

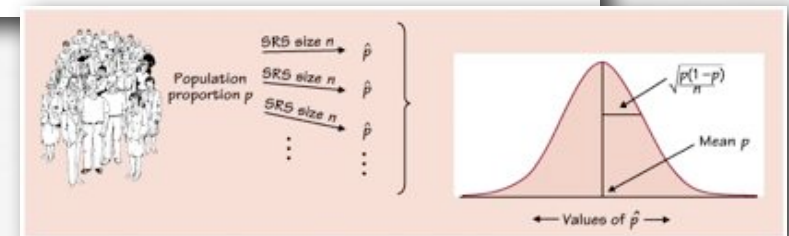
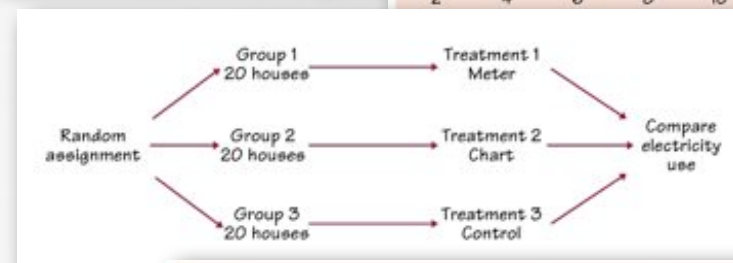
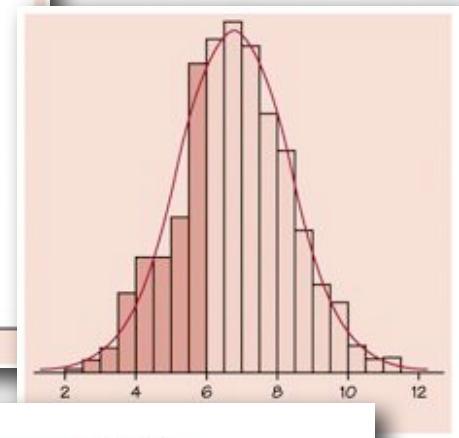
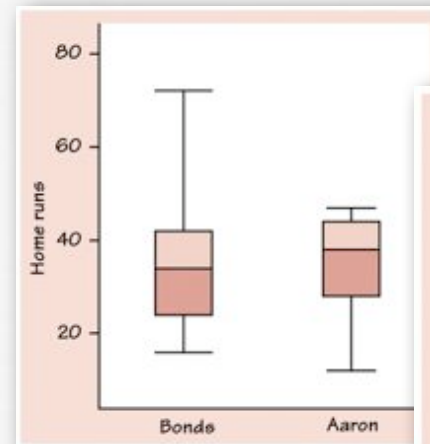
AP Statistics

Semester One
Review
Part 1
Chapters 1-5



AP Statistics Topics

- Describing Data
- Producing Data
- Probability
- Statistical Inference



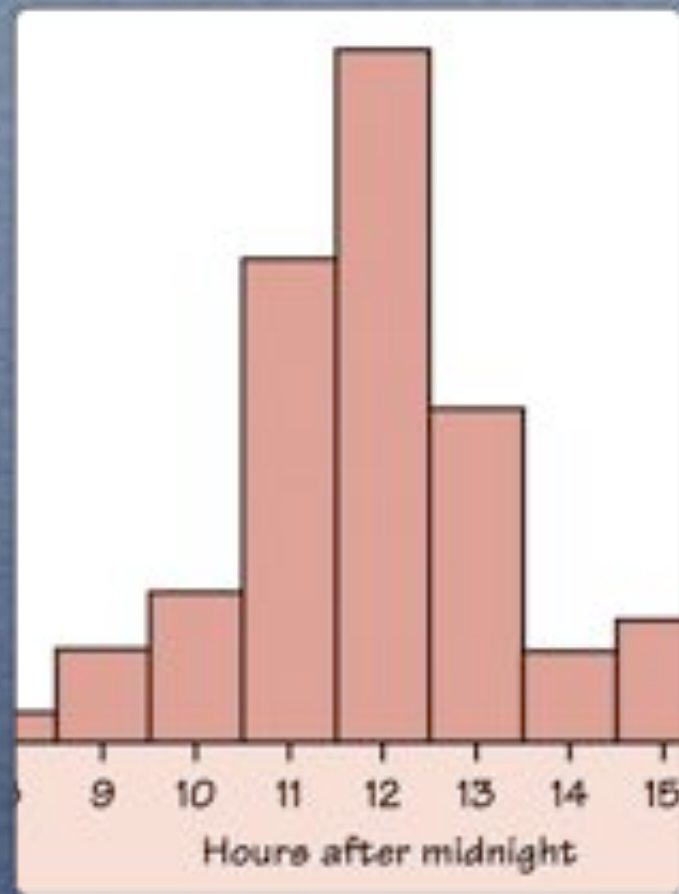
Describing Data

- ❑ Ch 1: Describing Data: Graphically and Numerically
- ❑ Ch 2: The Normal Distributions
- ❑ Ch 3: Describing BiVariate Relationships
- ❑ Ch 4: More BiVariate Relationships

Chapter 1: Describing Data

Our Introductory Chapter taught us how to describe a set of data graphically and numerically.

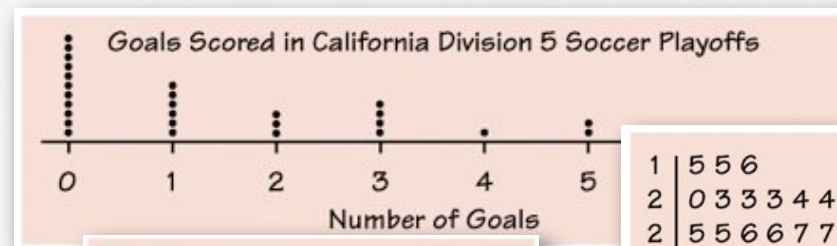
Our focus in this chapter was describing the Shape, Outliers, Center, and Spread of a dataset.



