**Lectures 8-10**

# Linear algebra: Dot products and Projections

**DSC 40A, Fall 2024**

# Announcements

- Homework 2 was released Friday. Remember that using the Overleaf template is required for Homework 2 (and only Homework 2).

- Groupwork 3 is due **tonight**.

- Check out FAQs page and the tutor-created supplemental resources on the course website.

# Agenda

- Recap: Simple linear regression and correlation.

- Connections to related models.
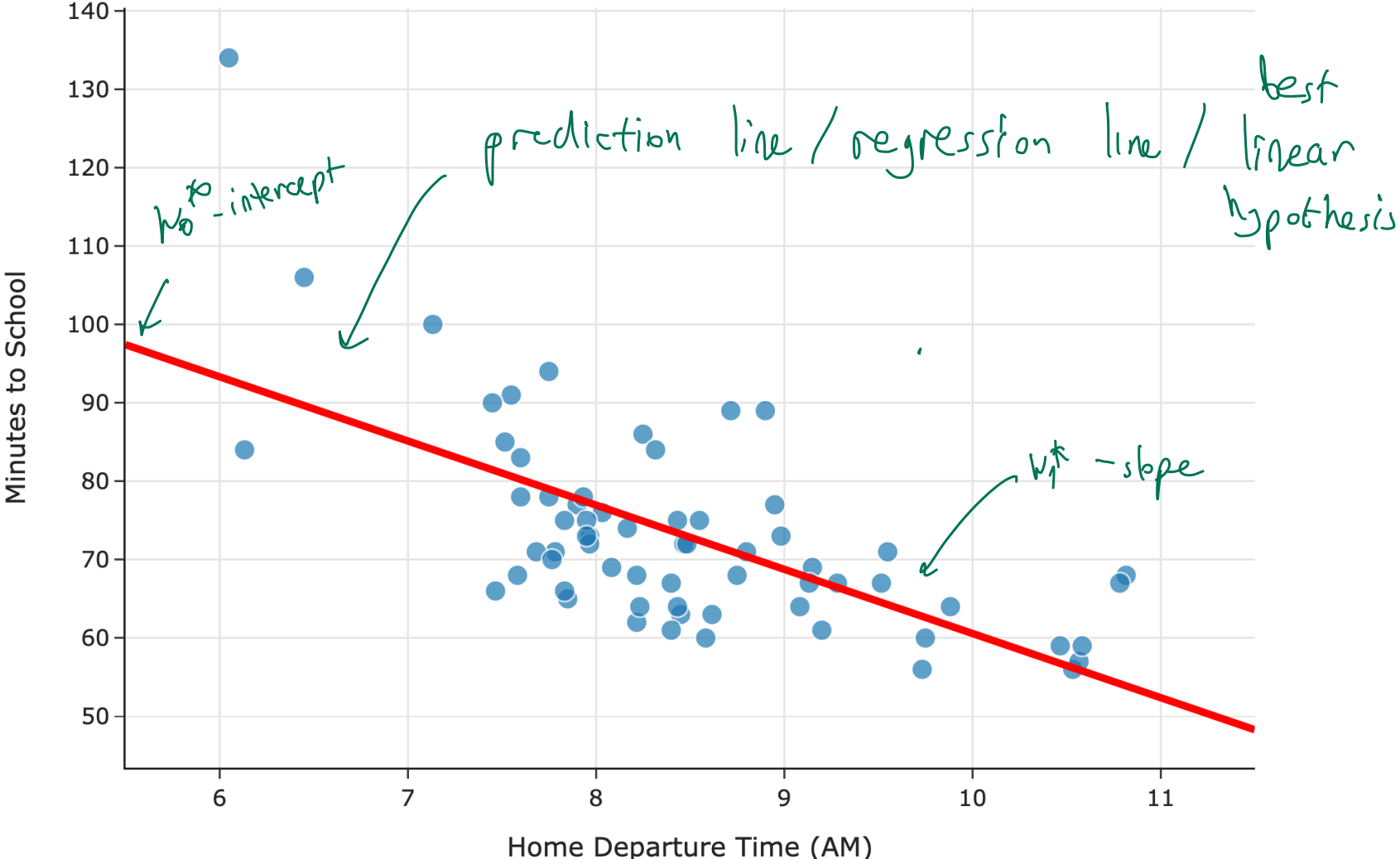
- Dot products.

