

1  
2  
3

4  
5  
6

**Mathematics Framework**  
**Chapter 5: Mathematical Foundations for Data Science**

7	Mathematics Framework Chapter 5: Mathematical Foundations for Data Science .....	1
8	Introduction.....	2
9	What Are Data Literacy and Data Science? .....	3
10	Why a Chapter on the Foundations of Data Science in the Math Framework? .....	5
11	Using Statistics and Data Science in a Problem-Solving Process to Support the	
12	Standards for Mathematical Practice and to Make Connections to Other Domains ....	8
13	Thematic Topics Within the CA CCSSM That Directly Support Data Science .....	12
14	Understanding and Describing Variability in Data and Data Distributions .....	15
15	Data Collection, Sampling, and Random Processes .....	16
16	Comparing Distributions and Identifying Associations Between Variables.....	17
17	Data Science in Each Grade Band.....	18
18	Transitioning From Pre-Kindergarten.....	18
19	Kindergarten Through Grade Five .....	18
20	Grades Six Through Eight.....	37
21	High School.....	51
22	Equitable and Inclusive Instruction.....	61
23	Connecting to the Drivers of Investigation and Content Connections .....	62
24	Conclusion.....	65
25	Additional Resources .....	67
26	Long Descriptions of Graphics for Chapter 5 .....	68

## 27 **Introduction**

28 The ability to work with and understand data is an essential life skill in a world  
29 continually inundated with data. As data become omnipresent in all sectors of life—  
30 personal, business, academic, and education—community members and citizens need  
31 the attitudes, skills, and practices to use data to make informed decisions in  
32 professional and personal matters. Data drive students’ lives, whether they see it or not.  
33 Making sense of data, being able to identify data that are misleading, and using data to  
34 make decisions are all important skills for students in their roles as global citizens.  
35 Almost all occupations now require that employees collect feedback data and use this  
36 information to adjust their practice. Stories about the world are illuminated by massive

37 quantities of data, and community members telling and listening to those stories need to  
38 be able to make sense of data to understand their health, their finances, and news.

## 39 **What Are Data Literacy and Data Science?**

40 Many groups have used different terms to describe the ability to work with and derive  
41 meaning from data. These terms include statistical literacy (Bargagliotti et al., 2020),  
42 data literacy (Education Development Center, 2016), data fluency, and data acumen  
43 (National Academies of Sciences, Engineering, and Medicine, 2018). Because a full  
44 discussion of the academic and industry differences between these terms is outside the  
45 scope of this framework, the rest of this chapter uses the terms “data literacy” and “data  
46 science.”

47 Wolff and colleagues (2016, 10) describe “data literacy” as

48 the ability to ask and answer real-world questions from large and small  
49 data sets through an inquiry process, with consideration of ethical use of  
50 data. It is based on core practical and creative skills, with the ability to  
51 extend knowledge of specialist data handling skills according to goals.  
52 These include the abilities to select, clean, analyze, visualize, critique,  
53 and interpret data, as well as to communicate stories from data and use  
54 data as part of a design process.

55 Within most mathematics courses, students at the prekindergarten through grade twelve  
56 level will be building strong mathematical foundations and be engaged in work that  
57 supports data literacy; all California students should graduate from high school with data  
58 literacy. However, students should also have access to experiences that extend beyond  
59 what many currently experience in their mathematics classrooms and that prepare them  
60 for future work in an emerging field called data science.

61 “Content expectations across multiple school subjects in U.S. primary and  
62 secondary education already incorporate at least some learning about data  
63 collection, utilization, and analysis. Data-related concepts consistently appear in

64 mathematics, science, computer science, and social studies across states.  
65 These existing standards may provide the building blocks or even partially  
66 comprise a data science education.” (Drozda, Johnstone, and Van Horne, n.d., 1)

67 At the broadest level, data science is “the science of learning from data” (IDSSP, 2019,  
68 9) or “the processes and systems that enable the extraction of knowledge or insights  
69 from data in various forms, either structured or unstructured” (National Science  
70 Foundation Computer and Information Science and Engineering Advisory Committee  
71 Data Science Working Group, 2016, 2). There has been an expansion in computing and  
72 visualization tools that have made many more techniques available for finding meaning  
73 in data—often relying on innovative visualizations of complex data that enable major  
74 features to be identified and explored further.

75 Data science is an emerging discipline, and there is yet to be a consensus on exactly  
76 what content constitutes “data science.” Data science education organizations have  
77 conceptualized it as a cyclical process that includes problem formulation, data  
78 collection, data analysis, and interpretation and communication of findings (GAISE II,  
79 IDSSP, 2019). While the data science field continues to be shaped by emerging  
80 technologies and techniques deeply influenced by academia and business sectors,  
81 most definitions of data science recognize an intersection of mathematics, statistics,  
82 and computer science. Data science may also include domain knowledge, ethics, and  
83 communication skills as well as specific approaches such as data mining (especially for  
84 data collected through the internet and electronic devices) and machine learning  
85 (Bargagliotti et al., 2020; NSF CISE Advisory Committee Data Science Working Group,  
86 2016; IDSSP, 2019; Rawlings-Goss et al., 2019). Although students may not encounter  
87 topics such as machine learning until after high school, kindergarten through grade  
88 twelve (K–12) mathematics provides an essential foundation in statistical concepts  
89 necessary for future learning in data science.

