

[illegible]

Statistics Notebook

by Tommy Yang

Updated: September 11, 2016

A word cloud featuring various Java-related terms. The largest word is 'CON' in purple. Other prominent words include 'Message' (green), 'JSP' (red), 'Servlet' (blue), 'Endpoint' (purple), 'Protocol' (green), 'Authenticator' (blue), 'File' (blue), 'Ast' (brown), 'XAccessor' (brown), 'Object ID' (blue), 'Coyote' (blue), 'Char' (brown), 'URL' (brown), 'Group' (blue), 'Application' (green), 'Callback' (green), 'Annotation' (brown), 'Helper' (purple), 'Sign' (green), 'Compiler' (green), 'Output' (green), 'Data' (green), 'Descriptor' (green), 'Service' (purple), 'Naming' (blue), 'Entry' (brown), 'Env' (blue), and 'F' (red). The words are arranged in a dense, overlapping manner with varying orientations.

Contents

I	Chapter 1: Exploring Data	
1	Exploring Data	7
1.1	Analyzing Categorical Data	7
1.2	Displaying Quantitative Data with Graphs	9
1.3	Describing Quantitative Data with Numbers	11
II	Chapter 2: Modeling Distributions of Data	
2	Modeling Distributions of Data	15
2.1	Describing Location in a Distribution	15
2.2	Density Curves and Normal Distributions	16
III	Chapter 3: Describing Relationships	
3	Describing Relationships	21
3.1	Scatterplots and Correlation	21
3.2	Least-Squares Regression	23
IV	Chapter 4: Designing Studies	
4	Designing Studies	29
4.1	Sampling and Surveys	29
4.2	Experiments	30
4.3	Using Studies Wisely	32
V	Chapter 5: Probability	
5	Probability	35
5.1	Randomness, Probability, and Simulation	35
5.2	Probability Rules	35
5.3	Conditional Probability and Independence	36

6	Random Variables	41
6.1	Discrete and Continuous Random Variables	41
6.2	Transforming and Combining Random Variables	43
6.3	Binomial and Geometric Random Variables	44
	Index	47



Chapter 1: Exploring Data

1	Exploring Data	7
1.1	Analyzing Categorical Data	
1.2	Displaying Quantitative Data with Graphs	
1.3	Describing Quantitative Data with Numbers	

1. Exploring Data

1.1 Analyzing Categorical Data

Definition 1.1 – Individuals and variables.

- **Individuals** are the objects described by a set of data. Individuals may be people, animals, or things.
- A **variable** is any characteristic of an individual. A variable can take different values for different individuals.

Definition 1.2 – Categorical variable and quantitative variable.

- A **categorical variable** places an individual into one of several groups or categories.
- A **quantitative variable** takes numerical values for which it makes sense to find an average.

Definition 1.3 – Distribution.

- The **distribution** of a variable tells us what values the variable takes and how often it takes these values.

Definition 1.4 – Distribution of a categorical variable.

- The distribution of a categorical variable lists the categories and gives the count (**frequency**) or percent (**relative frequency**) of individuals that fall within each category.
- **Pie charts** and **bar graphs** display the distribution of a categorical variable.
(All bars in a bar graph should have the same width; a change in area could be **misleading**)
- A **two-way table** of counts organizes data about two categorical variables measured for the same set of individuals.

- The **marginal distribution** of one of the categorical variables in a two-way table of counts is the distribution of values of that variable among all individuals described by the table.
- A **conditional distribution** of a variable describes the values of that variable among individuals who have a specific value of another variable.
- An **association** is a relationship between two variables if knowing the value of one variable helps predict the value of the other.

1.2 Displaying Quantitative Data with Graphs

Definition 1.1 – Describing Shape.

$$\begin{array}{l}
 \text{Describing Shape} \left\{ \begin{array}{l}
 \text{Type} \left\{ \begin{array}{l} \text{Unimodal} \\ \text{Bimodal} \\ \text{Multimodal} \end{array} \right. \\
 \text{Shape} \left\{ \begin{array}{l} \text{Symmetric} \\ \text{Skewed to the right} \\ \text{Skewed to the left} \end{array} \right. \\
 \text{Center} \Rightarrow \text{Midpoint}(\text{median}) \\
 \text{Spread} \Rightarrow \text{The range from Maximum to Minimum} \\
 \text{Outliers} \Rightarrow \left\{ \begin{array}{l} > Q_3 + (1.5 \times IQR) \\ < Q_1 - (1.5 \times IQR) \end{array} \right.
 \end{array} \right. \quad (1.1a)
 \end{array}$$

Definition 1.2 – Histogram.

- **Histogram** is an estimate of the probability distribution of a continuous variable (quantitative variable).