- Spans and projections.

# Question 🤔

Answer at q.dsc40a.com

**Remember, you can always ask questions at q.dsc40a.com!**

If the direct link doesn't work, click the "🤔 Lecture Questions"

link in the top right corner of dsc40a.com.

Predicted Commute Time = 142.25 - 8.19 * Departure Hour

$w_0^*$ -intercept

prediction line / regression line / best linear hypothesis

$w_1^*$ - slope

# Simple linear regression

- Model: $H(x) = w_0 + w_1 x$.

- Loss function: squared loss, i.e. $L_{\text{sq}}(y_i, H(x_i)) = (y_i - H(x_i))^2$.

- Average loss, i.e. empirical risk:

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^{n} (y_i - (w_0 + w_1 x_i))^2$$

- Optimal model parameters, found by minimizing empirical risk:

$$w_1^* = \frac{\displaystyle\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\displaystyle\sum_{i=1}^{n}(x_i - \bar{x})^2} = r \frac{\sigma_y}{\sigma_x} \qquad w_0^* = \bar{y} - w_1^* \bar{x}$$

$$= \frac{\sum_i (y_i - \bar{y}) x_i}{\sum_i (x_i - \bar{x}) x_i}$$

slope

intercept

5

# The correlation coefficient

- The correlation coefficient, $r$, is defined as the **average of the product of** $x$ **and** $y$, **when both are in standard units**.

- Let $\sigma_x$ be the standard deviation of the $x_i$s, and $\bar{x}$ be the mean of the $x_i$s.

- $x_i$ in standard units is $\frac{x_i - \bar{x}}{\sigma_x}$.

- The correlation coefficient, then, is:

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

# Correlation and mean squared error

- **Claim**: Suppose that $w_0^*$ and $w_1^*$ are the optimal intercept and slope for the regression line. Then,

$$R_{\mathrm{sq}}(w_0^*, w_1^*) = \sigma_y^2(1 - r^2)$$

- That is, the mean squared error of the regression line's predictions and the correlation coefficient, $r$, always satisfy the relationship above.

- Even if it's true, why do we care?

  - In machine learning, we often use both the mean squared error and $r^2$ to compare the performances of different models.

  - If we can prove the above statement, we can show that **finding models that minimize mean squared error** is equivalent to **finding models that maximize** $r^2$.

**Proof that** $R_{\mathrm{sq}}(w_0^*, w_1^*) = \boxed{\sigma_y^2(1 - r^2)}$

$$R_{sq}(w_0^*, w_1^*) = \frac{1}{n}\sum_i \left(y_i - (w_0^* + w_1^* x_i)\right)^2 = \frac{1}{n}\sum_i \left(y_i - (\bar{y} - w_1^* \bar{x} + w_1^* x_i)\right)^2$$

$$= \frac{1}{n}\sum_i \left((y_i - \bar{y}) - r\frac{\sigma_y}{\sigma_x}(x_i - \bar{x})\right)^2$$

$$= \frac{1}{n}\sum_i (y_i - \bar{y})^2 + \frac{1}{n}\sum_i r^2 \frac{\sigma_y^2}{\sigma_x^2}(x_i - \bar{x})^2 - 2\frac{1}{n}\sum_i (y_i - \bar{y})(x_i - \bar{x})\, r\frac{\sigma_y}{\sigma_x}$$

$$= \sigma_y^2 + r^2 \frac{\sigma_y^2}{\sigma_x^2}\sigma_x^2 - 2r\frac{\sigma_y}{\sigma_x}\boxed{\frac{1}{n}\sum_i (y_i - \bar{y})(x_i - \bar{x})}$$

$$\underbrace{\phantom{xxxx}}_{r \cdot \sigma_x \sigma_y}$$

$$= \sigma_y^2 + r^2 \sigma_y^2 - 2r^2 \sigma_y^2 = \sigma_y^2(1 + r^2 - 2r^2) = \boxed{\sigma_y^2(1 - r^2)}$$

8

# Connections to related models

# Exercise $(no\ slope)$

Suppose we choose the model $H(x) = w_0$ and squared loss.
What is the optimal model parameter, $w_0^*$?

