# Lecture 9 – Regression in Action and Linear Algebra Review



**DSC 40A, Winter 2024**

# Announcements

- ▶ Homework 3 is due **Wed at 11:59pm**.
  - ▶ Come to office hours. See `dsc40a.com/calendar` for the schedule.

- ▶ Discusssion session today
  - ▶ We modify the scope of discussion session/groupwork so that it aligns with the course better.

## Agenda

- ▶ Recap of Lecture 8.

- ▶ Connection with correlation.

- ▶ Interpretation of formulas.

- ▶ Regression demo.

- ▶ Linear algebra review.

**Recap of Lecture 8**

## Strategy

$$-\frac{2}{n} \sum_{i=1}^{n} \left( y_i - (w_0 + w_1 x_i) \right) = 0 \qquad -\frac{2}{n} \sum_{i=1}^{n} \left( y_i - (w_0 + w_1 x_i) \right) x_i = 0$$

1. Solve for $w_0$ in first equation.
   ▶ The result becomes $w_0^*$, since it is the "best intercept".

2. Plug $w_0^*$ into second equation, solve for $w_1$.
   ▶ The result becomes $w_1^*$, since it is the "best slope".

## Solve for $w_0^*$

$$-\frac{2}{n} \sum_{i=1}^{n} \left( y_i - (w_0 + w_1 x_i) \right) = 0$$

## Solve for $w_1^*$

$$-\frac{2}{n} \sum_{i=1}^{n} \left(y_i - (w_0 + w_1 x_i)\right) x_i = 0$$

## Least squares solutions

▶ We've found that the values $w_0^*$ and $w_1^*$ that minimize the function $R_{sq}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - (w_0 + w_1 x_i) \right)^2$ are

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

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i \qquad \bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$$

▶ Let's re-write the slope $w_1^*$ to be a bit more symmetric.

## Key fact

The **sum of deviations from the mean** for any dataset is 0.

$$\sum_{i=1}^{n} (x_i - \bar{x}) = 0 \qquad \sum_{i=1}^{n} (y_i - \bar{y}) = 0$$

Proof:

# Equivalent formula for $w_1^*$

Claim

$$w_1^* = \frac{\sum\limits_{i=1}^{n}(y_i - \bar{y})x_i}{\sum\limits_{i=1}^{n}(x_i - \bar{x})x_i} = \frac{\sum\limits_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2}$$

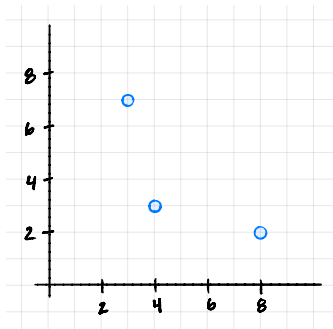Proof:

# Least squares solutions

▶ The **least squares solutions** for the slope $w_1^*$ and intercept $w_0^*$ are:

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

▶ We also say that $w_0^*$ and $w_1^*$ are **optimal parameters**.

▶ To make predictions about the future, we use the prediction rule

$$H^*(x) = w_0^* + w_1^*x$$

# Example



$$\bar{x} =$$

$$\bar{y} =$$

$$w_1^* = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} =$$

$$w_0^* = \bar{y} - w_1\bar{x} =$$

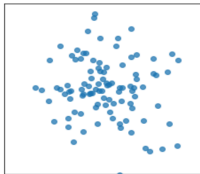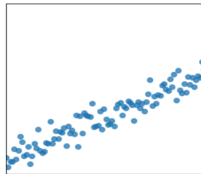| $x_i$ | $y_i$ | $(x_i - \bar{x})$ | $(y_i - \bar{y})$ | $(x_i - \bar{x})(y_i - \bar{y})$ | $(x_i - \bar{x})^2$ |
|-------|-------|-------------------|-------------------|----------------------------------|---------------------|
| 3 | 7 | | | | |
| 4 | 3 | | | | |
| 8 | 2 | | | | |

# Connection with correlation

# Correlation coefficient

▶ In DSC 10, you were introduced to the idea of correlation.

