

INTR 5057 Research Design & Methods

Juraj Medzihorsky



Day 11, 2016-12-02

Homework #1

- ▶ Almost done grading. Overall OK.
- ▶ Keep perspective when you see your grade.
- ▶ Common issues:
 - ▶ Unfocused Qs.
 - ▶ Focus on single cases even for “effects of causes” Qs.
 - ▶ Mixing up Qs and Hs.
 - ▶ Causes vs. conditions.

Summarizing a Single Variable

How to summarize the following information?

<hr/>
x
<hr/>
1
4
0
3
3
2
<hr/>

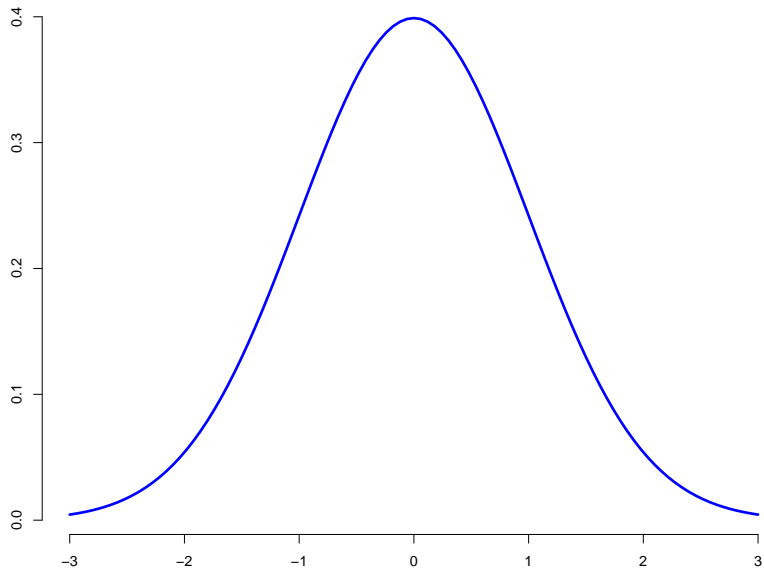
Summarizing a Single Variable

- ▶ Average (mean)
- ▶ Mode
- ▶ Median
- ▶ Midrange

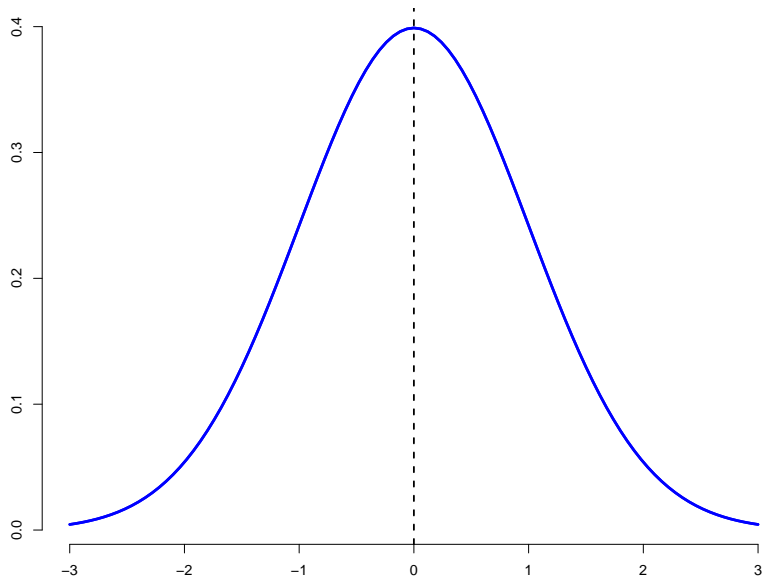
Summarizing a Single Variable

- ▶ Average: Sum divided by the number of values.
- ▶ Mode: The most common value.
- ▶ Median: Half of the observations have less, half more.
- ▶ Midrange: Midpoint between maximum and minimum.