90 Data literacy and data science should be thought of more as a continuum than as  
91 distinct concepts. The focus of data literacy is the ability to use an inquiry process to  
92 extract answers to real-world questions from data sets. All K–12 graduates should have

93 developed data literacy through rich data experiences in each grade level. Data literacy  
94 is part of data science, but data science also includes advanced mathematics, statistics,  
95 and computational skills that build upon—and go far beyond—the content contained in  
96 the K–12 California mathematics standards.

## 97 **Why a Chapter on the Foundations of Data Science in the** 98 **Math Framework?**

99 California is not alone in giving attention to this growing field. State departments of  
100 education in Georgia, Ohio, Oregon, Utah, and Virginia are also exploring how to  
101 increase access to data science concepts for their students, whether by revising  
102 mathematics standards, creating new mathematics course pathways or frameworks, or  
103 providing micro-credentials in data science for teachers (National Center for Education  
104 Research, Institute of Education Sciences, 2021). Other countries, including China, New  
105 Zealand, and South Korea, are also investing in more data science instruction for all K–  
106 12 students (DataScience4everyone, 2022).

107 This dedicated chapter on the foundations of data science is included in the framework  
108 for two primary reasons. The first has to do with relevance: Data literacy is increasingly  
109 central to being able to understand and fully participate in modern life. Discussions  
110 around healthcare, finance, electoral politics, and other major challenges of the current  
111 era all revolve around data. Since students frequently encounter claims made from such  
112 large data sets, it is crucial that all students have experiences in which they explore the  
113 ways in which such claims are made. The second reason has to do with the importance  
114 of data science in preparing students for the future. Data science is and will continue to  
115 be a fast-growing and lucrative field (Bureau of Labor Statistics, 2022). Both of these  
116 facts lead to the same conclusion: Students should have equitable access to data  
117 literacy and introductory data science at the K–12 level to facilitate equitable  
118 participation in a data-driven world as adults.

119 Studies by Walton and colleagues (2015) show that many students, particularly girls  
120 and students of color, do not feel that they belong in certain disciplines. These feelings

121 are often due to a history of negative and off-putting messages (Chestnut et al., 2018).  
122 Other studies have shown that different topics and teaching approaches can lead to  
123 feelings of belonging or not belonging (Boaler, Cordero, and Dieckmann, 2019). Data  
124 investigation can support teachers as they seek to create climates of belonging for  
125 students, inviting them to investigate real data that is likely relevant to their lives. This  
126 meaningful engagement can create opportunities for students to develop self-  
127 confidence and self-efficacy with mathematics.

128 Unequal access to data education in K–12 leads to unequal participation in occupations,  
129 systems, and outcomes that create or are heavily impacted by data-related products  
130 and endeavors. Consequently, historically marginalized communities have fewer  
131 opportunities to reap the economic benefits of data-driven industry. Unequal  
132 participation in data science can also lead to bias, not only in the solutions that are  
133 developed (e.g., résumé-screening tools that overlook historically underrepresented  
134 groups in the industry) but also in the selection of problems that get attention and  
135 resources.

136 Particular aspects of the California Common Core State Standards for Mathematics (CA  
137 CCSSM), especially the standards related to statistics instruction, help build the data  
138 understanding and skills that high school graduates require. However, the progression  
139 of these ideas—from counting, categorizing, and simple picture graphs studied at  
140 younger grade levels, to the complex skills and understanding that older students may  
141 develop—requires careful thought and considerably more focus throughout the K–12  
142 curriculum than most students have historically experienced. This chapter is a first step  
143 in identifying how the current standards can support data literacy across the grade  
144 levels and help K–12 students develop foundational knowledge and skills for data  
145 science. Subsequent chapters (6–8) will provide additional grade band–specific  
146 examples of how data can be integrated into mathematics. Learning about the  
147 mathematical and statistical concepts and practices associated with data science may  
148 help teachers energize their mathematics instruction and extend their students’  
149 mathematical experiences in new ways.

150 In the past, statistics instruction focused on just a few key ideas and procedures (mean,  
151 median, standard deviation, interquartile range, correlation, and linear regression, along  
152 with a few data visualizations such as line plots and scatter plots) or was overlooked  
153 altogether. As students progress through school, they should learn different approaches  
154 to statistical analysis, culminating in the investigation of large data sets using  
155 appropriate technological tools. Elevating statistics as a way to understand the world  
156 and solve problems at the K–12 level is an important step in supporting data literacy for  
157 all students and building a pathway for an introduction to data science in the third or  
158 fourth year of high school.