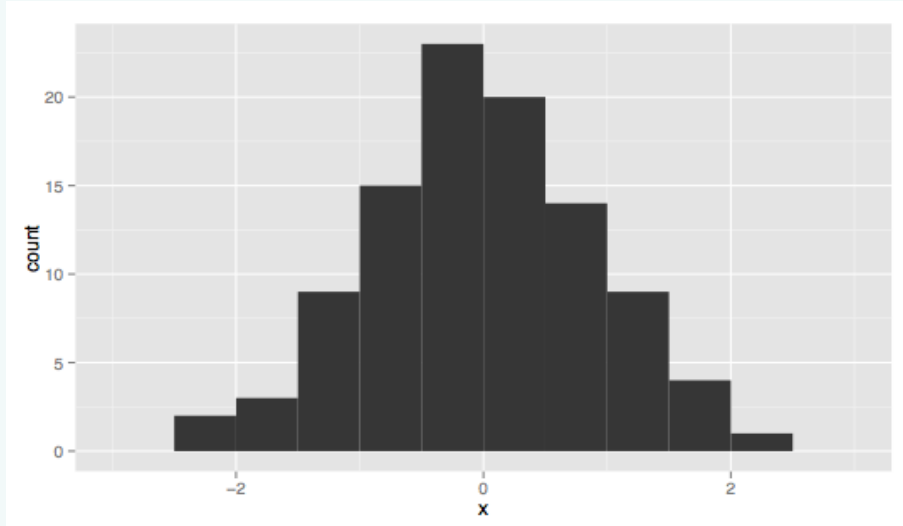


Figure 1.1: A symmetric, unimodal histogram

- **Remember:** histogram are for quantitative data; bar graphs are for categorical data. Also, be sure to use relative frequency histograms when comparing data sets of different sizes.

Definition 1.3 – Dotplot.

- A **Dotplot** is a representation of a distribution consists of group of data points plotted on a simple scale. Dotplots are used for continuous, quantitative, univariate data.

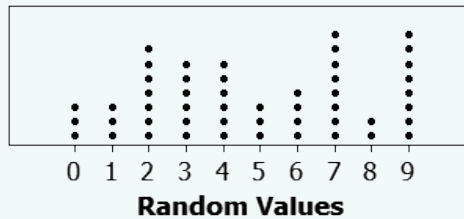
Dotplot of Random Values

Figure 1.2: A dotplot of 50 random values from 0 to 9.

Definition 1.4 – Stemplot.

- A **Stemplot** is a complicated device for presenting quantitative data in a graphical format, similar to a histogram, to assist in visualizing the shape of a distribution.

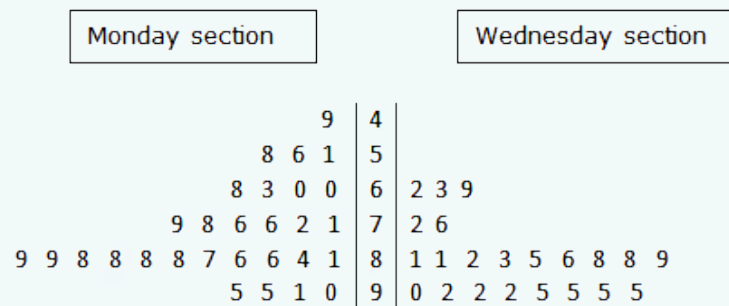


Figure 1.3: A back-to-back stemplot

1.3 Describing Quantitative Data with Numbers

Definition 1.5 – Mean and median.

- The **Mean** is the arithmetic average of a set of number.

Definition 1.3.1 – Mean.

$$\bar{x} = \frac{\text{sum of observations}}{n} = \frac{x_1 + x_2 + \cdots + x_n}{n} = \frac{\sum x_i}{n} \quad (1.2a)$$

- The **Median** is the midpoint of a distribution, the number such that half the observations are smaller and half are larger.
- Remember:** the median is a **resistant** measure of center because it is relatively unaffected by extreme observations. The mean is **nonresistant**. Among the measures of spread, the IQR is **resistant**, but the standard deviation and range are **nonresistant**.

Definition 1.6 – Measuring Spread.

- Interquartile range (IQR)

Definition 1.3.2 – Interquartile range (IQR).

$$IQR = Q_3 - Q_1 \quad (1.3a)$$

- The five-number summary

$$\text{The five-number summary} \left\{ \begin{array}{l} \text{Minimum} \\ Q_1 \Rightarrow \text{The first quartile} \\ \text{Median} \\ Q_3 \Rightarrow \text{The third quartile} \\ \text{Maximum} \end{array} \right. \quad (1.4a)$$

- Outliers

Definition 1.3.3 – Outlier.

$$\text{Outlier} \left\{ \begin{array}{l} > Q_3 + (1.5 \times IQR) \\ < Q_1 - (1.5 \times IQR) \end{array} \right. \quad (1.5a)$$

- **Standard Deviation** s_x measures the typical distance of the values in a distribution from the mean.

Definition 1.3.4 – Standard Deviation.

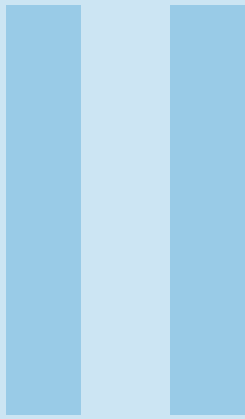
$$S_x = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} \quad (1.6a)$$

- **Variance** s_x^2 is the average squared deviation of a set of number.