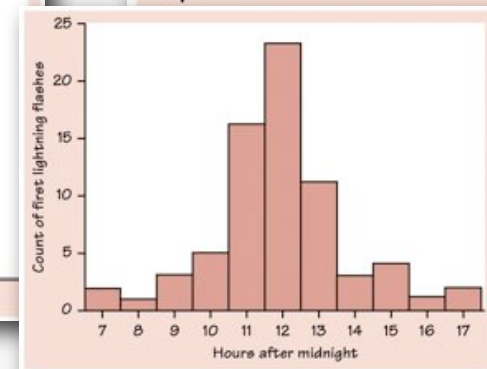
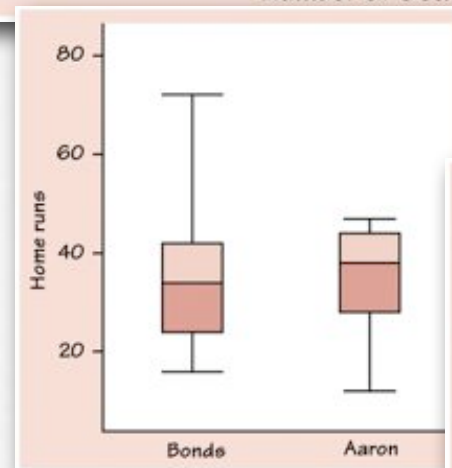
Describing Data

- When starting any data analysis, you should first PLOT your data and describe what you see...

- Dotplot
- Stemplot
- Box-n-Whisker Plot
- Histogram



1	5 5 6
2	0 3 3 3 4 4
2	5 5 6 6 7 7 7 8 8 8 8 8 9 9
3	1 1 3
3	5 5 5 6 7 7 7 8
4	3 3
4	7 7



Describe the SOCS

- After plotting the data, note the SOCS:
 - 📌 **Shape:** Skewed, Mound, Uniform, Bimodal
 - 📌 **Outliers:** Any “extreme” observations
 - 📌 **Center:** Typical “representative” value
 - 📌 **Spread:** Amount of variability

Numeric Descriptions

- While a plot provides a nice visual description of a dataset, we often want a more detailed numeric summary of the center and spread.

DataDesk

```
Summary of spending
No Selector
Percentile 25
Count 50
Mean 34.7022
Median 27.8550
StdDev 21.6974
Min 3.11000
Max 93.3400
Lower ith %tile 19.2700
Upper ith %tile 45.4000
```

Minitab

```
Descriptive Statistics
Variable N Mean Median TrMean StDev SEMean
spending 50 34.70 27.85 32.92 21.70 3.07
Variable Min Max Q1 Q3
spending 3.11 93.34 19.06 45.72
```

```
1-Var Stats
x̄=35.4375
Σx=567
Σx²=22881
Sx=13.63313977
σx=13.20023082
↓n=16
```

Measures of Center

□ When describing the “center” of a set of data, we can use the mean or the median.

□ **Mean:** “Average” value $\bar{x} = \frac{\sum x}{n}$

□ **Median:** “Center” value Q2

Measures of Variability

- When describing the “spread” of a set of data, we can use:
 - **Range:** Max-Min
 - **InterQuartile Range:** $IQR = Q3 - Q1$
 - **Standard Deviation:** $\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$

Numeric Descriptions

- When describing the center and spread of a set of data, be sure to provide a numeric description of each:
 - Mean and Standard Deviation
 - 5-Number Summary: *Min, Q1, Med, Q3, Max* {Box-n-Whisker Plot}

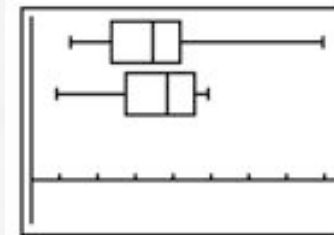
Determining Outliers

- When an observation appears to be an outlier, we will want to provide numeric evidence that it is or isn't "extreme"
- We will consider observations outliers if:
 - More than 3 standard deviations from the mean.

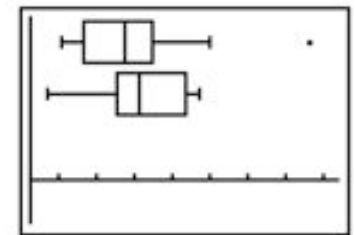
Or

- More than 1.5 IQR's outside the "box"

1	69
2	455
3	3344
3	77
4	02
4	69
5	
5	
6	
6	
7	3



(a)



(b)

