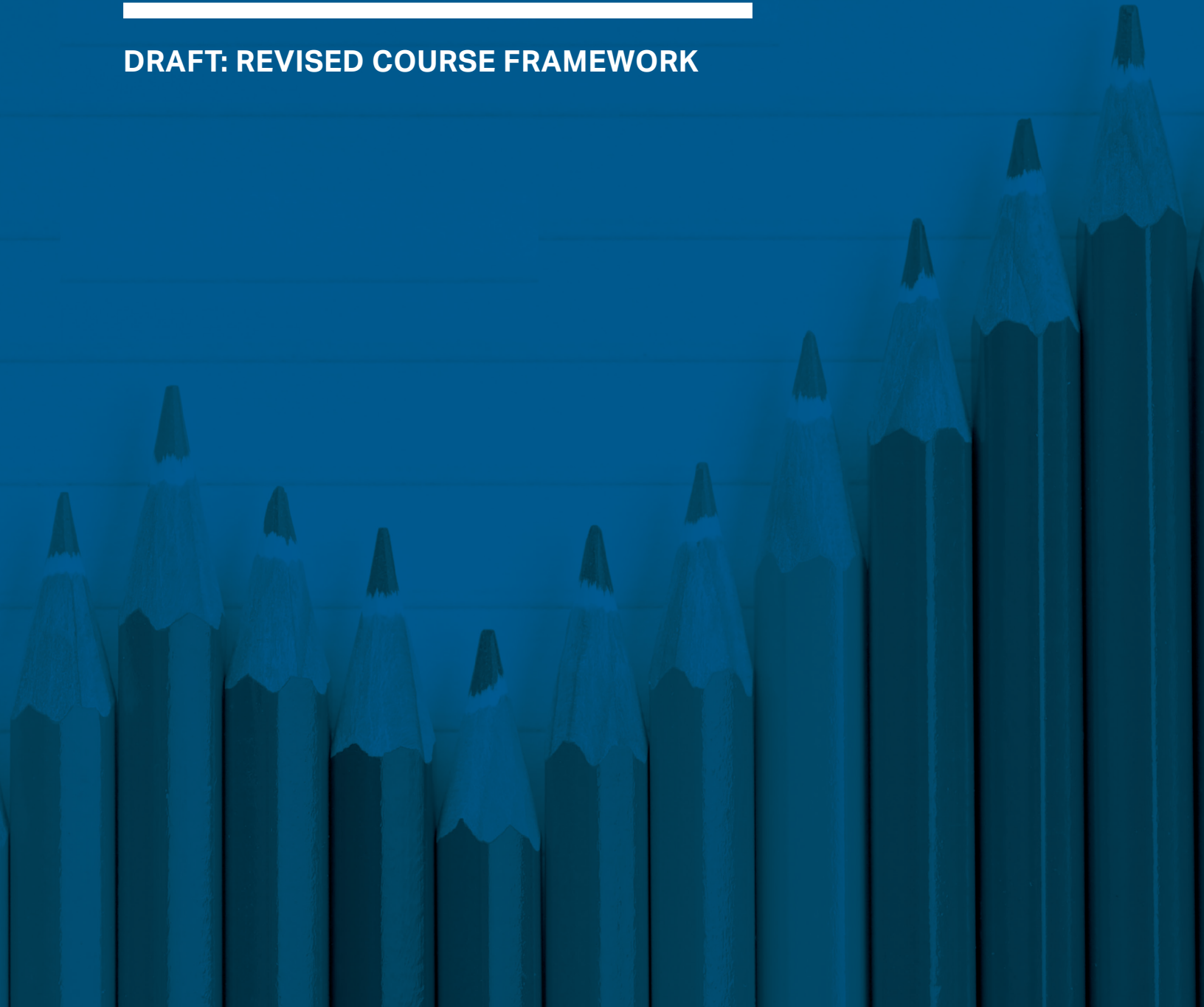




# AP Statistics

---

DRAFT: REVISED COURSE FRAMEWORK



THIS PAGE IS INTENTIONALLY LEFT BLANK.

# Contents

---

## 4 Acknowledgments

---

### **COURSE FRAMEWORK**

## 7 Course Framework Components

## 8 Statistical Practices

9 **UNIT 1:** Exploring One-Variable Data and Collecting Data

29 **UNIT 2:** Probability, Random Variables, and Probability  
Distributions

45 **UNIT 3:** Inference for Categorical Data: Proportions

69 **UNIT 4:** Inference for Quantitative Data: Means

87 **UNIT 5:** Regression Analysis

# Acknowledgments

---

**Prince Afriyie**, *University of Virginia, VA*  
**Javier Alvarez**, *Phillips Exeter Academy, NH*  
**Barb Barnet**, *University of Wisconsin-Platteville, WI*  
**Tim Bator**, *Sayre School, KY*  
**Ellen Breazel**, *Clemson University, SC*  
**Carol Chamberlain Hebert**, *Barnstable High School, MA*  
**Holly Deal**, *Kennesaw State University, GA*  
**Christine Franklin**, *University of Georgia, GA*  
**Donna LaLonde**, *American Statistical Association, VA*  
**Laura Marshall**, *Phillips Exeter Academy, NH*  
**Bridget Matamoros-Mota**, *John H. Guyer High School, TX*  
**Karen McGaughey**, *California Polytechnic State University, CA*  
**Chris Olsen**, *Grinnell College, IA*  
**Michael Posner**, *Villanova University, PA*  
**Al Reiff**, *The Taft School, CT*  
**Paul Roback**, *St. Olaf College, MN*  
**Tom Rohnkohl**, *Silverado High School, NV*  
**Allen Rossman**, *California Polytechnic State University, CA*  
**Joshua Sawyer**, *Camden County Schools, NC*

## College Board Staff

**Camille Pace**, *Director, AP Statistics Curriculum and Assessment*  
**Jason VanBilliard**, *Senior Director, AP Math and Computer Science*  
*Department Head*  
**Dana Kopelman**, *Executive Director, AP Instruction and Assessment Production*  
**Claire Lorenz**, *Senior Director, AP Classroom Instruction Products*  
**Jason Manoharan**, *Vice President, AP Program Development*  
**Daniel McDonough**, *Senior Director, AP Content and Assessment Publications*  
**Trevor Packer**, *Senior Vice President, AP and Instruction*  
**Allison Thurber**, *Vice President, AP Curriculum and Assessment*

**AP STATISTICS**

---

# Course Framework



# Course Framework Components

---

## Course Units

**Unit 1:** Exploring One-Variable Data and Collecting Data

**Unit 2:** Probability, Random Variables, and Probability Distributions

**Unit 3:** Inference for Categorical Data: Proportions

**Unit 4:** Inference for Quantitative Data: Means

**Unit 5:** Regression Analysis

## Curriculum Framework Overview

This curriculum framework provides a clear and detailed description of the course requirements necessary for student success. The framework specifies what students must know, be able to do, and understand to qualify for college credit or placement.

The curriculum framework includes two essential components:

- **AP Statistical Practices** (p. 8)  
The statistical practices are central to the study and practice of statistics. Students should develop and apply the described practices on a regular basis over the span of the course.
- **Course Content** (p. 9)  
The course content is organized into commonly taught units of study that provide a suggested sequence for the course and detail required content and conceptual understandings that colleges and universities typically expect students to master to qualify for college credit and/or placement.

# Statistical Practices

Practice 1	Practice 2	Practice 3	Practice 4
<i>Formulate Questions</i>	<i>Collect Data</i>	<i>Analyze Data</i>	<i>Interpret Results</i>
<i>Determine a research question for a statistical study.</i>	<i>Identify and justify methods for collecting data and conducting statistical inference.</i>	<i>Construct representations of data and calculate numerical statistical outputs.</i>	<i>Interpret results and justify conclusions and methods.</i>
<p><b>1.A:</b> Determine a valid research question that requires a statistical investigation.</p>	<p><b>2.A:</b> Identify information to answer a question or solve a problem.</p> <p><b>2.B:</b> Justify an appropriate method for ethically gathering and representing data.</p> <p><b>2.C:</b> Identify appropriate statistical inference methods.</p> <p><b>2.D:</b> Identify relationships among components in statistical inference methods.</p> <p><b>2.E:</b> Identify the null and alternative hypotheses.</p>	<p><b>3.A:</b> Construct tabular and graphical representations of data and distributions.</p> <p><b>3.B:</b> Calculate summary statistics, relative positions of points within a distribution, and predicted responses.</p> <p><b>3.C:</b> Calculate and estimate expected counts, percentages, probabilities, and intervals.</p> <p><b>3.D:</b> Calculate parameters for probability distributions.</p> <p><b>3.E:</b> Calculate appropriate statistical inference method results.</p>	<p><b>4.A:</b> Describe and compare tabular and graphical representations of data.</p> <p><b>4.B:</b> Justify a claim based on statistical calculations and results.</p> <p><b>4.C:</b> Describe distributions and compare relative positions of points within a distribution.</p> <p><b>4.D:</b> Interpret statistical calculations and results to assess meaning or a claim.</p> <p><b>4.E:</b> Justify the use of a chosen statistical inference method by verifying conditions.</p> <p><b>4.F:</b> Interpret results of statistical inference methods.</p> <p><b>4.G:</b> Justify a claim based on statistical inference method results.</p>



## AP STATISTICS

# UNIT 1

# Exploring One-Variable Data and Collecting Data

THIS PAGE IS INTENTIONALLY LEFT BLANK.

## TOPIC 1.1

## Introducing Statistics: What Can We Learn from Data?

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.A	<p data-bbox="342 688 396 709"><b>1.1.A</b></p> <p data-bbox="342 722 553 806">Identify components within a statistical study.</p>	<p data-bbox="605 688 675 709"><b>1.1.A.1</b></p> <p data-bbox="605 722 1468 779">A statistical study is a study in which data are collected from a sample to answer a research question about a larger population.</p> <p data-bbox="605 793 675 814"><b>1.1.A.2</b></p> <p data-bbox="605 827 1419 884">Statistical studies are necessary when the population is too large or it is too difficult to collect data from every item or individual in the population.</p> <p data-bbox="605 898 675 919"><b>1.1.A.3</b></p> <p data-bbox="605 932 1446 989">A datum (singular for data) is a piece of information about an item or individual. A collection of data is called a data set.</p> <p data-bbox="605 1003 675 1024"><b>1.1.A.4</b></p> <p data-bbox="605 1037 1446 1094">A population consists of all items or individuals of interest. The population size is represented by the symbol <math>N</math>.</p> <p data-bbox="605 1108 675 1129"><b>1.1.A.5</b></p> <p data-bbox="605 1142 1398 1220">A sample selected for study is a subset of the population from which data are obtained. The number of items in the sample, called the sample size, is represented by the symbol <math>n</math>.</p> <p data-bbox="605 1234 675 1255"><b>1.1.A.6</b></p> <p data-bbox="605 1268 1419 1388">Each component of a statistical study and the resulting calculations can be related to an aspect of the corresponding real-world context from which the components were derived. This identification of a statistical result with the corresponding contextual component is what is meant by “in context.”</p>
1.A	<p data-bbox="342 1444 396 1465"><b>1.1.B</b></p> <p data-bbox="342 1478 537 1583">Determine a research question within a statistical study.</p>	<p data-bbox="605 1444 675 1465"><b>1.1.B.1</b></p> <p data-bbox="605 1478 1419 1535">A research question for a specific study should have a defined purpose and should not be changed based on the data analysis or results.</p> <p data-bbox="605 1549 675 1570"><b>1.1.B.2</b></p> <p data-bbox="605 1583 1468 1625">A research question should be posed so that the required data can be collected and analyzed.</p>

## TOPIC 1.2

# Variables

Instructional Periods: 1

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.A	<p><b>1.2.A</b> Identify observational units, variables, parameters, and statistics from a statistical study or data set.</p>	<p><b>1.2.A.1</b> An observational unit is an item or individual from which a datum is collected.</p> <p><b>1.2.A.2</b> A variable is a characteristic that may change from one observational unit to another.</p> <p><b>1.2.A.3</b> Data collected on numerical and categorical variables measured on observational units, which could include photographs, sounds, videos, and text, can be used to convey meaningful information.</p> <p><b>1.2.A.4</b> A parameter is a numerical attribute or summary of the variable of interest for a population.</p> <p><b>1.2.A.5</b> A statistic is a numerical attribute or summary of the variable of interest for a sample. The value of a statistic from a certain sample is often not equal to the unknown value of the population parameter but may provide the basis for making inferences about the population parameter.</p>
2.A	<p><b>1.2.B</b> Identify types of variables.</p>	<p><b>1.2.B.1</b> A qualitative variable, also called a categorical variable, takes on values that are category names or group labels.</p> <p><b>1.2.B.2</b> A quantitative variable, also called a numerical variable, takes on numerical values for a measured or counted quantity and generally has units of measure.</p>
2.A	<p><b>1.2.C</b> Identify types of quantitative variables.</p>	<p><b>1.2.C.1</b> A discrete quantitative variable can take on a countable number of values. The number of values may be finite or countably infinite, as with the whole numbers.</p> <p><b>1.2.C.2</b> A continuous quantitative variable can take on an infinite number of possible values within a given interval. The number of values the variable can take on is measurable but not countable. This variable can take on all possible values between any pair of values.</p>

## TOPIC 1.3

# Tabular Representation and Summary Statistics for One Categorical Variable

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.A	<b>1.3.A</b> Construct categorical one-variable tabular representations.	<b>1.3.A.1</b> A frequency table gives the number of observational units falling into each category of a categorical variable.  <b>1.3.A.2</b> A relative frequency table gives the proportion of observational units falling into each category of a categorical variable.
4.A	<b>1.3.B</b> Describe categorical one-variable tabular representations.	<b>1.3.B.1</b> Percentages, relative frequencies, and ratios all provide the same information as proportions.  <b>1.3.B.2</b> Counts and relative frequencies of categorical variables reveal information that can be used to justify claims about the variables in context.