Definition 1.3.5 – Variance.

$$S_x^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n-1} = \frac{1}{n-1} \sum (x_i - \bar{x})^2 \quad (1.7a)$$

- **Remember:** The mean and standard deviation are good descriptions for roughly **symmetric** distributions without outliers. The median and IQR are a better description for **skewed** distributions.



Chapter 2: Modeling Distributions of Data

2	Modeling Distributions of Data	15
2.1	Describing Location in a Distribution	
2.2	Density Curves and Normal Distributions	

2. Modeling Distributions of Data

2.1 Describing Location in a Distribution

Definition 2.1 – Measuring Positions.

- **Percentile:** the p th percentile of a distribution is the value with p percent of the observations less than it.
- **Standardized Score (Z-Score):** if x is an observation from a distribution that has known mean and standard deviation, the **standardized score** for x is

Definition 2.1.1 – Normal Distribution.

$$Z = \frac{x_i - \mu}{\sigma} \quad (2.1a)$$

- A **cumulative relative frequency graph** allows us to examine location within a distribution, beginning by grouping the observations into equal-width classes.
- For a common **transform data** like changing units of measurement:

$$\text{Add a constant } a \begin{cases} \text{Median, mean, quartiles, and percentiles} \Rightarrow \text{increase by } a. \\ \text{Spread} \Rightarrow \text{do not change.} \end{cases} \quad (2.2a)$$

$$\text{Multiply a constant } b \begin{cases} \text{Median, mean, quartiles, percentiles} \Rightarrow \text{multiply by } b. \\ \text{Spread} \Rightarrow \text{also multiply by } b. \end{cases} \quad (2.3a)$$

- **Neither** of these transformations changes the shape of the distribution. ■

2.2 Density Curves and Normal Distributions

Definition 2.1 – Density Curves.

- A **density curve** is that

$$\left\{ \begin{array}{l} \text{is always on or above the horizontal axis, and} \\ \text{has area exactly 1 underneath it.} \end{array} \right. \quad (2.4a)$$

Definition 2.2 – Normal Distribution.

- A **Normal Distribution** is described by a Normal density curve. The **mean** of a Normal distribution μ is at the center of the symmetric **Normal curve**. The **standard deviation** σ is the distance from the center to the change-of-curvature points on either side.
- We **abbreviate** the Normal distribution with mean μ and standard deviation σ as $N(\mu, \sigma)$.

Definition 2.2.1 – Normal Distribution.

$$f(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (2.5a)$$

- **The 68-95-99.7 rule:**
For a Normal distribution with mean μ and standard deviation σ :
68% of the observations fall within σ of the mean μ .
95% of the observations fall within 2σ of the mean μ .
99.7% of the observations fall within 3σ of the mean μ .

Definition 2.2.2 – The 68-95-99.7 rule.

$$\text{Three Sigma} \left\{ \begin{array}{ll} \int_{-1}^1 \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx = \mathbf{68.2\%} & \text{(within 1 SD)} \\ \int_{-2}^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx = \mathbf{95.4\%} & \text{(within 2 SD)} \\ \int_{-3}^3 \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx = \mathbf{99.7\%} & \text{(within 3 SD)} \end{array} \right. \quad (2.6a)$$

- The **standard Normal distribution** is the Normal distribution with mean 0 and standard deviation 1.

Definition 2.2.3 – Standard Normal Distribution.

$$Z = \frac{x_i - \mu}{\sigma} \Rightarrow \text{Standardization} = \frac{\text{Obs} - \text{Mean}}{\text{SD}} \quad (2.7a)$$

$$f(x|0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \quad (2.7b)$$

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx = 1 \quad (2.7c)$$



Chapter 3: Describing Relationships

3	Describing Relationships	21
3.1	Scatterplots and Correlation	
3.2	Least-Squares Regression	

3. Describing Relationships

3.1 Scatterplots and Correlation

Definition 3.1 – Scatterplots.

- A **scatterplot** displays the relationship between two quantitative variables measured on the same individuals.
- Explanatory and response variable:

$$\begin{cases} x \Rightarrow \text{Explanatory variable} \\ y \Rightarrow \text{Response variable} \end{cases} \quad (3.1a)$$

Definition 3.2 – Describing Scatterplots.

- Describing Scatterplots

$$\text{Scatterplots} \left\{ \begin{array}{l} \text{Direction (general point trend)} \left\{ \begin{array}{l} \text{up to right: positive association} \\ \text{down: negative association} \end{array} \right. \\ \text{Form (shape of direction)} \left\{ \begin{array}{l} \text{Linear} \\ \text{Curved} \end{array} \right. \\ \text{Strength (how closely graph fits form)} \left\{ \begin{array}{l} \text{Strong} \\ \text{Weak} \end{array} \right. \\ \text{Outlier (a departure falls outside the overall relationship)} \end{array} \right. \quad (3.2a)$$

Definition 3.3 – Measuring Linear Association.

