

# Introduction to Regression



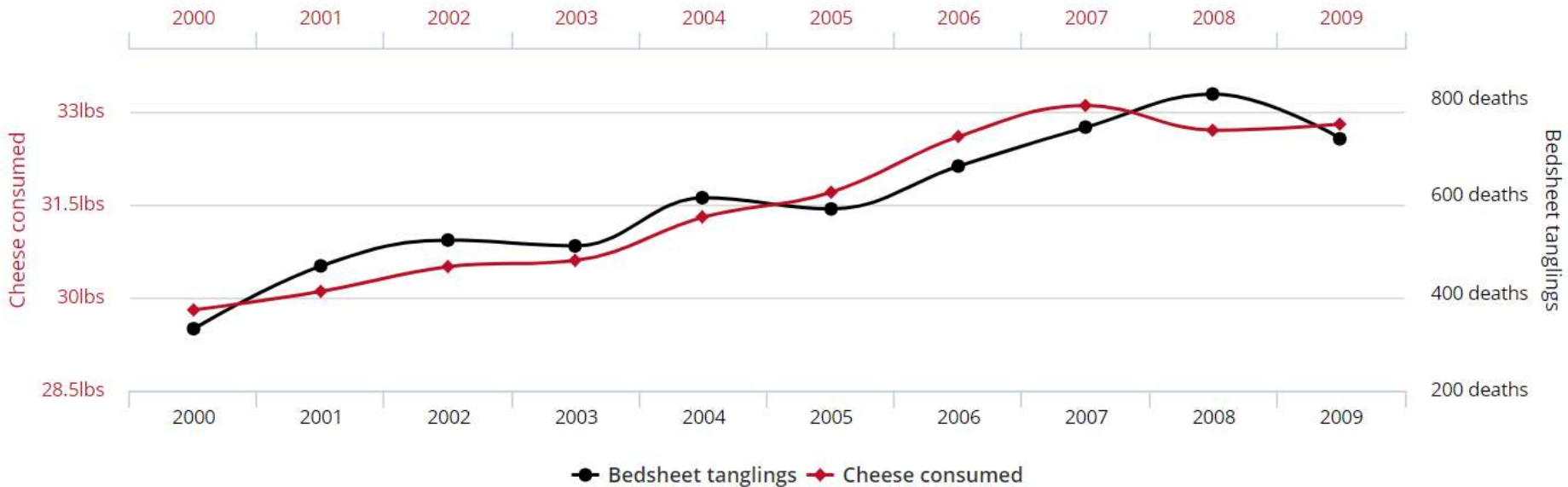
# Coming up...

- ▶ Lab Assignment 8 due Friday 4/5 11:55pm (extension)
- ▶ Peer Evaluation 2 due Tuesday 4/9 11:55 pm
- ▶ Problem Set 6 due Wednesday 4/10
- ▶ Readiness Assessment 7 Wednesday 4/10
- ▶ Don't forget the Project Stage 2 due in ~2 weeks



# Per capita cheese consumption correlates with Number of people who died by becoming tangled in their bedsheets

Correlation: 94.71% (r=0.947091)



Data sources: U.S. Department of Agriculture and Centers for Disease Control & Prevention

tylervigen.com

# What's new in Unit 6?

		Analyses and Ways to Conduct Them (that we know so far)	
Types of Variable(s) Involved	Population Parameter	Confidence Interval for the Population Parameter	Hypothesis Test for the Population Parameter <i>H<sub>0</sub>: Pop. Param = #</i> <i>H<sub>a</sub>: Pop. Param (≠ or &lt; or &gt;) #</i>
Single <b>Numerical</b> Variable	$\mu$	CLT Confidence Interval (Unit 3+4) Bootstrap Confidence Interval (Unit 4)	CLT Hypothesis Test (Unit 3+4) Bootstrap Hypothesis Test (Unit 4)
	$\mu_{diff}$	CLT Confidence Interval (Unit 4) Bootstrap Confidence Interval (Unit 4)	CLT Hypothesis Test (Unit 4) Bootstrap Hypothesis Test (Unit 4)
	<b>Median</b>	Bootstrap Confidence Interval (Unit 4)	Bootstrap Hypothesis Test (Unit 4)
Single <b>Categorical</b> Variable (2 levels)	<b>p</b>	CLT Confidence Interval (Unit 5) Bootstrap Confidence Interval (Unit 4)	CLT Hypothesis Test (Unit 5) Bootstrap Hypothesis Test (Unit 4) Randomization Testing (Unit 5-Selecting ball/chips out of bag, rolling dice)
<b>Numerical Response</b> Variable <b>Categorical Explanatory</b> Variable (2 levels)	$\mu_1 - \mu_2$	CLT Confidence Interval (Unit 4) Bootstrap Confidence Interval (Unit 5)	CLT Hypothesis Test (Unit 4) Randomization Testing (Unit 1-Shuffling Cards)
	<b>Median1-Median2</b>	Bootstrap Confidence Interval (Unit 5)	Randomization Testing (Unit 1-Shuffling Cards)
<b>Categorical Response</b> Variable <b>Categorical Explanatory</b> Variable (both have 2 levels)	<b>p1-p2</b>	CLT Confidence Interval (Unit 5) Bootstrap Confidence Interval (Unit 5)	CLT Hypothesis Test (Unit 5) Randomization Testing (Unit 1-Shuffling Cards)
<b>Numerical Response</b> Variable <b>Numerical and/or Categorical Explanatory</b> Variable(s)	$\beta_i$ (i=0,1,...,k)	Regression Coefficient Confidence Interval (Units 6 +7)	Regression Coefficient Hypothesis Test (Units 6 +7)

# How could we answer this research question?

- “Is there an association between **number of dependents** and **income**?”