159 This emphasis is not meant to suggest that statistics should replace other math content.  
160 Statistics and other math domains are mutually reinforcing. For example, understanding  
161 of linear regression is closely related to understanding of functions and polynomials. A  
162 comprehensive understanding of all domains in K–12 mathematics is necessary for  
163 successful postsecondary work in data science. If students are intending to pursue  
164 STEM majors in college (including data science), they should take courses that, at a  
165 minimum, allow them to enter college having completed the prerequisites for calculus.  
166 As of this writing, undergraduate data science programs typically require a core math  
167 sequence that includes calculus and linear algebra.

168 Nor are the points made above meant to suggest that learning statistics is equivalent to  
169 learning data science. The types of data being collected are vast and the types of  
170 techniques used to extract insights from the data depend on a strong understanding of  
171 multiple areas. Students will need to learn how to use computational tools to store,  
172 transform, and analyze the quantity of data being generated. Students will also need to  
173 have domain knowledge to identify questions that can be investigated through data  
174 science and to interpret the results of their analyses in a thoughtful way.

175 Statistics has become increasingly relevant in applications of mathematics and provides  
176 many contemporary illustrations of the significance of the CA CCSSM. Accordingly, this  
177 chapter focuses on ways in which teachers can create rich statistics and data  
178 experiences across PreK–12 that can help (a) modernize the teaching of Statistics and

179 Probability standards, (b) engage students in the kind of authentic problem solving that  
180 broadens participation in STEM fields generally and data science specifically, and (c)  
181 prepare students for life and work in the data age.

182 To avoid confusion about terminology, the remainder of this chapter will use the term  
183 “data science” to encompass K–12 work that serves these ends.

## 184 **Using Statistics and Data Science in a Problem-Solving** 185 **Process to Support the Standards for Mathematical Practice** 186 **and to Make Connections to Other Domains**

187 With data serving as the basis of large-scale decisions and predictions, all California  
188 high school graduates need skills in interpreting and visualizing data, making and  
189 critiquing data-based arguments, and some facility with spreadsheets and other tools  
190 used to store and analyze data. It is crucial for students to develop the ability to identify  
191 types of questions that are subject to exploration through data. Just as crucial is their  
192 understanding of some misuses of data. Students must ultimately approach data  
193 science and statistics as a problem-solving process that consists of formulating  
194 statistical investigative questions, collecting and interrogating existing data, analyzing  
195 data, and interpreting and communicating findings (Bargagliotti et al., 2020). Across the  
196 K–12 grade levels, students should have opportunities to experience data in different  
197 ways, such as the following:

- 198 ● **Encountering and understanding the role of data in the world:** How do we  
199 explore and interpret data and make ethical decisions about how it is used?  
200 Students should experience working with data from a context that is meaningful  
201 to them personally. They should have opportunities to solve problems of value to  
202 them and to their schools and communities.
  
- 203 ● **Collecting and exploring data:** How can we collect data? Data explorations  
204 should be investigative and collaborative, with students working together to ask  
205 investigative questions or engage in statistics as a problem-solving process.



206 Students should have multiple opportunities to become familiar with a variety of  
207 technology and modern tools to access, collect, explore, and make sense of  
208 data.

209 ● **Considering variability and engaging in multivariate thinking:** How can we  
210 describe, display, and compare data effectively? How can we determine the  
211 relationship between different variables or quantities? Students should learn to  
212 engage with real data that include multiple variables. At first, students can learn  
213 to understand the relationship between two variables with bivariate data; as they  
214 progress through the grades, they can learn to handle multivariable data and  
215 multivariate thinking. Multivariable data often include three or more variables. For  
216 example, students could categorize their favorite stuffed animals by considering  
217 their size, fluffiness, and type of animal (three variables).

218 ● **Considering data sampling and probability:** How can we use random  
219 sampling to help understand a population? How can we determine the chances  
220 that an event or events will occur? Technology and tool use should become more  
221 complex as students progress through the grades and can help them explore the  
222 role of sampling and probability.

223 ● **Interpreting and communicating findings:** What do our data mean? Does our  
224 analysis address any of our questions? What are the best ways to communicate  
225 our findings? What impacts might the findings have? As students learn to  
226 interpret data in increasingly sophisticated ways, they should also have  
227 opportunities to make statements about the data and to practice using data  
228 visualizations to communicate results. Especially in middle and high school,  
229 students' encounters with data should revisit the context from which the data  
230 originated, interpreting results in that context.