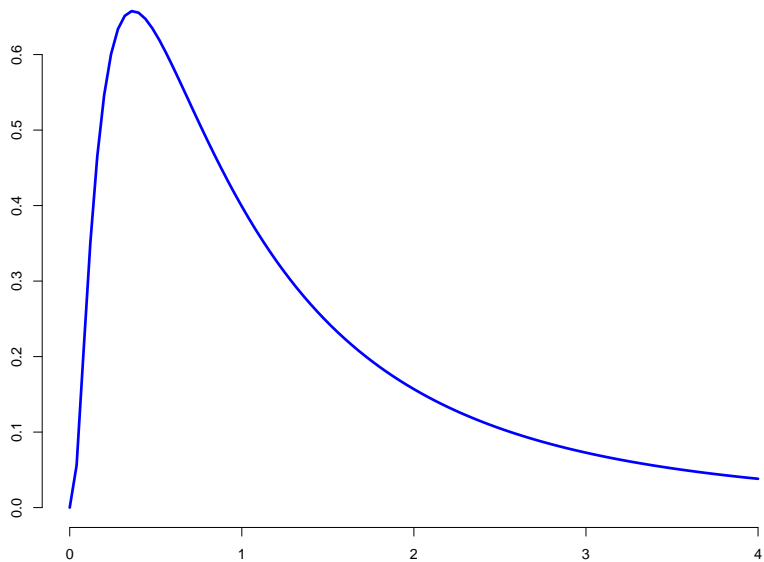
A Continuous Variable



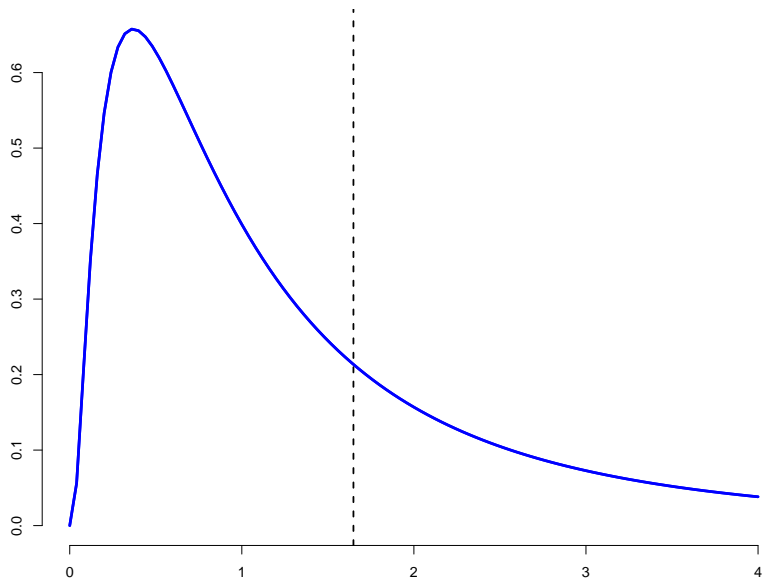
Mean, Median, Mode, Midrange



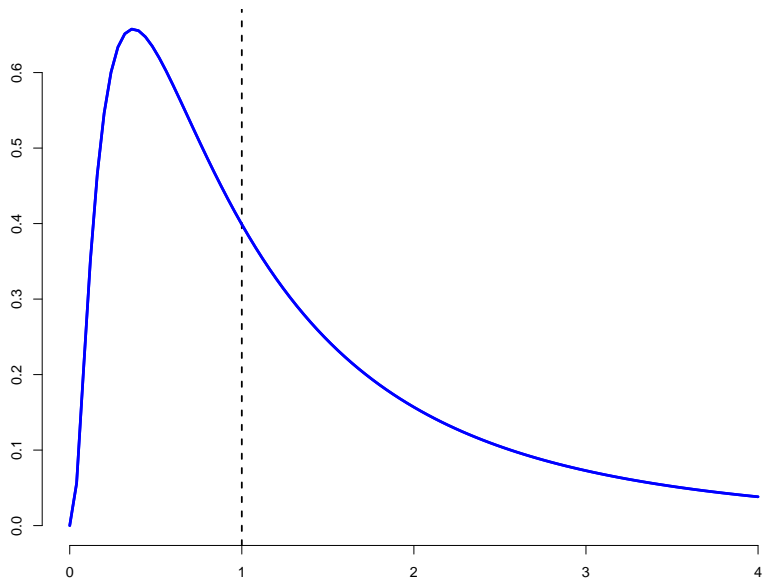
A Continuous Variable



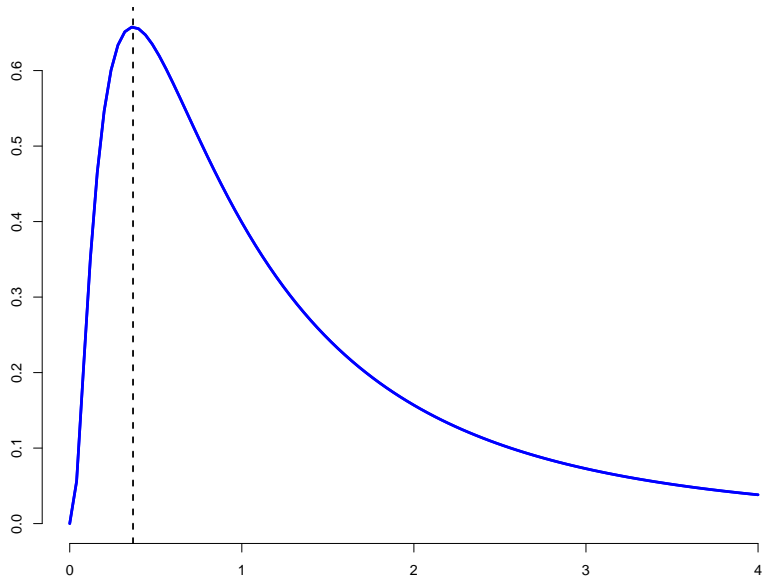
Mean



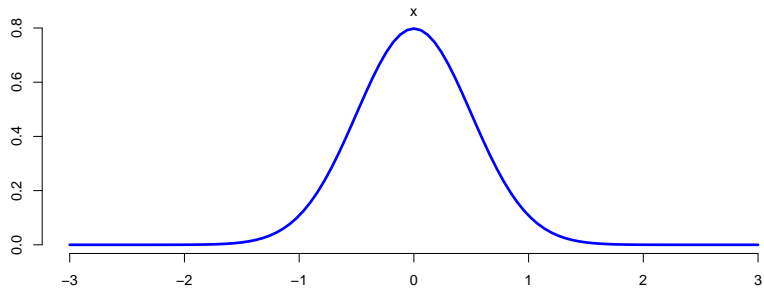
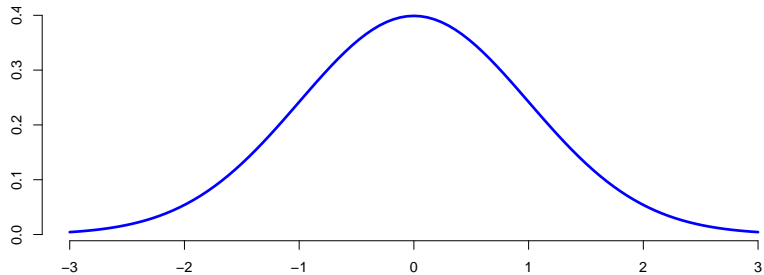
Median