- The **correlation** r measures the direction and strength of the linear relationship between two quantitative variables.

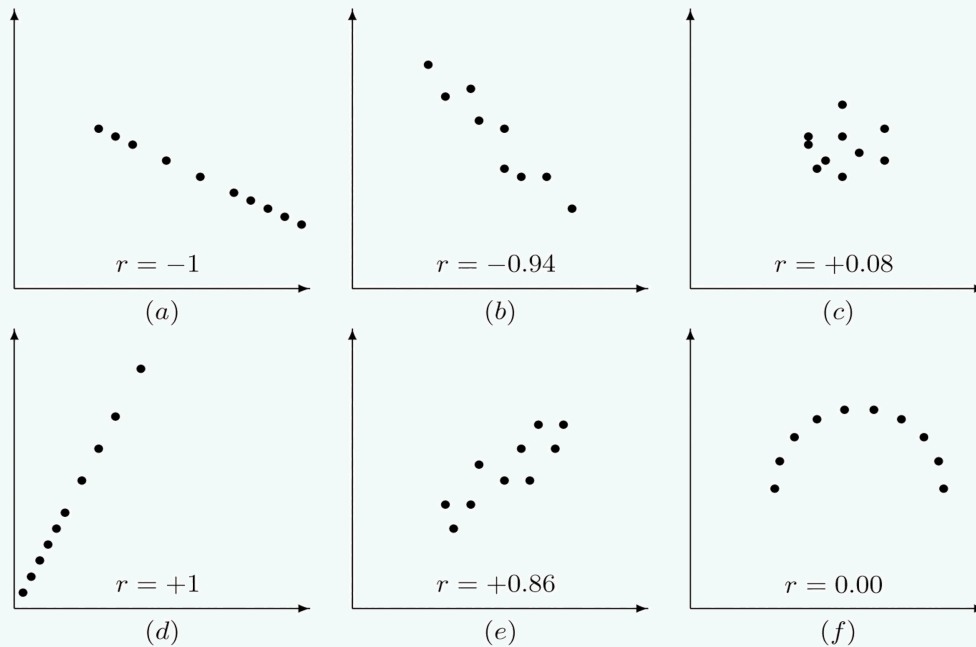


Figure 3.1: Plots with different correlations

Definition 3.1.1 – Correlation.

$$r = \frac{Z_{x_1}Z_{y_1} + Z_{x_2}Z_{y_2} + \cdots + Z_{x_n}Z_{y_n}}{n - 1} \quad (3.3a)$$

$$= \frac{1}{n - 1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{S_x} \right) \left(\frac{y_i - \bar{y}}{S_y} \right) \quad (3.3b)$$

$$= \frac{1}{n - 1} \sum Z_x Z_y \quad (Z = \frac{x_i - \bar{x}}{S_x}) \quad (3.3c)$$

- Features of Correlation:**

- $r = 1 \Rightarrow$ **linear**, perfect positive correlation.
- $r > 0 \Rightarrow$ **positive** association.
- $r < 0 \Rightarrow$ **negative** association.
- $r = -1 \Rightarrow$ **linear**, perfect negative correlation.
- $r = 0 \Rightarrow$ **doesn't** guarantee there's **no** relationships between two variables, just **No** linear relationship.

- **only** measures the strength of a **linear relationship**.
- can **Not** conclude that change in one variable cause in the other.
- is the **same** when you **inverse** x and y .
- is the **same** when you change the unit.
e.g: height in inches or meters, weight in pounds or kilograms.
- is **Not** resistant to outliers.
- is **Not** a complete summary of two-variable data.
- both variables must be **quantitative**.

3.2 Least-Squares Regression

Definition 3.1 – Regressions.

- A **Regression line** is a line that describes how a response variable y changes as an explanatory variable x changes.

Definition 3.2.1 – Regression line, predicted value, slope, y intercept.

A **regression line** relating y to x has an equation of the form:

$$\hat{y} = a + bx \quad (3.4a)$$

- \hat{y} is the **predicted value**.
- b is the **slop**.
- a is the **y intercept** (y value when $x = 0$).
- **Extrapolation** is the use of regression line for prediction far outside the interval of values of the explanatory variable x used to obtain the line.

Definition 3.2 – Residuals and the Least-Squares Regression Line.

- A **residual** is the difference between an observed value of the response variable and the value predicted by the regression line.

Definition 3.2.2 – Residual.

$$\begin{aligned} \text{residual} &= \text{observed } y - \text{predicted } y \\ &= y - \hat{y} \end{aligned} \quad (3.5a)$$

- The **least-squares regression line** of y on x is the line that makes the sum of the squared residuals as small as possible.

Definition 3.2.3 – Least-squares regression line.