## TOPIC 1.4

# Graphical Representations, Descriptions, and Comparisons for One Categorical Variable

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.A	<b>1.4.A</b> Construct categorical one-variable graphical representations.	<b>1.4.A.1</b> Bar charts, also called bar graphs, display frequencies (counts) or relative frequencies (proportions) for the categories of a single categorical variable. Each bar on a bar graph is associated with a category for the categorical variable of interest. The height or length of each bar corresponds to either the frequency or relative frequency of the observational units falling within each category. <b>1.4.A.2</b> Pie charts are used to display frequencies (counts) or relative frequencies (proportions) for categorical data. Each slice on a pie chart is associated with a category for the categorical variable of interest. The area of each slice, as a fraction of the total area, corresponds to the relative frequency of observational units falling within each category. The sum of the slices together will equal 1, or 100% of the total area.
4.A	<b>1.4.B</b> Describe categorical one-variable graphical representations.	<b>1.4.B.1</b> Graphical representations of a categorical variable reveal information that can be used to justify claims about the variable in context.
4.A	<b>1.4.C</b> Compare multiple categorical one-variable graphical and tabular representations.	<b>1.4.C.1</b> Frequency tables, bar charts, and pie charts can be used to compare two or more data sets in terms of the same categorical variable.

## TOPIC 1.5

# Graphical Representations for One Quantitative Variable

Instructional Periods: 2

**SKILL****3.A****LEARNING OBJECTIVE****1.5.A**

Construct quantitative one-variable graphical representations.

**ESSENTIAL KNOWLEDGE****1.5.A.1**

Histograms, stem-and-leaf plots, and dotplots provide a visual representation of the distribution of the values of a quantitative variable. These graphs show the frequency or relative frequency of the quantitative variable values or intervals of values and maintain the natural ordering, smallest to largest, of the quantitative variable.

**1.5.A.2**

A histogram puts the observed values of the quantitative variable into ordered intervals, or bins, along the horizontal axis. A bar is associated with each interval or bin, and the height of each bar shows the frequency or relative frequency of the observations that fall within the interval of the quantitative variable values corresponding to that bar. Altering the interval widths, or bin widths, can change the appearance of the histogram.

**1.5.A.3**

A stem-and-leaf plot splits each value of the quantitative variable into two parts: a “stem” (the first digit or digits) and a “leaf” (usually the single digit after the stem digit(s)). Both stems and leaves are ordered from smallest to largest.

**1.5.A.4**

A dotplot represents each value of the quantitative variable by a dot. Each dot is placed above the horizontal or beside the vertical axis corresponding to the value of that observation, with nearly identical values stacked on top of each other.

## TOPIC 1.6

# Descriptions for One Quantitative Variable

Instructional Periods: 2

**SKILL****4.A****LEARNING OBJECTIVE****1.6.A**

Describe quantitative one-variable graphical representations.

**ESSENTIAL KNOWLEDGE****1.6.A.1**

Descriptions of the distribution of one quantitative variable include shape, center, and variability (spread) as well as any unusual features such as outliers, gaps, or clusters in context.

**1.6.A.2**

The shape of the distribution of one quantitative variable is skewed to the right (positively skewed) if the right tail (toward larger values) is longer than the left. The shape of the distribution is skewed to the left (negatively skewed) if the left tail (toward smaller values) is longer than the right. The shape of the distribution is symmetric if the left half is the mirror image of the right half.

**1.6.A.3**

Distributions of one quantitative variable with one main peak are called unimodal. Distributions with two prominent peaks are called bimodal. A distribution in which each frequency or each relative frequency is approximately the same with no prominent peaks is approximately uniform.

**1.6.A.4**

Outliers for one quantitative variable are data points that are unusually small or large relative to the rest of the data.

**1.6.A.5**

A gap is a region in a distribution between two values in which there are no observed data.

**1.6.A.6**

Clusters are concentrations of values usually separated by gaps.



## TOPIC 1.7

# Summary Statistics for One Quantitative Variable

Instructional Periods: 2

## SKILL

3.B

## LEARNING OBJECTIVE

1.7.A

Calculate measures of center and position for quantitative data.

## ESSENTIAL KNOWLEDGE

1.7.A.1

Two commonly used measures of center in the distribution of a quantitative variable are the mean and median.

1.7.A.2

The mean is the sum of all the values divided by the number of values and can be found using technology. For a sample, the mean is denoted by  $\bar{x}$ :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \text{ where } x_i \text{ represents the } i^{\text{th}} \text{ data point in the sample and } n$$

represents the number of data values in the sample.

1.7.A.3

The median is the middle value when the data set is ordered from smallest to largest and can be found using technology. One common method for determining the median of a data set with an even number of values is to use the mean of the two middle values. A common method for determining the median of a data set with an odd number of values is to use the value in the middle of all the values.

1.7.A.4

In an ordered data set, the smallest value is the minimum value, and the largest value is the maximum value.

1.7.A.5

The first quartile, denoted by Q1, is the median value of the lower half of the ordered data set from the minimum value to the position of the median. Approximately 25% of the values in the data set are less than or equal to Q1. The third quartile, denoted by Q3, is the median value of the upper half of the ordered data set from the position of the median to the maximum value. Approximately 75% of the values in the data set are less than or equal to Q3. The second quartile, Q2, is also the median of the data set. Q1 and Q3 form the boundaries for the middle 50% of values in an ordered data set.

1.7.A.6

The  $p$ th percentile is the value that has  $p\%$  of the data less than or equal to it when the data set is ordered from smallest to largest. The first and third quartiles are the 25th and 75th percentiles, respectively.

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.B	<p><b>1.7.B</b> Calculate measures of variability for quantitative data.</p>	<p><b>1.7.B.1</b> Three commonly used measures of variability (or spread) in the distribution of a quantitative variable are the range, interquartile range, and standard deviation.</p> <p><b>1.7.B.2</b> The range is the difference between the maximum data value and the minimum data value.</p> <p><b>1.7.B.3</b> The interquartile range (IQR) is the difference between the third and first quartiles: <math>Q3 - Q1</math>.</p> <p><b>1.7.B.4</b> The standard deviation is a typical deviation of the data values from their mean and can be found using technology. The sample standard deviation is denoted by <math>s_x</math> and calculated as follows: <math>s_x = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}</math>, where <math>x_i</math> is the data value, <math>\bar{x}</math> is the mean, and <math>n</math> is the number of data values in the sample. The square of the sample standard deviation, <math>s^2</math>, is called the sample variance.</p>
3.B	<p><b>1.7.C</b> Calculate different units of measurement for summary statistics.</p>	<p><b>1.7.C.1</b> Changing units of measurement affects the values of the calculated statistics.</p>
3.B	<p><b>1.7.D</b> Calculate outliers for quantitative data.</p>	<p><b>1.7.D.1</b> There are many methods for determining potential outliers. Two methods frequently used are as follows:</p> <ol style="list-style-type: none"> <li>In a skewed distribution, an outlier is a value located more than <math>1.5 \times IQR</math> above the third quartile or more than <math>1.5 \times IQR</math> below the first quartile.</li> <li>In a symmetric distribution, an outlier is a value located 2 or more standard deviations above, or below, the mean.</li> </ol>
4.B	<p><b>1.7.E</b> Justify the selection of a particular measure of center and variability for describing quantitative data.</p>	<p><b>1.7.E.1</b> The median and <i>IQR</i> are considered resistant (or robust) measures of center and variability, respectively, because outliers do not greatly (if at all) affect their values. The mean, standard deviation, and range are considered nonresistant (or non-robust) measures because outliers can affect their values greatly.</p>

## TOPIC 1.8

# Graphical Representations of Summary Statistics for One Quantitative Variable

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.A	<p><b>1.8.A</b></p> <p>Construct quantitative one-variable graphical representations of summary statistics.</p>	<p><b>1.8.A.1</b></p> <p>A five-number summary is made up of the minimum data value, the first quartile (Q1), the median, the third quartile (Q3), and the maximum data value.</p> <p><b>1.8.A.2</b></p> <p>A boxplot is a graphical representation of the five-number summary (minimum, first quartile, median, third quartile, maximum). The box represents the middle 50% of data, with a line at the median and the ends of the box corresponding to the quartiles. Lines (“whiskers”) that represent 25% of the data extend from the first quartile to the minimum and from the third quartile to the maximum. If there are outliers in the data, the whiskers extend to the most extreme data values that are not outliers, and outliers are usually denoted with an asterisk or other symbol.</p>
4.A	<p><b>1.8.B</b></p> <p>Describe quantitative one-variable graphical representations of summary statistics based on the relationship of the mean and the median.</p>	<p><b>1.8.B.1</b></p> <p>If a distribution is relatively symmetric, then the values of the mean and median are relatively close to each other. If a distribution is skewed right, then the value of the mean is usually larger than the median. If the distribution is skewed left, then the value of the mean is usually smaller than the median.</p>

TOPIC 1.9

# Comparisons of the Distribution for One Quantitative Variable

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.A	<p><b>1.9.A</b></p> <p>Compare multiple quantitative one-variable graphical representations.</p>	<p><b>1.9.A.1</b></p> <p>Graphical representations of a quantitative variable can be used to compare important features between two or more distributions of the same quantitative variable. Histograms, back-to-back stem-and-leaf plots, and dotplots may be used to compare center, variability, shape, outliers, clusters, or gaps in two or more distributions. Boxplots may be used to compare center, variability, outliers, and skewness (or symmetry).</p>
4.A	<p><b>1.9.B</b></p> <p>Compare multiple quantitative one-variable graphical representations of summary statistics.</p>	<p><b>1.9.B.1</b></p> <p>A comparison of graphical representations for two or more distributions can include any of the numerical summaries (e.g., mean, standard deviation, etc.).</p>
3.B	<p><b>1.9.C</b></p> <p>Calculate <math>z</math>-scores with population parameters.</p>	<p><b>1.9.C.1</b></p> <p>A standardized score measures the number of standard deviations a data value falls above or below the mean.</p> <p><b>1.9.C.2</b></p> <p>A <math>z</math>-score is calculated as <math>\frac{x_i - \mu}{\sigma}</math>, where <math>x_i</math> is the data value, <math>\mu</math> is the population mean, and <math>\sigma</math> is the population standard deviation. A <math>z</math>-score measures how many standard deviations a data value is above (positive <math>z</math>-score) or below (negative <math>z</math>-score) the mean. When the population mean and standard deviation are unknown, the sample mean and standard deviation may be used to determine a <math>z</math>-score.</p>
4.C	<p><b>1.9.D</b></p> <p>Compare <math>z</math>-scores as measures of relative position for distributions.</p>	<p><b>1.9.D.1</b></p> <p><math>z</math>-scores may be used to compare relative positions of individual values within a distribution or between distributions.</p>

## TOPIC 1.10

# Data Collection

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.A	<b>1.10.A</b> Identify a census.	<b>1.10.A.1</b> A census consists of recording information from all items or individuals in a population.
2.A	<b>1.10.B</b> Identify an observational study.	<b>1.10.B.1</b> An observational study is one in which the researcher is a passive observer, recording the values of variables of interest in order to investigate a research question of interest.  <b>1.10.B.2</b> A prospective study is one in which the observational units of study are selected at a point in time, and data are gathered both at that time and into the future.  <b>1.10.B.3</b> A retrospective study is one in which the observational units of study are selected at a point in time and data from the past are gathered.  <b>1.10.B.4</b> A survey is an observational study in which the data are collected from humans using a standard set of questions.
2.A	<b>1.10.C</b> Identify an experiment.	<b>1.10.C.1</b> An experiment is a study in which a researcher assigns conditions, or treatments, to experimental units to investigate a research question of interest about the population.  <b>1.10.C.2</b> The experimental unit is the observational unit to which the treatment is assigned. When experimental units consist of people, they are sometimes referred to as subjects or participants.  <b>1.10.C.3</b> An explanatory variable, or factor, is a variable whose different categories, or levels, are imposed on the experimental units. The different categories, or levels, of the explanatory variable are called treatments. When there is more than one explanatory variable, the combinations of the categories, or levels, of the explanatory variables are called treatments.  <b>1.10.C.4</b> A response variable is an outcome measured on each experimental unit after the treatment has been administered.

**SKILL****2.B****LEARNING OBJECTIVE****1.10.D**

Justify the appropriateness of generalizations for a statistical study.

**ESSENTIAL KNOWLEDGE****1.10.D.1**

A sample is considered random when all observational units in the sample have an equal chance of being selected from the population. A random mechanism (e.g., random number table, lottery, etc.) is used to select the observational units to be included in the sample.

**1.10.D.2**

When observational units, or experimental units, in a sample are randomly selected from a population, it is appropriate to make generalizations about the entire population of individuals from which the sample was selected.

**1.10.D.3**

A sample is not randomly selected when observational units are deliberately chosen or volunteer themselves to be in the sample.

**1.10.D.4**

When observational units, or experimental units, in a sample are not randomly selected from a population, it is appropriate to make generalizations only about a population of individuals that are similar to those used in the study.

## TOPIC 1.11

# Random Sampling

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.A	<p><b>1.11.A</b> Identify a sampling method given a description of a study.</p>	<p><b>1.11.A.1</b> Sampling without replacement is a sampling strategy in which an observational unit from a population can be selected only once. The observational unit is not returned to the population before subsequent selections of observational units are made, so there is no chance that observational unit can be selected again.</p> <p><b>1.11.A.2</b> Sampling with replacement is a sampling strategy in which an observational unit from the population can be selected more than once. The observational unit is returned to the population before subsequent selections of observational units are made, so it is possible that observational unit could be selected again.</p> <p><b>1.11.A.3</b> In a simple random sample (SRS) of size <math>n</math>, every sample of the size <math>n</math> has the same chance of being selected. This method is the basis for many types of sampling mechanisms. There are several procedures to obtain a simple random sample, for example, using a random number generator or a lottery method.</p> <p><b>1.11.A.4</b> A stratified random sample involves the division of all individuals in a population into non-overlapping groups, called strata, based on one or more shared attributes or characteristics (homogeneous grouping). Within each stratum a simple random sample is selected, and the selected individuals are combined to form one sample.</p> <p><b>1.11.A.5</b> A cluster sample involves the division of a population into smaller groups, called clusters. Ideally, each cluster mirrors the heterogeneity of the population, with clusters similar to one another. A simple random sample of clusters is selected from the population to form the sample of clusters. Data are collected from all observational units in each of the selected clusters.</p> <p><b>1.11.A.6</b> A systematic random sample is a method in which sample members from a population are selected according to a random starting point and a fixed, periodic interval between successive sampling units.</p>
2.B	<p><b>1.11.B</b> Justify the appropriateness of a sampling method.</p>	<p><b>1.11.B.1</b> Sampling methods can be used to mitigate the consequences of using simple random samples of certain populations for certain research questions.</p>

## TOPIC 1.12

# Potential Problems with Sampling

Instructional Periods: 2

**SKILL****2.A****LEARNING OBJECTIVE****1.12.A**

Identify potential sources of bias in sampling methods.