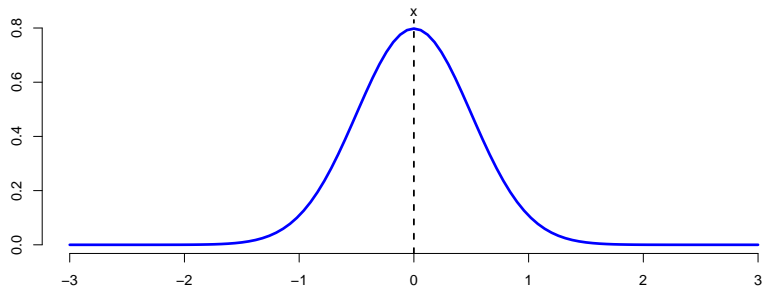
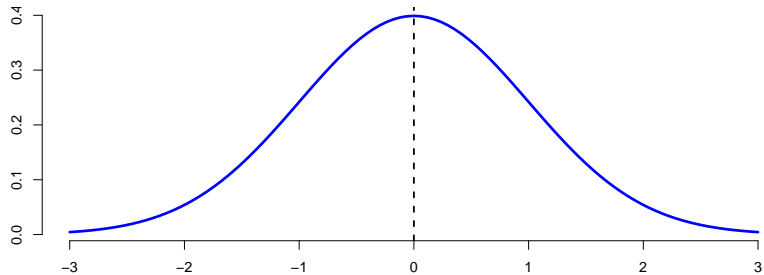
Mode



Two Continuous Variables



Same Mean



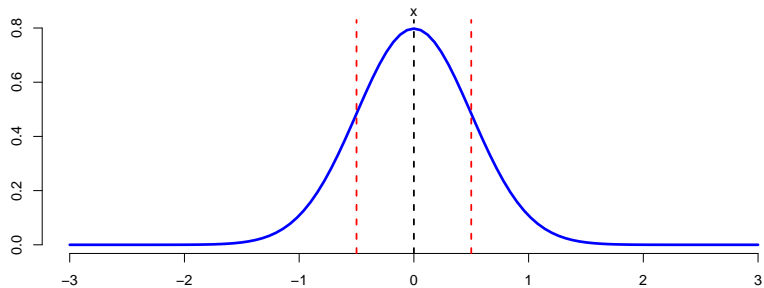
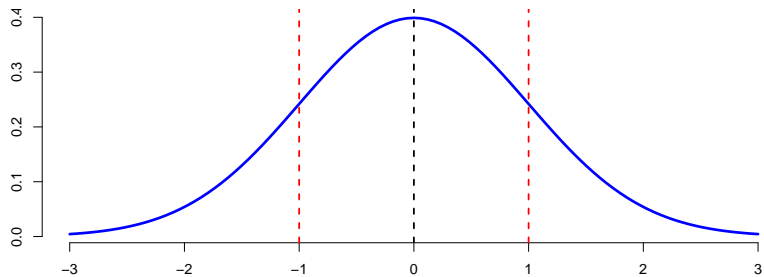
Standard Deviation

A measure of dispersion from the average.

Square root of the average squared distance from the average.

$$\sigma = \sqrt{\frac{\sum(\mu - x_i)^2}{N}}$$

Standard Deviation



Summarizing Two Variables

How to summarize the following information?

<hr/>	
x	y
<hr/>	
1	1
0	1
1	0
1	0
0	0
1	0
<hr/>	

Summarizing Two Variables

- ▶ Summarize each of them separately.
- ▶ Capture information about their **association**.

Cross Table

		y	
		0	1
x	1	3	1
	0	1	1

Cross Table

		y	
		0	1
x	1	a	b
	0	c	d

Cross Product Ratio

A measure of association between two binary variables also known as odds ratio.

If cross product ratio = 1 then the variables are independent.

$$cpr = \frac{a \times d}{b \times c}$$

Cross Sum Ratio

An alternative to the cross product ratio.

$$csr = \frac{a + d}{b + c}$$

Relative Risk Ratio

A measure of association between two binary variables.