Chapter 1 Summary

Plot your data
Dotplot, Stemplot, Histogram



Interpret what you see
Shape, Center, Spread, Outliers

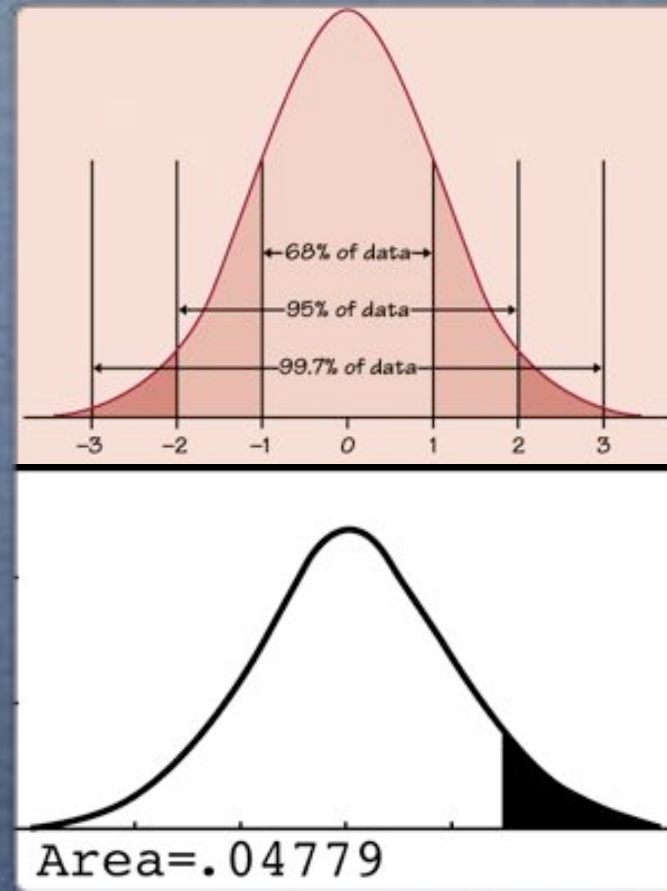


Choose numerical summary
 \bar{x} and s , Five-Number Summary

Chapter 2: Normal Distributions

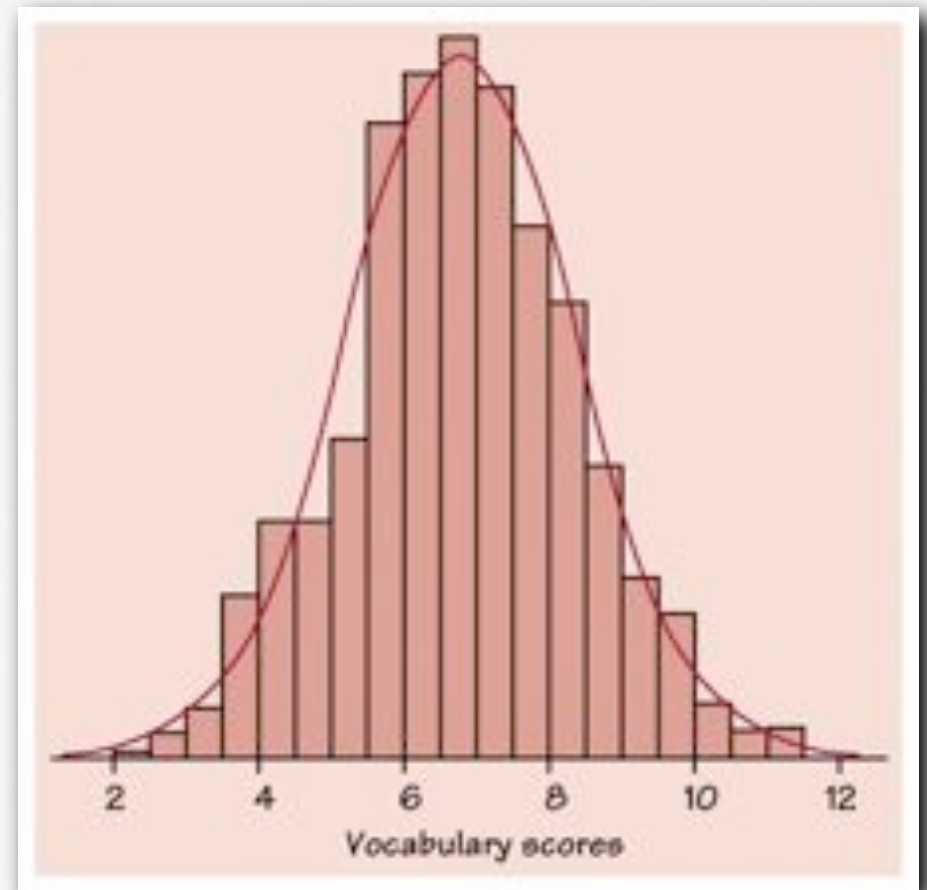
Many distributions in statistics can be described as approximately Normal.

In this chapter, we learned how to identify and describe normal distributions and how to do Standard Normal Calculations.



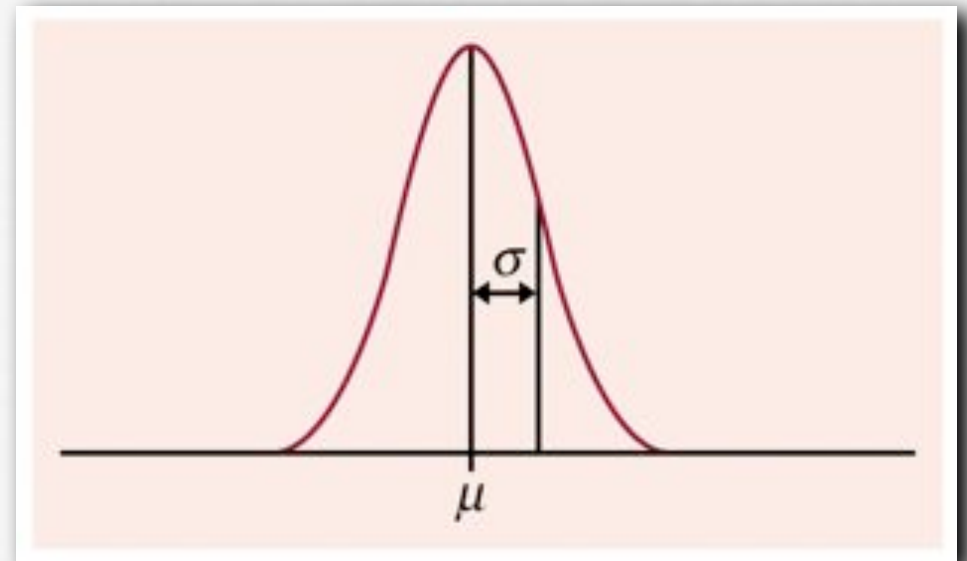
Density Curves

- A Density Curve is a smooth, idealized mathematical model of a distribution.
- The area under every density curve is 1.