**ESSENTIAL KNOWLEDGE****1.12.A.1**

Bias in a sampling method is a systematic error in the sampling procedure that results in a statistic being consistently larger or consistently smaller than the parameter the statistic is used to estimate.

**1.12.A.2**

Voluntary response bias is a bias that may occur when a sample consists entirely of volunteers or is a convenience sample.

**1.12.A.3**

Undercoverage bias may occur when the sampling method fails to include part of the population.

**1.12.A.4**

Nonresponse bias may occur because of a failure to obtain responses from some individuals chosen to be sampled. The respondents and nonrespondents would differ significantly in ways that are important for the study.

**1.12.A.5**

Response bias may occur when responses to a survey or measurements of entities tend to differ from the “true” value in one direction. Examples include questions that are confusing or leading (question wording bias) or self-reported responses.

**1.12.A.6**

Nonrandom sampling methods (for example, samples chosen by convenience or voluntary response) introduce potential bias because they do not use random chance to select the individuals.



## TOPIC 1.13

# Experimental Design

Instructional Periods: 3

**SKILL****2.A****LEARNING OBJECTIVE****1.13.A**

Identify elements of a well-designed experiment.

**ESSENTIAL KNOWLEDGE****1.13.A.1**

A well-planned experiment should include the following:

- i. Comparisons of at least two treatment groups, one of which is a control group
- ii. Random assignment of treatments to experimental units
- iii. Replication
- iv. Direct control of potential extraneous sources of variation in the response

**1.13.A.2**

A control group is a collection of experimental units that are constructed for comparison. A control group may be given a treatment different from the treatment of interest to determine if the treatment of interest has an effect (i.e., a treatment with an inactive substance, a placebo, may be given).

**1.13.A.3**

The placebo effect is the difference between the average response to a placebo and the average response to no treatment.

**1.13.A.4**

In a single-blind, also called single-masked, experiment, participants do not know which treatment they are receiving, but members of the research team who interact with them know which treatment each participant is receiving, or vice versa.

**1.13.A.5**

In a double-blind, also called double-masked, experiment, neither the participants nor the members of the research team who interact with them know which treatment each participant is receiving.

**1.13.A.6**

An extraneous source of variation, also referred to as an extraneous variable, in the response variable is a variable that is known (or believed) to affect the response but is not an explanatory variable being studied.

**SKILL**

**LEARNING OBJECTIVE**

**ESSENTIAL KNOWLEDGE**

2.A

**1.13.A**

Identify elements of a well-designed experiment.

**1.13.A.7**

The purpose of random assignment is to create treatment groups that are as similar as possible with respect to extraneous sources of variation. If random assignment is successful, the respective distributions of each extraneous variable will be approximately the same for all the treatment groups.

**1.13.A.8**

A confounding variable is a variable that is distributed differently among treatment groups and affects the response variable. A confounding variable provides an alternative explanation for the observed relationship between the response and explanatory variables determined in the study, thereby lowering the credibility of the assertion of a causal relationship between the explanatory and response variables of interest. To be a confounding variable, a variable must be associated with both the explanatory variable and the response variable.

**1.13.A.9**

Replication within an experiment means more than one experimental unit is assigned to each treatment.

**1.13.A.10**

Direct control in an experiment means keeping the settings of certain potential extraneous sources of variation in the response variable the same from experimental unit to experimental unit.

2.A

**1.13.B**

Identify experimental designs.

**1.13.B.1**

In a completely randomized design, treatments are assigned to experimental units completely at random. Often the number of experimental units assigned to each treatment will be the same, but the sample sizes in each treatment do not have to be the same.

**1.13.B.2**

A blocking variable is a source of extraneous variation in the response variable. In a randomized block design, the experimental units are first grouped according to similar values of a blocking variable. These groups are called blocks. Units within the same block are homogeneous with respect to the blocking variable. After the blocks are formed, the treatments are randomly assigned to experimental units within each block so that all treatments occur within every block.

**1.13.B.3**

The purpose of blocking is to separate the variation in the response caused by the blocking variable from the rest of the extraneous variation in the response. Blocking allows for more precise comparisons of the response across the treatments. Within a block, the treatments can be compared without having to worry about variation in the response caused by changes in the blocking variable.

**1.13.B.4**

In randomized complete block designs, one experimental unit is assigned to each treatment within each block.

**1.13.B.5**

A matched pairs design is a randomized block design with only two treatments. Experimental units are arranged in pairs by matching on one or more extraneous sources of variation in the response variable. Each pair receives both treatments by randomly assigning one treatment to one member of the pair and the other treatment to the second member of the pair. Alternatively, each experimental unit may get both treatments while the order of the treatments is randomized.

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.B	<p><b>1.13.C</b> Justify the appropriateness of a particular experimental design.</p>	<p><b>1.13.C.1</b> One experimental design may be more appropriate than another experimental design based on the goals of the research study, the characteristics of the population, and the sample and variables involved.</p>
2.B	<p><b>1.13.D</b> Justify the appropriateness of the conclusions based on a well-designed experiment.</p>	<p><b>1.13.D.1</b> Using random assignment of treatments to experimental units allows for cause-and-effect conclusions between the explanatory and the response variables if the data show a statistically significant result, because the potential for confounding variables is reduced.</p> <p><b>1.13.D.2</b> Depending on the experimental unit, it may be unethical or difficult to randomly select experimental units to participate in an experiment. In that case, the study's experimental units are obtained from volunteers and will represent the population of experimental units similar to those who participated in the study.</p>

THIS PAGE IS INTENTIONALLY LEFT BLANK.

## AP STATISTICS

# UNIT 2

# Probability, Random Variables, and Probability Distributions

THIS PAGE IS INTENTIONALLY LEFT BLANK.

## TOPIC 2.1

# Graphical and Tabular Representations for the Distributions of Two Categorical Variables

Instructional Periods: 2

**SKILL****4.A****LEARNING OBJECTIVE****2.1.A**

Compare graphical and tabular representations for the distributions of two categorical variables.

**ESSENTIAL KNOWLEDGE****2.1.A.1**

Side-by-side bar charts, segmented bar charts, and mosaic plots are examples of graphs used to display and compare two categorical variables. In these graphs, the frequency or relative frequency of each category, or level, of one of the categorical variables is displayed for each category of the other categorical variable.

**2.1.A.2**

Graphical representations of two categorical variables can be used to compare the distributions of one categorical variable across the levels of the other categorical variable and determine whether the two variables are associated.

**2.1.A.3**

A two-way table, also called a contingency table, can be used to summarize and compare data for two categorical variables. The entries in the cells of the table can be frequencies (i.e., counts) or relative frequencies (i.e., proportions).

## TOPIC 2.2

# Summary Statistics for Two Categorical Variables

Instructional Periods: 1

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.B	<b>2.2.A</b> Calculate summary statistics from two-way tables.	<b>2.2.A.1</b> A joint relative frequency in a two-way table is a cell frequency divided by the total for the entire table. <b>2.2.A.2</b> A marginal relative frequency in a two-way table is a row total divided by the total for the entire table or a column total divided by the total for the entire table. <b>2.2.A.3</b> A conditional relative frequency is a relative frequency computed by restricting to a particular level, or category of interest. A conditional relative frequency can be a cell frequency in a row divided by the total for that row or it can be a cell frequency in a column divided by the total for that column.
4.A	<b>2.2.B</b> Compare summary statistics for two categorical variables.	<b>2.2.B.1</b> Summary statistics for two categorical variables can be used to compare distributions for evidence of association between the two variables.



## TOPIC 2.3

# Estimating Probabilities

Instructional Periods: 1

**SKILL****3.C****LEARNING OBJECTIVE****2.3.A**

Estimate probabilities.

**ESSENTIAL KNOWLEDGE****2.3.A.1**

A random process generates results that are determined by chance.

**2.3.A.2**

An outcome is the result of one trial of a random process.

**2.3.A.3**

An event is a collection of outcomes.

**2.3.A.4**

The probability of an outcome or event is its long-run relative frequency, that is, its relative frequency over a large number of trials.

**2.3.A.5**

The relative frequency of an outcome or event determined from empirical data can be used to estimate the actual, or true, probability of that outcome or event.

**2.3.A.6**

The law of large numbers states that for independent events, as the number of trials increases, the long-run relative frequency of the outcome or event gets closer and closer to a single value.

## TOPIC 2.4

# Introduction to Probability

Instructional Periods: 2

## SKILL

3.C

LEARNING  
OBJECTIVE

2.4.A

Calculate probabilities for events and their complements.

## ESSENTIAL KNOWLEDGE

2.4.A.1

The sample space of a random process is the set of all possible nonoverlapping outcomes. The probability of the sample space is 1.

2.4.A.2

If all outcomes in the sample space are equally likely, then the theoretical probability an event  $E$  will occur is:

$$\frac{\text{number of outcomes in event } E}{\text{total number of outcomes in the sample space}}.$$

The probability of event  $E$  occurring is written as  $P(E)$ .

2.4.A.3

The probability of an event is a number between 0 and 1, inclusive.

2.4.A.4

The probability of the complement of an event  $E$ , which can be written as  $E'$ ,  $\bar{E}$ , or  $E^C$  (i.e., the probability of "not  $E$ "), is equal to  $1 - P(E)$ .

## TOPIC 2.5

**Mutually Exclusive Events**

Instructional Periods: 2

**SKILL****4.B****LEARNING  
OBJECTIVE****2.5.A**

Justify why two events are mutually exclusive (or disjoint) using joint probability.

**ESSENTIAL KNOWLEDGE****2.5.A.1**

The probability that events  $A$  and  $B$  both will occur, sometimes called the joint probability, is the probability of the intersection of  $A$  and  $B$ . Joint probability is defined as  $P(A \cap B)$ .

**2.5.A.2**

Two events are mutually exclusive, or disjoint, if they cannot occur at the same time. This means that if two events are mutually exclusive, then  $P(A \text{ intersect } B) = 0$  or  $P(A \cap B) = 0$ .

## TOPIC 2.6

# Conditional Probability

Instructional Periods: 2

## SKILL

3.C

LEARNING  
OBJECTIVE

2.6.A

Calculate conditional probabilities.

## ESSENTIAL KNOWLEDGE

2.6.A.1

The probability that event  $A$  will occur given that event  $B$  has occurred is called a conditional probability and is written as  $P(A|B)$ . Conditional probability is

defined as 
$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$
.

2.6.A.2

The general multiplication rule states that the probability that events  $A$  and  $B$  both will occur is equal to the probability that event  $A$  will occur multiplied by the conditional probability that event  $B$  will occur given that  $A$  has occurred. The multiplication rule is defined as  $P(A \cap B) = P(A) \cdot P(B|A)$ .

## TOPIC 2.7

# Independent Events and Unions of Events

Instructional Periods: 2

**SKILL****3.C****LEARNING OBJECTIVE****2.7.A**

Calculate probabilities for independent events and for the union of two events.

**ESSENTIAL KNOWLEDGE****2.7.A.1**

Events  $A$  and  $B$  are independent if, and only if, knowing whether event  $A$  has occurred (or will occur) does not change the probability that event  $B$  will occur. When events  $A$  and  $B$  are independent, then  $P(A|B) = P(A)$ ,  $P(B|A) = P(B)$ , and  $P(A \cap B) = P(A) \cdot P(B)$ .

**2.7.A.2**

The probability that event  $A$  or event  $B$  (or both) will occur is the probability of  $A$  union  $B$ . The probability of the union is defined as  $P(A \cup B)$ .

**2.7.A.3**

The probability that  $P(A \text{ union } B) = P(A) + P(B) - P(A \text{ intersect } B)$ , or  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

## TOPIC 2.8

# Introduction to Random Variables and Probability Distributions

Instructional Periods: 2

**SKILL****3.A****LEARNING OBJECTIVE****2.8.A**

Construct a probability distribution for a discrete random variable.

**ESSENTIAL KNOWLEDGE****2.8.A.1**

A random variable is a variable whose values have numerical outcomes that result from a random phenomenon.

**2.8.A.2**

A probability distribution for a discrete random variable shows the probability associated with every possible value of the random variable. The sum of the probabilities over all possible values of a discrete random variable is 1.

**2.8.A.3**

A discrete probability distribution can be determined using the rules of probability or estimated with a simulation.

**2.8.A.4**

A discrete probability distribution can be represented as a graph, table, or function showing the probabilities associated with values of a random variable.

**2.8.A.5**

A cumulative probability distribution can be represented as a table or function and shows the probability of being less than or equal to each value of the discrete random variable.