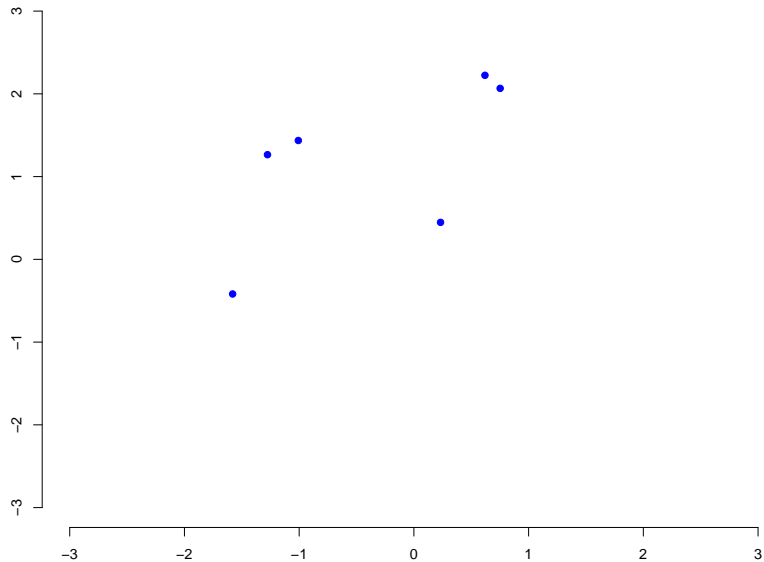
$$rr = \frac{\frac{a}{a+b}}{\frac{c}{c+d}}$$

Summarizing Two Variables

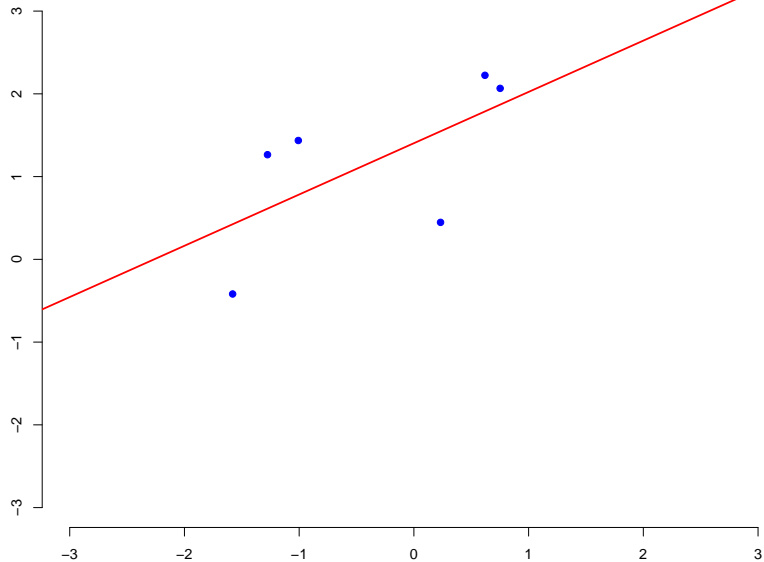
How to summarize the following information?

y	x
0.8	2.1
0.6	2.2
-1.6	-0.4
0.2	0.5
-1.3	1.3
-1.0	1.4

Scatter Plot



Correlation



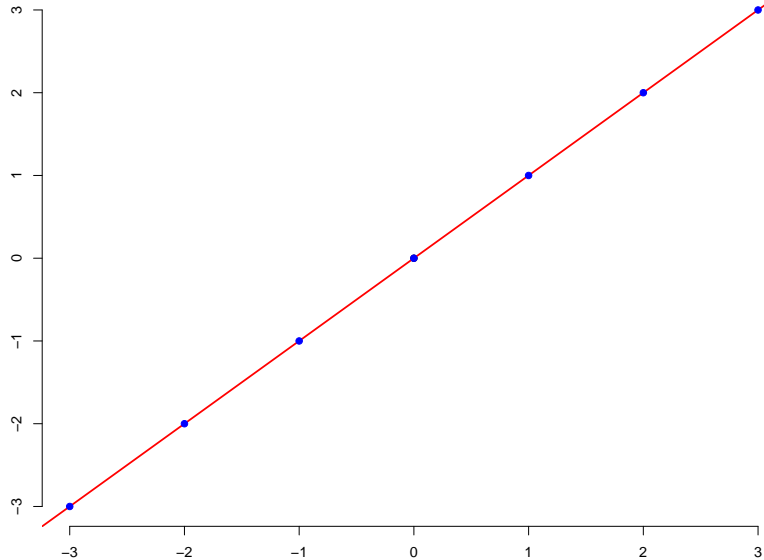
Correlation

A measure of association between two continuous variables.

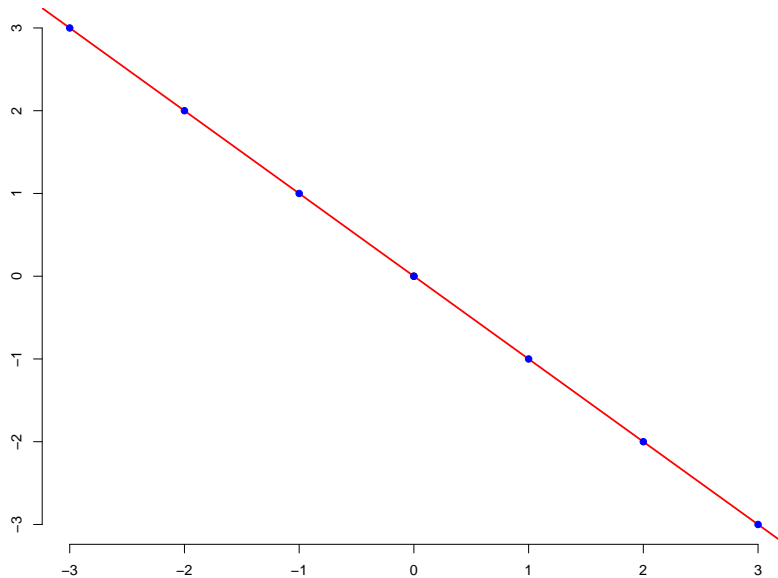
$$\rho_{x,y} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

Ranges from -1 to 1 on a closed interval.

