

Quantitative Methods for Political Science

Descriptive Statistics
September 14, 2022

Three Parts of Research

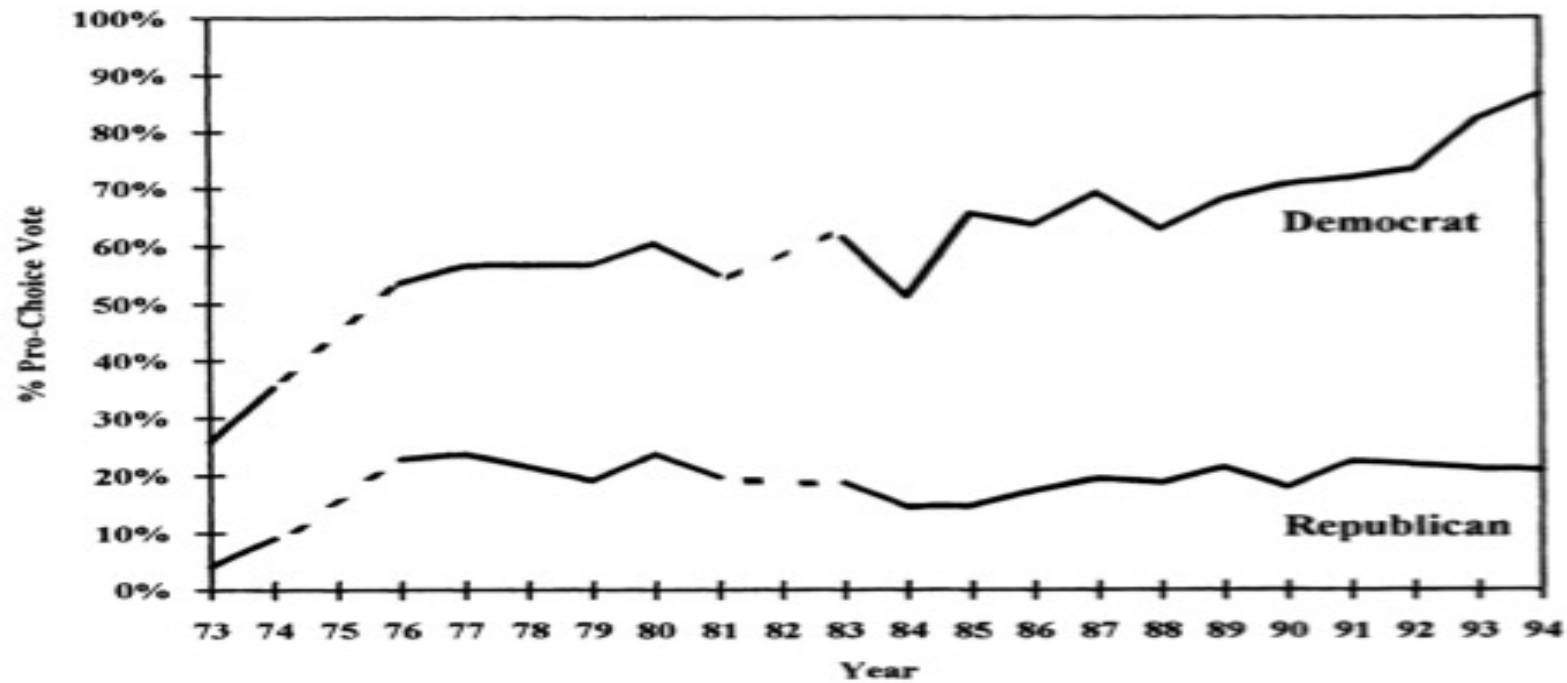
- Design—what is the question? How will we try to answer it?
- Description—what do the data say?
- Inference—What can we learn about the question from the evidence? (And, technically, what can we infer about the population given the sample)

Partisan Attitudes toward Abortion

- Have partisan positions on abortion shifted since Roe vs. Wade (1973)?
- Two levels:
 - Elites—members of Congress
 - Masses—individuals identifying with each political party
- From Greg D. Adams (1997), “Abortion: Evidence of an Issue Evolution.” *American Journal of Political Science* 41(3): 718-737.