Constant model

loss: squared loss

week 1 ?    $w_0^* = $ mean $\{y_1, \ldots, y_n\}$

(this was h is week 1)

# Exercise  $(no\ intercept)$

Suppose we choose the model $H(x) = w_1 x$ and squared loss.
What is the optimal model parameter, $w_1^*$?

Groupwork 3!

# Comparing mean squared errors

- With both:

  - the constant model, $H(x) = h$, and

  - the simple linear regression model, $H(x) = w_0 + w_1 x$,

  when we chose squared loss, we minimized mean squared error to find optimal parameters:
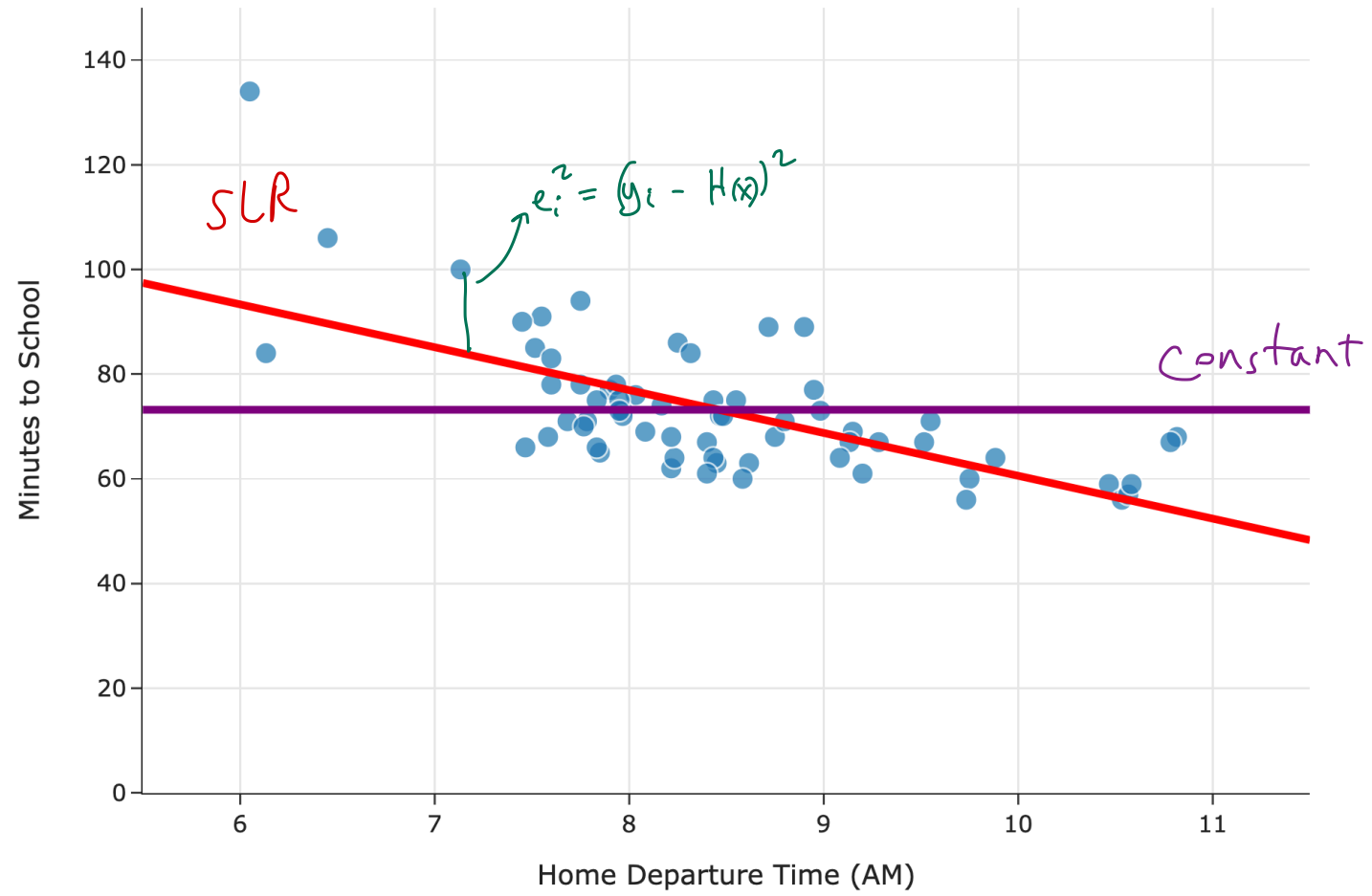
  $$R_{\text{sq}}(H) = \frac{1}{n} \sum_{i=1}^{n} (y_i - H(x_i))^2$$