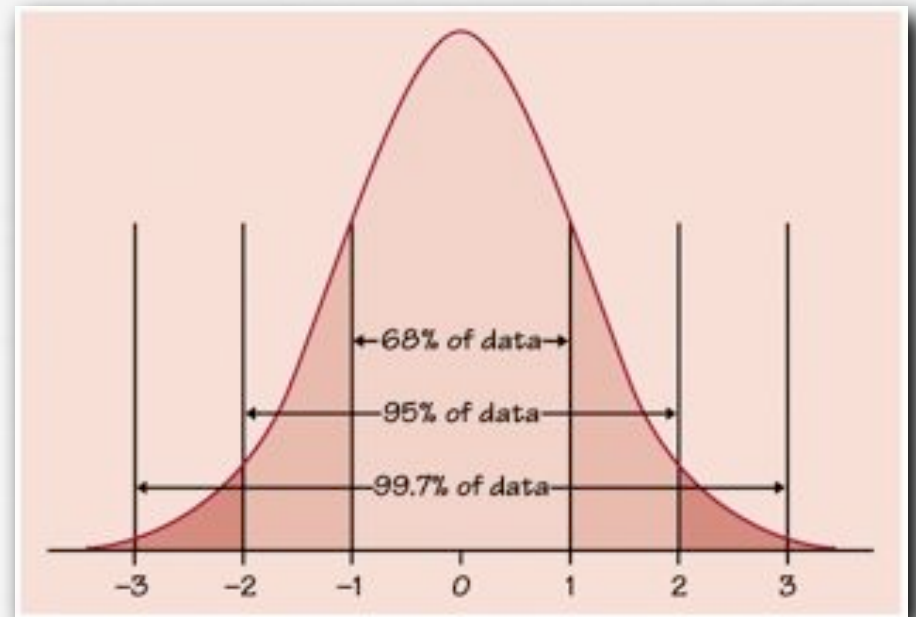
The Normal Distribution

- Many distributions of data and many statistical applications can be described by an approximately normal distribution.
- Symmetric, Bell-shaped Curve
- Centered at Mean μ
- Described as $N(\mu, \sigma)$



Empirical Rule

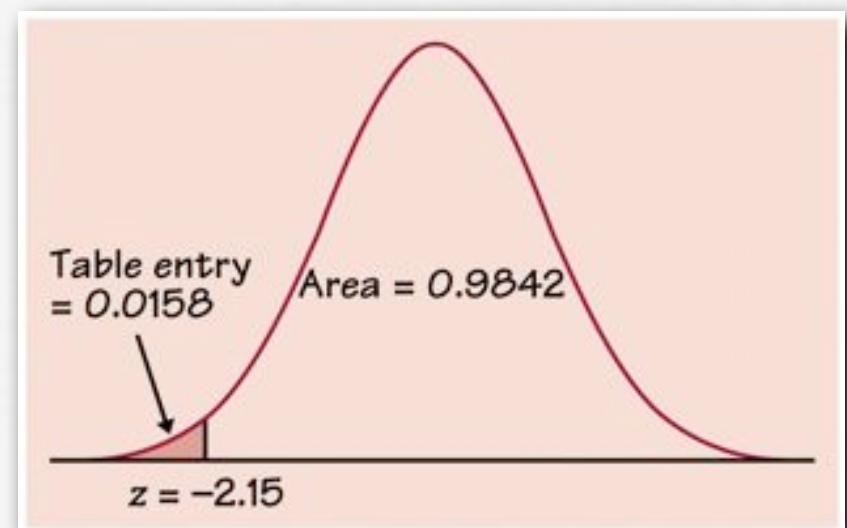
- One particularly useful fact about approximately Normal distributions is that
 - 68% of observations fall within one standard deviation of μ
 - 95% fall within 2 standard deviations of μ
 - 99.7% fall within 3 standard deviations of μ



Standard Normal Calculations

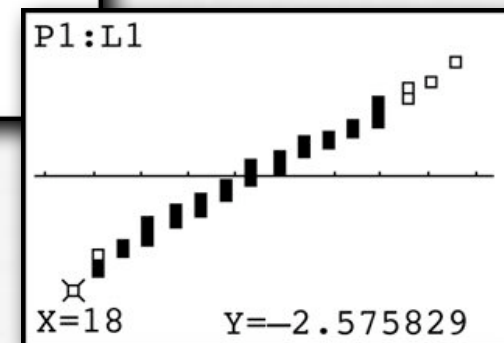
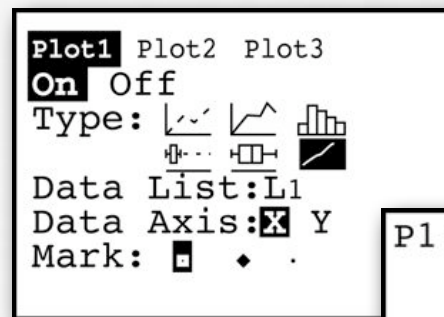
- The empirical rule is useful when an observation falls exactly 1,2,or 3 standard deviations from μ . When it doesn't, we must standardize the value {z-score} and use a table to calculate percentiles, etc.

$$z = \frac{x - \mu}{\sigma}$$



Assessing Normality

- ❑ To assess the normality of a set of data, we can't rely on the naked eye alone - not all mound shaped distributions are normal.
- ❑ Instead, we should make a *Normal Quantile Plot* and look for linearity.
- ❑ Linearity \rightarrow Normality