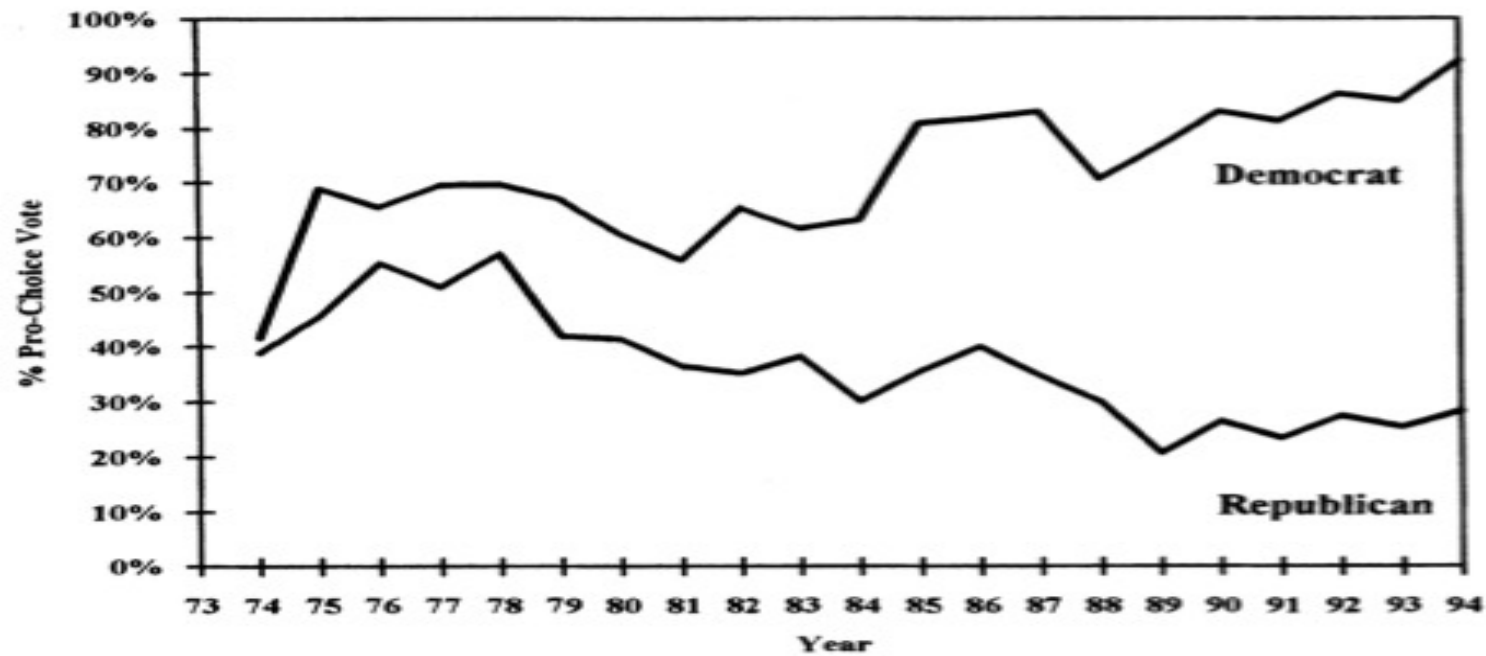
Evidence for Elites

Figure 1A. Percentage of House Abortion Votes That Are Pro-Choice



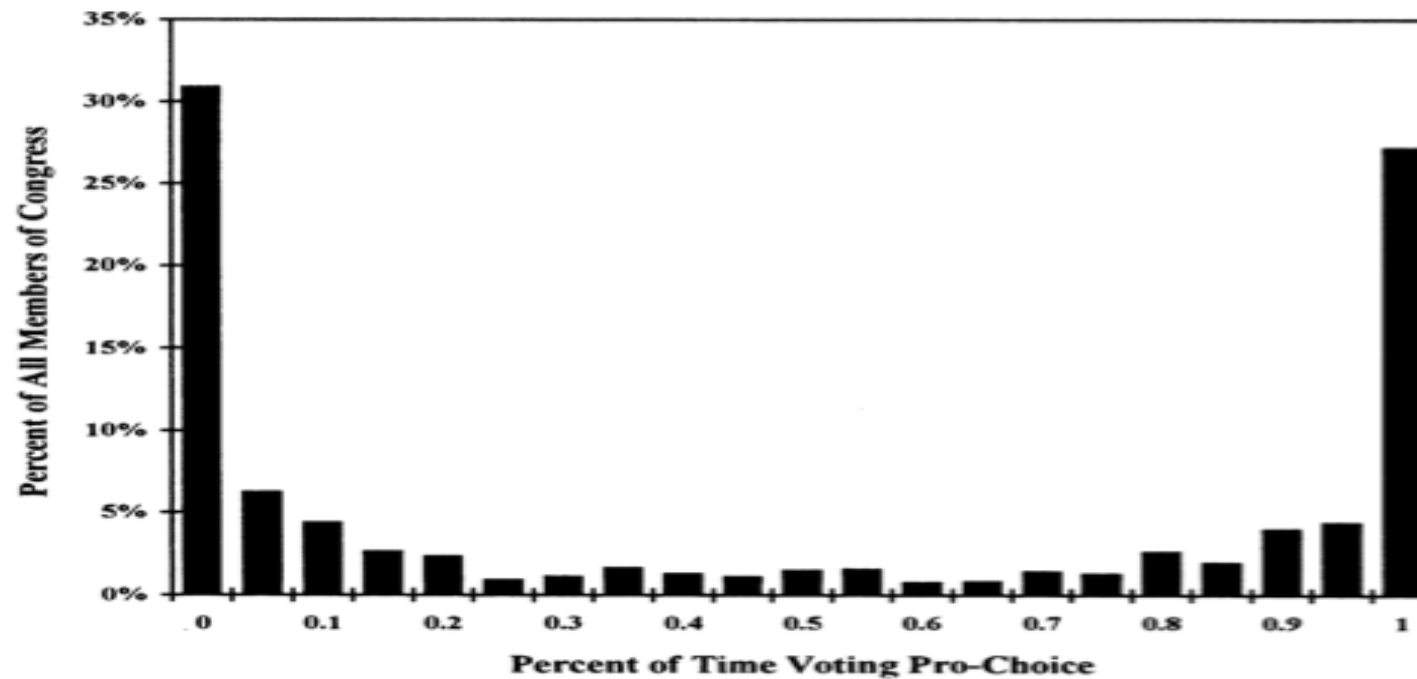
Evidence for Elites

Figure 1B. Percentage of Senate Abortion Votes That Are Pro-Choice



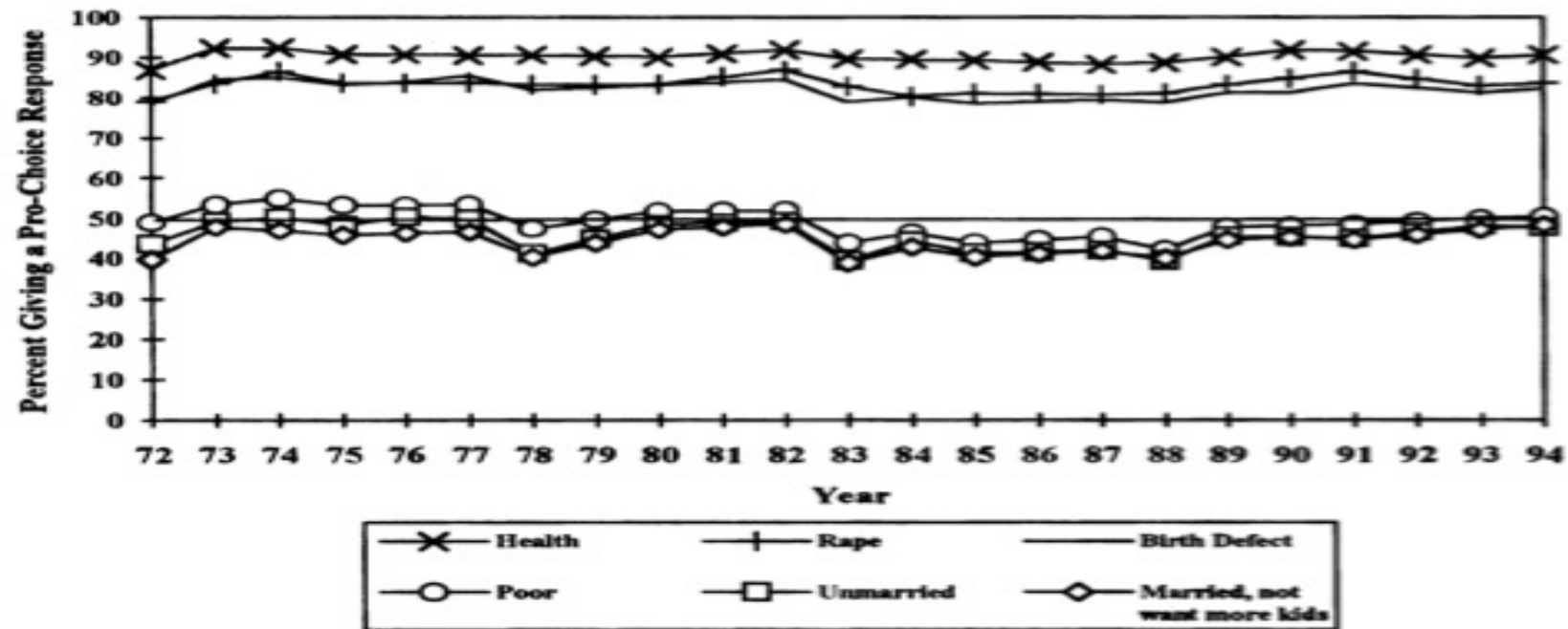
Why? Do representatives switch?

Figure 3. Distribution of Individual Legislators' Abortion Votes Over Entire Career



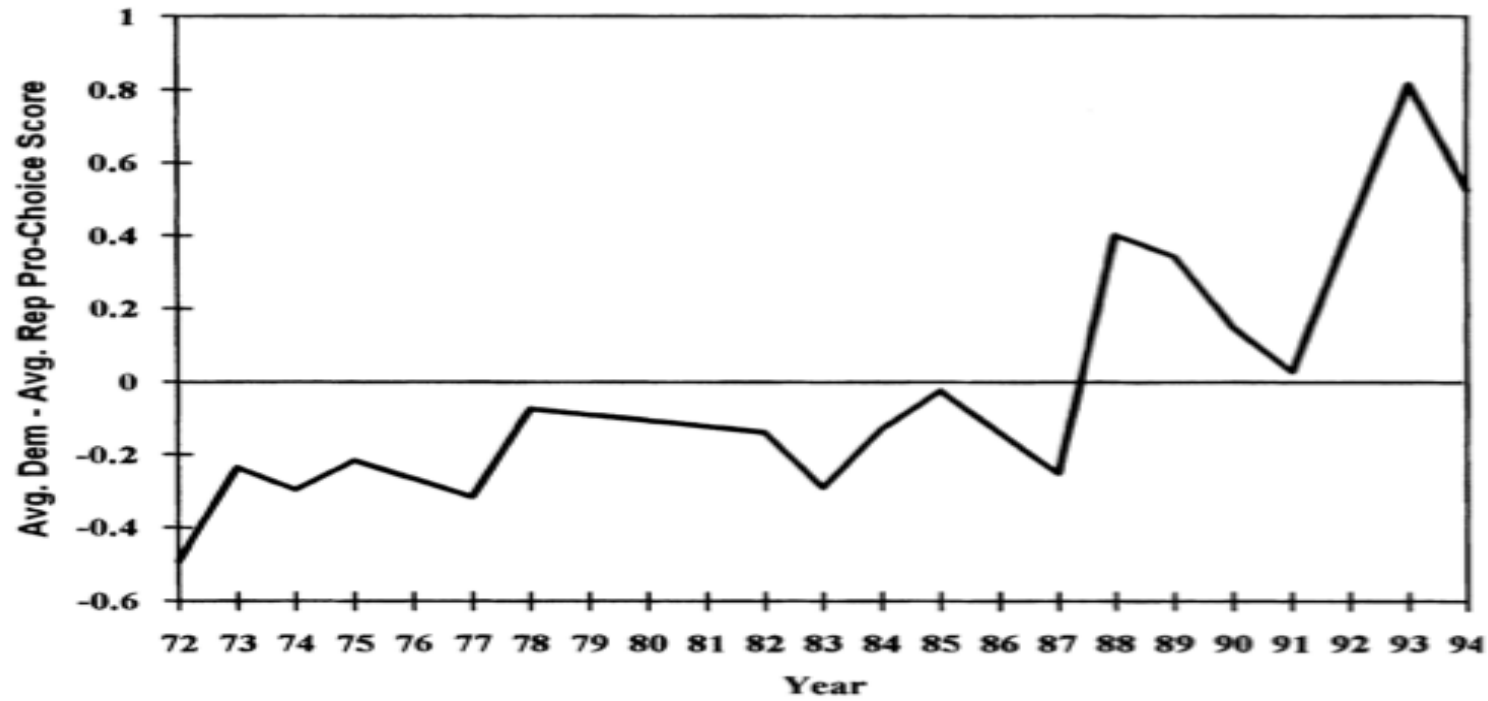
What about the masses?

Figure 4. Support for Abortion Rights Among Survey Respondents
Source: General Social Surveys, 1972–94.



Partisan differences among the masses

Figure 5. Difference Between Average Mass Republican and Democrat Pro-Choice Scores



Research Design

- First step—ask an interesting question
- Then, develop a theory to answer that question
- Think about how to design a study to answer that question
 - Who should participate in the study?
 - How can we isolate the relationship between variables?
 - How can we get the data we need to answer the question?
 - How will you measure concepts?