- **Which model minimizes mean squared error more?**

$$\text{MSE}(\text{SLR}) \leq \text{MSR}(\text{constant})$$

# Comparing mean squared errors

Predicted Commute Time = 142.25 - 8.19 * Departure Hour
Predicted Commute Time = 73.18



$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (y_i - H(x_i))^2$$

- The MSE of the best simple linear regression model is $\approx 97$

- The MSE of the best constant model is $\approx 167$

- The simple linear regression model is a more flexible version of the constant model.

# Linear algebra

# Wait... why do we need linear algebra?

- Soon, we'll want to make predictions using more than one feature.
  - Example: Predicting commute times using departure hour and temperature.
- Thinking about linear regression in terms of **matrices and vectors** will allow us to find hypothesis functions that:
  - Use multiple features (input variables).
  - Are nonlinear in the features, e.g. $H(x) = w_0 + w_1 x + w_2 x^2$.

# Warning ⚠️

- We're **not** going to cover every single detail from your linear algebra course.
- There will be facts that you're expected to remember that we won't explicitly say.
  - For example, if $A$ and $B$ are two matrices, then $AB \neq BA$.
  - This is the kind of fact that we will only mention explicitly if it's directly relevant to what we're studying.
  - But you still need to know it, and it may come up in homework questions.
- We **will** review the topics that you really need to know well.

# Dot Products

# Vectors

*\mathbb{R}^n (latex)*

- A **vector** in $\mathbb{R}^n$ is an **ordered collection of $n$ numbers**.

- We use lower-case letters with an arrow on top to represent vectors, and we usually write vectors as **columns**.

$$\vec{v} = \begin{bmatrix} 8 \\ 3 \\ -2 \\ 5 \end{bmatrix}$$

*size* $\longrightarrow n \times 1$

- Another way of writing the above vector is $\vec{v} = [8, 3, -2, 5]^{\mathsf{T}}$.

*transpose*

- Since $\vec{v}$ has four **components**, we say $\vec{v} \in \mathbb{R}^4$.

*element*

*\in (latex)*

# The geometric interpretation of a vector

$$\vec{v} = \begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

$$\|\vec{v}\| = \sqrt{5^2 + 3^2} = \sqrt{34}$$

- A vector $\vec{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$ is an arrow to the point $(v_1, v_2, \ldots, v_n)$ from the origin.
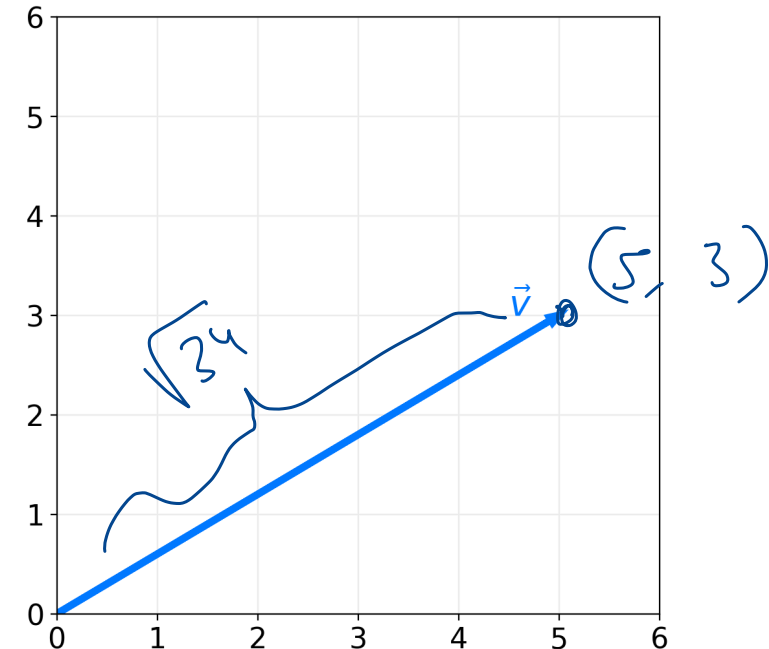