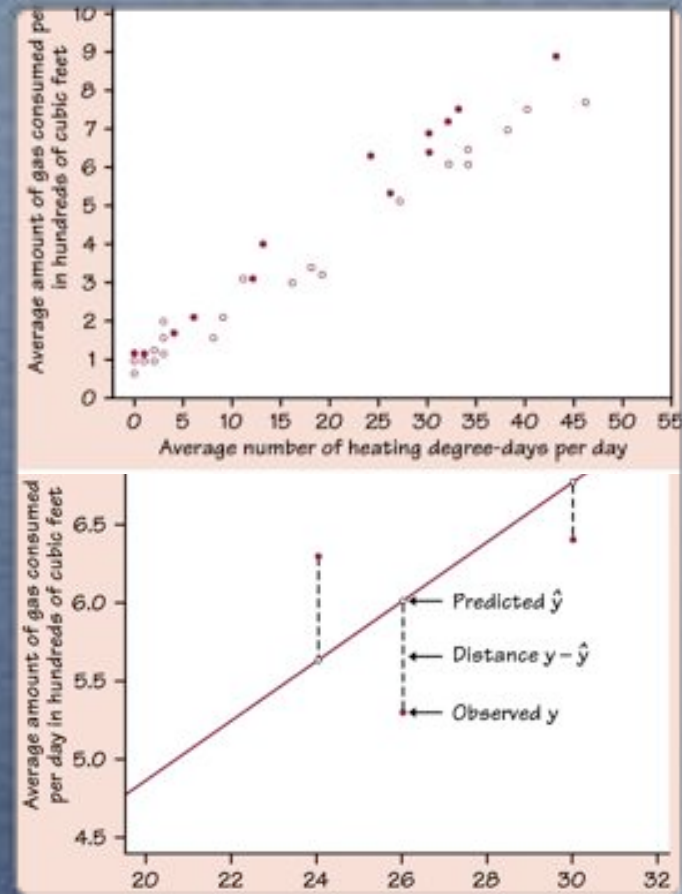


Chapter 3

Describing BiVariate Relationships

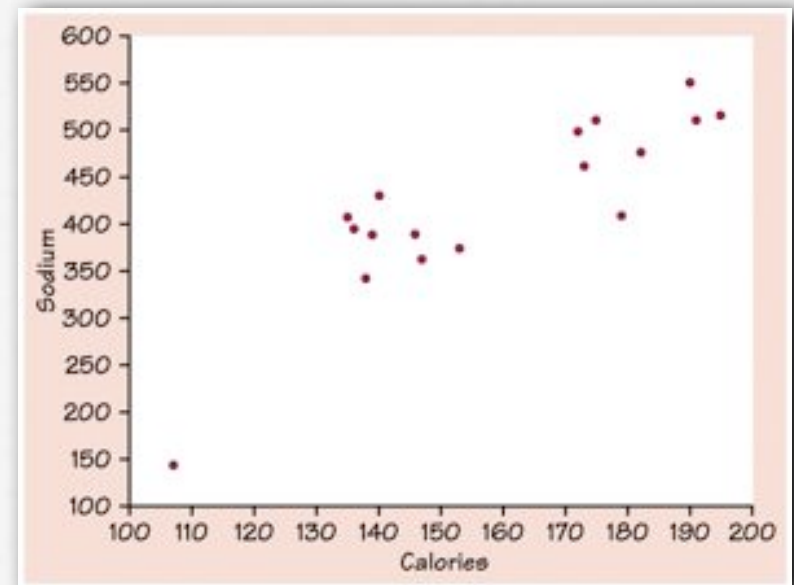
In this chapter, we learned how to describe bivariate relationships.

We focused on quantitative data and learned how to perform least squares regression.

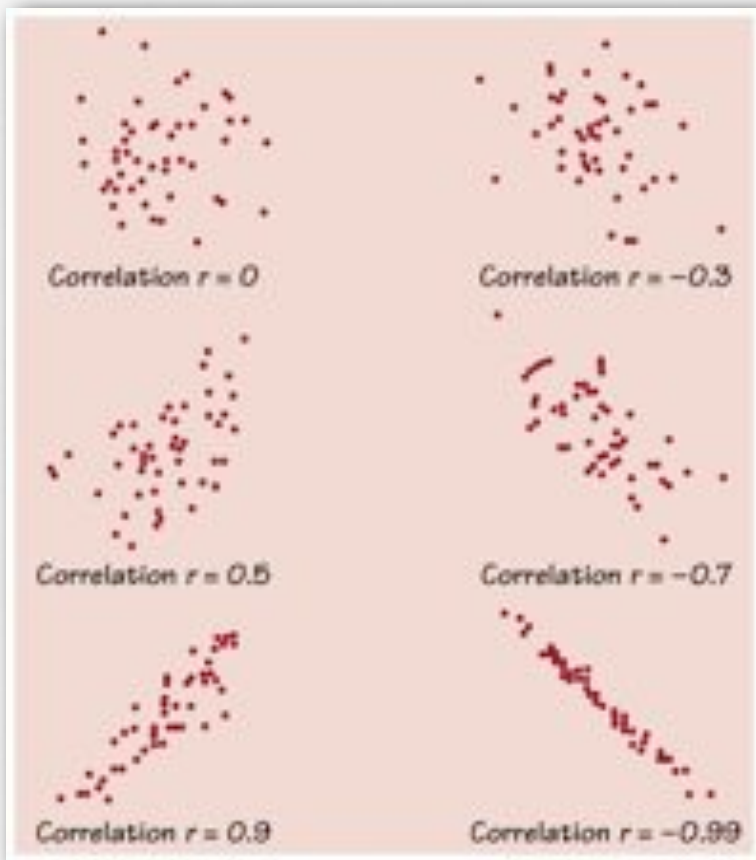


Bivariate Relationships

- Like describing univariate data, the first thing you should do with bivariate data is make a plot.
- Scatterplot
- Note Strength, Direction, Form



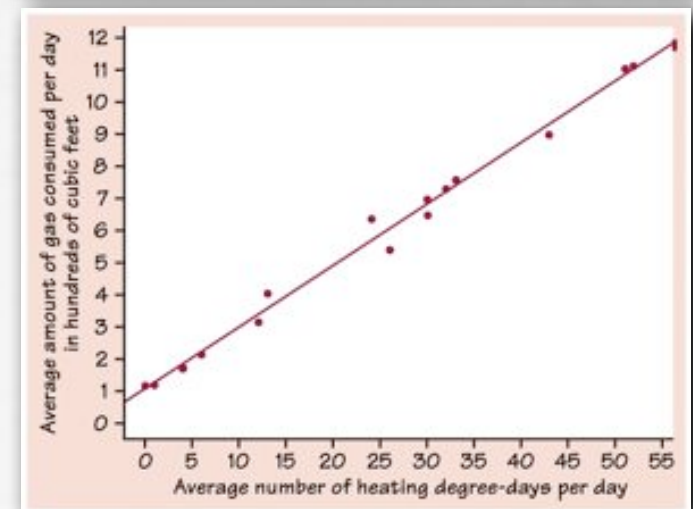
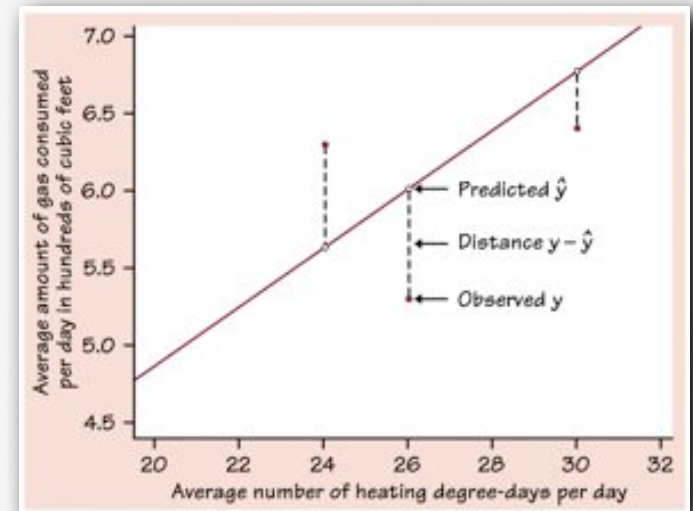
Correlation “r”



- We can describe the strength of a linear relationship with the Correlation Coefficient, r
- $-1 \leq r \leq 1$
- The closer r is to 1 or -1, the stronger the linear relationship between x and y .