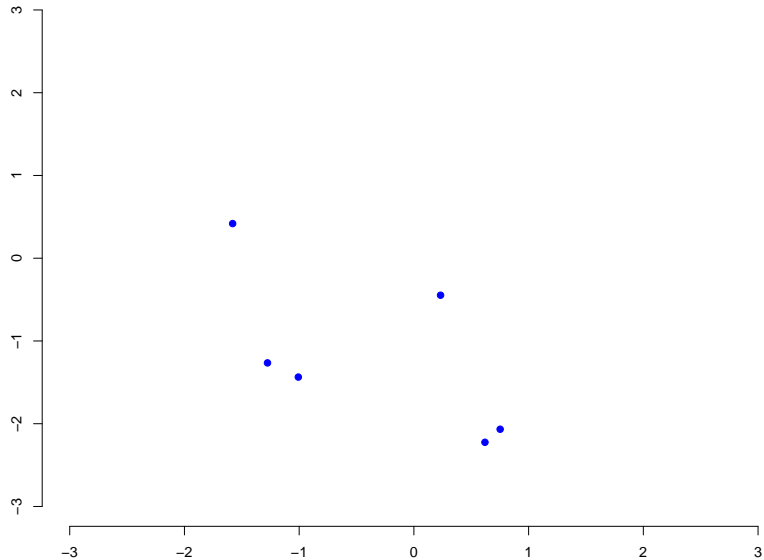
$$\rho = 1$$



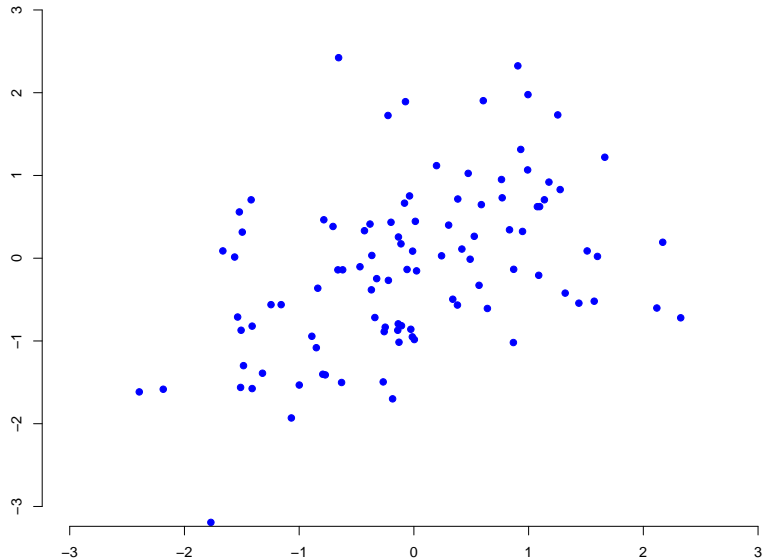
$$\rho = -1$$



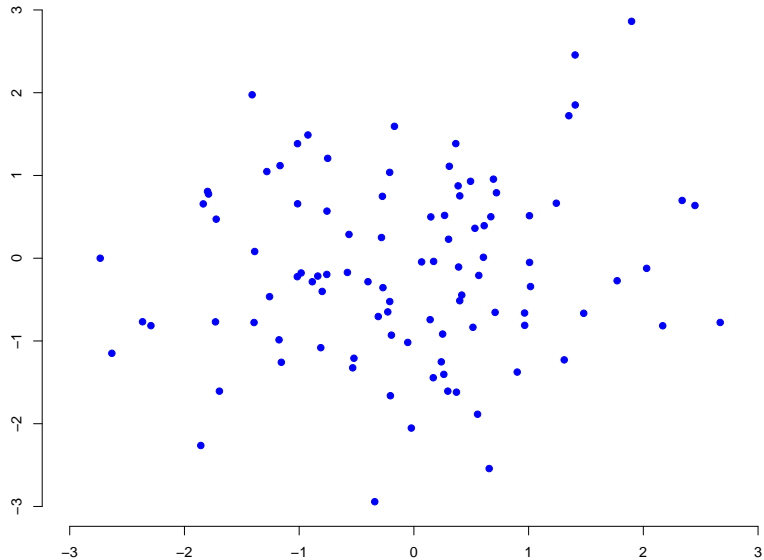
$$\rho = -0.63$$



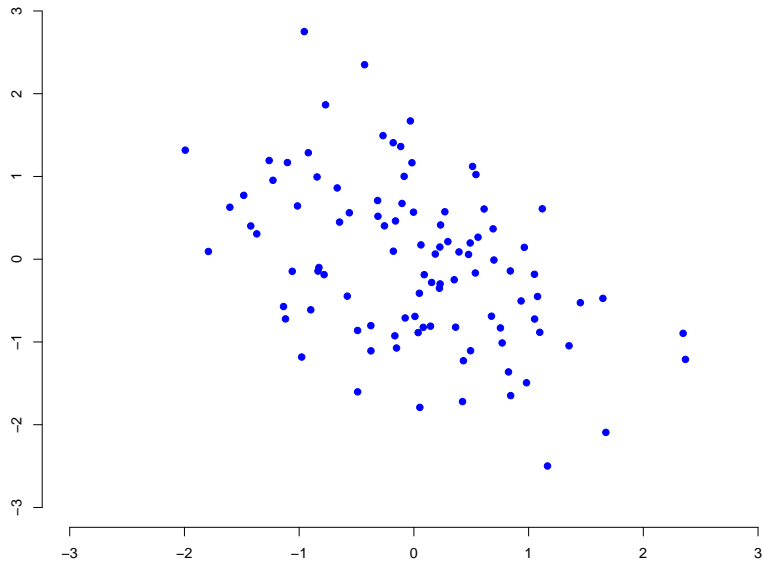
$$\rho = 0.44$$



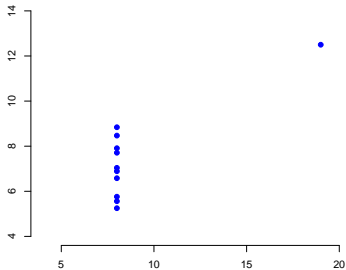
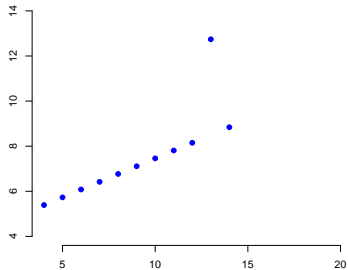
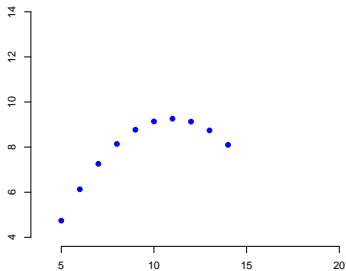
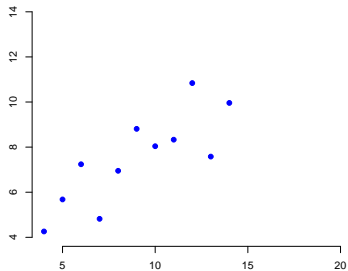
$$\rho = 0.09$$



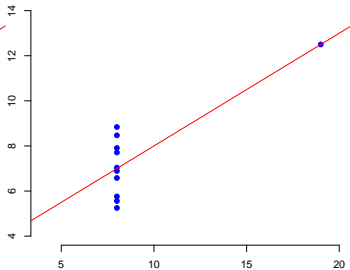
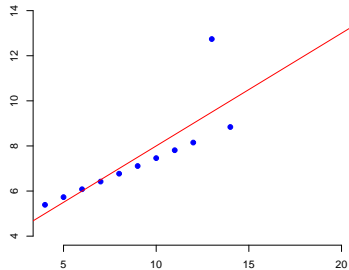
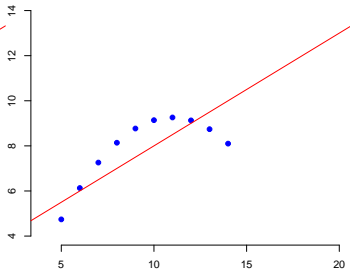
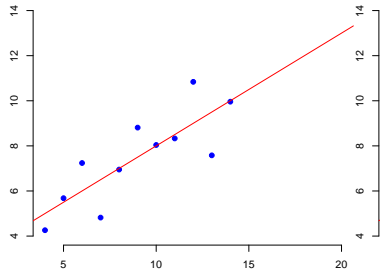
$$\rho = -0.46$$