## Sample Data (Registered Voters)

Gender <i>(Categorical with 2 Levels)</i>	Political Affiliation <i>(Categorical with 3 Levels)</i>	Voted in 2018? <i>(Categorical with 2 Levels)</i>	Age <i>(Numerical - Continuous)</i>	Income <i>(Numerical - Continuous)</i>	Number of Dependents <i>(Numerical - Discrete)</i>
Male	Democrat	Yes	22	\$20,000	0
Female	Republican	Yes	39	\$40,000	3
Male	Independent	Yes	47	\$400,000	0
Male	Republican	Yes	18	\$129,000	0
Female	Republican	No	29	\$85,000	2
Female	Democrat	No	80	\$72,000	1
Female	Democrat	Yes	56	\$55,000	0
Male	Independent	Yes	72	\$34,000	0
...	...	...	...	...	...

# “Is there a **linear** association between **number of dependents** and **income**?”

Numerical response variable

Numerical explanatory variable



## Simple Linear Regression

\*provided necessary conditions are met

$$\widehat{Income} = \beta_0 + \beta_1 (\text{Num. of Dependents})$$

$$H_0: \beta_1 = 0$$

→ No **linear** association

$$H_a: \beta_1 \neq 0$$

→ **Linear** Association

# Is there more than one way to answer this research question?

- “Is there an association between **gender** and **income**?”



<https://www.wsj.com/articles/parsing-the-gender-pay-gap-1542917969>

## Sample Data (Registered Voters)

Gender <i>(Categorical with 2 Levels)</i>	Political Affiliation <i>(Categorical with 3 Levels)</i>	Voted in 2018? <i>(Categorical with 2 Levels)</i>	Age <i>(Numerical - Continuous)</i>	Income <i>(Numerical - Continuous)</i>	Number of Dependents <i>(Numerical - Discrete)</i>
Male	Democrat	Yes	22	\$20,000	0
Female	Republican	Yes	39	\$40,000	3
Male	Independent	Yes	47	\$400,000	0
Male	Republican	Yes	18	\$129,000	0
Female	Republican	No	29	\$85,000	2
Female	Democrat	No	80	\$72,000	1
Female	Democrat	Yes	56	\$55,000	0
Male	Independent	Yes	72	\$34,000	0
...	...	...	...	...	...

# “Is there an association between **gender** and **income?**”

Numerical response variable

Categorical explanatory variable with 2 levels



<https://www.wsj.com/articles/parsing-the-gender-pay-gap-1542917969>

**Independent Means Hypothesis Test**  
\*provided necessary conditions are met

$$H_0: \mu_{male} - \mu_{female} = 0$$

→ No association

$$H_a: \mu_{male} - \mu_{female} \neq 0$$

→ Association

# “Is there a linear association between gender and income?”

Numerical response variable

Categorical explanatory variable with 2 levels



<https://www.wsj.com/articles/parsing-the-gender-pay-gap-1542917969>

## Simple Linear Regression

\*provided necessary conditions are met

$$\widehat{Income} = \beta_0 + \beta_1(\text{Gender: Female})$$

$$H_0: \beta_1 = 0$$

$$H_a: \beta_1 \neq 0$$

→ No linear association

→ Linear Association

If  $X \sim \text{Bin}(n, p)$  and when certain conditions are met...

$$X \sim N(\text{mean} = np, \text{standard dev.} = \sqrt{np(1-p)})$$

When other certain conditions are met...

$$\bar{x} \sim N(\text{mean} = \mu, \text{standard dev./error} = \frac{\sigma}{\sqrt{n}})$$

When other certain conditions are met...

$$\bar{x}_1 - \bar{x}_2 \sim N(\text{mean} = \mu_1 - \mu_2, \text{standard dev./error} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}})$$

When other certain conditions are met...

$$\hat{p} \sim N(\text{mean} = p, \text{standard dev./error} = \sqrt{\frac{p(1-p)}{n}})$$

When other certain conditions are met...

$$\hat{p}_1 - \hat{p}_2 \sim N\left(\text{mean} = p_1 - p_2, \text{standard dev./error} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}\right)$$

When other certain conditions are met...

$$b_i \sim N(\text{mean} = \beta_i, \text{standard dev./error} = \underline{\hspace{2cm}})$$

**What's new in Unit 6? New "observation" that follows a normal distribution under certain conditions.**

# Outliers