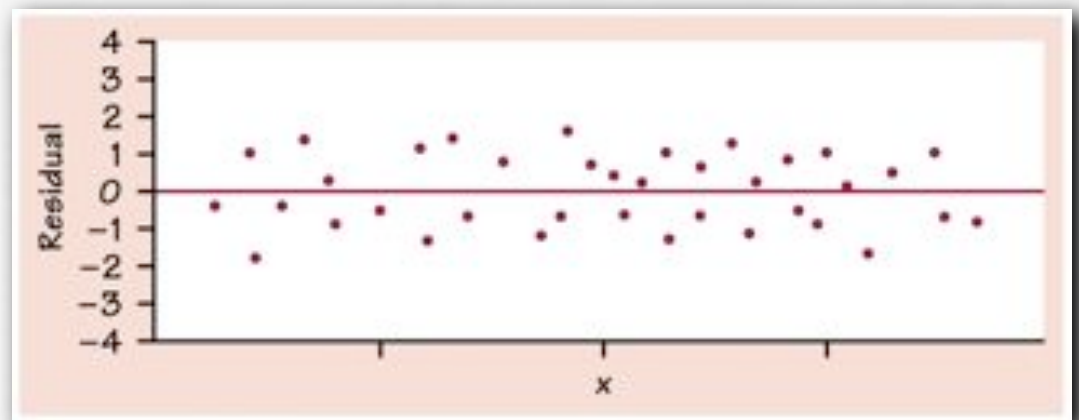
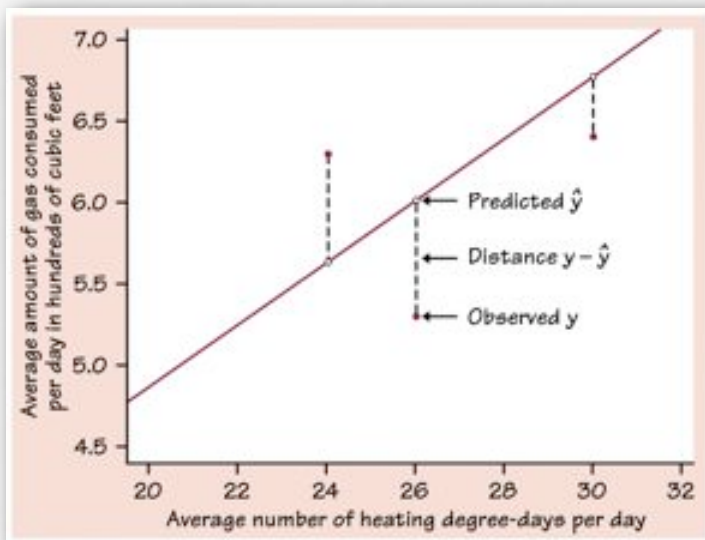
Least Squares Regression

- When we observe a linear relationship between x and y , we often want to describe it with a “line of best fit” $y=a+bx$.
- We can find this line by performing least-squares regression.
- We can use the resulting equation to predict y -values for given x -values.



Assessing the Fit

- If we hope to make useful predictions of y we must assess whether or not the LSRL is indeed the best fit. If not, we may need to find a different model.
- Residual Plot



Making Predictions

- If you are satisfied that the LSRL provides an appropriate model for predictions, you can use it to predict a \hat{y} for x 's within the observed range of x -values.

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

- Predictions for observed x -values can be assessed by noting the residual.

- Residual = observed y - predicted y

Chapter 3 Summary

Plot your data.
Scatterplot

Interpret what you see:
direction, form, strength, Linear?

Numerical summary?
 \bar{x} , \bar{y} , s_x , s_y , and r ?

Mathematical model?
Regression line?

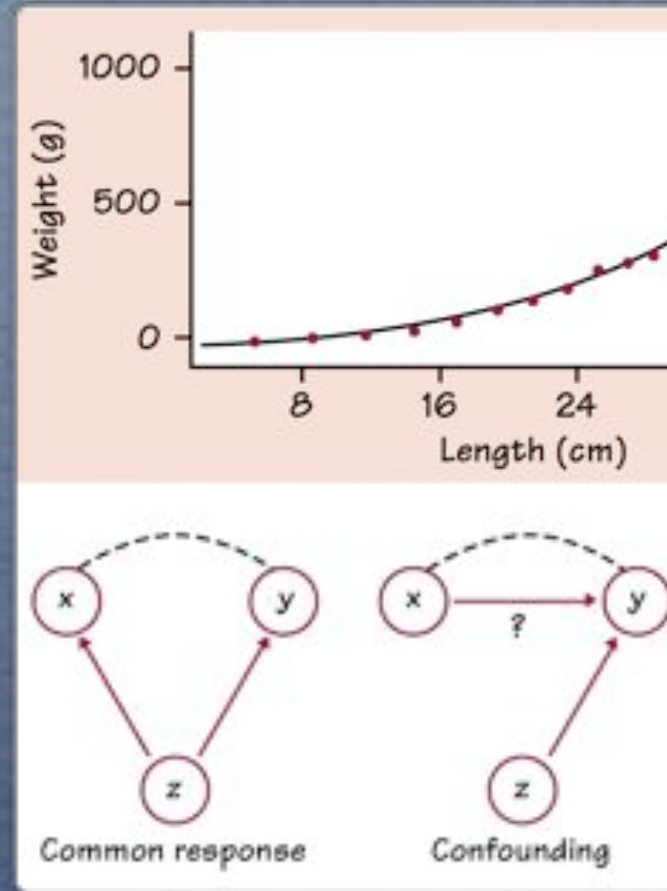


Chapter 4

More BiVariate Relationships

In this chapter, we learned how to find models that fit some nonlinear relationships.