231 Throughout the CA CCSSM, there are multiple opportunities to support such data-rich  
232 experiences and integrate the five components of equitable and engaging teaching  
233 described in chapter two, even if the standards domains do not appear explicitly within a  
234 grade or grade band. When approaching the grade band chapters in this framework

235 from a data science lens, educators can find additional moments to integrate data into  
236 students' mathematical experiences. For example, chapter six includes a vignette  
237 describing Mrs. Verners' fourth grade lessons supporting Number and Operations in  
238 Base Ten and Operations and Algebraic Thinking ([Comparing Numbers and Place](#)  
239 [Value Relationships in Grade Four, with Integrated English Language Development](#)). In  
240 her class, students explore population data by making estimates based on prior  
241 knowledge, exploring data both in written and standard form, considering place value,  
242 and making multiplicative comparisons. The lessons help students focus on changing  
243 mathematical quantities while also connecting to social studies content and integrating  
244 English language arts/English language development standards in a meaningful way.  
245 Similarly, as described in chapter eight, alternative third- and fourth-year high school  
246 courses can also provide valuable opportunities to explore important data science topics  
247 beyond statistics, such as ethics, data modeling and simulations, and data cleaning.  
248 Two important sources for contexts in which to explore statistics and data science are:

- 249           • The California Next Generation Science Standards (CA NGSS) (California  
250           Department of Education, 2013a)
- 251           • The California Environmental Principles and Concepts (EP&Cs) (California  
252           Department of Education, 2013b)

253 In addition to connecting data-rich experiences to other content, data investigations can  
254 support students to draw on the Standards for Mathematical Practices (SMPs) as they  
255 engage in this problem-solving process across every grade band. For example,  
256 students should reason abstractly and quantitatively by engaging in statistical thinking  
257 while considering where data come from (SMP.2); apply statistical models to include  
258 descriptions of the variability present in data (SMP.4); and consider available tools such  
259 as calculators, spreadsheets, applets, statistical packages, and graphical displays to  
260 help facilitate the statistical problem-solving process (SMP.5). When students  
261 participate in the analysis of large data sets, they should be able to decide which  
262 questions matter and identify which ones can be answered with a given data set  
263 (SMP.4). Figure 5.1 illustrates an example of how the SMPs can be highlighted within  
264 an elementary data investigation.

Brief Description of the Learning Activity (Level A Ladybugs Example)	Connections to the Standards for Mathematical Practice
<p>Students formulate statistical investigative questions:                      “How many spots do ladybugs typically have?” Or “Do red-bodied ladybugs tend to have more spots than black-bodied ladybugs?”</p>	<p><b>SMP.1: Make sense of problems</b>                      Consider which of our questions can be answered with data.                      Interesting investigations anticipate that data collected will vary or are not the same for every observation.</p>
<p>Students collect data:                      Students recognize that they can use photographs to help answer data collection questions—e.g., “What is the body color?” or “How many spots are on each ladybug?” These questions generate data that are both numeric and categorical.</p>	<p><b>SMP.5: Use appropriate tools strategically</b>                      Students use a nontraditional data source—photographs of ladybugs—and develop a data collection plan for the class to use. Tables are helpful tools for organizing individual data or for collecting and organizing data from multiple students.  <b>SMP.6: Attend to precision</b>                      Humans make decisions that impact data collection and resulting analyses or interpretations.</p>
<p>Students analyze data by making plots and describing them:                      Students make dot plots for their ladybug data and describe the number of spots that were most common, visually estimate the median, and identify how these values compared for ladybugs of different colors.                      Students can use a probability chart (0 = not probable, 1/2 = equal chance, 1 = very likely) to express whether they think something will happen.</p>	<p><b>SMP.2: Reason abstractly and quantitatively</b>                      Students can compare the distribution of the number of spots for ladybugs of different colors.                      Students make informal associations between ladybug color and their numbers of spots, for example black-bodied ladybugs have fewer spots than red-bodied and orange-bodied ladybugs.                      Students use data to consider the likelihood that something will happen, such as finding a black-bodied ladybug with 14 spots.</p>
<p>Students interpret and use the plots to answer their initial questions:                      Students describe the distribution of spot numbers and range of spot numbers for ladybugs.</p>	<p><b>SMP.3: Construct viable arguments</b>                      Students consider the limitations of their data—for example, that their information probably is not sufficient to describe all the ladybugs in the world (make population inferences).</p>

266 Source: Adapted from GAISE II (Bargagliotti et al., 2020)

## 267 **Thematic Topics Within the CA CCSSM That Directly Support** 268 **Data Science**

269 Mathematics and statistics make up a significant portion of a data science education.  
270 Topics related to data primarily are found within two domains of the standards:  
271 a) Measurement and Data and b) Statistics and Probability. This section describes three  
272 thematic topics developed from the CA CCSSM that support data science and  
273 demonstrate how the topics progress across the grade bands. Thematic topics were  
274 created by applying the “Essential Understandings” of statistics described by the  
275 National Council of Teachers of Mathematics (NCTM) (Kader and Jacobbe, 2013; Peck  
276 et al., 2013) and aligned to the CCSSM and the statistics developmental levels from  
277 GAISE II (Bargagliotti et al., 2020).