Description

- Summarize the raw data
- Important rule—the techniques you can use to describe (and analyze) data differ depending on type of variables you have
 - Categorical variables
 - Quantitative variables
- Present the data in a useful format
 - Can serve as exploratory data analysis

Variables

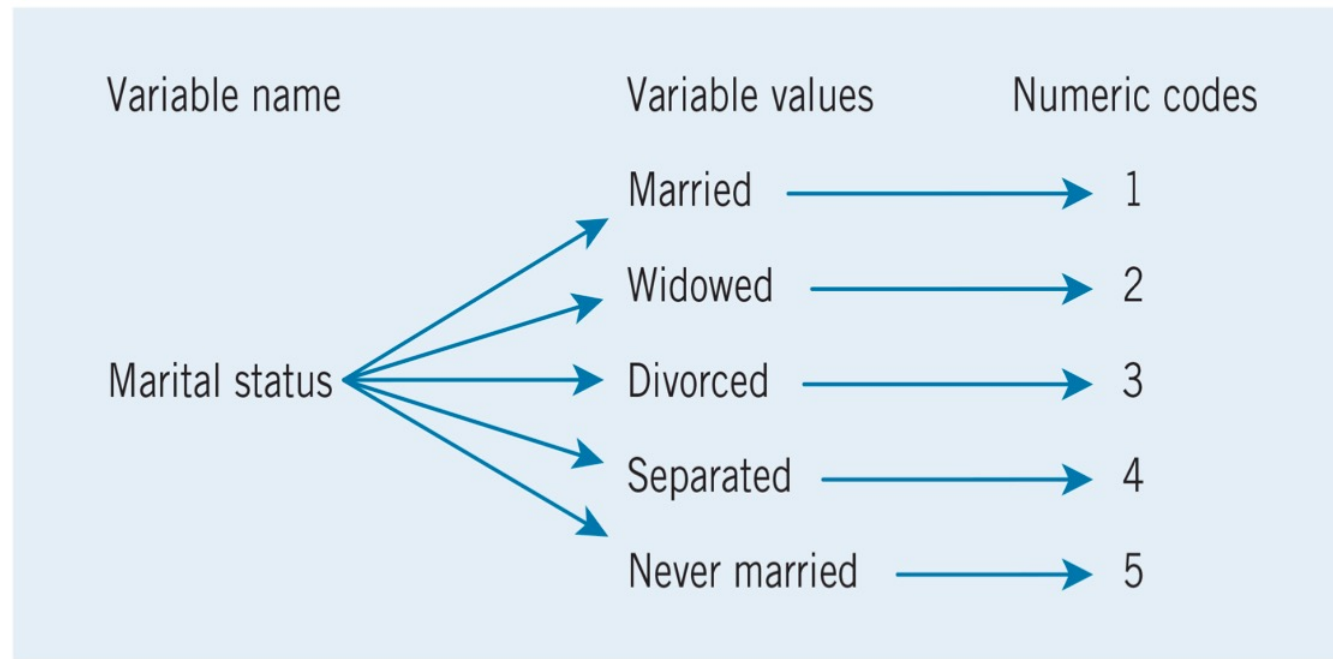
- Research pertains to some unit of analysis:
 - Individual, Household, Congressional District, State, Country, International System
- Observe characteristics of these units (observations, or cases)
- A variable is an empirical measurement of a characteristic
- Key rule—variables vary across observations
- Different types of variables based on how they are measured and what numbers mean

Types of Variables

- Quantitative/Interval/Continuous-Observations take on numerical values
- Categorical-Observations belong to one of a set of categories
 - Nominal-Categories are named
 - Ordinal-Categories are ranked
- Dichotomous variables—two values, yes/no
 - War/not war, win/lose, etc.
- Important—all variables can be measured with numbers, these numbers mean different things for quantitative vs. categorical variables
- Examples—Age, Race, Percentage of vote candidate receives, Gender
- Many concepts can be measured at different levels

Coding a Nominal Variable

Figure 2-1 Anatomy of a Variable



Transforming Variables

- Can collapse quantitative variables into ordinal (or nominal) variables
 - Ex—income to categories (low/middle/high)
- Cannot go the other way
- Generally want to retain as much information as possible
- Level of measurement for variable should be driven by theory

Describing Variables

- Distribution of a variable tells us what values it takes and how often it takes these values
- Techniques used to describe distribution of variables differ depending on type of variable
 - Quantitative—center, spread and shape of distribution
 - Categorical—percentage in each category
- Can describe variables individually and also examine relationship between them descriptively
 - But, need more rigorous analysis for actual hypothesis testing

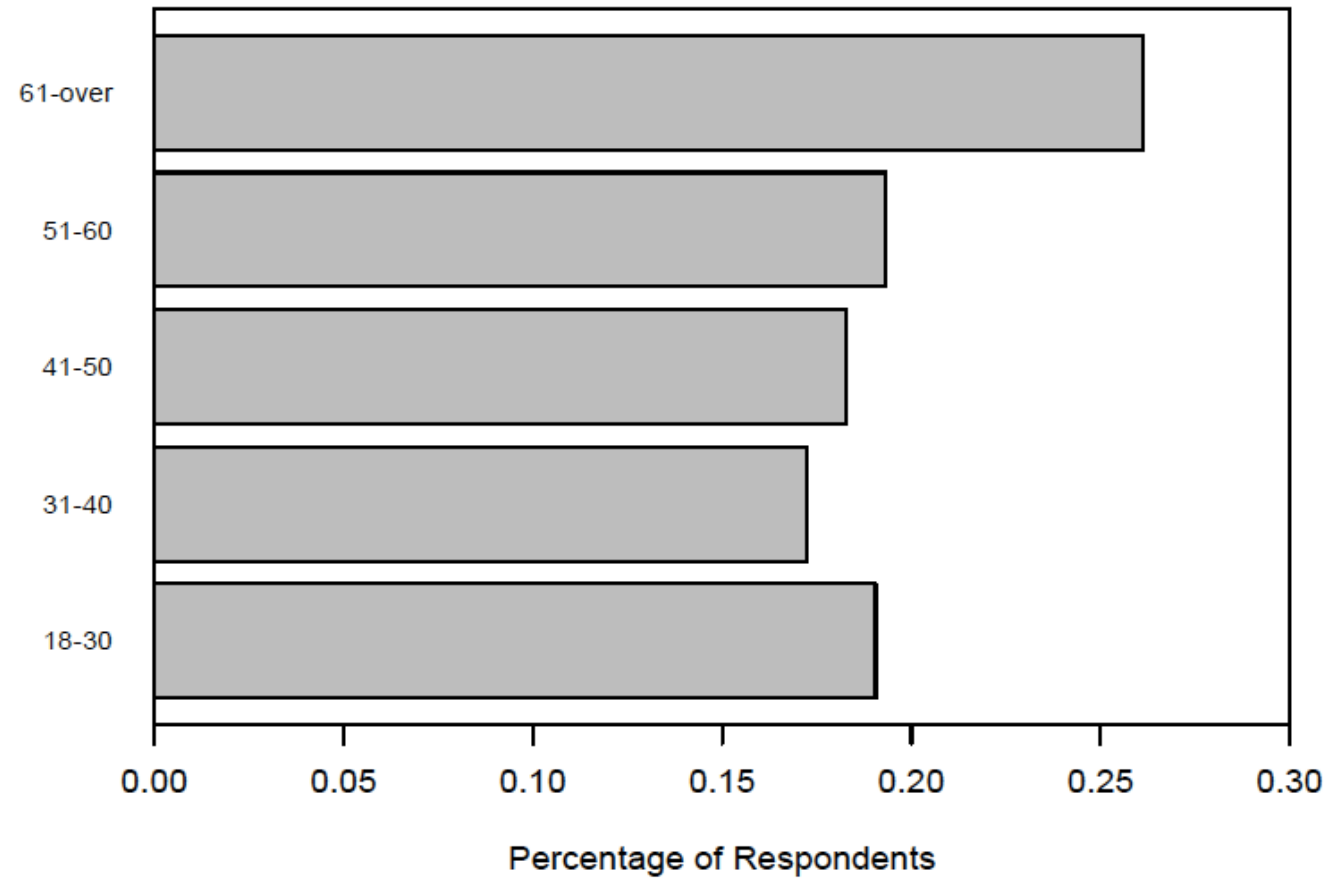
Describing Categorical Variables

- Interested in count and/or percentage of cases that fall in each category of variable
- Ways to display this:
 - Frequency Distribution
 - Bar Chart

Frequency Distribution

Age Group	Frequency	Percent (%)	Cumulative (%)
18-30	73	19.06	19.06
31-40	66	17.23	36.29
41-50	70	18.28	54.57
51-60	74	19.32	73.89
61-older	100	26.11	100
Total	383	100	100

Bar Chart



Describing Quantitative Variables

- Several graphical options for quantitative variables:
 - Dot Plot-useful for smaller data files
 - Stem-and-Leaf Plot-also for smaller data files
 - Histogram-larger data files