We also explored how to describe categorical relationships.



NonLinear Relationships

□ If data is not best described by a LSRL, we may be able to find a Power or Exponential model that can be used for more accurate predictions.

□ Power Model: $\hat{y} = 10^a x^b$

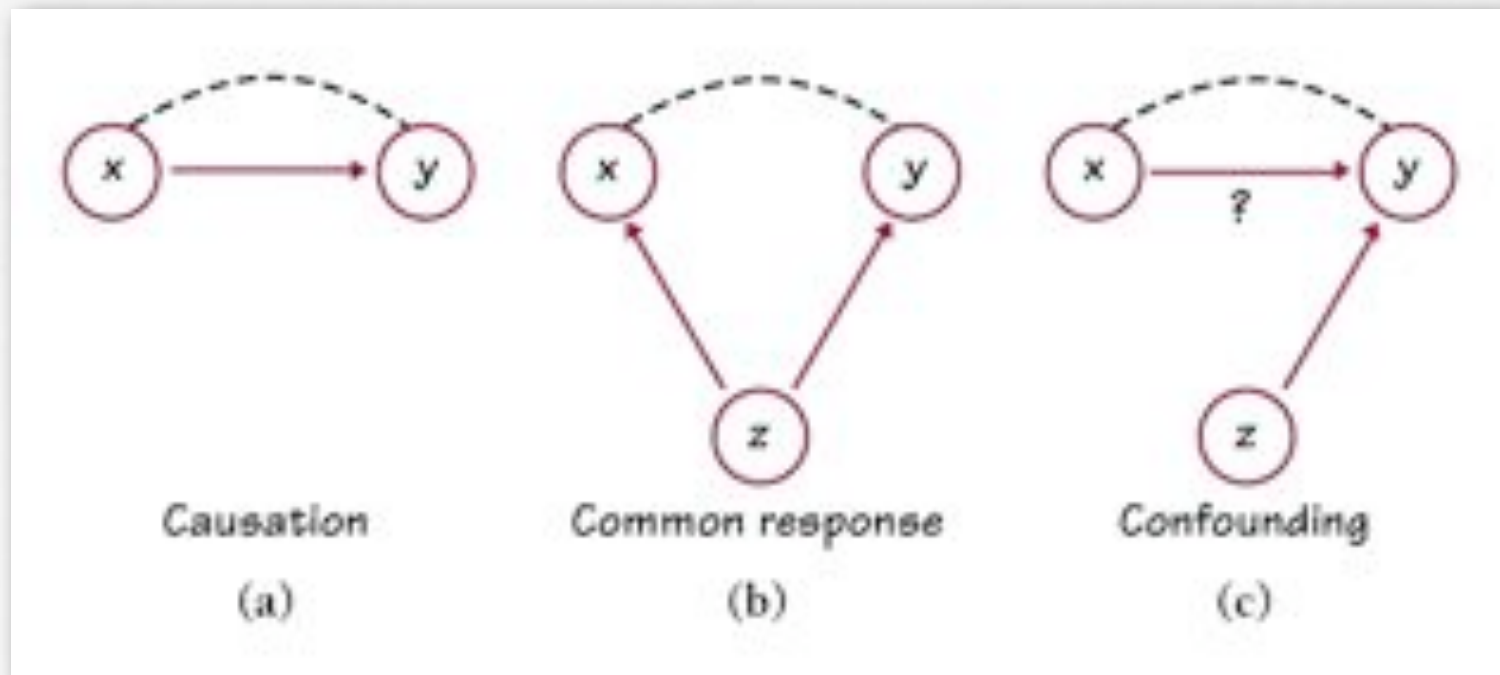
□ Exponential Model: $\hat{y} = 10^a 10^{bx}$

Transforming Data

- If (x,y) is non-linear, we can transform it to try to achieve a linear relationship.
 - If transformed data appears linear, we can find a LSRL and then transform back to the original terms of the data
- $(x, \log y)$ LSRL \rightarrow Exponential Model
- $(\log x, \log y)$ LSRL \rightarrow Power Model

The Question of Causation

- Just because we observe a strong relationship or strong correlation between x and y , we can not assume it is a causal relationship.

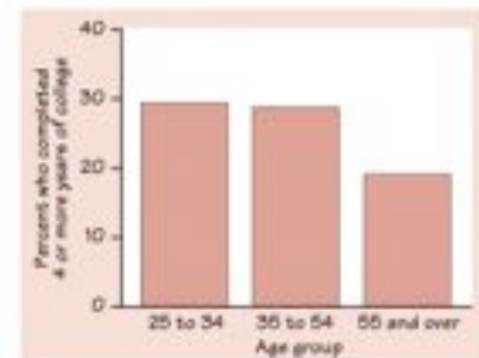
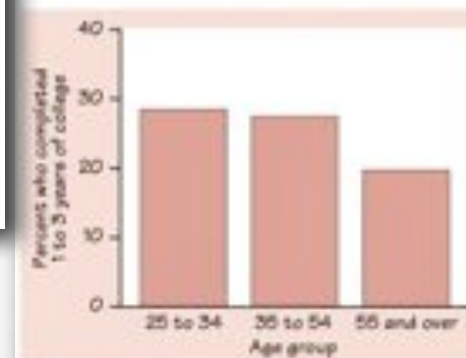
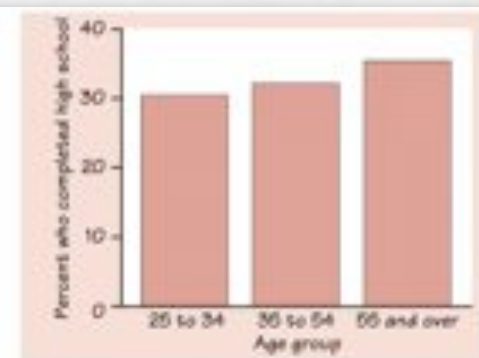
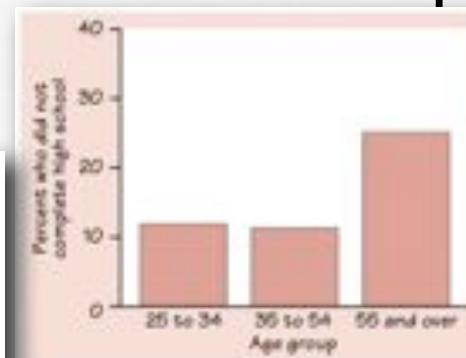


Relations in Categorical Data

- When categorical data is presented in a two-way table, we can explore the marginal and conditional distributions to describe the relationship between the variables.