278 The thematic topics are:

- 279 ● Understanding and describing variability in data and data distributions
- 280 ● Data collection, sampling, and random processes
- 281 ● Comparing distributions and identifying associations between variables

282 The subsequent sections of this chapter show how the thematic topics connect across  
283 grade bands to create the foundational knowledge for data science, beginning with  
284 simple counting and categorizing activities in kindergarten and culminating in high  
285 school where students integrate all these concepts during sophisticated investigations  
286 involving linear models and inferential statistics (see figure 5.2). In figure 5.2, the bullet  
287 points represent the CA CCSSM clusters from the Measurement and Data and  
288 Statistics and Probability domains. Additional details on how the individual standards  
289 can be implemented appear in the sections that follow, which are specific to each grade  
290 band.

291 Figure 5.2 Data-Focused CA CCSSM Content Clusters, Organized Into Thematic  
292 Topics That Span K–12

Grade Levels	Understanding and describing variability in data and data distributions	Data collection, sampling, and random processes	Comparing distributions and identifying associations between variables
K–5 Measurement and data	<ul style="list-style-type: none"> <li>Describe and compare measurable attributes</li> <li>Represent and interpret data</li> </ul>	<ul style="list-style-type: none"> <li>Classify objects and count the number of objects in categories</li> </ul>	n/a
6–8 Statistics and probability	<ul style="list-style-type: none"> <li>Develop an understanding of statistical variability</li> <li>Summarize and describe distributions</li> </ul>	<ul style="list-style-type: none"> <li>Use random sampling to draw inferences about a population</li> <li>Investigate chance processes and develop use and evaluate probability models</li> </ul>	<ul style="list-style-type: none"> <li>Draw informal comparative inferences about two populations</li> <li>Investigate patterns of associations in bivariate data</li> </ul>

**High School  
Statistics and  
probability**

- Summarize, represent and interpret data on a single count of measurement variable
- Summarize, represent and interpret data on two categorical and quantitative variables
- Understand and evaluate random processes underlying statistical investigation
- Interpret linear models
- Make inferences and justify conclusions from sample surveys, experiments, and observational studies
- Understand independence and conditional probability and use them to interpret data
- Use the rules of probability to compute probabilities of compound

- Summarize, represent and interpret data on a single count of measurement variable
- Summarize, represent and interpret data on two categorical and quantitative variables
- Understand and evaluate random processes underlying statistical investigation
- Interpret linear models
- Make inferences and justify conclusions from sample surveys, experiments, and observational studies
- Understand independence and conditional probability and use them to interpret data
- Use the rules of probability to compute probabilities of compound

- Summarize, represent and interpret data on a single count of measurement variable
- Summarize, represent and interpret data on two categorical and quantitative variables
- Understand and evaluate random processes underlying statistical investigation
- Interpret linear models
- Make inferences and justify conclusions from sample surveys, experiments, and observational studies
- Understand independence and conditional probability and use them to interpret data
- Use the rules of probability to compute probabilities of compound

Grade Levels	Understanding and describing variability in data and data distributions	Data collection, sampling, and random processes	Comparing distributions and identifying associations between variables
	<ul style="list-style-type: none"> <li>events in a uniform probability model</li> <li>• Use probability to evaluate outcomes of decisions</li> </ul>	<ul style="list-style-type: none"> <li>events in a uniform probability model</li> <li>• Use probability to evaluate outcomes of decisions</li> </ul>	<ul style="list-style-type: none"> <li>events in a uniform probability model</li> <li>• Use probability to evaluate outcomes of decisions</li> </ul>

293 Source: Adapted from the CA CCSSM

294 **Understanding and Describing Variability in Data and Data**  
 295 **Distributions**

296 Many important outcomes (e.g., health, wealth, education) vary in the world. Gathering  
 297 data provides a way to capture how these outcomes vary in order to understand the  
 298 causes of the variation. The patterns of variation that are seen in the data are called  
 299 distributions. Across the curriculum, students consider how their observed, counted, or  
 300 measured values and data characteristics might not be the same—that is to say, they  
 301 vary. The statistical work of understanding and describing variability provides a strong  
 302 footing for students to engage in the work of data science. In kindergarten through  
 303 grade five, it is essential that students encounter variation in a variety of ways, including  
 304 by counting, measuring, and observing quantities and characteristics that vary in order  
 305 to be prepared for more sophisticated work with statistics later. Elementary students  
 306 develop visualizations to show variability in data. Early elementary students begin with  
 307 creating picture graphs, showing one or more categories of data in whole units. By the  
 308 end of elementary school, students should have had experience with line plots and bar  
 309 graphs for data with three or four categories and experience with plots with scales in  
 310 fractions of units.

311 From sixth grade, students begin learning more formal methods to understand data and  
 312 to create models of variation. Students continue to produce data visualizations. They  
 313 learn to describe distributions by their overall shape (e.g., symmetric versus skewed) as  
 314 well as measures of center (mean, median, mode) and spread. This foundational work  
 315 is important for being able to compare distributions and identify associations between  
 316 variables beginning in seventh grade.