## TOPIC 2.9

## Parameters of Random Variables

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.B	<p><b>2.9.A</b></p> <p>Calculate parameters for a discrete random variable.</p>	<p><b>2.9.A.1</b></p> <p>A numerical value measuring a characteristic of a probability distribution of a random variable, or a population, is a parameter. The value of a parameter is a single, fixed value.</p> <p><b>2.9.A.2</b></p> <p>The expected value (or mean) of a probability distribution is a parameter and is denoted by <math>E(X)</math> or <math>\mu_X</math>. For a discrete random variable <math>X</math>, the expected value is calculated as <math>\mu_X = \sum x_i \cdot P(x_i)</math>, where <math>x_i</math> is the possible value of the random variable and <math>P(x_i)</math> is the probability of the possible value of the random variable. The expected value can be interpreted as the long-run average outcome of the random variable. The discrete random variable can only take on values that are countable or finite.</p> <p><b>2.9.A.3</b></p> <p>The standard deviation of a probability distribution is a parameter represented by <math>SD(X)</math> or <math>\sigma_X</math>. For a discrete random variable <math>X</math>, the standard deviation is calculated as <math>\sigma_X = \sqrt{\sum (x_i - \mu_X)^2 \cdot P(x_i)}</math>, where <math>x_i</math> is the possible value of the random variable, <math>\mu_X</math> is the mean, and <math>P(x_i)</math> is the probability of the possible value of the random variable. The standard deviation can be interpreted as the typical deviation of the values of the random variable from the mean value (or expected value) of the random variable over the long run. The square of the standard deviation of a random variable is called the variance of the random variable and is denoted as <math>V(X)</math> or <math>\sigma_X^2</math>.</p>
4.D	<p><b>2.9.B</b></p> <p>Interpret parameters and probabilities for a discrete random variable.</p>	<p><b>2.9.B.1</b></p> <p>Parameters and probabilities for the probability distribution of a discrete random variable should be interpreted in the context of a specific population.</p>

## TOPIC 2.10

## The Binomial Distribution

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.B	<p><b>2.10.A</b></p> <p>Justify why a random variable is or is not a binomial random variable.</p>	<p><b>2.10.A.1</b></p> <p>A binomial random variable, <math>X</math>, is a discrete random variable that counts the number of successes in repeated independent trials that have only two possible outcomes (success or failure), with the probability of success <math>p</math> and the probability of failure <math>1 - p</math>.</p>
3.D	<p><b>2.10.B</b></p> <p>Calculate parameters for a binomial distribution.</p>	<p><b>2.10.B.1</b></p> <p>If a random variable is binomial, its mean, <math>\mu_x</math>, is <math>np</math> and its standard deviation, <math>\sigma_x</math>, is <math>\sqrt{np(1-p)}</math>.</p>
3.C	<p><b>2.10.C</b></p> <p>Calculate probabilities for a binomial distribution.</p>	<p><b>2.10.C.1</b></p> <p>The probability that a binomial random variable, <math>X</math>, has exactly <math>x</math> successes for <math>n</math> independent trials, when the probability of success is <math>p</math>, is calculated as <math>P(X = x) = \binom{n}{x} p^x (1-p)^{n-x}</math>, <math>x = 0, 1, 2, \dots, n</math>. This is called the binomial probability function.</p>
4.D	<p><b>2.10.D</b></p> <p>Interpret parameters and probabilities for a binomial distribution.</p>	<p><b>2.10.D.1</b></p> <p>Parameters and probabilities for a binomial distribution should be interpreted in context.</p>



## TOPIC 2.11

**The Normal Distribution**

Instructional Periods: 3

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.B	<p><b>2.11.A</b> Justify why a random variable is or is not a normal random variable.</p>	<p><b>2.11.A.1</b> A continuous random variable is a variable that can take on any value within a specified domain. Every interval within the domain has a probability associated with it.</p> <p><b>2.11.A.2</b> Many continuous random variables are well-modeled by a normal distribution.</p> <p><b>2.11.A.3</b> A normal distribution can be described as a continuous, unimodal, bell-shaped, and symmetric curve.</p>
4.C	<p><b>2.11.B</b> Describe a normal distribution.</p>	<p><b>2.11.B.1</b> The normal distribution, or the normal curve, is identified by two parameters, the mean, <math>\mu</math>, and the standard deviation, <math>\sigma</math>. The smaller the standard deviation, the taller and more concentrated the normal curve is around its mean. The larger the standard deviation, the shorter and less concentrated the normal curve is around its mean.</p>
3.D	<p><b>2.11.C</b> Calculate parameters for a normal distribution.</p>	<p><b>2.11.C.1</b> A standard normal distribution is a normal distribution with mean <math>\mu = 0</math> and standard deviation <math>\sigma = 1</math>.</p>
3.C	<p><b>2.11.D</b> Calculate percentages from a normal distribution using the empirical rule.</p>	<p><b>2.11.D.1</b> The empirical rule can be used to estimate the area of a region under the graph of the normal distribution curve. For a normal distribution, approximately 68% of the observations are within 1 standard deviation of the mean, approximately 95% of observations are within 2 standard deviations of the mean, and approximately 99.7% of observations are within 3 standard deviations of the mean. This is called the empirical rule, or the 68–95–99.7 rule.</p>

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.C	<p><b>2.11.E</b> Calculate the probability that a particular value lies within a given interval of a normal distribution.</p>	<p><b>2.11.E.1</b> If the distribution of a random variable is approximately normal, the probability that the random variable takes on values within a particular interval of the random variable is determined by the area under the normal curve within that interval. The total probability or area under the normal curve is 1.</p>
3.C	<p><b>2.11.F</b> Calculate the associated intervals and areas of a normal distribution.</p>	<p><b>2.11.F.1</b> The boundaries of an interval associated with a given area in a normal distribution can be determined using <math>z</math>-scores and a standard normal table or technology.</p> <p><b>2.11.F.2</b> Intervals associated with a given area in a normal distribution can be determined by assigning appropriate inequalities to the boundaries of the intervals. To determine the intervals, <math>p</math> is defined as a number between 0 and 1, <math>x_a</math> is the lower bound, and <math>x_b</math> is the upper bound on a normal distribution.</p> <p>i. <math>P(X &lt; x_a) = \frac{p}{100}</math> means that the lowest <math>p\%</math> of the values lie to the left of <math>x_a</math>.</p> <p>ii. <math>P(x_a &lt; X &lt; x_b) = \frac{p}{100}</math> means that <math>p\%</math> of the values lie between <math>x_a</math> and <math>x_b</math>.</p> <p>iii. <math>P(X &gt; x_b) = \frac{p}{100}</math> means that the highest <math>p\%</math> of the values lie to the left of <math>x_b</math>.</p> <p>iv. To determine the most extreme <math>p\%</math> of values on both sides requires dividing the area associated with <math>p\%</math> into two equal areas on either extreme of the distribution: <math>P(X &lt; x_a) = \frac{1}{2} \left( \frac{p}{100} \right)</math> and <math>P(X &gt; x_b) = \frac{1}{2} \left( \frac{p}{100} \right)</math> mean that half of the <math>p\%</math> most extreme values lie to the left of <math>x_a</math> and half of the <math>p\%</math> most extreme values lie to the right of <math>x_b</math>.</p>
4.C	<p><b>2.11.G</b> Compare measures of relative position for distributions.</p>	<p><b>2.11.G.1</b> Percentiles and proportions may be used to compare relative positions of individual values within a normal distribution or between normal distributions.</p>

## TOPIC 2.12

# The Central Limit Theorem

Instructional Periods: 3

**SKILL****4.C****LEARNING OBJECTIVE****2.12.A**

Describe sampling distributions.

**ESSENTIAL KNOWLEDGE****2.12.A.1**

A sampling distribution of a statistic is the distribution of values of the statistic for all possible samples of a given size from a given population.

**2.12.A.2**

The sampling distribution of a statistic can be simulated by repeatedly generating a large number of random samples from the population assuming known value(s) for the parameter(s). The value of the statistic is determined and recorded for each sample. The resulting distribution of the sample statistic values approximates the sampling distribution of the statistic.

**2.12.A.3**

A randomization distribution is the distribution of a statistic generated from repeatedly randomly reallocating, or reassigning, the response values to treatment groups. The value of the statistic is determined and recorded for each reallocation, or reassignment. The resulting distribution of the statistic values approximates the sampling distribution of the statistic.

**2.12.A.4**

The Central Limit Theorem (CLT) states that the sampling distribution of a mean of a random sample has a shape that can be approximated by a normal distribution. The larger the sample is, the better the approximation will be.

THIS PAGE IS INTENTIONALLY LEFT BLANK.

**AP STATISTICS**

**UNIT 3**

**Inference for  
Categorical  
Data:  
Proportions**

THIS PAGE IS INTENTIONALLY LEFT BLANK.

# TOPIC 3.1

## Estimators

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.B	<p><b>3.1.A</b> Justify why an estimator is or is not unbiased.</p>	<p><b>3.1.A.1</b> When estimating a population parameter, an estimator is unbiased if, on average, the value of the estimator does not underestimate or overestimate the population parameter.</p>
3.D	<p><b>3.1.B</b> Calculate estimates for a population parameter.</p>	<p><b>3.1.B.1</b> A sample statistic is a point estimator of the corresponding population parameter and can be thought of as the estimate of the population parameter. For example, the sample proportion <math>\hat{p}</math> is a point estimator for the population proportion <math>p</math>.</p>

TOPIC 3.2

# Sampling Distributions for Sample Proportions

Instructional Periods: 1

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.D	<p><b>3.2.A</b> Calculate parameters of a sampling distribution for a sample proportion.</p>	<p><b>3.2.A.1</b> For a population with population proportion <math>p</math>, when the sampled values are independent, the sampling distribution of a sample proportion <math>\hat{p}</math> has a mean <math>\mu_{\hat{p}} = p</math> and a standard deviation <math>\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}</math>.</p>
4.E	<p><b>3.2.B</b> Justify why a sampling distribution of a sample proportion can or cannot be described as approximately normal.</p>	<p><b>3.2.B.1</b> The sampling distribution of the sample proportion <math>\hat{p}</math> is approximately normal provided the sample size is large enough. To ensure the sample size is large enough, the following condition must be met: <math>np \geq 10</math> and <math>n(1-p) \geq 10</math>, where <math>np</math> is the number of successes and <math>n(1-p)</math> is the number of failures.</p>
4.D	<p><b>3.2.C</b> Interpret parameters and probabilities for a sampling distribution of a sample proportion.</p>	<p><b>3.2.C.1</b> Parameters and probabilities for a sampling distribution of a sample proportion should be interpreted in the context of a specific population.</p>



TOPIC 3.3

# Constructing a Confidence Interval for a Population Proportion

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.C	<p><b>3.3.A</b></p> <p>Identify an appropriate confidence interval procedure including the parameter for a population proportion.</p>	<p><b>3.3.A.1</b></p> <p>A confidence interval is an interval estimate for a population parameter. Based on the sample proportion, a confidence interval can be calculated to estimate the value of a single population proportion. The appropriate confidence interval procedure is a one-sample <math>z</math>-interval for a population proportion.</p> <p><b>3.3.A.2</b></p> <p>The parameter for a confidence interval for a population proportion should reference the proportion, the response variable, and the population in context.</p>
4.E	<p><b>3.3.B</b></p> <p>Justify the appropriateness of constructing a confidence interval for a population proportion by verifying conditions.</p>	<p><b>3.3.B.1</b></p> <p>A one-sample <math>z</math>-interval for a population proportion requires that three conditions be met as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: the data should be collected using a random sample.</li> <li>The 10% condition: when sampling without replacement, the population size must be at least 10 times larger than the sample size (<math>n &lt; 10\%N</math>), where <math>N</math> is the size of the population and <math>n</math> is the sample size.</li> <li>The normality condition: The number of successes, <math>n\hat{p}</math>, and the number of failures, <math>n(1-\hat{p})</math>, should be at least 10.</li> </ol>
3.E	<p><b>3.3.C</b></p> <p>Calculate an appropriate confidence interval for a population proportion.</p>	<p><b>3.3.C.1</b></p> <p>An interval estimate can be constructed as point estimate <math>\pm</math> (margin of error). For a population proportion, the one-sample <math>z</math>-interval estimate is</p> $\hat{p} \pm z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

**SKILL**

**3.E**

**LEARNING OBJECTIVE**

**3.3.D**

Calculate the standard error and margin of error of a sample statistic for a confidence interval for a population proportion.

**ESSENTIAL KNOWLEDGE**

**3.3.D.1**

$z^*$  denotes a “critical value”, such that  $-z^*$  and  $+z^*$  represent the boundaries enclosing the middle  $C\%$  of the standard normal distribution, in which  $C\%$  is an approximate confidence level with which the population proportion is estimated.

**3.3.D.2**

The standard error of a statistic is an estimate of the standard deviation of the sampling distribution of the statistic. The standard error of the sample

proportion  $\hat{p}$  is  $SE_{\hat{p}} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ .

**3.3.D.3**

The standard error quantifies the “typical” amount that a statistic will vary from the value of the corresponding population parameter.

**3.3.D.4**

The margin of error is half the width of the confidence interval and is calculated as the critical value ( $z^*$ ) times the standard error (SE) of  $\hat{p}$ , which

equals  $z^* \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ .