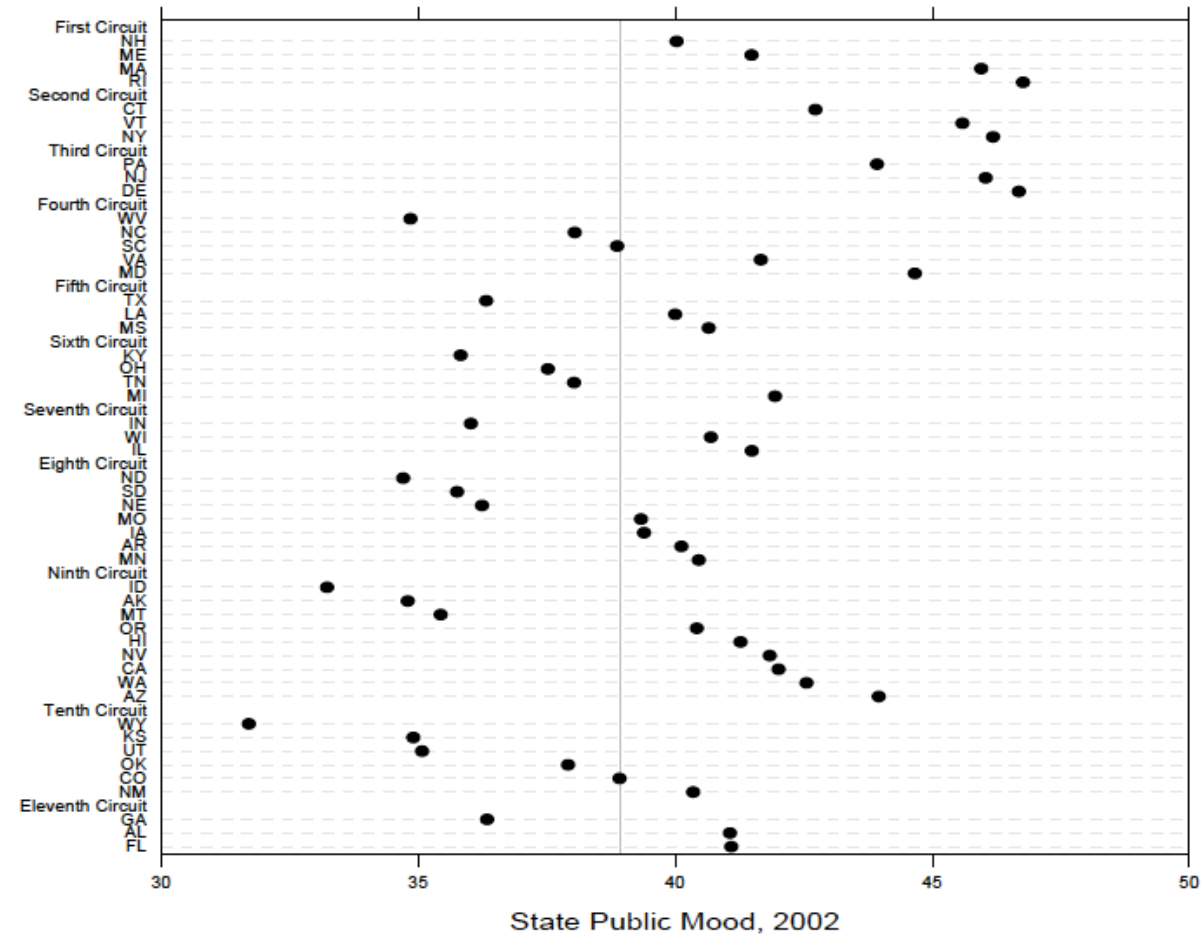
Stem-and-Leaf Plot

- Stem—all but the final (rightmost) digit
- Leaf—the final digit
- Stems go on a vertical column
- Write the leafs for each stem
- Gives you a sense of the distribution of the data
- I have never done this, nor have I ever seen anyone else present one of these

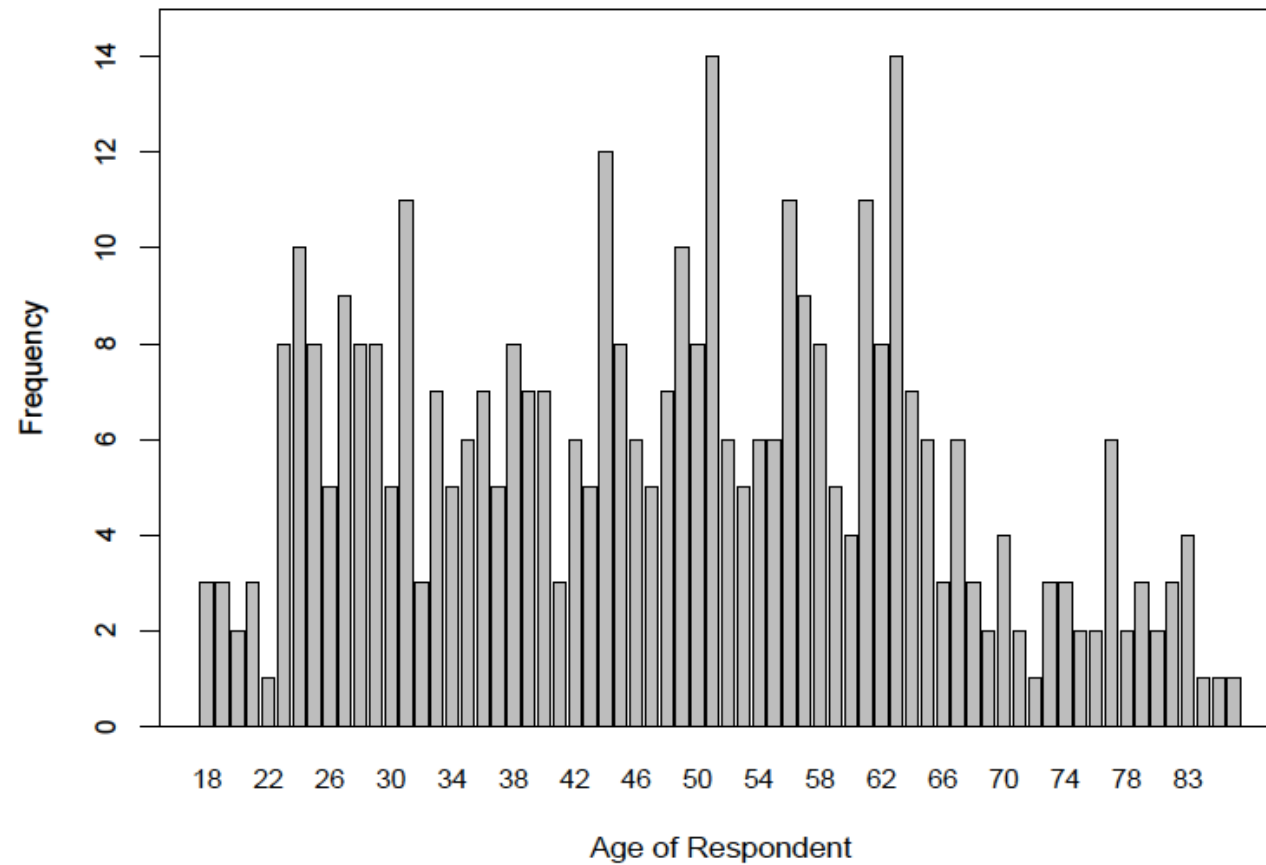
Stemplot

- These are the scores on an exam in an undergrad statistics class:
- 55, 63, 68, 69, 71, 74, 77, 79, 81, 81, 82, 83, 84, 84, 85, 86, 87, 87, 88, 88, 89, 90, 91, 93, 93, 95, 97, 98, 100
- Can also use stemplots for comparison
- Men—55, 63, 68, 69, 74, 77, 81, 81, 83, 85, 86, 87, 88, 93, 95
- Women—71, 79, 82, 84, 84, 87, 88, 89, 90, 91, 93, 97, 98, 100