317 **Data Collection, Sampling, and Random Processes**

318 Data collection can underpin data science activities. As students look at the world  
 319 through a data lens, they might notice that data can take different forms. Data collected  
 320 and represented fall into two categories: categorical (non-numerical, or qualitative) data  
 321 and numerical or quantitative data. For instance, consider a set of colored blocks in the  
 322 classroom. Color is a categorical or qualitative variable that students could observe  
 323 about each block, while length is quantitative, a data point generated through  
 324 measuring. The standards focus on students developing understanding of categorical  
 325 data in kindergarten through grade three; in grade two, students begin to also learn  
 326 about measurement data. Figure 5.3 illustrates several examples of categorical and  
 327 quantitative data.

328 Figure 5.3 Examples of Categorical and Quantitative Data

Categorical Data	Quantitative (or Measurement) Data
<ul style="list-style-type: none"> <li>• Temperature (hot, room temperature, cold)</li> <li>• Color (red, green, blue, yellow) of blocks in the classroom</li> <li>• Species of trees at the school</li> <li>• Identification of schools in the district as “elementary school,” “middle school,” or “high school”</li> </ul>	<ul style="list-style-type: none"> <li>• Temperature (80F)</li> <li>• Pixel or RGB color values</li> <li>• Height (or circumference of trunk, or biomass) of trees at the school</li> <li>• Number of pages (or weight, or height) of books in the classroom</li> <li>• Annual income for households in a census tract</li> </ul>

329 As students pose statistical investigative questions, they should also encounter  
 330 opportunities to help determine how data might be produced to answer those questions,  
 331 and what forms of data would be best to use for producing the answers. In addition to



332 producing data directly through their own observations, students should gain exposure  
333 to designing and using surveys and simple experiments. By producing their own data  
334 from their classroom or community, students recognize data as having context and  
335 deriving from observation and measurement, and they come to see data (and  
336 mathematics more broadly) as a tool to help think about their worlds.

337 In seventh grade, students are introduced to the idea of random sampling and the idea  
338 that data collected from a subset of a population can help them understand the whole  
339 population. Students are also introduced to probability and chance processes in seventh  
340 grade, building theoretical probability models and conducting experiments to calculate  
341 long-run probabilities of chance events. Students should continue to develop their  
342 understanding of sampling and random sampling and probability models through eighth  
343 grade to prepare for work in high school.

## 344 **Comparing Distributions and Identifying Associations Between** 345 **Variables**

346 In elementary school and early in middle school, students are primarily working with  
347 data sets that include a single variable measured in a single population in mathematics  
348 and one to three variables in science. In seventh grade, students continue working with  
349 univariate data but begin informally comparing a single variable measured across two  
350 populations or at two points in time.

351 In eighth grade, students begin working with bivariate data: two variables measured in  
352 the same population. Students are introduced to the use of scatter plots to visualize  
353 bivariate data and depict how the two variables are associated. Students also begin  
354 working with informally fitting linear models to scatter plots that suggest a linear  
355 association and using those linear models to solve problems and make predictions.  
356 Although students' statistical explorations of linear models are informal in middle school,  
357 middle school work with expressions, equations, and functions is critical to preparing  
358 students for more formal use of linear models in high school.

359 In high school, students use technological tools to create a line of best fit and compute a  
360 correlation coefficient. At the high school level, students should be able to interpret the  
361 slope and intercept of a linear model and distinguish between correlation and causation.

362 In high school, students integrate their knowledge of random sampling, comparing  
363 distributions, and identifying associations into more complex statistical investigations in  
364 which they make inferences from data. Students begin asking whether observed  
365 differences between two samples could happen through random chance.

## 366 **Data Science in Each Grade Band**

### 367 **Transitioning From Pre-Kindergarten**

368 Before kindergarten, children begin to describe their world in language, identifying  
369 characteristics of objects, places, people, and events: *The ball is red. My classroom is*  
370 *warm. My teacher is old or young. Our trip to the park was too short.* Identifying  
371 characteristics is the beginning of using data and wondering about characteristics—  
372 including countable characteristics—is the beginning of asking questions that data can  
373 help to answer. In the California Preschool Learning Foundations, this content is located  
374 under the heading “Algebra and Functions (Classification and Patterning),” in which  
375 children “sort and classify objects in their everyday environment” (by one attribute at  
376 around 48 months and by more than one attribute at around 60 months of age); and in  
377 “Measurement,” in which students compare and order objects directly at around 48  
378 months of age and may use an intermediate object for comparison at around 60 months  
379 of age (Preschool Learning Foundations, Volume 1). These preschool activities directly  
380 enable the types of kindergarten through grade five learning trajectories described  
381 below.

### 382 **Kindergarten Through Grade Five**

383 Within the kindergarten through fifth grade band, a sense of curiosity about the world is  
384 a crucial first step in building an understanding of what data are and how they can be  
385 generated. All work with data should begin with noticing and wondering: “I notice that...”