TOPIC 3.4

# Justifying a Claim Based on a Confidence Interval for a Population Proportion

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.F	<p><b>3.4.A</b></p> <p>Interpret a confidence interval in context for a population proportion.</p>	<p><b>3.4.A.1</b></p> <p>Since the confidence interval for a population proportion is calculated based on a sample from a population, the computed interval may or may not contain the value of the population proportion.</p> <p><b>3.4.A.2</b></p> <p>The interpretation of the confidence level is that in repeated random sampling with the same sample size, approximately <math>C\%</math> of confidence intervals calculated will capture the population proportion, with <math>C</math> representing the numerical value of the confidence level used.</p> <p><b>3.4.A.3</b></p> <p>When interpreting a <math>C\%</math> confidence interval for a population proportion, we say we are <math>C\%</math> confident that the interval <math>(a, b)</math> contains the true value of the parameter for the population. An interpretation of a confidence interval for a population proportion includes a reference to the parameter with details about the population it represents in the context of the study.</p>
4.G	<p><b>3.4.B</b></p> <p>Justify a claim based on a confidence interval for a population proportion.</p>	<p><b>3.4.B.1</b></p> <p>A confidence interval for a population proportion provides a range of plausible values that may serve as convincing evidence to support a particular claim about the population proportion.</p>
2.D	<p><b>3.4.C</b></p> <p>Identify the relationships among sample size, confidence interval width, confidence level, and margin of error for a population proportion.</p>	<p><b>3.4.C.1</b></p> <p>For a given sample, increasing the confidence level will result in the following:</p> <ol style="list-style-type: none"> <li>The critical value will increase.</li> <li>The margin of error will increase.</li> <li>The width of the confidence interval will increase.</li> </ol> <p><b>3.4.C.2</b></p> <p>Increasing the sample size decreases the standard error. Thus, when all other things remain the same, the width of the confidence interval for a population proportion tends to decrease as the sample size increases. For a confidence interval for a population proportion with a given confidence level, the width of the interval is approximately proportional to <math>\frac{1}{\sqrt{n}}</math>.</p>

TOPIC 3.5

# Setting Up a Test for a Population Proportion

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.C	<p><b>3.5.A</b> Identify an appropriate testing method for a population proportion including the parameter for the population proportion.</p>	<p><b>3.5.A.1</b> A hypothesis test is a statistical inference procedure that is used to make a decision about the value of a population parameter. The appropriate hypothesis testing procedure is a one-sample <math>z</math>-test for a population proportion.</p> <p><b>3.5.A.2</b> The parameter for a hypothesis test for a population proportion should reference the proportion, the response variable, and the population in context.</p>
2.E	<p><b>3.5.B</b> Identify the null and alternative hypotheses for a population proportion.</p>	<p><b>3.5.B.1</b> In the hypothesis testing procedure, the null hypothesis, <math>H_0</math>, is the statement about a parameter that is correct unless evidence suggests otherwise. It is the status quo condition. The alternative hypothesis, <math>H_a</math>, is the claim or belief about a parameter for which evidence is being collected. A researcher's claim or belief about the population parameter is represented by the alternative hypothesis.</p> <p><b>3.5.B.2</b> The null hypothesis contains an equality reference (<math>=</math>, <math>\geq</math>, or <math>\leq</math>). Although the null hypothesis for a one-sided test may include an inequality symbol, in AP Statistics, it is tested at the boundary of equality. The alternative hypothesis with <math>&lt;</math> or <math>&gt;</math> is called one-sided, and the alternative hypothesis with <math>\neq</math> is called two-sided.</p> <p><b>3.5.B.3</b> The null hypothesis for a one-sample <math>z</math>-test for a population proportion is as follows: <math>H_0 : p = p_0</math>, where <math>p_0</math> is the null hypothesized value for the population proportion. A one-sided alternative hypothesis for a one-sample <math>z</math>-test for a population proportion is either <math>H_a : p &lt; p_0</math> or <math>H_a : p &gt; p_0</math>. A two-sided alternative hypothesis is <math>H_a : p \neq p_0</math>.</p>
4.E	<p><b>3.5.C</b> Justify the appropriateness of a hypothesis test for a population proportion by verifying conditions.</p>	<p><b>3.5.C.1</b> A one-sample <math>z</math>-test for a population proportion requires that three conditions be met as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: The data should be collected using a random sample.</li> <li>The 10% condition: when sampling without replacement, the population size must be at least 10 times larger than the sample size (<math>n &lt; 10\%N</math>), where <math>N</math> is the size of the population and <math>n</math> is the sample size.</li> <li>The normality condition: The number of successes, <math>np_0</math>, and the number of failures, <math>n(1 - p_0)</math>, should be at least 10.</li> </ol>

## TOPIC 3.6

# *p*-Values

Instructional Periods: 2

**SKILL****4.F****LEARNING OBJECTIVE****3.6.A**

Interpret the *p*-value of a hypothesis test for a population proportion.

**ESSENTIAL KNOWLEDGE****3.6.A.1**

The *p*-value is the probability of obtaining a test statistic as extreme or more extreme (i.e., in the direction of the alternative hypothesis) than the test statistic that is observed given that the null hypothesis is true. That is, when  $x$  is the test statistic, the *p*-value is determined by finding the following:

- The probability at or above the observed value of the test statistic ( $P(z \geq x)$ ), if the alternative is  $>$
- The probability at or below the observed value of the test statistic ( $P(z \leq x)$ ), if the alternative is  $<$
- The probability less than or equal to the negative of the absolute value of the test statistic plus the probability greater than or equal to the absolute value of the test statistic, ( $P(z \leq -|x|) + P(z \geq |x|)$ ), if the alternative is  $\neq$

**3.6.A.2**

If the distribution of the test statistic has been simulated, the *p*-value is the proportion of values in the null distribution that are as extreme or more extreme than the observed value of the test statistic. This is as follows:

- The proportion at or above the observed value of the test statistic, if the alternative is  $>$
- The proportion at or below the observed value of the test statistic, if the alternative is  $<$
- The proportion less than or equal to the negative of the absolute value of the test statistic plus the proportion greater than or equal to the absolute value of the test statistic, if the alternative is  $\neq$

**3.6.A.3**

An interpretation of the *p*-value of a hypothesis test for a population proportion should include a statement that the *p*-value is computed by assuming the null hypothesis is true (i.e., by assuming the true population proportion is equal to the particular value stated in the null hypothesis in context).

**3.6.A.4**

Small *p*-values indicate that the observed value of the test statistic would be unusual if the null hypothesis were true and therefore provide evidence for the alternative hypothesis. The lower the *p*-value, the more convincing the statistical evidence for the alternative hypothesis.

**3.6.A.5**

*p*-values that are not small indicate that the observed value of the test statistic would not be unusual if the null hypothesis were true and therefore do not provide convincing statistical evidence for the alternative hypothesis, nor do they provide evidence that the null hypothesis is true.

## TOPIC 3.7

# Carrying Out a Test for a Population Proportion

Instructional Periods: 2

## SKILL

3.E

## LEARNING OBJECTIVE

3.7.A

Calculate an appropriate test statistic and  $p$ -value for testing a hypothesis about a population proportion.

## ESSENTIAL KNOWLEDGE

3.7.A.1

The test statistic for testing a population proportion is  $z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$ .

The  $z$ -statistic has a standard normal distribution when the null hypothesis is true.

3.7.A.2

The distribution of the test statistic assuming the null hypothesis is true (null distribution) can be approximated by a probability model (e.g., a theoretical distribution such as the standard normal distribution).

3.7.A.3

The  $p$ -value of a one-sample  $z$ -test for a population proportion is found from the standard normal distribution table or from the standard normal distribution using technology.

## SKILL

4.G

## LEARNING OBJECTIVE

3.7.B

Justify a claim about the population based on the results of a hypothesis test for a population proportion.

## ESSENTIAL KNOWLEDGE

3.7.B.1

The significance level of a hypothesis test, denoted by  $\alpha$ , is the predetermined probability of rejecting the null hypothesis given that it is true. The significance level may be given or determined by the researcher.

3.7.B.2

A formal decision in a hypothesis test explicitly compares the  $p$ -value to the significance level,  $\alpha$ . If the  $p$ -value  $\leq \alpha$  then reject the null hypothesis,  $H_0 : p = p_0$ . If the  $p$ -value  $> \alpha$  then fail to reject the null hypothesis.

3.7.B.3

Rejecting the null hypothesis means there is convincing statistical evidence to support the alternative hypothesis. Failing to reject the null hypothesis means there is not convincing statistical evidence to support the alternative hypothesis.

3.7.B.4

A hypothesis test can lead to rejecting or not rejecting the null hypothesis but can never lead to concluding or proving that the null hypothesis is true. Lack of statistical evidence for the alternative hypothesis is not the same as evidence for the null hypothesis.

3.7.B.5

The results of a hypothesis test for a population proportion can serve as the statistical reasoning to support the answer to a research question about the population that was sampled.

3.7.B.6

A conclusion for the hypothesis test for a population proportion is stated in context consistent with, and in terms of, the alternative hypothesis using non-definitive language. The conclusion should contain a reference to the parameter and the population.

TOPIC 3.8

# Sampling Distributions for the Difference Between Sample Proportions

Instructional Periods: 1

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.D	<p><b>3.8.A</b> Calculate parameters of a sampling distribution for the difference between two sample proportions.</p>	<p><b>3.8.A.1</b> For two independent populations, with population proportions <math>p_1</math> and <math>p_2</math>, when the sampled values are independent, the sampling distribution for the difference in sample proportions, <math>\hat{p}_1 - \hat{p}_2</math>, has a mean, <math>\mu_{\hat{p}_1 - \hat{p}_2} = p_1 - p_2</math>, and standard deviation, <math>\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}</math>.</p> <p><b>3.8.A.2</b> When sampling without replacement, two conditions must be met as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: the data should be collected using two independent random samples.</li> <li>The 10% condition: both samples must be less than 10% of the size of their respective populations.</li> </ol> <p><b>3.8.A.3</b> If the data comes from an experiment, the data only needs to meet the randomization condition. The treatments must be randomly assigned to participants or experimental units to meet the randomization condition.</p>
4.E	<p><b>3.8.B</b> Justify why a sampling distribution for the difference between sample proportions can or cannot be described as approximately normal.</p>	<p><b>3.8.B.1</b> The sampling distribution for the difference between sample proportions, <math>\hat{p}_1 - \hat{p}_2</math>, will have an approximately normal distribution provided both sample sizes are large enough. To ensure that both samples are large enough, the data must meet the following conditions: <math>n_1 p_1 \geq 10</math>, <math>n_1(1-p_1) \geq 10</math>, <math>n_2 p_2 \geq 10</math>, and <math>n_2(1-p_2) \geq 10</math>, where <math>n_1 p_1</math> and <math>n_2 p_2</math> are the number of successes and <math>n_1(1-p_1)</math> and <math>n_2(1-p_2)</math> are the number of failures.</p>
4.D	<p><b>3.8.C</b> Interpret parameters and probabilities for a sampling distribution for the difference between proportions.</p>	<p><b>3.8.C.1</b> Parameters and probabilities for a sampling distribution of a sample proportion should be interpreted within the context of a specific population.</p>



## TOPIC 3.9

# Constructing a Confidence Interval for the Difference Between Two Population Proportions

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.C	<p><b>3.9.A</b></p> <p>Identify an appropriate confidence interval procedure including the parameters for the difference between two population proportions.</p>	<p><b>3.9.A.1</b></p> <p>Based on the sample data, a confidence interval can be calculated to estimate the difference between two independent population proportions. The appropriate confidence interval procedure is a two-sample z-interval for a difference between population proportions.</p> <p><b>3.9.A.2</b></p> <p>The parameters of a confidence interval for the difference between two population proportions should refer to the difference in the proportions, the response variable, and the populations in context.</p>
4.E	<p><b>3.9.B</b></p> <p>Justify the appropriateness of constructing a confidence interval for the difference between two population proportions by verifying conditions.</p>	<p><b>3.9.B.1</b></p> <p>A two-sample z-interval for a difference between two population proportions requires three conditions to be met as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: the data should be collected using two independent random samples or a randomized experiment.</li> <li>The 10% condition: When sampling without replacement, check that <math>n_1 \leq 10\%N_1</math> and <math>n_2 \leq 10\%N_2</math>, where <math>N_1</math> is the size of population 1 and <math>N_2</math> is the size of population 2. The sample sizes are represented as <math>n_1</math> and <math>n_2</math>. (Note: This condition is not necessary when the data are from a randomized experiment.)</li> <li>The normality condition: the number of successes, <math>n_1\hat{p}_1</math> and <math>n_2\hat{p}_2</math>, and number of failures, <math>n_1(1-\hat{p}_1)</math> and <math>n_2(1-\hat{p}_2)</math>, for both samples are all at least 10.</li> </ol>

**SKILL**

**LEARNING OBJECTIVE**

**ESSENTIAL KNOWLEDGE**

3.E

3.9.C

Calculate an appropriate confidence interval for the difference between two population proportions.

3.9.C.1

The point estimate for the difference between two population proportions  $(\hat{p}_1 - \hat{p}_2)$  is the difference in sample proportions.

3.9.C.2

For the difference between two population proportions, the interval estimate can be constructed as point estimate  $\pm$  (margin of error). The interval estimate for the difference between two population proportions is

$$(\hat{p}_1 - \hat{p}_2) \pm z^* \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$$

3.E

3.9.D

Calculate the standard error and margin of error for the difference between two population proportions.

3.9.D.1

The standard error for the difference between two population proportions is

$$\sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$$

3.9.D.2

For the difference between two population proportions, the margin of error is the critical value ( $z^*$ ) times the standard error (SE) of the difference

between the two proportions, which equals  $z^* \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$ .

TOPIC 3.10

# Justifying a Claim Based on a Confidence Interval for the Difference Between Two Population Proportions

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.F	<p><b>3.10.A</b></p> <p>Interpret a confidence interval in context for the difference between two population proportions.</p>	<p><b>3.10.A.1</b></p> <p>Since the confidence interval for the difference between two population proportions is calculated based on samples from two populations, the computed interval may or may not contain the value for the difference between those two population proportions.</p> <p><b>3.10.A.2</b></p> <p>The interpretation of the confidence level is as follows: In repeated random sampling with the same sample sizes, approximately <math>C\%</math> of confidence intervals created will capture the difference between the two population proportions, where <math>C</math> represents the numerical value of the confidence level used.</p> <p><b>3.10.A.3</b></p> <p>When interpreting a <math>C\%</math> confidence interval for the difference between two population proportions, we say we are <math>C\%</math> confident that the interval <math>(a, b)</math> contains the parameter for the population. An interpretation of a confidence interval for the difference between two population proportions includes a reference to the parameter with the details about the populations it represents in the context of the study.</p>
4.G	<p><b>3.10.B</b></p> <p>Justify a claim based on a confidence interval for the difference between two population proportions.</p>	<p><b>3.10.B.1</b></p> <p>A confidence interval for the difference between two population proportions provides an interval of values that may provide convincing evidence to support a particular claim about the difference between the two population proportions. For example, if the interval contains 0, then there is insufficient evidence to conclude there is a difference between the two population proportions. If the interval does not contain 0, there is sufficient evidence to conclude there is a difference between the two population proportions.</p>

TOPIC 3.11

# Setting Up a Test for the Difference Between Two Population Proportions

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.C	<p><b>3.11.A</b> Identify an appropriate testing method for the difference between population proportions including the parameters.</p>	<p><b>3.11.A.1</b> The appropriate testing method for the difference between two population proportions is a two-sample <math>z</math>-test for the difference between two population proportions.</p> <p><b>3.11.A.2</b> The parameters for a hypothesis test for the difference between two population proportions should refer to the difference in the proportions, the response variable, and the populations in context.</p>
2.E	<p><b>3.11.B</b> Identify the null and alternative hypotheses for the difference between population proportions.</p>	<p><b>3.11.B.1</b> For a two-sample <math>z</math>-test for the difference between two population proportions, the null hypothesis indicates no difference. The null hypothesis for the difference between two population proportions can be written as either <math>H_0 : p_1 = p_2</math> or <math>H_0 : p_1 - p_2 = 0</math>. A one-sided alternative hypothesis for the difference between two population proportions can be written as either:  <math>H_a : p_1 &lt; p_2</math>, or equivalently <math>H_a : p_1 - p_2 &lt; 0</math>, or  <math>H_a : p_1 &gt; p_2</math>, or equivalently <math>H_a : p_1 - p_2 &gt; 0</math>.                      A two-sided alternative hypothesis for the difference between two population proportions can be written as either <math>H_a : p_1 \neq p_2</math>, or equivalently <math>H_a : p_1 - p_2 \neq 0</math>.</p>

**SKILL**

**4.E**

**LEARNING OBJECTIVE**

**3.11.C**

Justify the appropriateness of a hypothesis test for the difference between two population proportions by verifying conditions.