Anscombe's Quartet



$$\rho = 0.82$$



Simpson's Paradox

The whole sample:

	heal	didn't
drug	20	20
no drug	16	24

Simpson's Paradox

Females:

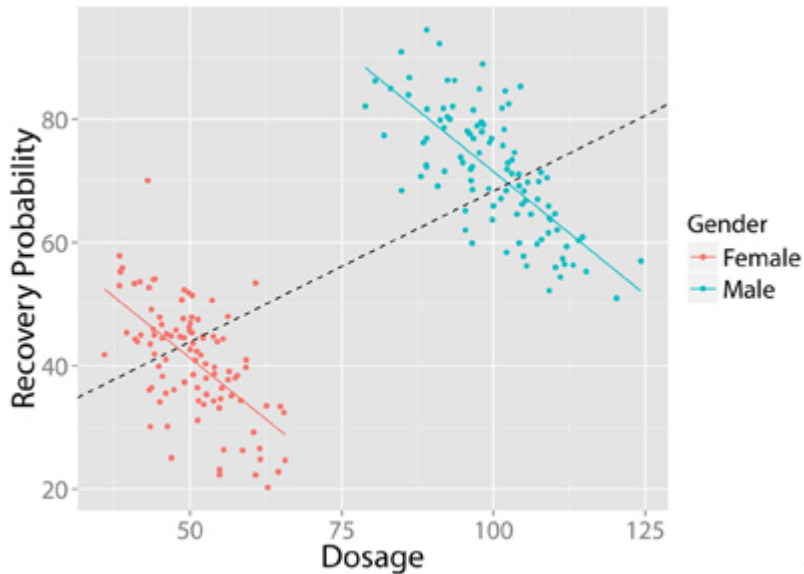
	heal	didn't
drug	2	8
no drug	9	21

Simpson's Paradox

Males:

	heal	didn't
drug	18	12
no drug	7	3

Simpson's Paradox



Simpson's Paradox

- ▶ In the whole population association in one direction.
- ▶ In subsets of the population association in the opposite direction.
- ▶ Not really a paradox when you think about it.
- ▶ A serious problem is that people rush ahead with causal interpretations.