386 or “I wonder what...” or “I wonder how many....” To prompt wonder, teachers can ask:  
387 “What do you notice or wonder about here [in this context] that we could  
388 [count/measure/keep track of] to figure out or explore further?” To establish effective  
389 routines, and to support language development in “I wonder” activities, it can be  
390 effective to provide these examples as sentence starters. Early data explorations might  
391 begin with students asking questions that can be answered with a single value, “How  
392 many students are there in our class?” or “How long is recess?” (SMP.1). With support  
393 from their teachers, students can also start to pose or explore statistical investigative  
394 questions that involve multiple variables, such as “I wonder if plants grow more with  
395 more sunlight?” Or “I wonder if age affects which color people like?” Questions guide  
396 much of students’ work with mathematics and include “those used to frame an  
397 investigation, those used to collect data, and those used to guide analysis and  
398 interpretation” (GAISEII, p. 44).

399 At the lower elementary level, students encounter data through exposure to small data  
400 sets or numbers that were collected manually through counting, classifying, comparing,  
401 or possibly measuring objects. Simple peer-questioning activities such as gathering  
402 answers to questions like “How many siblings do you have?” or “Was the sky clear or  
403 cloudy today?” engage students with basic data collection. Activities focused on a single  
404 attribute (e.g., number of siblings) also provide an opportunity for students to engage in  
405 the SMPs while the students represent data graphically, and these activities support  
406 looking for patterns in data, which is crucial work for any statistical investigation. Upper  
407 elementary students explore various types of data representations and use those data  
408 representations, generated by themselves or others, to describe the world around them  
409 (SMP.2). These experiences lay a critical foundation for mathematical and statistical  
410 thinking within middle school and are crucial steps in supporting data literacy.

411 As students gain confidence in their ability to communicate the mathematical ideas, the  
412 teacher should encourage students to generate questions themselves to build their  
413 agency in using mathematics to make sense of their worlds or to use their data to  
414 develop claims in response to their questions (SMP.3). For example, a weekly whole-  
415 class “I wonder” routine—in which students propose questions to investigate by

416 collecting data—contributes to students’ development of modeling with mathematics  
417 (SMP.4). Exploring data is an opportunity to help students see how mathematics can be  
418 used to make sense of problems or answer questions that are relevant to them  
419 (SMP.1).

420 Because the mathematical experiences that support data science increase in  
421 sophistication substantially between the early and later elementary grades, thematic  
422 topics are discussed separately for the kindergarten through grade two and grades  
423 three through five bands in the following sections.

## 424 ***Kindergarten Through Second Grade***

### 425 **Understanding and Describing Variability in Data and Data Distributions**

426 Students naturally encounter simple variation in their everyday lives through tasks that  
427 invite qualitative descriptions or comparisons, such as “The same kinds of plant are  
428 different sizes” or “The contents of our lunch boxes differ.”

429 Once students are invited to notice things in a context and wonder about a question,  
430 they begin to describe measurable, countable, and observable attributes of objects or  
431 situations (K.MD.1, K.G.1, K.G.4) and classify objects and count the number in each  
432 category (K.MD.3), such as categorizing a set of cubes by color. Sorting objects into two  
433 or three categories and representing these categories by their count (K.MD.3, 1.MD.4)  
434 are early examples of students representing data to help make sense of their worlds  
435 (SMP.4). Basic summary statements of objects pulled from a bag—“The shape is  
436 square” and “This cube is red” (categorical data) or “There are 13 red cubes in the set”  
437 (numeric data)—represent early work that builds toward an understanding of variability.  
438 Notably, most of the focus on *number* in kindergarten and first grade should be with  
439 numbers representing quantities of counted objects (SMP.2). Sorting and categorizing  
440 activities help students recognize that objects can naturally vary and help students  
441 develop language to express this variability, such as “We have three different shapes—  
442 squares, triangles and circles” or “Our red cubes come in two different sizes.” Students

443 also begin to use total counts to describe their observation—e.g., “When we pulled  
444 shapes from a bag one at a time, 12 were square and 6 were circles.”

445 Many opportunities to explore variation arise as students measure time to the nearest  
446 five minutes (2.MD.7) and measure length to the nearest whole unit (2.MD.9), using  
447 different standard units (centimeters, meters, inches, feet) (2.MD.3) and several tools  
448 (2.MD.1). They might recognize that their measurements are not always the same.  
449 Working with data collected in tables and in visualizations provides students an  
450 opportunity to explore questions such as, “For the objects measured, what was the most  
451 common length in inches? What was the smallest (minimum) or largest (maximum)  
452 object? What is the difference between our largest and smallest object (range)?” By  
453 second grade, students begin to expand their focus on data representation, being  
454 introduced to line plots (whole number units only; 2.MD.9), picture graphs, and bar  
455 graphs. These graphs can be used to answer put-together, take-apart, and compare  
456 questions (2.MD.10).