**ESSENTIAL KNOWLEDGE**

**3.11.C.1**

A two-sample z-test for a difference between two population proportions requires three conditions to be met as follows:

- i. The randomization condition: the data should be collected using two independent random samples or a randomized experiment.
- ii. The 10% condition: When sampling without replacement, check that  $n_1 \leq 10\%N_1$  and  $n_2 \leq 10\%N_2$ , where  $N_1$  is the size of population 1 and  $N_2$  is the size of population 2. The sample sizes are represented as  $n_1$  and  $n_2$ . (Note: This condition is unnecessary when the data are from a randomized experiment.)
- iii. The normality condition:  $n_1\hat{p}_c$ ,  $n_1(1-\hat{p}_c)$ ,  $n_2\hat{p}_c$ , and  $n_2(1-\hat{p}_c)$  must all be at least 10, with  $\hat{p}_c = \frac{n_1\hat{p}_1 + n_2\hat{p}_2}{n_1 + n_2}$  being the combined (or pooled) proportion assuming that  $H_0$  is true ( $H_0 : p_1 = p_2$  or  $H_0 : p_1 - p_2 = 0$ ).

## TOPIC 3.12

# Carrying Out a Test for the Difference Between Two Population Proportions

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.E	<b>3.12.A</b> Calculate an appropriate test statistic and $p$ -value for testing a hypothesis for the difference between two population proportions.	<b>3.12.A.1</b> The test statistic for the difference between two population proportions is as follows: $z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}_c(1 - \hat{p}_c)} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$ , where $\hat{p}_c = \frac{n_1\hat{p}_1 + n_2\hat{p}_2}{n_1 + n_2}$ is the proportion of successes for the two groups combined. The $z$ -statistic has a standard normal distribution when the null hypothesis is true. <b>3.12.A.2</b> The $p$ -value for a two-sample $z$ -test for the difference between two population proportions can be found using the standard normal table or from the standard normal distribution using technology.
4.F	<b>3.12.B</b> Interpret the $p$ -value of a hypothesis test for the difference between two population proportions.	<b>3.12.B.1</b> The $p$ -value is the probability of obtaining a test statistic as extreme or more extreme than the test statistic that was observed (i.e., in the direction of the alternative hypothesis) given that the null hypothesis is true. An interpretation of the $p$ -value of a hypothesis test for a difference between two population proportions should include a statement that the $p$ -value is computed assuming the null hypothesis is true (i.e., by assuming that the true population proportions are equal to each other in context).

## SKILL

4.G

## LEARNING OBJECTIVE

3.12.C

Justify a claim about the population based on the results of a hypothesis test for the difference between two population proportions.

## ESSENTIAL KNOWLEDGE

3.12.C.1

A formal decision in a hypothesis test for the difference between two population proportions explicitly compares the  $p$ -value to the significance level,  $\alpha$ . If the  $p$ -value  $\leq \alpha$  then reject the null hypothesis,  $H_0 : p_1 = p_2$  or  $H_0 : p_1 - p_2 = 0$ . If the  $p$ -value  $> \alpha$  then fail to reject the null hypothesis.

3.12.C.2

The results of a hypothesis test for the difference between two population proportions can serve as the statistical reasoning to support the answer to a research question about the two populations that were sampled.

3.12.C.3

A conclusion for the hypothesis test for the difference between two population proportions is stated in context consistent with, and in terms of, the alternative hypothesis using non-definitive language. The conclusion should contain a reference to the parameters and the populations.

TOPIC 3.13

# Setting Up a Chi-Square Test for Homogeneity or Independence

Instructional Periods: 3

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.C	<p><b>3.13.A</b> Describe chi-square distributions.</p>	<p><b>3.13.A.1</b> The chi-square statistic measures the distance between observed and expected counts relative to expected counts.</p> <p><b>3.13.A.2</b> Chi-square distributions have positive values and are skewed right. Within this family of density curves, the skew becomes less pronounced with increasing degrees of freedom.</p>
2.C	<p><b>3.13.B</b> Identify an appropriate testing method for comparing distributions in two-way tables of categorical data.</p>	<p><b>3.13.B.1</b> To determine whether the distributions of a categorical variable for two or more populations are different, the appropriate test is the chi-square test for homogeneity.</p> <p><b>3.13.B.2</b> To determine whether row and column variables in a two-way table of categorical data might be associated in the single population from which the data were sampled, the appropriate test is the chi-square test for independence.</p>
2.E	<p><b>3.13.C</b> Identify the null and alternative hypotheses for a chi-square test for homogeneity or independence.</p>	<p><b>3.13.C.1</b> The appropriate null hypothesis for a chi-square test for homogeneity is <math>H_0</math>: there is no difference in the distributions of the categorical variable across populations or treatments. The appropriate alternative hypothesis for a chi-square test for homogeneity is <math>H_a</math>: there is a difference in the distributions of the categorical variable across populations or treatments.</p> <p><b>3.13.C.2</b> The appropriate null hypothesis for a chi-square test for independence is <math>H_0</math>: there is no association between two categorical variables in a given population or the two categorical variables in a given population are independent of each other. The appropriate alternative hypothesis for a chi-square test for independence is <math>H_a</math>: there is an association between two categorical variables in a given population or the two categorical variables in a given population are dependent on each other.</p>



## SKILL

4.E

## LEARNING OBJECTIVE

3.13.D

Justify the appropriateness of a hypothesis test for a chi-square distribution for independence or homogeneity by verifying conditions.

## ESSENTIAL KNOWLEDGE

3.13.D.1

For a chi-square test for homogeneity or independence, three conditions must be met as follows:

- i. The randomization condition: the test of independence states that the data should be collected using an independent random sample. The test for homogeneity states that the data should be collected using independent random samples or a randomized experiment.
- ii. The 10% condition: When sampling without replacement, check that  $n < 10\%N$ , where  $N$  is the size of the population and  $n$  is the sample size. (Note: this condition is unnecessary when the data are from a randomized experiment.)
- iii. The expected values condition: all expected values should be greater than 5.

TOPIC 3.14

# Carrying Out a Chi-Square Test for Homogeneity or Independence

Instructional Periods: 3

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.C	<p><b>3.14.A</b> Calculate expected counts for two-way tables of categorical data.</p>	<p><b>3.14.A.1</b> The expected values (under the null hypothesis) in a particular cell of a two-way table of categorical data can be calculated using the following formula:  <math display="block">\text{expected value} = \frac{(\text{row total})(\text{column total})}{\text{total table}}</math></p>
3.E	<p><b>3.14.B</b> Calculate the appropriate test statistic and <math>p</math>-value for a chi-square test for homogeneity or independence.</p>	<p><b>3.14.B.1</b> The appropriate test statistic for a chi-square test for homogeneity or independence is the chi-square statistic:  <math display="block">\chi^2 = \sum \frac{(\text{Observed value} - \text{Expected value})^2}{\text{Expected value}}</math>                     , where the sum is taken over all cells of the two-way table. The chi-square statistics have a chi-square distribution with degrees of freedom equal to <math>(\text{number of rows} - 1) \cdot (\text{number of columns} - 1)</math> when the null hypothesis is true.</p> <p><b>3.14.B.2</b> The <math>p</math>-value for a chi-square test for independence or homogeneity is found using a chi-square critical values table or from a chi-square distribution using technology.</p>
4.F	<p><b>3.14.C</b> Interpret the <math>p</math>-value for the chi-square test for homogeneity or independence.</p>	<p><b>3.14.C.1</b> The <math>p</math>-value is the probability of obtaining a test statistic as extreme or more extreme than the test statistic that was observed (i.e., in the direction of the alternative hypothesis) given that the null hypothesis is true. An interpretation of the <math>p</math>-value for the chi-square test for homogeneity or independence should include a statement that the <math>p</math>-value is computed by assuming the null hypothesis is true.</p>

## SKILL

4.G

## LEARNING OBJECTIVE

3.14.D

Justify a claim about the population based on the results of a chi-square test for homogeneity or independence.

## ESSENTIAL KNOWLEDGE

3.14.D.1

A formal decision in a hypothesis test explicitly compares the  $p$ -value to the significance level,  $\alpha$ . If the  $p$ -value  $\leq \alpha$  then reject the null hypothesis for the appropriate chi-square test. If the  $p$ -value  $> \alpha$  then fail to reject the null hypothesis.

3.14.D.2

The results of a chi-square test for homogeneity or independence can serve as the statistical reasoning to support the answer to a research question about the population that was sampled (independence) or the populations that were sampled (homogeneity).

3.14.D.3

A conclusion for a chi-square test for homogeneity or independence is stated in context consistent with, and in terms of, the alternative hypothesis using non-definitive language. The conclusion should contain a reference to the population(s).

THIS PAGE IS INTENTIONALLY LEFT BLANK.

**AP STATISTICS**

**UNIT 4**

**Inference for  
Quantitative  
Data: Means**

THIS PAGE IS INTENTIONALLY LEFT BLANK.

TOPIC 4.1

# Sampling Distributions for Sample Means

Instructional Periods: 1

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.D	<p><b>4.1.A</b></p> <p>Calculate parameters of a sampling distribution of a sample mean.</p>	<p><b>4.1.A.1</b></p> <p>For a population with population mean <math>\mu</math> and population standard deviation <math>\sigma</math>, when the sampled values are independent, the sampling distribution of the sample mean has mean <math>\mu_{\bar{x}} = \mu</math> and standard deviation <math>\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}</math>.</p> <p><b>4.1.A.2</b></p> <p>When sampling without replacement, two conditions must be met as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: the data should be collected using a random sample.</li> <li>The 10% condition: The population size must be at least 10 times larger than the sample size (<math>n &lt; 10\%N</math>), where <math>N</math> is the size of the population and <math>n</math> is the sample size.</li> </ol>
4.E	<p><b>4.1.B</b></p> <p>Justify why the sampling distribution of a sample mean can or cannot be described as approximately normal.</p>	<p><b>4.1.B.1</b></p> <p>For a quantitative variable, if the population distribution can be modeled with a normal distribution, the sampling distribution of the sample mean, <math>\bar{x}</math>, can be modeled with a normal distribution regardless of the sample size.</p> <p><b>4.1.B.2</b></p> <p>For a quantitative variable, if the population distribution is not normal, the sampling distribution of the sample mean, <math>\bar{x}</math>, can be modeled approximately by a normal distribution, provided <math>n \geq 30</math>. If the population distribution is extremely skewed, a sample size much larger than 30 may be needed to ensure the sampling distribution is approximately normal.</p>
4.D	<p><b>4.1.C</b></p> <p>Interpret parameters and probabilities for the sampling distribution of a sample mean.</p>	<p><b>4.1.C.1</b></p> <p>Parameters and probabilities for a sampling distribution of a sample mean should be interpreted within the context of a specific population.</p>

TOPIC 4.2

# Constructing a Confidence Interval for a Population Mean

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.C	<p><b>4.2.A</b> Describe <math>t</math>-distributions.</p>	<p><b>4.2.A.1</b> <math>t</math>-distributions, also called Student's <math>t</math>-distributions, make up a family of symmetric, bell-shaped, standardized distributions with wider tails than that of the standard normal distribution. Specific <math>t</math>-distributions are identified using a parameter known as the number of degrees of freedom (<math>df</math>), which is based on the sample size(s). When the degrees of freedom are small, the <math>t</math>-distribution has a much narrower peak and fatter tails than a normal distribution. As the degrees of freedom increase, the <math>t</math>-distribution more closely resembles the standard normal distribution (mean <math>\mu = 0</math> and standard deviation <math>\sigma = 1</math>).</p> <p><b>4.2.A.2</b> <math>t</math>-distributions are used for finding critical values and test statistics for inferences about a population mean, <math>\mu</math>, when the population standard deviation, <math>\sigma</math>, is unknown and the sample standard deviation, <math>s</math>, must be used instead.</p>
2.C	<p><b>4.2.B</b> Identify an appropriate confidence interval procedure including the parameter for a population mean, including matched pairs.</p>	<p><b>4.2.B.1</b> The appropriate confidence interval procedure for estimating the population mean of a quantitative variable for one sample is a one-sample <math>t</math>-interval for a population mean. (The population standard deviation, <math>\sigma</math>, is not typically known for distributions for quantitative variables).</p> <p><b>4.2.B.2</b> For a matched pairs design with two dependent samples, the appropriate analysis calculates differences between pairs of values to produce one sample of differences. The confidence interval procedure for the matched pairs design is a one-sample <math>t</math>-interval for a population mean difference.</p> <p><b>4.2.B.3</b> The parameter for a confidence interval for a population mean, including matched pairs, should reference the population mean or population mean difference and the response variable, in context. For the population mean difference, it is important to state the order of subtraction for the difference.</p>



SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.E	<p><b>4.2.C</b> Justify the appropriateness of constructing a confidence interval for a population mean, including the mean difference between values in matched pairs, by verifying conditions.</p>	<p><b>4.2.C.1</b> A one-sample <math>t</math>-interval for a population mean or population mean difference requires that three conditions be met as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: the data should be collected using a random sample or a randomized experiment.</li> <li>The 10% condition: when sampling without replacement, check that <math>n &lt; 10\%N</math>, where <math>N</math> is the size of the population and <math>n</math> is the sample size.</li> <li>The sample data condition: <math>n \geq 30</math>, or if <math>n &lt; 30</math>, the sample data distribution should be free from strong skewness and outliers. For matched pairs, the number of differences should be greater than or equal to 30. If the number of differences is less than 30, the sample of differences should be free from strong skewness and outliers.</li> </ol>
3.E	<p><b>4.2.D</b> Calculate an appropriate confidence interval for a population mean, including the mean difference between values in matched pairs.</p>	<p><b>4.2.D.1</b> A point estimate for a population mean is the sample mean, <math>\bar{x}</math>, or <math>\bar{x}_d</math> for the sample mean difference.</p> <p><b>4.2.D.2</b> To estimate the population mean for one sample or the population mean difference between values in matched pairs, when the population standard deviation is unknown, the confidence interval is <math>\bar{x} \pm t^* \frac{s}{\sqrt{n}}</math>, where <math>\pm t^*</math> is the critical value for the central C% of a <math>t</math>-distribution with degrees of freedom, <math>n - 1</math>.</p>
3.E	<p><b>4.2.E</b> Calculate the standard error and margin of error for a sample size for a one-sample <math>t</math>-interval for a population mean.</p>	<p><b>4.2.E.1</b> The standard error for a sample mean is given by <math>s_{\bar{x}} = \frac{s}{\sqrt{n}}</math>.</p> <p><b>4.2.E.2</b> For a one-sample <math>t</math>-interval for a population mean, the margin of error is the critical value (<math>t^*</math>) times the standard error (SE), which equals <math>t^* \left( \frac{s}{\sqrt{n}} \right)</math>.</p>

## TOPIC 4.3

# Justifying a Claim Based on a Confidence Interval for a Population Mean

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.F	<b>4.3.A</b> Interpret a confidence interval in context for a population mean, including the mean difference between values in matched pairs.	<b>4.3.A.1</b> Since the confidence interval for a population mean or population mean difference is calculated based on a sample from a population, the computed interval may or may not contain the value of the population mean or population mean difference. <b>4.3.A.2</b> The interpretation of the confidence level is as follows: In repeated random sampling with the same sample size, approximately $C\%$ of confidence intervals created will capture the population mean or population mean difference, where $C$ represents the numerical value of the confidence level used. <b>4.3.A.3</b> When interpreting a $C\%$ confidence interval for a population mean or population mean difference, we say we are $C\%$ confident the interval $(a, b)$ contains the value of the population mean or population mean difference. An interpretation of a confidence interval for a population mean includes a reference to the parameter.
4.G	<b>4.3.B</b> Justify a claim based on a confidence interval for a population mean, including the mean difference between values in matched pairs.	<b>4.3.B.1</b> A confidence interval for a population mean or population mean difference provides an interval of values that may serve as convincing evidence to support a particular claim about the population mean.

**SKILL**

**2.D**

**LEARNING OBJECTIVE**

**4.3.C**

Identify the relationships among sample size, confidence interval width, confidence level, and margin of error for a population mean.

**ESSENTIAL KNOWLEDGE**

**4.3.C.1**

For a given sample, increasing the confidence level will result in the following:

- i. The critical value will increase.
- ii. The margin of error will increase.
- iii. The width of the confidence interval will increase.

**4.3.C.2**

Increasing the sample size decreases the standard error. Thus, when all other things remain the same, the width of a confidence interval for a population mean tends to decrease as the sample size increases. For a confidence interval for a population mean with a given confidence level, the width of the interval is approximately proportional to  $\frac{1}{\sqrt{n}}$ .

TOPIC 4.4

# Setting Up a Test for a Population Mean

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.C	<p><b>4.4.A</b></p> <p>Identify an appropriate testing method and parameter for a population mean with unknown <math>\sigma</math>, including the mean difference between values in matched pairs.</p>	<p><b>4.4.A.1</b></p> <p>The appropriate test for a population mean with unknown population standard deviation <math>\sigma</math> is a one-sample <math>t</math>-test for a population mean.</p> <p><b>4.4.A.2</b></p> <p>For a matched pairs design with two dependent samples, the appropriate analysis calculates differences between pairs of values to produce one sample of differences. The hypothesis testing procedure for the matched pairs design is a one-sample <math>t</math>-test for the mean difference.</p> <p><b>4.4.A.3</b></p> <p>The parameter for a hypothesis test for a population mean, including matched pairs, should reference the population mean or population mean difference and the response variable in context.</p>
2.E	<p><b>4.4.B</b></p> <p>Identify the null and alternative hypotheses for a population mean with unknown <math>\sigma</math>, including the mean difference between values in matched pairs.</p>	<p><b>4.4.B.1</b></p> <p>The null hypothesis for a one-sample <math>t</math>-test for a population mean is <math>H_0 : \mu = \mu_0</math>, in which <math>\mu_0</math> is the null hypothesized value for the population mean. A one-sided alternative hypothesis for a one-sample <math>t</math>-test for a population mean is either <math>H_a : \mu &lt; \mu_0</math> or <math>H_a : \mu &gt; \mu_0</math>. A two-sided alternative hypothesis is <math>H_a : \mu \neq \mu_0</math>.</p> <p><b>4.4.B.2</b></p> <p>The null hypothesis for a matched pairs design is <math>H_0 : \mu_d = 0</math>. A one-sided alternative hypothesis for a matched pair design is either <math>H_a : \mu_d &lt; 0</math> or <math>H_a : \mu_d &gt; 0</math>. A two-sided alternative hypothesis is <math>H_a : \mu_d \neq 0</math>.</p>
4.E	<p><b>4.4.C</b></p> <p>Justify the appropriateness of a hypothesis test for a population mean by verifying conditions.</p>	<p><b>4.4.C.1</b></p> <p>A one-sample <math>t</math>-test for a population mean requires that three conditions be met as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: the data should be collected using a random sample or a randomized experiment.</li> <li>The 10% condition: when sampling without replacement, check that <math>n &lt; 10\%N</math>, where <math>N</math> is the size of the population and <math>n</math> is the sample size.</li> <li>The sample data condition: <math>n \geq 30</math>, or if <math>n &lt; 30</math>, the sample data distribution should be free from strong skewness and outliers. For matched pairs, the number of differences should be greater than or equal to 30. If the number of differences is less than 30, the sample of differences should be free from strong skewness and outliers.</li> </ol>

TOPIC 4.5

# Carrying Out a Test for a Population Mean

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.E	<p><b>4.5.A</b></p> <p>Calculate an appropriate test statistic and <math>p</math>-value for a population mean, including the mean difference between values in matched pairs.</p>	<p><b>4.5.A.1</b></p> <p>The test statistic for a one-sample <math>t</math>-test for a population mean or population mean difference is <math>t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}</math>, where <math>t</math> has degrees of freedom <math>n - 1</math>. The <math>t</math>-statistic has a <math>t</math>-distribution with degrees of freedom <math>n - 1</math> when the null hypothesis is true.</p> <p><b>4.5.A.2</b></p> <p>The <math>p</math>-value for a one-sample <math>t</math>-test for a population mean is found using the appropriate <math>t</math>-distribution table or from the appropriate <math>t</math>-distribution using technology.</p>
4.F	<p><b>4.5.B</b></p> <p>Interpret the <math>p</math>-value of a hypothesis test for a population mean, including the mean difference between values in matched pairs.</p>	<p><b>4.5.B.1</b></p> <p>The <math>p</math>-value is the probability of obtaining a test statistic as extreme or more extreme than the test statistic that was observed (i.e., in the direction of the alternative hypothesis) given that the null hypothesis is true. An interpretation of the <math>p</math>-value of a hypothesis test for a population mean should include a statement that the <math>p</math>-value is computed by assuming that the null hypothesis is true (i.e., by assuming that the population mean is equal to the particular value stated in the null hypothesis).</p>
4.G	<p><b>4.5.C</b></p> <p>Justify a claim about the population based on the results of a hypothesis test for a population mean, including the mean difference between values in matched pairs.</p>	<p><b>4.5.C.1</b></p> <p>A formal decision explicitly compares the <math>p</math>-value to the significance level, <math>\alpha</math>. If the <math>p</math>-value <math>\leq \alpha</math> then reject the null hypothesis, <math>H_0 : \mu = \mu_0</math>. If the <math>p</math>-value <math>&lt; \alpha</math> then fail to reject the null hypothesis.</p> <p><b>4.5.C.2</b></p> <p>The results of a hypothesis test for a population mean can serve as the statistical reasoning to support the answer to a research question about the population that was sampled.</p> <p><b>4.5.C.3</b></p> <p>A conclusion for the hypothesis test for a population mean is stated in context consistent with, and in terms of, the alternative hypothesis using non-definitive language. The conclusion should contain a reference to the parameter and the population.</p>

## TOPIC 4.6

# Sampling Distributions for the Difference Between Two Sample Means

Instructional Periods: 1

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.D	<p><b>4.6.A</b></p> <p>Calculate parameters of a sampling distribution for the difference between two sample means.</p>	<p><b>4.6.A.1</b></p> <p>For two independent populations with population means <math>\mu_1</math> and <math>\mu_2</math> and population standard deviations <math>\sigma_1</math> and <math>\sigma_2</math>, when the sampled values are independent, the sampling distribution of the difference in sample means <math>\bar{x}_1 - \bar{x}_2</math> has a mean <math>\mu_{(\bar{x}_1 - \bar{x}_2)} = \mu_1 - \mu_2</math> and standard deviation</p> $\sigma_{(\bar{x}_1 - \bar{x}_2)} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}.$ <p><b>4.6.A.2</b></p> <p>When sampling without replacement, the data must meet two conditions as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: the data should be collected using two independent random samples.</li> <li>The 10% condition: both samples must be less than 10% of the size of their respective populations.</li> </ol> <p><b>4.6.A.3</b></p> <p>If the data come from an experiment, the data only need to meet the randomization condition. The treatments must be randomly assigned to participants or experimental units to meet the randomization condition.</p>
4.E	<p><b>4.6.B</b></p> <p>Justify why a sampling distribution for the difference between sample means can or cannot be described as approximately normal.</p>	<p><b>4.6.B.1</b></p> <p>The sampling distribution for the difference between sample means, <math>\bar{x}_1 - \bar{x}_2</math>, can be modeled with a normal distribution if the two population distributions can each be modeled by a normal distribution.</p> <p><b>4.6.B.2</b></p> <p>The sampling distribution for the difference between sample means, <math>\bar{x}_1 - \bar{x}_2</math>, can be modeled approximately by a normal distribution if the two population distributions cannot be modeled by a normal distribution but <math>n_1 \geq 30</math> and <math>n_2 \geq 30</math>.</p>

## SKILL

**4.D**

## LEARNING OBJECTIVE

**4.6.C**

Interpret parameters and probabilities for a sampling distribution for the difference between sample means.

## ESSENTIAL KNOWLEDGE

**4.6.C.1**

Parameters and probabilities for a sampling distribution for the difference between sample means should be interpreted within the context of specific populations.

## TOPIC 4.7

# Constructing a Confidence Interval for the Difference Between Two Sample Means

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.C	<b>4.7.A</b> Identify an appropriate confidence interval procedure including the parameter for the difference between two population means.	<b>4.7.A.1</b> The sample data can be used to calculate a confidence interval to estimate the difference between two independent population means. The appropriate confidence interval procedure for two independent samples is a two-sample $t$ -interval for the difference between population means. <b>4.7.A.2</b> The parameter for a confidence interval for a two-sample $t$ -interval for the difference between population means should reference the population mean and the response variable in context.
4.E	<b>4.7.B</b> Justify why a sampling distribution for the difference between sample means can or cannot be described as approximately normal.	<b>4.7.B.1</b> A two-sample $t$ -interval for a difference between population means requires that three conditions be met as follows: <ol style="list-style-type: none"><li>The randomization condition: the data should be collected using two random samples or a randomized experiment.</li><li>The 10% condition: When sampling without replacement, the size of each sample should be less than or equal to 10% of the respective population size: <math>n_1 &lt; 10\%N_1</math> and <math>n_2 &lt; 10\%N_2</math>, where <math>N_1</math> is the size of population 1 and <math>N_2</math> is the size of population 2. The sample sizes are represented as <math>n_1</math> and <math>n_2</math>. (Note: This condition is unnecessary when the data are from a randomized experiment.)</li><li>The sample data condition: Both samples should have a sample size greater than or equal to 30. If either sample size is less than 30, both sample data distributions should be free from strong skewness and outliers.</li></ol>



SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.E	<p><b>4.7.C</b></p> <p>Calculate an appropriate confidence interval for the difference between two population means.</p>	<p><b>4.7.C.1</b></p> <p>A point estimate for the difference between two population means is the difference in sample means, <math>\bar{x}_1 - \bar{x}_2</math>.</p> <p><b>4.7.C.2</b></p> <p>For the difference between population means when the population standard deviations are unknown, the confidence interval can be constructed as point estimate <math>\pm</math> (margin of error). The confidence interval for the difference between population means is <math>(\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}</math>, where <math>\pm t^*</math> are the critical values for the central C% of a <math>t</math>-distribution with appropriate degrees of freedom that can be found using technology. The degrees of freedom fall between <math>n_1 + n_2 - 2</math> and the smaller of <math>n_1 - 1</math> and <math>n_2 - 1</math>.</p>
3.E	<p><b>4.7.D</b></p> <p>Calculate the standard error and margin of error for the difference between two population means.</p>	<p><b>4.7.D.1</b></p> <p>The standard error for the difference between two sample means is <math>\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}</math>, where <math>s_1</math> and <math>s_2</math> are the sample standard deviations.</p> <p><b>4.7.D.2</b></p> <p>For the difference between two sample means, the margin of error is the critical value (<math>t^*</math>) times the standard error (<math>SE</math>) of the difference of two sample means, which equals <math>t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}</math>.</p>

## TOPIC 4.8

# Constructing a Confidence Interval for the Difference Between Two Sample Means

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
4.F	<b>4.8.A</b> Interpret a confidence interval in context for the difference between two population means.	<b>4.8.A.1</b> Since the confidence interval for the difference between two population means is calculated based on samples from two populations, the computed interval may or may not contain the value for the difference between the two population means. <b>4.8.A.2</b> The interpretation of the confidence level is as follows: In repeated random sampling with the same sample size, approximately $C\%$ of confidence intervals created will capture the difference between the two population means, where $C$ represents the numerical value of the confidence level used. <b>4.8.A.3</b> When interpreting a $C\%$ confidence interval for the difference between two population means, we say we are $C\%$ confident that the interval $(a, b)$ contains the value of the difference in the population means. An interpretation of a confidence interval for the difference between two population means includes a reference to the difference in the population means with the details about the populations it represents in the context of the study.
4.G	<b>4.8.B</b> Justify a claim based on a confidence interval for the difference between two population means.	<b>4.8.B.1</b> A confidence interval for the difference between two population means provides an interval of values that may serve as convincing evidence to support a particular claim about the difference in two population means. For example, if the interval contains 0, then there is insufficient evidence to conclude there is a difference between the two population means. If the interval does not contain 0, then there is sufficient evidence to conclude there is a difference between the two population means.