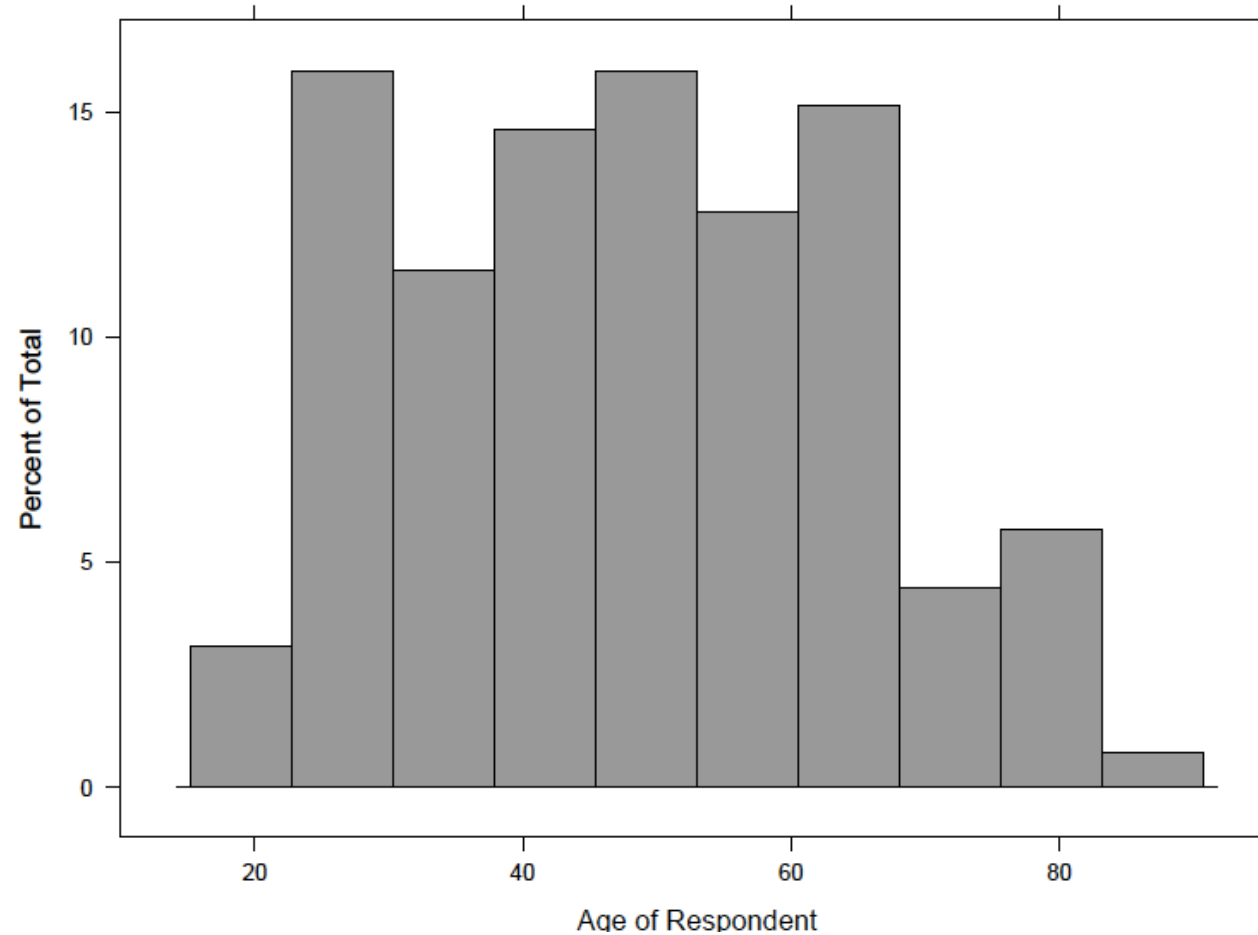
Dotplot



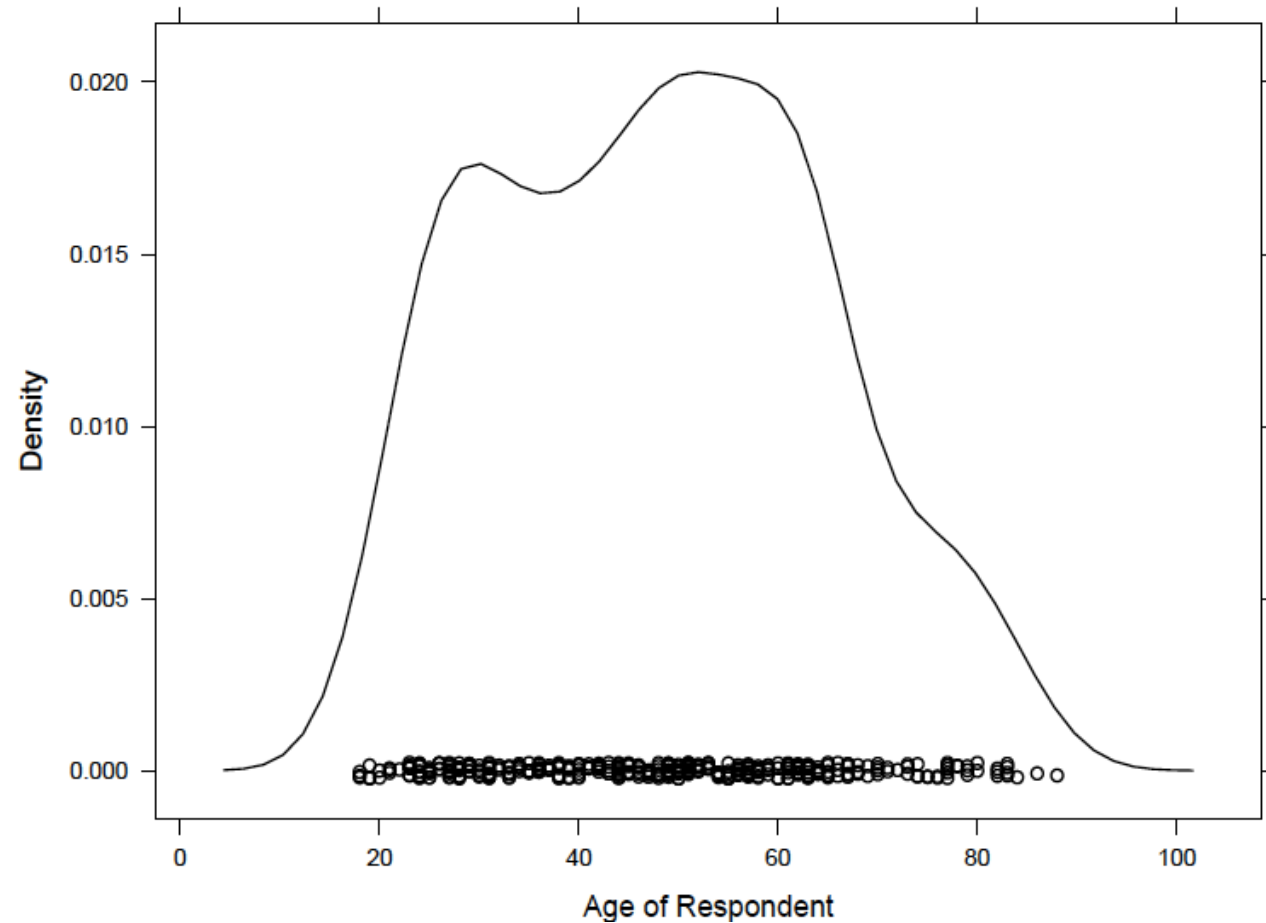
Histogram



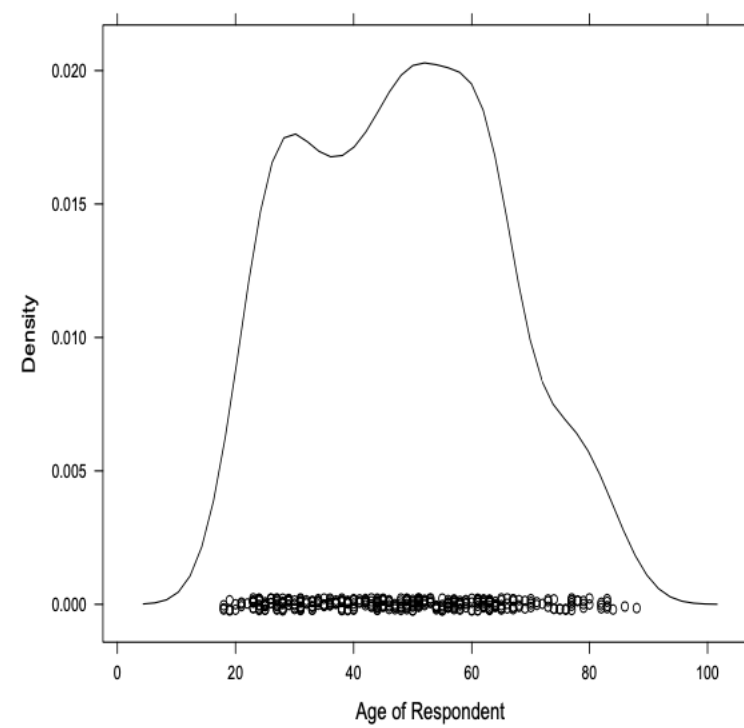
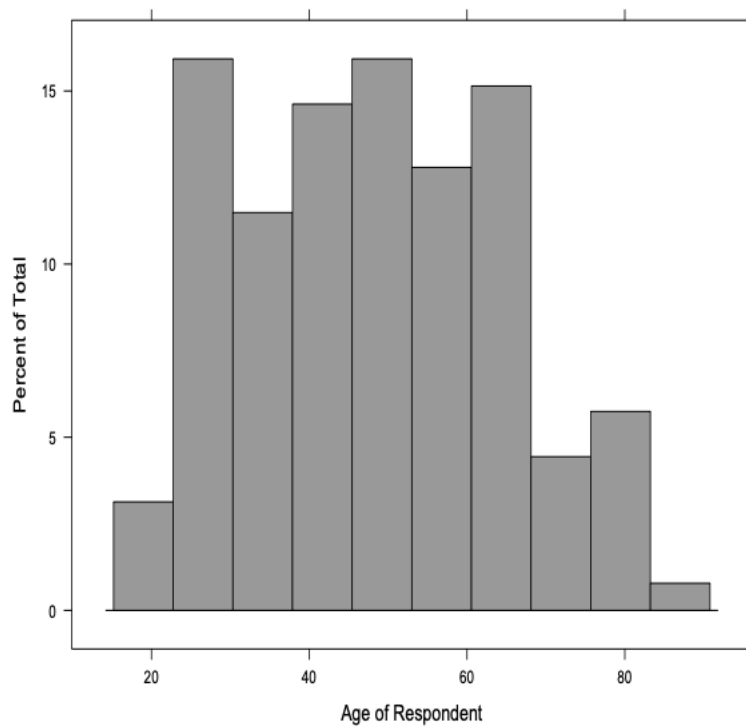
Histogram (with a larger “bin” width)