The least-squares regression line is the line $\hat{y} = a + bx$ with **slope**

$$b = r \frac{s_y}{s_x} \quad (3.6a)$$

and **y intercept**

$$a = \bar{y} - b\bar{x} \quad (3.6b)$$

The least-squares regression line always passes through the point (\bar{x}, \bar{y})

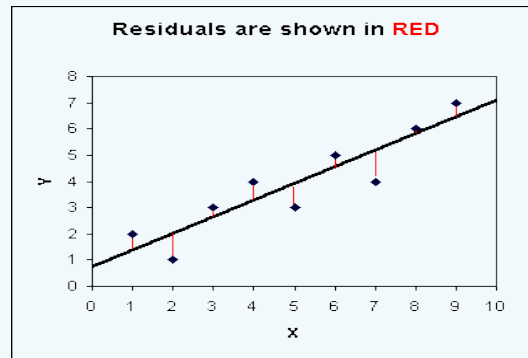


Figure 3.2: A Least-squares regression line

- A **residual plot** is a scatterplot of the residuals against the explanatory variable, helping us assess whether a linear model is appropriate.

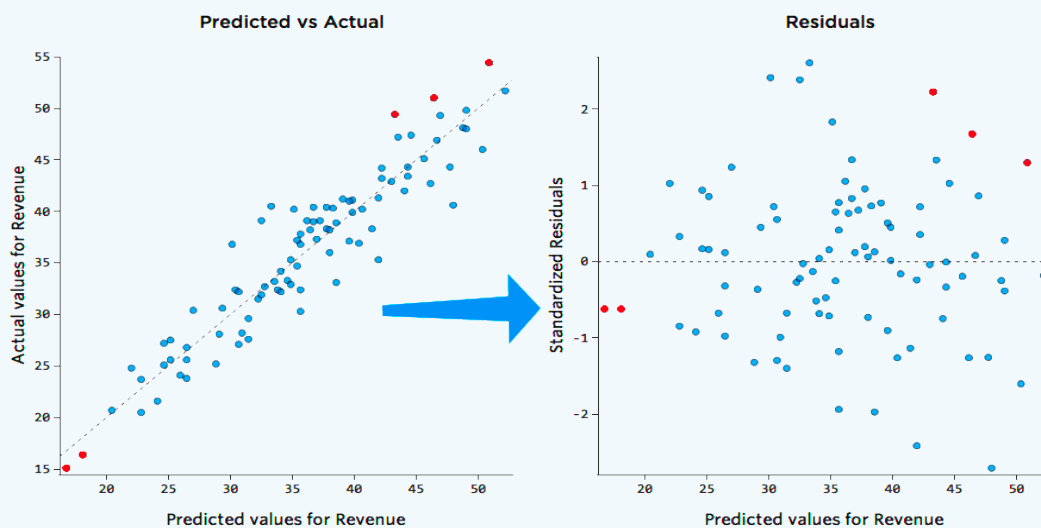


Figure 3.3: A residual plot

Definition 3.3 – The Role of s and r^2 in Regression.

- If we use a least-squares line to predict the values of a response variable y from an explanatory variable x , the **standard deviation of the residuals** s is

Definition 3.2.4 – Standard deviation of the residuals s .

$$s = \sqrt{\frac{\sum \text{residuals}^2}{n - 2}} = \sqrt{\frac{\sum (y_i - \hat{y})^2}{n - 2}} \quad (3.7a)$$

This value gives the approximate size of a typical **prediction error** (residual)

- The **coefficient of determination** r^2 is the fraction of the variation in the values of y that is accounted for by the least-squares regression line of y on x . (The **percentage** of how well the line fits those data)

Definition 3.2.5 – Coefficient of determination r^2 .

$$r^2 = 1 - \frac{\sum \text{residuals}^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\sum (\text{Obs} - \text{prediction})^2}{\sum (\text{Obs} - \text{mean Obs})^2} \quad (3.8a)$$

Definition 3.4 – Outliers and influential observations in regression.

- An **outlier** is an observation that lies outside the overall pattern of the other observation. Points that are outliers in the y direction but not the x direction of a scatterplot have large residuals. Other outliers may not have large residuals.
- An observation is **influential** for a statistical calculation if removing it would markedly change the result of the calculation. Outliers in x are often influential for the regression line.



Chapter 4: Designing Studies

4	Designing Studies	29
4.1	Sampling and Surveys	
4.2	Experiments	
4.3	Using Studies Wisely	

4. Designing Studies

4.1 Sampling and Surveys

Definition 4.1 – Population, census, and sample.

- The **population** in a statistical study is the entire group of individuals we want information about.
- A **census** collect data from every individual in the population.
- A **sample** is a subset of individuals in the population from which we actually collect data.
- A **survey** (sample survey) is used to infer statistics of a population.

Definition 4.2 – How to Sample Badly.

- **Convenience sample** is to choose individuals from the population who are easy to reach results.
- A design of a statistical study with **bias** would consistently underestimate or consistently overestimate the value you want to know.
- A **voluntary response sample** consist of people who choose themselves by responding to a general invitation.

Definition 4.3 – How to Sample Well.