$(5, 3)$

$\sqrt{34}$

- The **length**, or $L_2$ **norm**, of $\vec{v}$ is:

$$\|\vec{v}\| = \sqrt{v_1^2 + v_2^2 + \ldots + v_n^2} = \sqrt{\sum_{i=1}^n v_i^2} = \sqrt{\vec{v} \cdot \vec{v}} = \sqrt{\vec{v}^T \vec{v}}$$

- A vector is sometimes described as an object with a **magnitude/length** and **direction**.

$$\neq \sqrt{\vec{v}^2} = \sqrt{[\ ]\cdot[\ ]}$$

cannot
perform
this operation

$n \times 1$

$n \times 1$

27

# Dot product: coordinate definition

- The **dot product** of two vectors $\vec{u}$ and $\vec{v}$ in $\mathbb{R}^n$ is written as:

  $$\vec{u} \cdot \vec{v} = \vec{u}^\mathsf{T}\vec{v}$$

  *↖ \cdot (latex)*
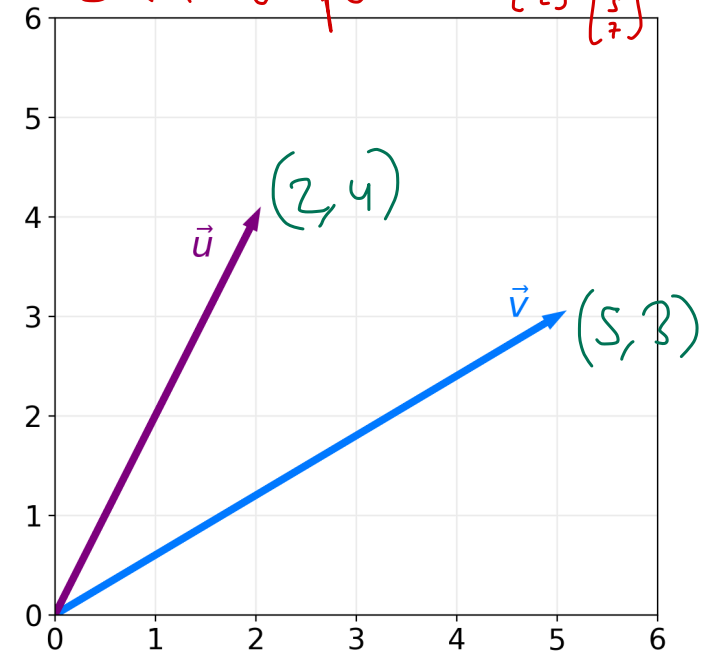
- The computational definition of the dot product:

  $$\vec{u} \cdot \vec{v} = \sum_{i=1}^{n} u_i v_i = u_1 v_1 + u_2 v_2 + \ldots + u_n v_n$$

- The result is a **scalar**, i.e. a single number.

both vectors need to have same number of elements

cannot perform $\begin{bmatrix} 1 \\ 2 \end{bmatrix} \cdot \begin{bmatrix} 4 \\ 5 \\ 7 \end{bmatrix}$

$\vec{u}$ (2,4)

$\vec{v}$ (5,3)

$\vec{u} \cdot \vec{v} = 2 \cdot 5 + 4 \cdot 3 = 10 + 12 = 22 \in \mathbb{R}$ ← scalar

$\vec{u}^\mathsf{T}\vec{v} = \begin{bmatrix} 2 & 4 \end{bmatrix} \begin{bmatrix} 5 \\ 3 \end{bmatrix} = 10 + 12 = 22$