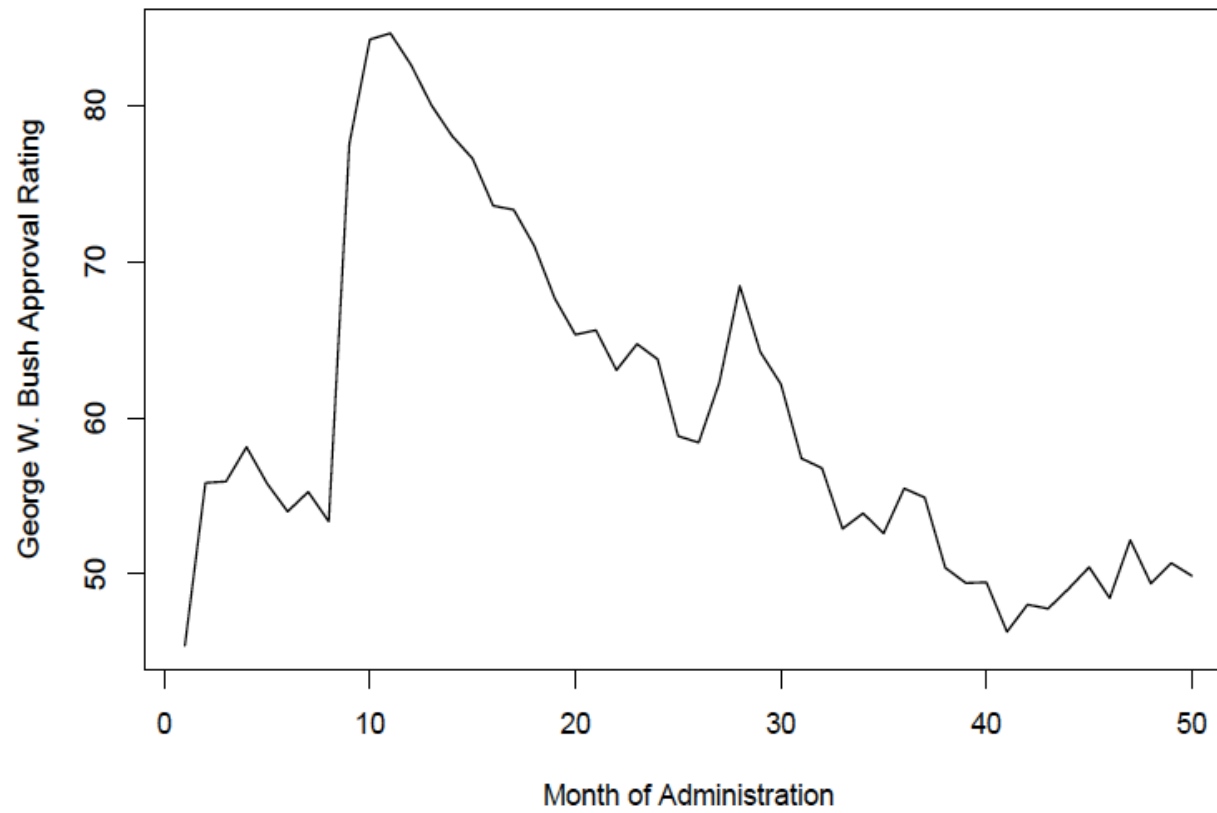
Density Curve



Histogram and Density Curve



Line/Time Plots

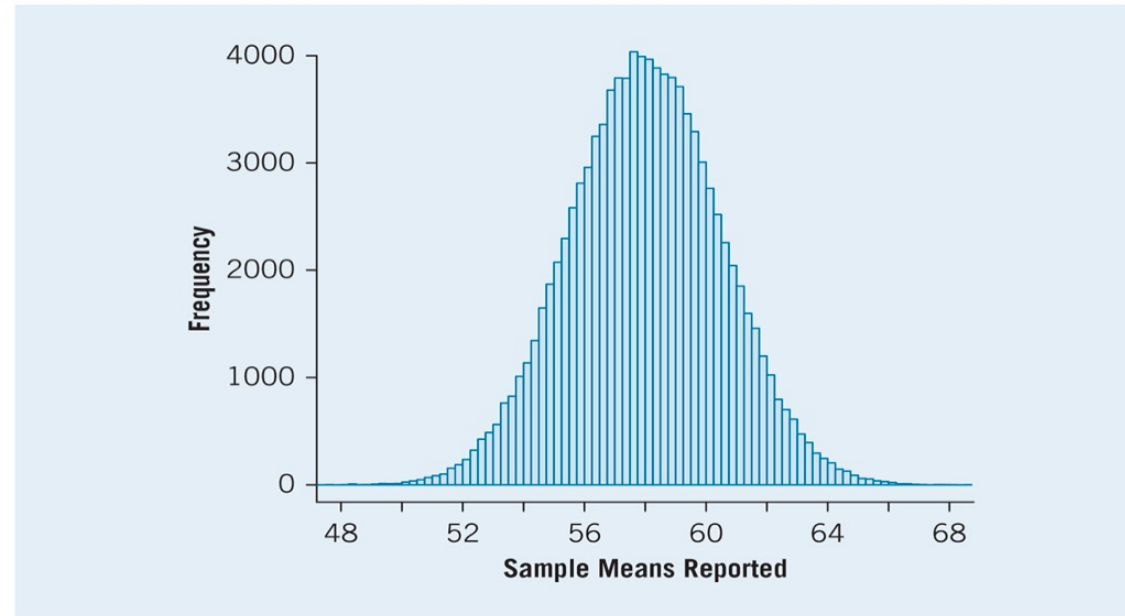


Shapes of Distribution

- Is the distribution symmetric or skewed?
- How many modes does the distribution have?
 - Unimodal
 - Bimodal
- Are there outliers or deviations from the overall shape?

Normal Distribution

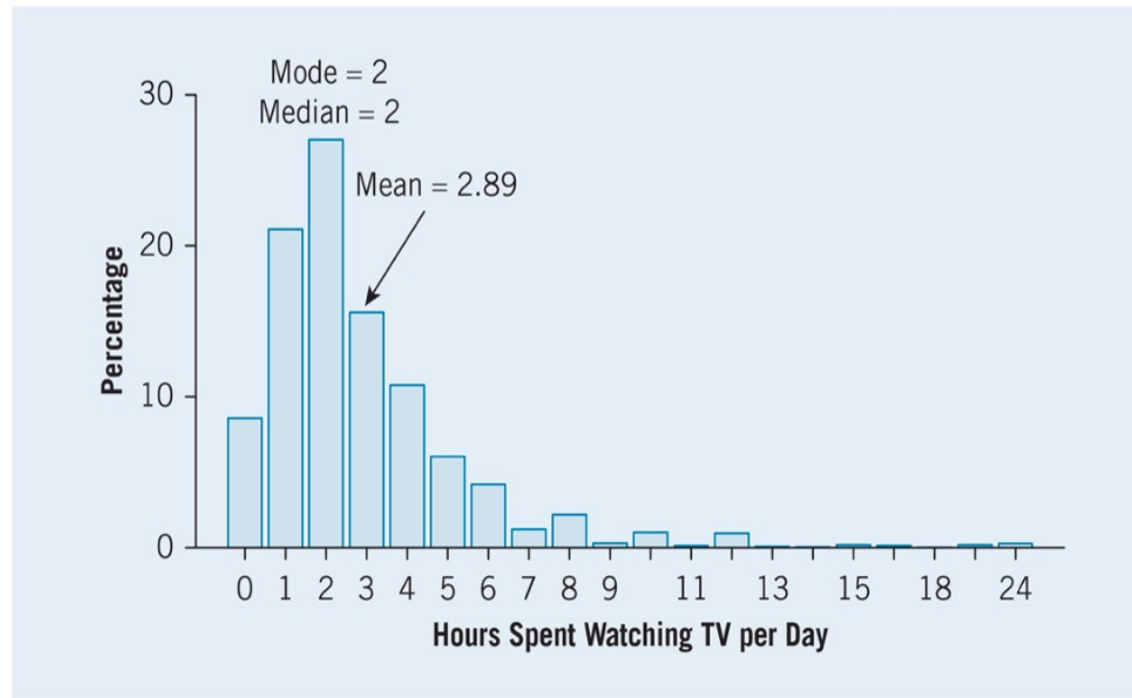
Figure 6-2 Distribution of Means from 100,000 Random Samples



Note: The figure shows means from 100,000 samples of $n = 100$. Population parameters: $\mu = 58$ and $\sigma = 24.8$.