- **Random sampling** involves using a chance process to determine which members of a population are included in the sample.
- A **simple random sample** (SRS) of size n is chosen in such a way that every group of n individuals in the population has an equal chance to be selected as the sample.

Definition 4.4 – Other Random Sampling Methods.

- To get a **stratified random sample**, start by classifying the population into groups of similar individuals, called **strata**. Then choose a separate SRS in each stratum and combine these SRSs to form the sample.
- To get a **cluster sample**, start by classifying the population into groups of individuals that are located near each other, called **clusters**. Then choose an SRS of the clusters. All individuals in the chosen clusters are included in the sample.

Definition 4.5 – Sample Surveys: What Can Go Wrong?.

- **Undercoverage** occurs when some members of the population cannot be chosen in a sample.
- **Nonresponse** occurs when an individual chosen for the sample can't be contacted or refuses to participate.

Incorrect answers \implies Response bias

Wording of questions has a big influence on the answer.

4.2 Experiments**Definition 4.1 – Observational Study VS Experiment.**

- An **observational study** observes individuals and measures variables of interest but does not attempt to influence the response.
- An **experiment** deliberately imposes some treatment on individuals to measure their responses.
- **Confounding** occurs when two variables are associated in such a way that their effects on a response variable cannot be distinguished from each other.

Definition 4.2 – The Language of Experiments.

- **Treatment** is the condition applied to subjects in an experiment.
- The **experimental units** are the smallest collection of individuals to which treatments are applied. When they are human beings, they often are called **subjects**.

Definition 4.3 – How to Experiment Badly.

- In an experiment, **random assignment** means that experimental units are assigned to treatments using a chance process.

Principles of Experimental Design:

$$\text{A well experiment} \left\{ \begin{array}{l} \text{Must make a } \boxed{\text{comparison}} \text{ between treatments.} \\ \text{Must } \underline{\text{randomly assign}} \text{ subjects to treatments.} \\ \text{Must } \underline{\text{control}} \text{ all other variables that affect the response to be the same for all groups.} \\ \text{Must select enough subjects} \Rightarrow \boxed{\text{replication}}. \end{array} \right. \quad (4.2a)$$

A good control minimizes confounding and reduces variability in the response. ■

Definition 4.4 – Completely Randomized Designs.

- In a **completely randomized design**, the experimental units are assigned to the treatments completely by chance.
All experiments need a control (no treatment or placebo group)

Definition 4.5 – Experiments: What Can Go Wrong?.

- In a **double-blind** experiment, neither the subjects nor those who interact with them and measure the response variable know which treatment a subject received. If one party knows and the other doesn't, then the experiment is **single-blind**. ■

Definition 4.6 – Inference for Experiments.

- **Statistically significant** is an observed effect so large that it would rarely occur by chance. ■

Definition 4.7 – Blocking.

- A **block** is a group of experimental units that are known before the experiment to be similar in some way that is expected to affect the response to the treatments.
- In a **randomized block design**, the random assignment of experimental units to treatments is carried out separately within each block. ■

4.3 Using Studies Wisely

Definition 4.6 – Inference about a population.

- **Inference about a population** requires that the individuals taking part in a study be randomly selected from the population. A well-designed experiment that randomly assigns experimental units to treatments allows **inference about cause and effect**.
- **Lack of realism** in an experiment can prevent us from generalizing its results.
- Any information about the individuals in the study must be kept **confidential** ■



Chapter 5: Probability

5	Probability	35
5.1	Randomness, Probability, and Simulation	
5.2	Probability Rules	
5.3	Conditional Probability and Independence	

5. Probability

5.1 Randomness, Probability, and Simulation

Definition 5.1 – The Idea of Probability.

- The **Law of Large Numbers** says that the proportion of times that a particular outcome occurs in many repetitions will approach a single number, which is the **possibility**.
- **Probability** is a number between 0 and 1 describing the proportion of the time the outcome would occur over the long run.

Definition 5.2 – Simulation.

- A **simulation** is an imitation of chance behavior, most often carried out with random number.

Four-step process of simulation:

Simulation { **State:** Ask a question of interest about some chance process.
Plan: Describe one repetition of process.
Do: Perform many repetitions of the simulation.
Conclude: Use the results of your simulation to answer the question of interest.

(5.1a)

5.2 Probability Rules

Definition 5.1 – Probability Models.

- The **sample space** S of chance process is the set of all possible outcomes.
- A **probability model** is a description of some chance process that consists of two parts: a sample space S and a probability for each outcome.
- An **event** is any collection of outcomes from some chance process.

Definition 5.2 – Basic Rules of Probability.

- For any event A , $0 \leq P(A) \leq 1$
- If S is the sample space in a probability model, $P(S) = 1$.
- $P(A) = \frac{\text{number of outcomes corresponding to event } A}{\text{total number of outcomes in sample space}}$
- **Complement rule:** $P(A') = 1 - P(A)$ $P(A) = P(A \cap B) + P(A \cap B')$.
- **Mutually exclusive (disjoint):** two events A and B have no outcomes in common and so can never occur together—that is, if $P(A \cap B) = 0$.
In other words: $P(A \cup B) = P(A) + P(B)$.