$f(\vec{u}, \vec{v}) \in \mathbb{R}$

$\mathbb{R}^n \nearrow \quad \uparrow \mathbb{R}^n$

28

# Dot product: geometric definition

- The computational definition of the dot product:

$$\vec{u} \cdot \vec{v} = \sum_{i=1}^{n} u_i v_i = u_1 v_1 + u_2 v_2 + \ldots + u_n v_n$$
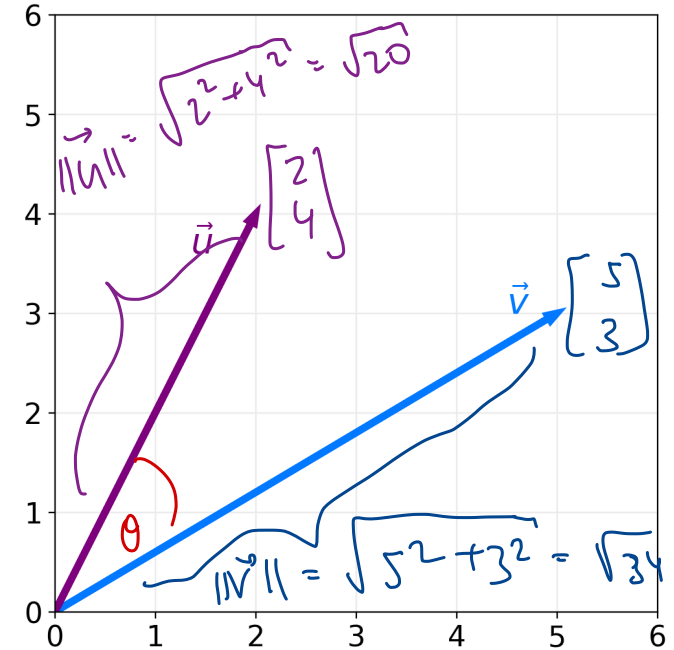
- The geometric definition of the dot product:

$$\vec{u} \cdot \vec{v} = \|\vec{u}\| \|\vec{v}\| \cos\theta$$

  where $\theta$ is the angle between $\vec{u}$ and $\vec{v}$.

- The two definitions are equivalent! This equivalence allows us to find the angle $\theta$ between two vectors.

$$\vec{u} \cdot \vec{v} = 22 \quad (\text{from previous slide})$$

$$\cos\theta = \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|} = \frac{22}{\sqrt{20}\sqrt{34}} = \frac{11}{\sqrt{5} \cdot \sqrt{34}} = \frac{11}{\sqrt{170}}$$



$$\|\vec{u}\| = \sqrt{2^2 + 4^2} = \sqrt{20}$$

$$\vec{u} \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$\vec{v} \begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

$$\|\vec{v}\| = \sqrt{5^2 + 3^2} = \sqrt{34}$$

# Question 4: $\cos\theta$

What is $\cos\theta$?

- A. $\dfrac{6}{\sqrt{85}}$
- B. $\dfrac{-6}{\sqrt{85}}$
- C. $\dfrac{-3}{85}$
- D. $\dfrac{-2}{3}$

$$\vec{u}\cdot\vec{v} = \|\vec{u}\|\cdot\|\vec{v}\|\cdot\cos\theta$$

$$\cos\theta = \frac{\vec{u}\cdot\vec{v}}{\|\vec{u}\|\,\|\vec{v}\|}$$

$$\vec{u}\cdot\vec{v} = \vec{u}^T\vec{v} = \begin{bmatrix} 4 & 1\end{bmatrix}\begin{bmatrix} -2 \\ -4 \end{bmatrix} = -8-4$$
$$= -12$$

$$\cos\theta = \frac{-12}{\sqrt{20}\,\sqrt{17}} = \frac{-6}{\sqrt{5}\,\sqrt{17}} = \frac{-6}{\sqrt{85}}$$

$$\frac{-2\cdot 6}{\sqrt{4\cdot 5}\cdot\sqrt{17}} = \frac{-6}{\sqrt{5}\cdot\sqrt{17}}$$

$$\|\vec{u}\| = \sqrt{4^2+1} = \sqrt{17}$$

$$\sqrt{4^2+2^2} = \sqrt{20}$$



$\vec{u}$ $\begin{bmatrix} 4 \\ 1 \end{bmatrix}$

$\vec{v}$ $\begin{bmatrix} -2 \\ -4 \end{bmatrix}$

$\theta$