TABLE 4.6 Years of school completed, by age, 2000 (thousands of persons)

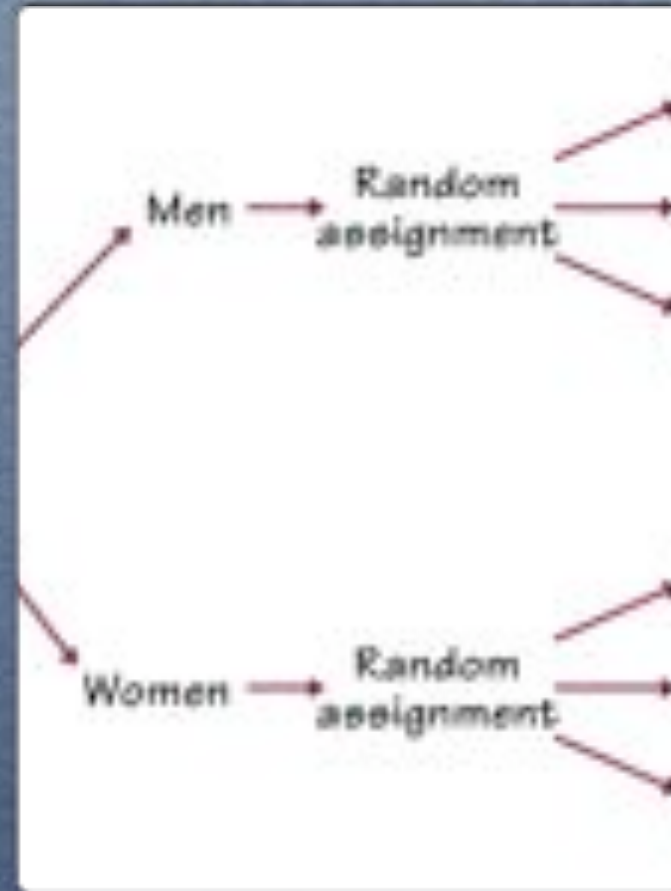
Education	Age group			Total
	25 to 34	35 to 54	55+	
Did not complete high school	4,474	9,155	14,224	27,853
Completed high school	11,546	26,481	20,060	58,087
1 to 3 years of college	10,700	22,618	11,127	44,445
4 or more years of college	11,066	23,183	10,596	44,845
Total	37,786	81,435	56,008	175,230



Chapter 5

Producing Data

In this chapter, we learned methods for collecting data through sampling and experimental design.



Sampling Design

- Our goal in statistics is often to answer a question about a population using information from a sample.
- Observational Study vs. Experiment
 - There are a number of ways to select a sample.
 - We must be sure the sample is representative of the population in question.

Sampling

- If you are performing an observational study, your sample can be obtained in a number of ways:
 - Convenience - Cluster
 - Systematic
 - Simple Random Sample
 - Stratified Random Sample

```
randInt(0,9,5)
      {5 6 5 7 1}
randInt(1,6,7)
      {5 6 5 5 3 4 1}
randInt(0,99,10)

{81 23 86 2 40...
```

Experimental Design

- ❑ In an experiment, we impose a treatment with the hopes of establishing a causal relationship.
- ❑ Experiments exhibit 3 Principles
 - ❑ Randomization
 - ❑ Control
 - ❑ Replication



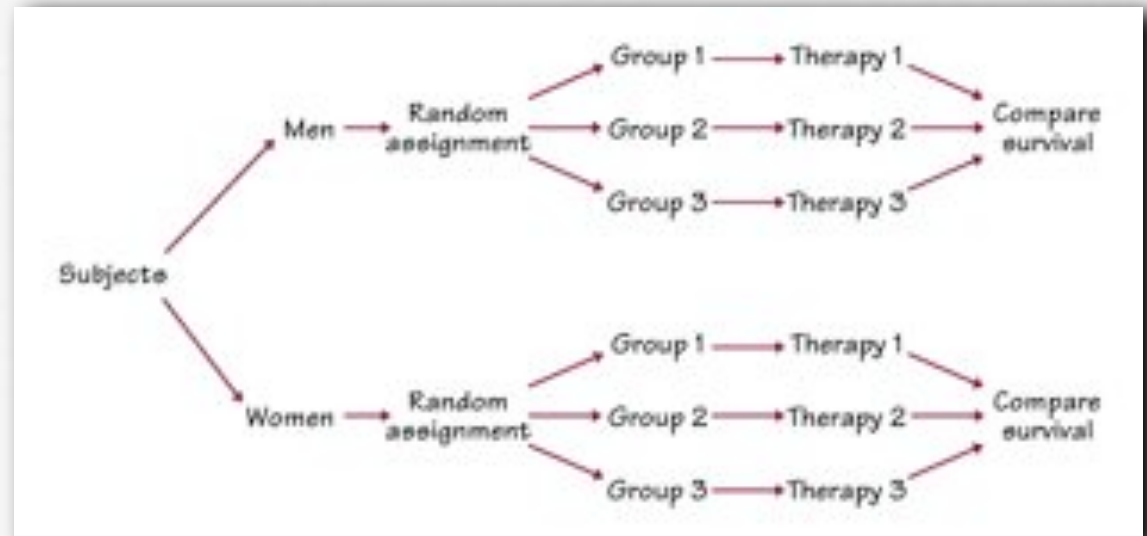
Experimental Designs

□ Like Observational Studies, Experiments can take a number of different forms:

□ Completely Controlled Randomized Comparative Experiment

□ Blocked

□ Matched Pairs





Chapters 6-9 Tomorrow