

# STAT 3011 Discussion 015

## Week 14: Regression Analysis

Talha Hamza

University of Minnesota  
College of Science and Engineering

Spring 2025

# Scatterplots: Visualizing Relationships

---

## Purpose

Display the relationship between two quantitative variables

## Construction

- X-axis: Explanatory variable
- Y-axis: Response variable
- Each point represents  $(x,y)$  pair

## What to Look For

- Overall pattern (linear, curved, etc.)
- Strength of relationship
- Direction (positive/negative)
- Outliers

# Scatterplot Example: Olympic Dash Times

```
1 # Load the dash data
2 dash <- read.table("http://users.stat.umn.edu/~wuxxx725/data/dash.txt",
3                     header=TRUE)
4
5 # Basic scatterplot
6 plot(dash$Year, dash$Time,
7       xlab = "Years after 1900",
8       ylab = "Winning Time (seconds)",
9       pch = 16, # Solid dots
10      col = "maroon") # UMN colors!
```

## Parameters Definition

xlab/ylab	Axis labels
pch=16	Solid points
col	Point color
cor()	Calculates correlation

# Correlation: Measuring Linear Association

## Definition

$$r = \frac{1}{n-1} \sum \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

Measures strength and direction of linear relationship (always between -1 to 1)

## Interpretation

Value Range	Interpretation
$0.8 \leq  r  \leq 1.0$	Very strong
$0.6 \leq  r  < 0.8$	Strong
$0.4 \leq  r  < 0.6$	Moderate
$0.2 \leq  r  < 0.4$	Weak
$0.0 \leq  r  < 0.2$	Very weak

# Least Squares Regression

---

## Regression Line

$$\hat{y} = a + bx$$

Where:

- $a = \bar{y} - b\bar{x}$  (y-intercept)
- $b = r \frac{s_y}{s_x}$  (slope)

## Key Concepts

- Minimizes sum of squared residuals
- $r^2$  = proportion of variation in y explained by x
- Always passes through  $(\bar{x}, \bar{y})$

# R Implementation

---

```
1 # Scatterplot
2 plot(x, y, xlab="Explanatory", ylab="Response", pch=16)
3 # Correlation
4 cor(x, y)
5 # Regression
6 model <- lm(y ~ x)
7 summary(model)
8 # Add regression line
9 abline(model)
```

## Output Interpretation

- Coefficients table shows intercept (a) and slope (b)
- Residual standard error measures typical prediction error
- Multiple R-squared shows proportion of variance explained

# Inference for Regression

---

## Hypothesis Test for Slope

- $H_0 : \beta = 0$  (no linear relationship)
- $H_a : \beta \neq 0$  (linear relationship exists)
- Test statistic:  $t = \frac{b}{SE(b)} \sim t_{n-2}$

## Confidence Interval for Slope

$$b \pm t_{n-2}^* \times SE(b)$$

Where  $t^*$  is critical value for desired confidence level

# Example: Olympic Dash Times

```
1 > model <- lm(Time ~ Year)
2 > summary(model)
3 Coefficients:
4             Estimate Std. Error t value Pr(>|t|)
5 (Intercept) 22.355087   0.143763  155.50  <2e-16 ***
6 Year        -0.030266   0.002354  -12.86  1.65e-10 ***
7 ---
8 Multiple R-squared:  0.9018
```

## Interpretation

- Slope = -0.0303: Times decrease by about 0.03 sec/year
- p-value < 0.05: Statistically significant relationship
- $r^2 = 0.902$ : 90.2% of variation is explained, rest is error



# Final Remarks

## Hypothesis Testing for $\beta$ Assumptions

1. **Random Sample:** Pairs  $(x_i, y_i)$  are independent observations
2. **Model Conditions:**
  - True relationship is linear:  $\mu_y = \alpha + \beta x$
  - Residuals are approximately normal  $\epsilon \sim N(0, \sigma)$
  - Constant variance

**Note:** The  $t$ -test for  $\beta$  is robust to minor violations of normality when  $n \geq 30$

## Cautions in Regression

- Correlation  $\neq$  Causation
- Beware of lurking variables
- Don't extrapolate beyond data range

# Questions?

Email me at [hamza050@umn.edu](mailto:hamza050@umn.edu) or attend my office hours