Skewed Distribution

Figure 2-5 Bar Chart of Hours Spent Watching Television Per Day

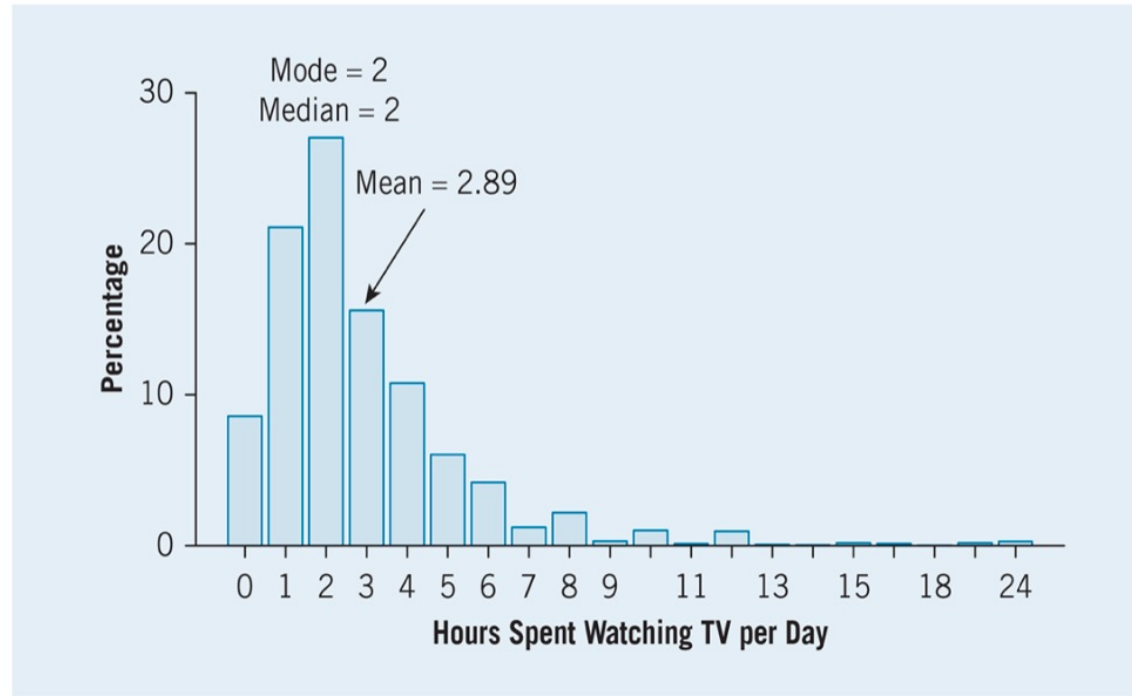


Measures of Central Tendency

- Mean: Average of all values
 - $\bar{x} = \frac{\sum x_i}{n}$
- Median: midpoint of the observations
- Mode: value that occurs most frequently (often used for categorical data)
- Outliers: An observation that falls well above or below the overall data
 - The mean is *sensitive* to outliers
 - The median is *resistant* to outliers

Graphical Display of Central Tendency

Figure 2-5 Bar Chart of Hours Spent Watching Television Per Day



Quartiles

- Splits the data into four parts
- The median is the second quartile, Q_2
- The first quartile, Q_1 , is the median of the lower half of the observations
- The third quartile, Q_3 , is the median of the upper half of the observations
- The interquartile range is the distance between the third quartile and first quartile
 - $IRQ = Q_3 - Q_1$

Five number summary

- Minimum value
- Q_1
- Median
- Q_3
- Maximum value

Five number summary

- ▶ Below is a random sample of 40 tree diameters, in centimeters, from the Wade Tract in Thomas County, Georgia. What is the five-number summary for these data?

#	D.	#	D.	#	D.	#	D.
1	2.2	11	11.4	21	29.1	31	43.3
2	2.2	12	11.4	22	31.5	32	43.6
3	2.3	13	13.3	23	31.8	33	44.2
4	2.7	14	16.9	24	32.6	34	44.4
5	4.3	15	17.6	25	35.7	35	44.6
6	4.9	16	18.3	26	37.5	36	47.2
7	5.4	17	22.3	27	38.1	37	51.5
8	7.8	18	26	28	39.7	38	51.8
9	9.2	19	26.1	29	40.3	39	52.2
10	10.5	20	27.9	30	40.5	40	69.3

Five number summary

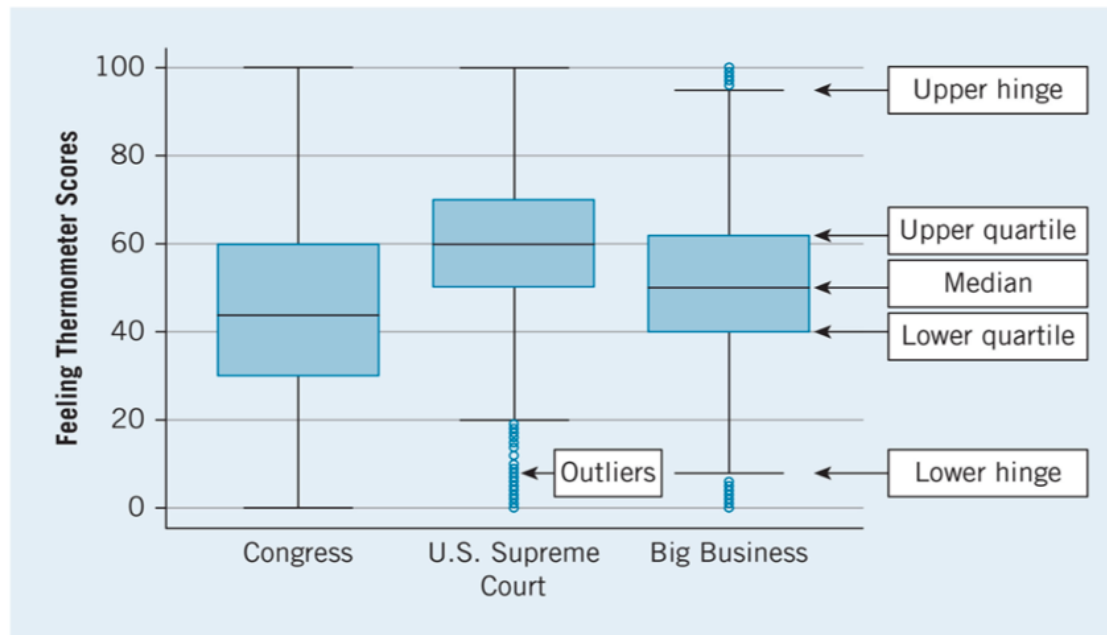
- Min=2.2cm
- $Q_1=10.95\text{cm}$
- M=28.5cm
- $Q_3=41.9\text{cm}$
- Max=69.3cm

Boxplot

- Construct a box from Q_1 to Q_3 (i.e., the IQR)
- Draw a line inside the box at the median value
- Draw a line out to the lowest value that is not an outlier
- Draw a line out to the highest value that is not an outlier
- Rule of thumb for outliers—an observation is a suspected outlier if it is more than 1.5 times the IQR above the third quartile (Q_3), or below the first quartile (Q_1).

Boxplot

Figure 2-7 Box Plots of Three Feeling Thermometer Variables



Source: 2016 American National Election Study.