  ▶ It is a measure of the strength of the **linear association** of two variables, *x* and *y*.

  ▶ Intuitively, it measures how tightly clustered a scatter plot is around a straight line.

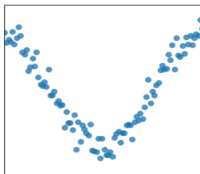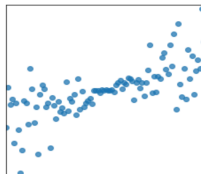  ▶ It ranges between -1 and 1.

# Patterns in scatter plots



r = -0.121

r = 0.949

r = 0.052

r = 0.704

# Definition of 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 $\dfrac{x_i - \bar{x}}{\sigma_x}$.

  - The correlation coefficient is

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

## Another way to express $w_1^*$

▶ It turns out that $w_1^*$, the optimal slope for the linear prediction rule, can be written in terms of $r$!

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

▶ It's not surprising that $r$ is related to $w_1^*$, since $r$ is a measure of linear association.

▶ Concise way of writing $w_0^*$ and $w_1^*$:

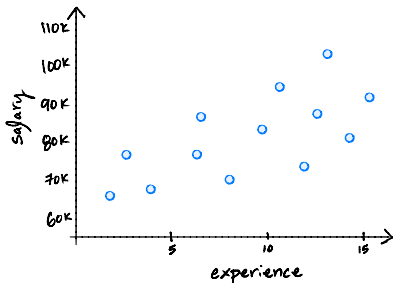$$w_1^* = r\frac{\sigma_y}{\sigma_x} \qquad w_0^* = \bar{y} - w_1^*\bar{x}$$

**Proof that** $w_1^* = r\dfrac{\sigma_y}{\sigma_x}$
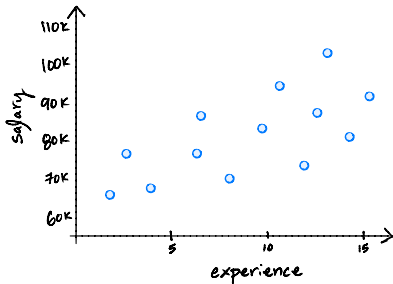
# Interpretation of formulas

# Interpreting the slope

$$w_1^* = r \frac{\sigma_y}{\sigma_x}$$



- $\sigma_y$ and $\sigma_x$ are always non-negative. As a result, the sign of the slope is determined by the sign of $r$.

- As the $y$ values get more spread out, $\sigma_y$ increases and so does the slope.

- As the $x$ values get more spread out, $\sigma_x$ increases and the slope decreases.

# Interpreting the intercept



$$w_0^* = \bar{y} - w_1^* \bar{x}$$

▶ What is $H^*(\bar{x})$?

### Discussion Question

We fit a linear prediction rule for salary given years of experience. Then everyone gets a $5,000 raise. Which of these happens?

a) slope increases, intercept increases

b) slope decreases, intercept increases

c) slope stays same, intercept increases

d) slope stays same, intercept stays same

# Regression demo

Let's see regression in action. Follow along here.

# Linear algebra review

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

- ▶ Soon, we'll want to make predictions using more than one feature (e.g. predicting salary using years of experience and GPA).

- ▶ Thinking about linear regression in terms of **linear algebra** will allow us to find prediction rules that
  - ▶ use multiple features.

  - ▶ are non-linear.

- ▶ Before we dive in, let's review.

# Matrices

- An $m \times n$ **matrix** is a table of numbers with $m$ rows and $n$ columns.

- We use upper-case letters for matrices.

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$$

- $A^T$ denotes the transpose of $A$:

$$A^T = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$$

# Matrix addition and scalar multiplication

▶ We can add two matrices only if they are the same size.

