = (X^T X)^{-1} X^T \vec{y}$

## Terminology for parameters

- With  $d$  features,  $\vec{w}$  has  $d + 1$  entries.
- $w_0$  is the **bias**, also known as the **intercept**.
- $w_1, w_2, \dots, w_d$  each give the **weight**, or **coefficient**, or **slope**, of a feature.

$$H(\vec{x}) = w_0 + w_1 x^{(1)} + w_2 x^{(2)} + \dots + w_d x^{(d)}$$

# Interpreting parameters

## Example: Predicting sales

27

- For each of 27 stores, we have:

- **net sales**,  $\hat{y}$
- square feet,
- inventory,
- advertising expenditure,
- district size, and
- number of competing stores.

X

27 x 6

- **Goal:** Predict net sales given the other five features.
- To begin, we'll start trying to fit the hypothesis function to predict sales:

$$H(\text{square feet, competitors}) = w_0 + w_1 \cdot \text{square feet} + w_2 \cdot \text{competitors}$$

$$d=2$$

## Question 🤔

Answer at [q.dsc40a.com](http://q.dsc40a.com)

$$H(\text{square feet, competitors}) = w_0 + w_1 \cdot \text{square feet} + w_2 \cdot \text{competitors}$$

What will be the signs of  $w_1^*$  and  $w_2^*$ ?

- A.  $w_1^* + w_2^* +$
- B.  $w_1^* + w_2^* -$
- C.  $w_1^* - w_2^* +$
- D.  $w_1^* - w_2^* -$

↳ bigger stores  
sell more

↳ More competitors  
sell less

Let's find out! Follow along in [this notebook](#).