Definition 5.3 – General Addition Rule For Two Events.

- If A and B are any two events resulting from some chance process, the

Definition 5.2.1 – General Addition Rule For Two Events.

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \quad (5.2a)$$

5.3 Conditional Probability and Independence**Definition 5.3 – Conditional Probability.**

- **Conditional Probability** is the possibility that one event happens given that another event is already known to have happened. Suppose we know that event A has happened. Then the probability that event B happens given that event A has happened is denoted by $P(B|A)$.

Definition 5.4 – Calculation Conditional Probabilities.

- To find the conditional probability $P(A|B)$, use formula

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (5.3a)$$

- The conditional probability $P(B|A)$ is given by

$$P(B|A) = \frac{P(B \cap A)}{P(A)} \quad (5.4a)$$

Definition 5.5 – General Multiplication Rule.

- The probability that event A and B both occur can be found using the **General Multiplication Rule**

$$P(A \cap B) = P(A) \cdot P(B|A) \quad (5.5a)$$

where $P(B|A)$ is the conditional probability that event B occurs given that event A has already occurred.

Definition 5.6 – Independent events.

- When events A and B are **independent**:

$$P(A|B) = P(A) \quad \text{and} \quad P(B|A) = P(B) \quad (5.6a)$$

where $P(B|A)$ is the conditional probability that event B occurs given that event A has already occurred.

Definition 5.3.1 – Multiplication rule for independent events.

$$P(A \cap B) = P(A) \cdot P(B) \quad (5.7a)$$



Chapter 6: Random Variables

6	Random Variables	41
6.1	Discrete and Continuous Random Variables	
6.2	Transforming and Combining Random Variables	
6.3	Binomial and Geometric Random Variables	
	Index	47

6. Random Variables

6.1 Discrete and Continuous Random Variables

Definition 6.1 – Random variable and probability distribution.

- A **random variable** takes numerical values that describe the outcomes of some chance process.
- The **probability distribution** of a random variable gives its possible values and their probability.

Definition 6.2 – Discrete Random Variables.

- A **Discrete random variable** X takes a **fixed** set of possible values with gaps between. The probability distribution of a discrete random variable X lists the values x_i and their probabilities p_i :

Value:	x_1	x_2	x_3	\cdots
Probability:	p_1	p_2	p_3	\cdots

The probabilities p_i must satisfy two requirements:

- Every probability p_i is a number between 0 and 1.
- The sum of the probabilities is 1: $p_1 + p_2 + p_3 + \cdots + p_n = 1$

Definition 6.3 – Mean (Expected Value) of a Discrete Random Variable.

- Suppose that X is a discrete random variable with probability distribution

Value:	x_1	x_2	x_3	\cdots
Probability:	p_1	p_2	p_3	\cdots

To find the **mean (expected value)** of X , multiply each possible value by its probability, then add all the products:

Definition 6.1.1 – Mean (Expected Value).

$$\begin{aligned}\mu_x = E(X) &= x_1p_1 + x_2p_2 + x_3p_3 + \cdots \\ &= \sum x_i p_i\end{aligned}\quad (6.1a)$$

Definition 6.4 – Variance and standard deviation of a discrete random variable.

- Suppose that X is a discrete random variable with probability distribution

Value:	x_1	x_2	x_3	\cdots
Probability:	p_1	p_2	p_3	\cdots

and that μ_x is the mean of X . The **variance** of X is

Definition 6.1.2 – Variance of a discrete random variable.

$$\begin{aligned}\text{Var}(X) = \sigma_x^2 &= (x_1 - \mu_x)^2 p_1 + (x_2 - \mu_x)^2 p_2 + (x_3 - \mu_x)^2 p_3 + \cdots \\ &= \sum (x_i - \mu_x)^2 p_i\end{aligned}\quad (6.2a)$$

The **standard deviation** of X , σ_x , is the square root of the variance.

Definition 6.1.3 – Standard deviation of a discrete random variable.

$$\sigma_x = \sqrt{\sum (x_i - \mu_x)^2 p_i}\quad (6.3a)$$

Definition 6.5 – Continuous Random Variables.

- A **Continuous random variable** X takes all values in an interval of numbers. The probability distribution of X is described by a density curve.

6.2 Transforming and Combining Random Variables

Definition 6.1 – Linear Transformations.

- Adding (or subtracting) each value of a random variable by a positive number a :

Definition 6.2.1 – Effect On A Random Variable of Adding (or subtracting) by A Constant.

- Mean, Medean, Quartiles, and percentiles $+ (-) a$
- Does not change Range, IQR, Standard deviation

- Multiplying (or dividing) each value of a random variable by a positive number b :

Definition 6.2.2 – Effect On A Random Variable of Multiplying (or Dividing) by A Constant.

- Mean, Medean, Quartiles, and percentiles $\times (\div) b$
- Range, IQR, and Standard deviation $\times (\div) b$
- Does not change the slop of the distribution

Definition 6.2 – Combining Random Variables.