## TOPIC 4.9

# Setting Up a Test for the Difference Between Two Population Means

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
2.C	<p><b>4.9.A</b></p> <p>Identify an appropriate testing method for the difference between two population means including the parameters for the difference between the two population means.</p>	<p><b>4.9.A.1</b></p> <p>For two independent samples, the appropriate test for the difference between two population means is a two-sample <math>t</math>-test for a difference between two population means.</p> <p><b>4.9.A.2</b></p> <p>The parameter for the two-sample <math>t</math>-test for the difference between two population means should reference the difference in means of the populations and the response variable in context.</p>
2.E	<p><b>4.9.B</b></p> <p>Identify the null and alternative hypotheses for the difference between two population means.</p>	<p><b>4.9.B.1</b></p> <p>The null hypothesis for a two-sample <math>t</math>-test for the difference between two population means, <math>\mu_1</math> and <math>\mu_2</math>, can be written as either: <math>H_0 : \mu_1 - \mu_2 = 0</math> or <math>H_0 : \mu_1 = \mu_2</math>. A one-sided alternative hypothesis for the difference between population means can be written as either <math>H_a : \mu_1 &lt; \mu_2</math> (or equivalently <math>H_a : \mu_1 - \mu_2 &lt; 0</math>) or <math>H_a : \mu_1 &gt; \mu_2</math> (or equivalently <math>H_a : \mu_1 - \mu_2 &gt; 0</math>). A two-sided alternative hypothesis for the difference between population means can be written as <math>H_a : \mu_1 \neq \mu_2</math> (or equivalently <math>H_a : \mu_1 - \mu_2 \neq 0</math>).</p>
4.E	<p><b>4.9.C</b></p> <p>Justify the appropriateness of a hypothesis test for the difference between two population means by verifying conditions.</p>	<p><b>4.9.C.1</b></p> <p>A two-sample <math>t</math>-test for a difference between population means requires that three conditions be met as follows:</p> <ol style="list-style-type: none"> <li>The randomization condition: the data should be collected using two random samples or a randomized experiment.</li> <li>The 10% condition: When sampling without replacement, the size of each sample should be less than or equal to 10% of the respective population size: <math>n_1 &lt; 10\%N_1</math> and <math>n_2 &lt; 10\%N_2</math>, where <math>N_1</math> is the size of population 1 and <math>N_2</math> is the size of population 2. The sample sizes are represented as <math>n_1</math> and <math>n_2</math>. (Note: This condition is unnecessary when the data are from a randomized experiment.)</li> <li>The sample data condition: Both samples should have a sample size greater than or equal to 30. If either sample size is less than 30, both sample data distributions should be free from strong skewness and outliers.</li> </ol>

## TOPIC 4.10

# Carrying Out a Test for the Difference Between Two Sample Means

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.E	<p><b>4.10.A</b></p> <p>Calculate an appropriate test statistic and <math>p</math>-value for testing a hypothesis for the difference between two population means.</p>	<p><b>4.10.A.1</b></p> <p>The test statistic for a two-sample <math>t</math>-test for the difference between two population means is <math>t = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}</math>. The <math>t</math>-statistic has a <math>t</math>-distribution when the null hypothesis is true. The <math>t</math>-statistics and the degrees of freedom, which fall between <math>n_1 + n_2 - 2</math> and the smaller of <math>n_1 - 1</math> and <math>n_2 - 1</math>, can be found with technology.</p>
4.F	<p><b>4.10.B</b></p> <p>Interpret the <math>p</math>-value of a hypothesis test for the difference between two population means.</p>	<p><b>4.10.B.1</b></p> <p>The <math>p</math>-value is the probability of obtaining a test statistic as extreme or more extreme than the test statistic that was observed (i.e., in the direction of the alternative hypothesis) given that the null hypothesis is true. An interpretation of the <math>p</math>-value of a hypothesis test for a two-sample test for the difference between two population means should include a statement that the <math>p</math>-value is computed by assuming that the null hypothesis is true (i.e., by assuming the population means are equal to each other).</p>

**SKILL**

**4.G**

**LEARNING OBJECTIVE**

**4.10.C**

Justify a claim about the population based on the results of a hypothesis test for the difference between two population means in context.

**ESSENTIAL KNOWLEDGE**

**4.10.C.1**

A formal decision explicitly compares the  $p$ -value to the significance level,  $\alpha$ . If the  $p$ -value  $\leq \alpha$  then reject the null hypothesis,  $H_0 : \mu_1 - \mu_2 = 0$  or  $H_0 : \mu_1 = \mu_2$ . If the  $p$ -value  $> \alpha$  then fail to reject the null hypothesis.

**4.10.C.2**

The results of a hypothesis test for a two-sample test for a difference between two population means can serve as the statistical reasoning to support the answer to a research question about the populations that were sampled.

**4.10.C.3**

A conclusion for the hypothesis test for the difference between two population means is stated in context consistent with, and in terms of, the alternative hypothesis using non-definitive language. The conclusion should contain a reference to the parameters and the populations.

THIS PAGE IS INTENTIONALLY LEFT BLANK.

**AP STATISTICS**

**UNIT 5**

# **Regression Analysis**

THIS PAGE IS INTENTIONALLY LEFT BLANK.



## TOPIC 5.1

# Graphical Representations Between Two Quantitative Variables

Instructional Periods: 1

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.A	<p><b>5.1.A</b></p> <p>Construct scatterplots depicting the distribution of two numeric variables.</p>	<p><b>5.1.A.1</b></p> <p>A bivariate quantitative data set consists of observations of ordered pairs from two quantitative variables, acquired from the same individuals in a sample or in a population, that can be used to construct a scatterplot.</p> <p><b>5.1.A.2</b></p> <p>A scatterplot shows the values of two quantitative variables for each observation, one corresponding to the value on the <math>x</math>-axis and one corresponding to the value on the <math>y</math>-axis. The explanatory variable is placed on the <math>x</math>-axis and is the variable whose values are used to explain or predict the corresponding values for the response variable, which is placed on the <math>y</math>-axis.</p>
4.A	<p><b>5.1.B</b></p> <p>Describe the characteristics of a scatterplot.</p>	<p><b>5.1.B.1</b></p> <p>A description of the association shown in a scatterplot includes form, direction, strength, and unusual features.</p> <p><b>5.1.B.2</b></p> <p>The form of the association shown in a scatterplot, if any, can be described as linear or nonlinear.</p> <p><b>5.1.B.3</b></p> <p>The direction of the association shown in a scatterplot, if any, can be described as positive or negative. A positive association means that as values of the explanatory variable increase, the values of the response variable tend to increase. A negative association means that as values of the explanatory variable increase, the values of the response variable tend to decrease.</p> <p><b>5.1.B.4</b></p> <p>The strength of the association shown in a scatterplot is how closely the points follow the general pattern. Strength can be described as strong, moderate, or weak.</p> <p><b>5.1.B.5</b></p> <p>Unusual features of a scatterplot include clusters of individual points or points that don't fit in the general pattern of association between the two variables.</p>

## TOPIC 5.2

**Correlation**

Instructional Periods: 2

**SKILL****4.D****LEARNING OBJECTIVE****5.2.A**

Interpret the correlation for a linear relationship.

**ESSENTIAL KNOWLEDGE****5.2.A.1**

The correlation coefficient,  $r$ , summarizes the strength and direction of the linear association between two quantitative variables. The correlation coefficient  $r$  is unit-free and always between  $-1$  and  $1$ , inclusive. A negative correlation coefficient value indicates a negative association, and a positive correlation coefficient value indicates a positive association.

**5.2.A.2**

The strength of the linear association is determined by how close the correlation coefficient is to  $1$  or  $-1$ . A value of  $r = 0$  indicates that there is no linear association. A value of  $r = 1$  or  $r = -1$  indicates that there is a perfect linear association.

**5.2.A.3**

A correlation coefficient close to  $-1$  or  $1$  does not necessarily mean that a linear model is appropriate.

**5.2.A.4**

A perceived or real relationship between two variables does not mean that changes in one variable cause changes in the other. That is, correlation does not necessarily imply causation.

## TOPIC 5.3

# Linear Regression Models

Instructional Periods: 2

**SKILL****3.B****LEARNING OBJECTIVE****5.3.A**

Calculate a predicted response value using a linear regression model.

**ESSENTIAL KNOWLEDGE****5.3.A.1**

If the form of the relationship between  $x$  and  $y$  appears linear, we can approximate the relationship between  $x$  and  $y$  using a simple linear regression model, which is a linear equation that uses an explanatory variable,  $x$ , to predict the response variable,  $y$ .

**5.3.A.2**

In a simple linear regression model, the predicted response value, denoted by  $\hat{y}$ , is calculated as  $\hat{y} = a + bx$ , where  $a$  is the  $y$ -intercept and  $b$  is the slope of the regression line and  $x$  is the explanatory variable.

**5.3.A.3**

Extrapolation is predicting a response value using a value for the explanatory variable that is beyond the interval of  $x$ -values used to determine the regression line. The predicted value is less reliable the further the estimate is extrapolated.

**5.3.A.4**

Interpolation is predicting a response value using a value for the explanatory variable that is within the interval of  $x$ -values used to determine the regression line.

## TOPIC 5.4

# Residuals

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.B	<p><b>5.4.A</b></p> <p>Calculate the differences between actual and predicted values.</p>	<p><b>5.4.A.1</b></p> <p>A residual is the difference between the observed response value and the predicted response value for the given value of the explanatory variable: <math>\text{residual} = y - \hat{y}</math> or <math>(\text{residual} = \text{observed } y - \text{predicted } \hat{y})</math>.</p>
4.D	<p><b>5.4.B</b></p> <p>Interpret the differences between actual and predicted values.</p>	<p><b>5.4.B.1</b></p> <p>If the residual is positive, the model underpredicts (underestimates) the value of the response variable. If the residual is negative, the model overpredicts (overestimates) the value of the response variable.</p>
4.A	<p><b>5.4.C</b></p> <p>Describe the form of association of bivariate data using residual plots.</p>	<p><b>5.4.C.1</b></p> <p>A residual plot is a scatterplot of the residuals versus the predicted response values (or the explanatory variable values).</p> <p><b>5.4.C.2</b></p> <p>Residual plots can be used to investigate the appropriateness of the simple linear regression model for the observed data.</p> <p><b>5.4.C.3</b></p> <p>The simple linear regression model should only be fit to the data if the data exhibits a linear trend. Apparent randomness in a residual plot for a simple linear regression model is confirmation of a linear form in the association between the two variables and indicates that the simple linear regression model is an appropriate model for the data.</p> <p><b>5.4.C.4</b></p> <p>Curvature in the residual plot for a simple linear regression model suggests that the linear model is not the most appropriate model for the data.</p>

## TOPIC 5.5

## Least-Squares Regression

Instructional Periods: 2

SKILL	LEARNING OBJECTIVE	ESSENTIAL KNOWLEDGE
3.B	<p><b>5.5.A</b></p> <p>Calculate the coefficients for the least-squares regression line model.</p>	<p><b>5.5.A.1</b></p> <p>The simple linear regression model is fit to the data by minimizing the sum of the squares of the residuals. Because of this, the resulting equation is often called the least-squares regression line (LSRL) and is calculated using technology. This regression line will pass through the point <math>(\bar{x}, \bar{y})</math>.</p> <p><b>5.5.A.2</b></p> <p>The slope of the regression line, <math>b</math>, is calculated using technology.</p> <p><b>5.5.A.3</b></p> <p>The <math>y</math>-intercept of the regression line, <math>a</math>, is calculated using technology.</p> <p><b>5.5.A.4</b></p> <p>In simple linear regression, the square of the correlation coefficient, <math>r^2</math>, is called the coefficient of determination. <math>r^2</math> is the proportion of variation in the response variable that is explained by the linear relationship with the explanatory variable.</p>
4.D	<p><b>5.5.B</b></p> <p>Interpret coefficients for the least-squares regression line model.</p>	<p><b>5.5.B.1</b></p> <p>The coefficients of the least-squares regression line model (line of best fit) are the estimated slope, <math>b</math>, and the estimated <math>y</math>-intercept, <math>a</math>, because they are based on a sample of values.</p> <p><b>5.5.B.2</b></p> <p>The slope coefficient of the least-squares regression line can be interpreted as the predicted increase or decrease in the response variable for a one-unit increase or decrease in the explanatory variable, and it should be interpreted in context.</p> <p><b>5.5.B.3</b></p> <p>The <math>y</math>-intercept coefficient in the least-squares regression line is the predicted value of the response variable when the explanatory variable is equal to 0, and it should be interpreted in context. Sometimes, the <math>y</math>-intercept of the line does not have a reasonable interpretation in context since <math>x = 0</math> might be beyond the interval of <math>x</math>-values used to determine the regression line (extrapolation). At other times, the <math>y</math>-intercept of the line does not have a logical interpretation in context since it might be a negative value for a response variable that has no negative values, such as height.</p>