Association & Causality

- ▶ Non-statisticians say “*correlation does not imply causation.*”
- ▶ Statisticians say “*association does not imply causation.*”
- ▶ Calling all association “correlation” is like calling all motor vehicles “cars.”

Goals

Goals

- ▶ Describe.
- ▶ Explain.
- ▶ Predict/Forecast.
- ▶ ...

Table 1: Electoral and Replacement Volatility in Post-Communist Europe

	Verification		Without Bosnia-Herzegovina		Corrected Bosnia-Herzegovina	
	Electoral Volatility	Replacement Volatility	Electoral Volatility	Replacement Volatility	Electoral Volatility	Replacement Volatility
GDP Change from 1989	0.639 (0.693)	-4.623*** (1.326)	0.116 (3.206)	-6.066 (7.178)	0.004 (3.233)	-6.002 (6.609)
GDP Change Between Elections	-2.059 (5.219)	9.019 (10.128)	-1.891 (5.898)	6.677 (10.704)	-1.076 (5.229)	4.576 (10.064)
Effective Number of Electoral Parties	0.446 (0.313)	-0.346 (0.533)	0.452 (0.316)	-0.264 (0.558)	0.471 (0.326)	-0.462 (0.546)
Log Weighted District Magnitude	-0.784 (0.887)	0.638 (2.931)	-0.789 (0.882)	0.603 (2.893)	-0.824 (0.886)	0.820 (2.872)
Presidential System	-4.631 (4.126)	6.784 (9.435)	-4.847 (4.606)	5.532 (10.241)	-4.928 (4.623)	6.659 (10.296)
Semi-Presidential System	-2.788 (2.211)	4.255 (5.897)	-2.813 (2.286)	4.017 (5.885)	-2.596 (2.266)	2.621 (5.887)
Proportional Representation	0.827 (2.228)	0.077 (6.004)	0.852 (2.265)	-0.146 (5.943)	0.987 (2.223)	-0.739 (5.948)
Ethnic Fractionalization	-6.163 (6.397)	-2.677 (18.978)	-6.716 (6.784)	-5.298 (22.939)	-5.713 (6.772)	-11.828 (22.931)
Years Since Collapse of Communism	0.848 (0.807)	-2.633 (2.153)	0.828 (0.863)	-1.989 (2.117)	0.732 (0.797)	-1.959 (1.976)
Years Since Collapse Squared	-0.031 (0.042)	0.070 (0.101)	-0.029 (0.044)	0.045 (0.097)	-0.026 (0.043)	0.049 (0.093)
Constant	13.059** (5.318)	41.941*** (13.329)	13.586*** (5.115)	43.661*** (14.057)	12.885** (5.034)	48.191*** (14.509)
Countries	21	21	20	20	21	21
Pairs of Elections	89	89	86	86	89	89
R ²	0.116	0.139	0.114	0.119	0.112	0.109

* $p < 0.10$; ** $p < 0.05$; *** $p < 0.01$ (two-tailed).

Regression

- ▶ Whether we like it or not, **regression** is the workhorse of quantitative social science.

Regression

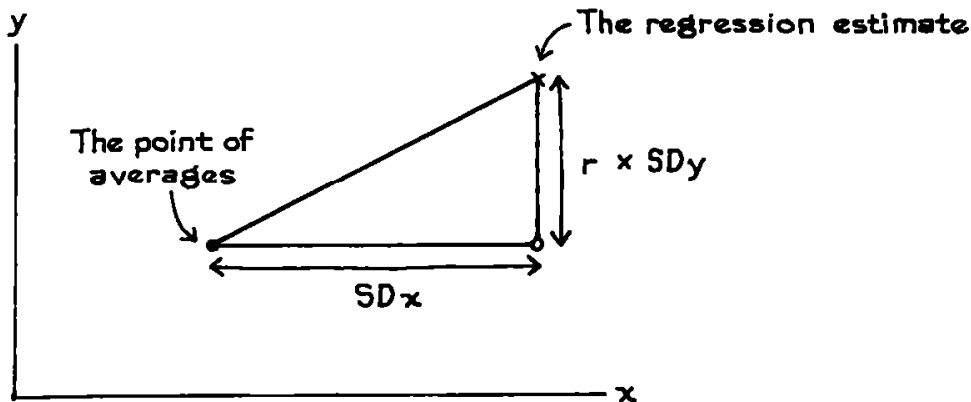
- ▶ Whether we like it or not, **regression** is the workhorse of quantitative social science.
- ▶ Any previous experiences with regression?

Regression

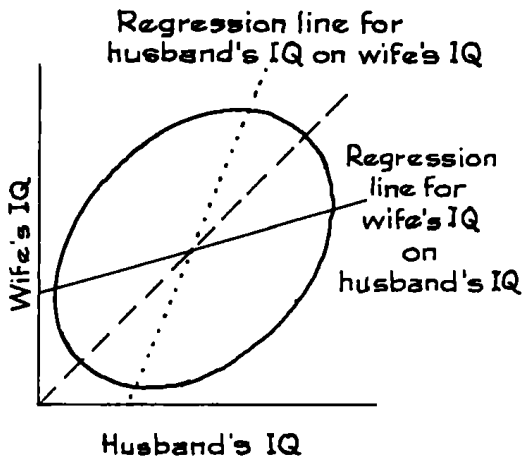
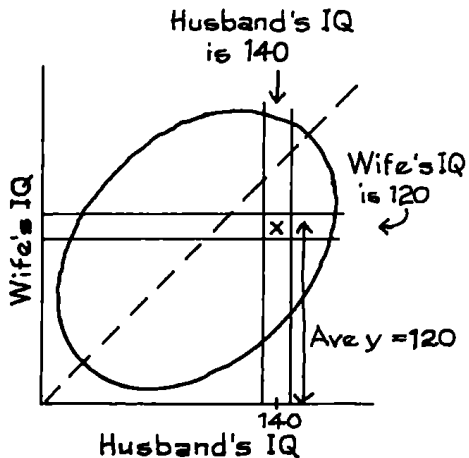
- ▶ One variable as a function of one or more other variables.
- ▶ **Conditional association.**
- ▶ Typically used to **explain** or **predict**.