- Mean of the Sum of Random Variables:

$$E(T) = \mu_T = \mu_x + \mu_y \quad (6.6a)$$

- Range of the Sum of Random Variables:

$$\text{range of } T = \text{range of } X + \text{range of } Y \quad (6.7a)$$

- Variance of the Sum of Random Variables:

$$\sigma_T^2 = \sigma_X^2 + \sigma_Y^2 \quad (6.8a)$$

- Mean of the Difference of Random Variables:

$$\mu_D = E(D) = \mu_x - \mu_y \quad (6.9a)$$

- Variance of the Difference of Random Variables:

$$\sigma_D^2 = \sigma_X^2 + \sigma_Y^2 \quad (6.10a)$$

6.3 Binomial and Geometric Random Variables

Definition 6.6 – Binomial Settings and Binomial Random Variables.

- A **Binomial setting** consists of n independent trials of the same chance process, each resulting in a success or a failure, with probability of success p on each trial.

$$\text{BINS} \left\{ \begin{array}{l} \text{Trails can be classified as "success" or "failure."} \\ \text{Trails must be independent.} \\ \text{The number of trials } n \text{ must be fixed.} \\ \text{There is the same probability } p \text{ of success on each trail.} \end{array} \right.$$

- The count X of successes is a **Binomial Random Variable**. Its probability distribution is a **Binomial Distribution**.

Definition 6.7 – Binomial Probabilities.

- The **Binomial Coefficient**

Definition 6.3.1 – Binomial Coefficient.

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \quad (6.12a)$$

counts the number of ways k successes can be arranged among n trials. The **factorial** of n is

$$n! = n(n-1)(n-2) \cdots (3)(2)(1) \quad (6.13a)$$

for positive whole numbers n , and $0! = 1$

- **Binomial Probability Formula**

Definition 6.3.2 – Binomial Probability.

$$P(X = K) = \binom{n}{k} p^k (1-p)^{n-k} \quad (6.14a)$$

If a count X of successes has the binomial distribution with number of trials n and probability of success p , the **mean** and **standard deviation** of X are

Definition 6.3.3 – Mean and Standard deviation of Binomial distribution.

$$\mu_x = np \quad (6.15a)$$

$$\sigma_x = \sqrt{np(1-p)} \quad (6.15b)$$

Definition 6.8 – Binomial Distributions in Statistical Sampling.

- **10% Condition**

When taking an simple random sample of size n from a population of size N , we can use a binomial distribution to model the count of successes in the sample as long as $n \leq \frac{1}{10}N$.

- **The Large Counts Condition**

Suppose that a count X of successes has the binomial distribution with n trials and success probability p . When n is large, the distribution of X is approximately Normal with

$$\text{mean: } \mu_x = np \text{ and standard deviation: } \sigma_X = \sqrt{np(1-p)} \quad (6.16a)$$

As an approximation, we will use the Normal approximation when n is no longer that

$$np \geq 10 \quad \text{and} \quad n(1-p) \geq 10 \quad (6.17a)$$

That is, the expected number of **successes** and **failures** are both at **least 10**. We refer to this as the **Large Counts condition**. ■

Definition 6.9 – Geometric Random Variables.

- A **Geometric setting** consists of repeated trials of the same chance process in which the probability p of successes is the same on each trial, and the goal is to count the number of trials it takes to get one success.
- If Y = the number of trials required to obtain the first success, then Y is a **Geometric probability** that Y takes any value is

Definition 6.3.4 – Geometric Probability.

$$P(Y = K) = (1-p)^{k-1}p \quad (6.18a)$$

The **mean** (expected value) of a geometric random variable Y is

Definition 6.3.5 – Mean of Geometric Probability.

$$\mu_Y = E(Y) = \frac{1}{p} \quad (6.19a)$$

which is the expected number of trials required to get the first success. ■

Index

- Analyzing Categorical Data, 5
- Binomial and Geometric Random Variables, 29
- Categorical variable and quantitative variable, 5
- Conditional Probability and Independence, 25
- Density Curves, 14
- Density Curves and Normal Distributions, 14
- Describing Location in a Distribution, 13
- Describing Quantitative Data with Numbers, 9
- Describing Shape, 7
- Discrete and Continuous Random Variables, 29
- Displaying Quantitative Data with Graphs, 7
- Distribution, 5
- Distribution of a categorical variable, 5
- Dotplot, 8
- Experiments, 21
- Histogram, 7
- Individuals and variables, 5
- Least-Squares Regression, 17
- Mean and median, 9
- Measuring Positions, 13
- Measuring Spread, 9
- Normal Distribution, 14
- Normal Distribution, 13, 14
- Probability Rules, 25
- Randomness, Probability, and Simulation, 25
- Sampling and Surveys, 21
- Scatterplots and Correlation, 17
- Stemplot, 8
- Transforming and Combining Random Variables, 29
- Using Studies Wisely, 21