▶ Addition occurs elementwise:

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 7 & 8 & 9 \\ -1 & -2 & -3 \end{bmatrix} = \begin{bmatrix} 8 & 10 & 12 \\ 3 & 3 & 3 \end{bmatrix}$$

▶ Scalar multiplication occurs elementwise, too:

$$2 \cdot \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} = \begin{bmatrix} 2 & 4 & 6 \\ 8 & 10 & 12 \end{bmatrix}$$

# Matrix-matrix multiplication

▶ We can multiply two matrices *A* and *B* only if

# columns in *A* = # rows in *B*.

▶ If *A* is *m* × *n* and *B* is *n* × *p*, the result is *m* × *p*.
  ▶ This is **very useful**.

▶ The *ij* entry of the product is:

$$(AB)_{ij} = \sum_{k=1}^{n} A_{ik} B_{kj}$$

# Some matrix properties

▶ Multiplication is Distributive:

$$A(B + C) = AB + AC$$

▶ Multiplication is Associative:

$$(AB)C = A(BC)$$

▶ Multiplication is **not commutative**:

$$AB \neq BA$$

▶ Transpose of sum:

$$(A + B)^T = A^T + B^T$$

▶ Transpose of product:

$$(AB)^T = B^T A^T$$

# Vectors

- ▶ An **vector** in $\mathbb{R}^n$ is an $n \times 1$ matrix.

- ▶ We use lower-case letters for vectors.

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

- ▶ Vector addition and scalar multiplication occur elementwise.

# Geometric meaning of vectors

▶ A vector $\vec{v} = (v_1, \ldots, v_n)^T$ is an arrow to the point $(v_1, \ldots, v_n)$ from the origin.

▶ The **length**, or **norm**, of $\vec{v}$ is $\|\vec{v}\| = \sqrt{v_1^2 + v_2^2 + \ldots + v_n^2}$.

# Dot products

▶ The **dot product** of two vectors $\vec{u}$ and $\vec{v}$ in $\mathbb{R}^n$ is denoted by:

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

▶ Definition:

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

▶ The result is a **scalar**!

**Discussion Question**

Which of these is another expression for the length of $\vec{u}$?

a) $\vec{u} \cdot \vec{u}$
b) $\sqrt{\vec{u}^2}$
c) $\sqrt{\vec{u} \cdot \vec{u}}$
d) $\vec{u}^2$

# Properties of the dot product

► Commutative:

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

► Distributive:

$$\vec{u} \cdot (\vec{v} + \vec{w}) = \vec{u} \cdot \vec{v} + \vec{u} \cdot \vec{w}$$

# Matrix-vector multiplication

▶ Special case of matrix-matrix multiplication.

▶ Result is always a vector with same number of rows as the matrix.

▶ One view: a "mixture" of the columns.

$$\begin{bmatrix} 1 & 2 & 1 \\ 3 & 4 & 5 \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = a_1 \begin{bmatrix} 1 \\ 3 \end{bmatrix} + a_2 \begin{bmatrix} 2 \\ 4 \end{bmatrix} + a_3 \begin{bmatrix} 1 \\ 5 \end{bmatrix}$$

▶ Another view: a dot product with the rows.

**Discussion Question**

If $A$ is an $m \times n$ matrix and $\vec{v}$ is a vector in $\mathbb{R}^n$, what are the dimensions of the product $\vec{v}^T A^T A \vec{v}$?

a) $m \times n$ (matrix)
b) $n \times 1$ (vector)
c) $1 \times 1$ (scalar)
d) The product is undefined.

**Summary**

## Summary, next time

▶ We can re-write the optimal parameters for the regression line

$$w_1^* = r\frac{\sigma_y}{\sigma_x} \qquad w_0^* = \bar{y} - w_1^*\bar{x}$$

▶ We can then make predictions using $H^*(x) = w_0^* + w_1^*x$.

▶ We will need linear algebra in order to generalize regression to work with multiple features.

▶ **Next time**: Continue linear algebra review. Formulate linear regression in terms of linear algebra.