Regression

Figure 2. Regression method. When x goes up by one SD, the average value of y only goes up by r SDs.



Regression



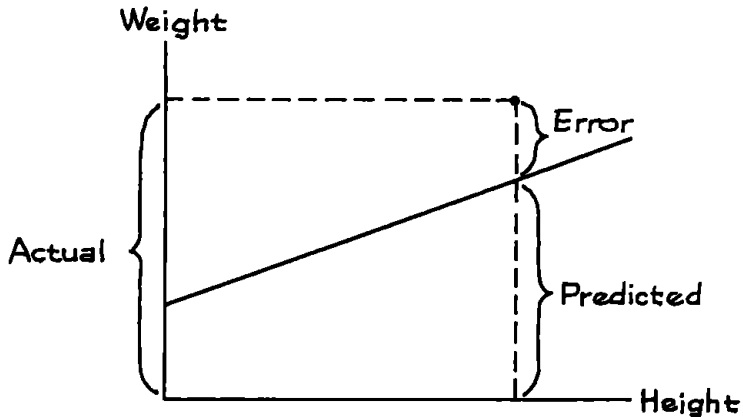
Linear Regression

$$y_i = \alpha + \beta \times x_i + \epsilon_i$$

- ▶ y : LHS, “dependent variable,” outcome
- ▶ x : RHS, “independent variable,” predictor, determinant
- ▶ α : intercept, “constant”
- ▶ β : slope, coefficient, “effect”
- ▶ ϵ : residual, “error”

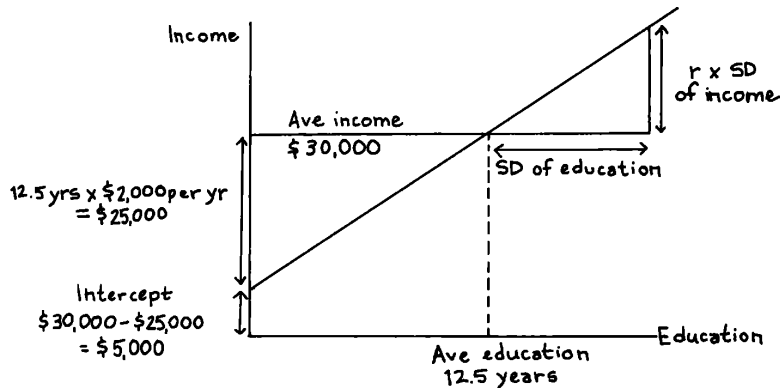
Regression

Figure 2. Prediction error equals vertical distance from the line.



Regression

Figure 3. Finding the slope and intercept of the regression line.



Linear Regression

$$y_i = \alpha + \beta x_i + \epsilon_i$$

- ▶ Ordinary Least Squares (OLS): $\sum_i \epsilon_i^2$.
- ▶ Probabilistic I.:

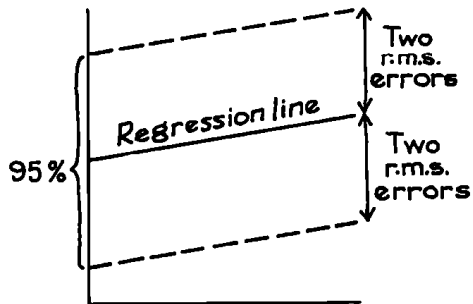
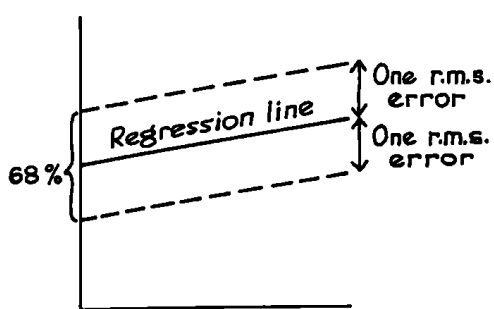
$$\epsilon_i \sim \text{Normal}(0, \sigma)$$

- ▶ Probabilistic II.:

$$y_i \sim \text{Normal}(\alpha + \beta x_i, \sigma)$$

Regression

Figure 3. Rule of thumb. About 68% of the points on a scatter diagram fall inside the strip whose edges are parallel to the regression line, and one r.m.s. error away (up or down). About 95% of the points are in the wider strip whose edges are parallel to the regression line, and twice the r.m.s. error away.



Regression

Figure 10. A football-shaped scatter diagram. Take the points inside a narrow vertical strip. Their y-values are a new data set. The new average is given by the regression method. The new SD is given by the r.m.s. error of the regression line. Inside the strip, a typical y-value is around the new average—give or take the new SD.