Measures of Spread

- Range: difference between the largest and smallest observations
 - Range = Max - Min
- Variance: average of squares of deviations of the observations from their mean
 - $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$
- Standard deviation: square root of the average squared deviation from the mean
 - $s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$

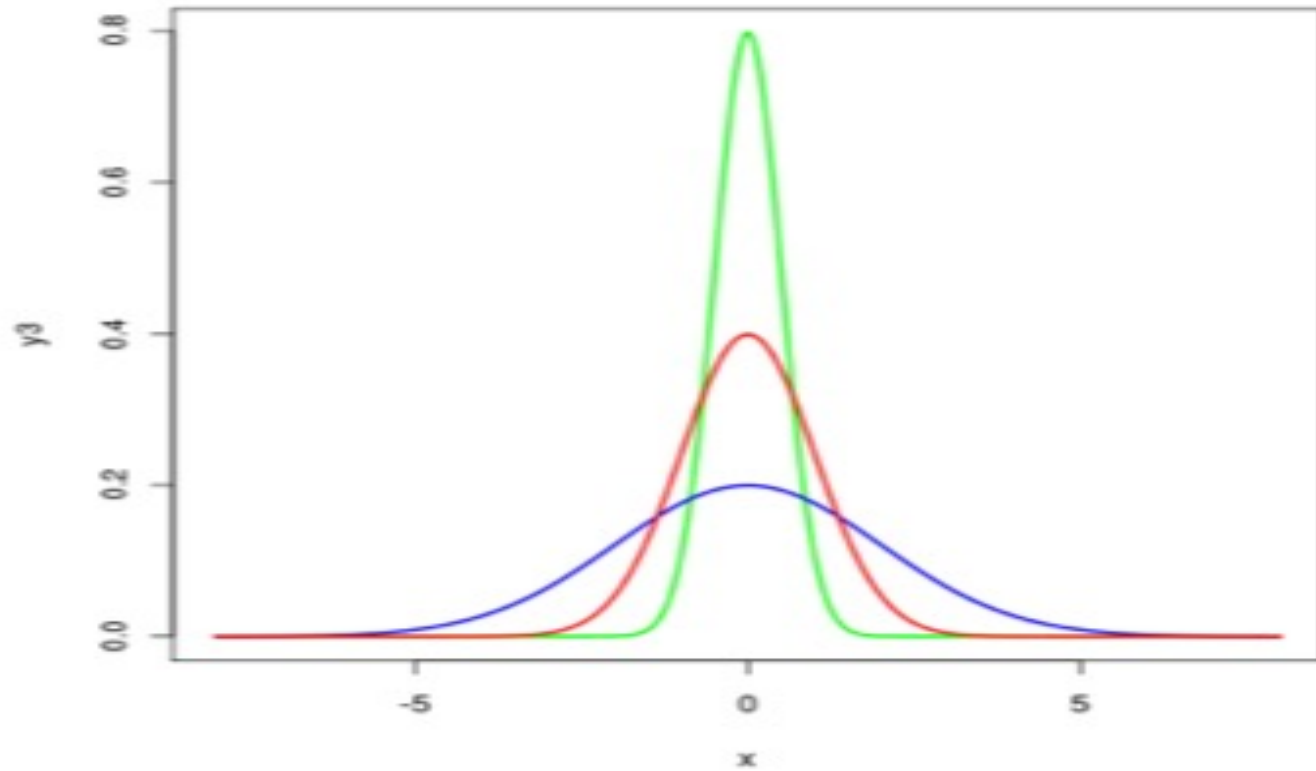
Standard Deviation

- S.D. should only be used with mean, not median, to describe spread
- If $SD=0$, there is no spread (i.e., all observations are the mean)
- SD not resistant to outliers
 - More sensitive than mean, because squared deviations used

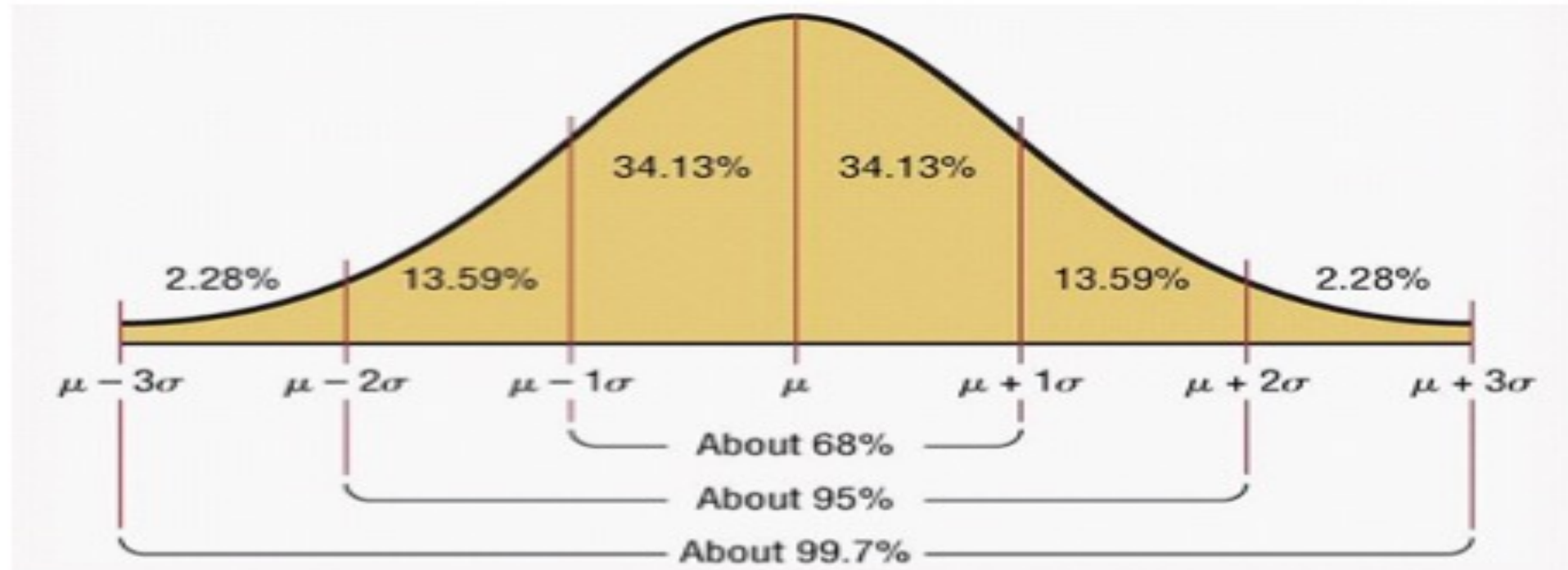
Linear Transformation

- Linear transformation applies the same linear equation to each observation of x
 - $x_{new} = a + b(x_{orig})$
 - Example—Convert Fahrenheit to Celsius
 - $C = (\frac{5}{9})(F - 32)$
- Linear transformations do not affect the shape of the distribution
- They do affect the measures of center and spread, but in a predictable way
 - Center—add “a” and multiply center by “b”
 - Spread—multiply spread by “b”, do not add “a”

Normal Distribution



Normal Distribution



Empirical Rule: for bell-shaped sets of data

- Approximately 68% of cases fall within 1 standard deviation of the mean
- Approximately 95% of cases fall within 2 standard deviations of the mean
- Approximately 99.7% of cases fall within 3 standard deviations of the mean

Normal distribution

- The length of human pregnancies from conception to birth varies according to a distribution that is approximately Normal with mean 266 days and a variance of 256 days. Use the empirical rule to answer the following questions.
 - Between what values do the lengths of the middle 95% of all pregnancies fall?
 - How short are the shortest 2.5% of all pregnancies?
 - How long do the longest 2.5% last?

Normal distribution

- The length of human pregnancies from conception to birth varies according to a distribution that is approximately Normal with mean 266 days and a variance of 256 days. Use the empirical rule to answer the following questions.
 - Between what values do the lengths of the middle 95% of all pregnancies fall?
 - The middle 95% fall within two standard deviations of the mean:
 - $266 \pm 2(16) = 236$ to 298 days
 - How short are the shortest 2.5% of all pregnancies?
 - How long do the longest 2.5% last?

Normal distribution

- The length of human pregnancies from conception to birth varies according to a distribution that is approximately Normal with mean 266 days and a variance of 256 days. Use the empirical rule to answer the following questions.
 - Between what values do the lengths of the middle 95% of all pregnancies fall?
 - The middle 95% fall within two standard deviations of the mean:
 - $266 \pm 2(16) = 236$ to 298 days
 - How short are the shortest 2.5% of all pregnancies?
 - The shortest 2.5% are shorter than 234 days
 - How long do the longest 2.5% last?

Normal distribution

- The length of human pregnancies from conception to birth varies according to a distribution that is approximately Normal with mean 266 days and a variance of 256 days. Use the empirical rule to answer the following questions.
 - Between what values do the lengths of the middle 95% of all pregnancies fall?
 - The middle 95% fall within two standard deviations of the mean:
 - $266 \pm 2(16) = 236$ to 298 days
 - How short are the shortest 2.5% of all pregnancies?
 - The shortest 2.5% are shorter than 234 days
 - How long do the longest 2.5% last?
 - The longest 2.5% are longer than 298 days.

Z-scores

- Because of the empirical rule, we can measure (in a standardized manner) how far away observations are from the mean along a normal distribution
- Z-score: how many s.d.s away from the mean the observation is
 - $Z_i = \frac{x_i - \mu_x}{\sigma_x}$
- Cumulative percentages
 - When we know z-score, can use a table to tell us percentage of cases that far from mean
 - Table A in back of book gives these

Conclusion

- Take home points:
 - Always plot your data first
 - Techniques for describing data differ based on type of variables
 - Important to consider outliers/skew when deciding which measures to use
 - Normal distributions have special properties, will be very important for inferential statistics
- Next week we will focus on techniques for examining relationships between variables