### 457 **Data Collection, Sampling, and Random Processes**

458 Data investigations should be investigative and collaborative, with students working  
459 together to learn about and describe the world around them. Collecting data through  
460 measurements, surveys, and experiments is an important part of the statistical  
461 investigative process and supports young learners in building their awareness of what  
462 data are and where data come from. Simple classroom polls provide opportunities for  
463 early elementary students to work with simple addition and subtraction equations to  
464 express relationships between the collected student responses. For example, the  
465 question “How many people took the bus today compared to yesterday?” requires  
466 students to collect data and consider how, and perhaps why, quantities might change  
467 from day to day. Simple surveys in the form of interviews help students practice  
468 expressing counts verbally and symbolically and provide opportunities for students to  
469 communicate with each other about their data.

470 Data collection tasks in the early elementary grades are usually constrained to the  
471 context of the classroom. When choosing data tasks, it is important to consider the

472 grade-level expectations for counting (up to 10 objects scattered or up to 20 if arranged  
473 in a line, array, or circle in kindergarten [K.CC.5]; 120 by the end of first grade  
474 [1.NBT.1]; and up to 1,000 by the end of second grade [2.NBT.2]). Counting tasks can  
475 also be structured to build understanding of place value.

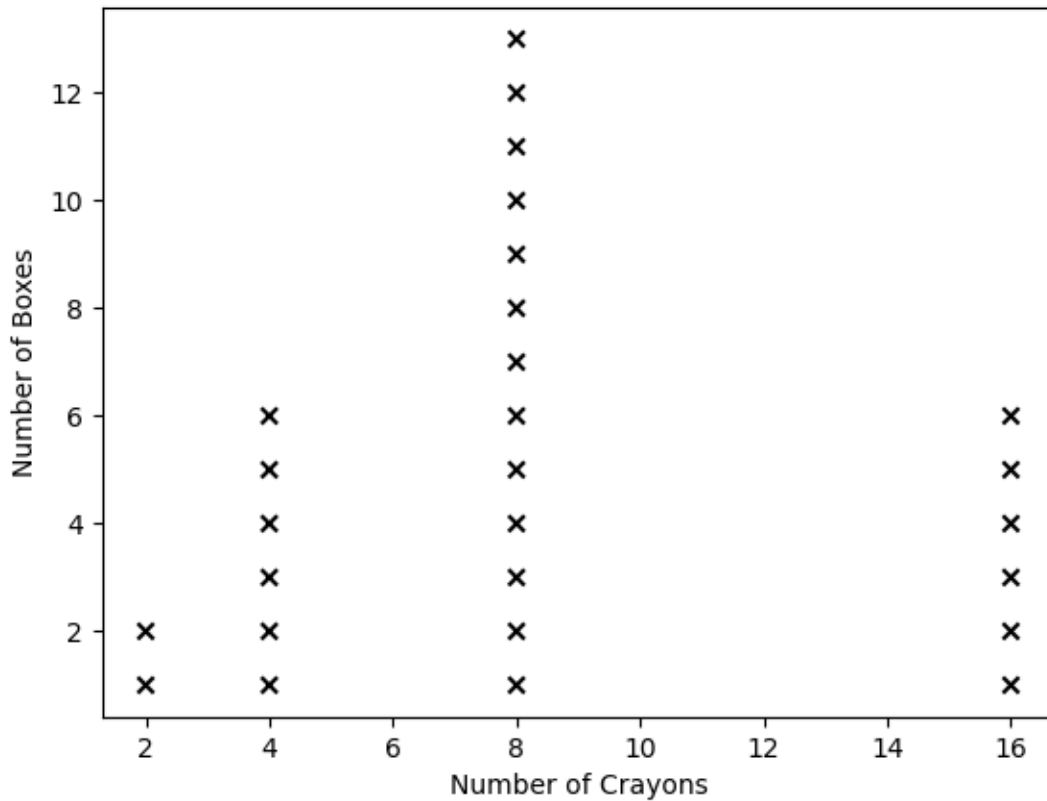
### 476 **Comparing Distributions and Identifying Associations Between Variables**

477 Within kindergarten through second grade, students use language, counts, and  
478 measures to describe and compare objects. Students compare the numbers of objects  
479 in different categories (K.CC.6) to answer “Which has more?” questions (e.g., “I wonder  
480 whether there are more square blocks or more triangular blocks on the desk?”). At first,  
481 the teacher suggests or specifies categories; eventually students generate ideas for  
482 classification. They also directly compare objects (rather than measuring each with a  
483 unit or an intermediate object) with common measurable/countable attributes to see  
484 which has more (K.MD.2, K.G.4) (“I wonder which shape has more sides? Which kind of  
485 block is heaviest?” [answering this question by using a balance or an informal one-in-  
486 each-hand comparison rather than a scale]). “I wonder...” questions should explore two-  
487 category “Which is more?” questions as well as comparisons of objects according to  
488 length, height, weight, and countable attributes like number of sides. Student-generated  
489 questions provide opportunities to work on precision of language as well; for example,  
490 by asking students to clarify what they mean by “bigger.” Mathematics discussions that  
491 are rooted in academic language can help students understand mathematical concepts  
492 more deeply and discover new ones.

493 To explore foundational ideas of distribution, students might explore the variety of colors  
494 offered in crayons or markers in their classroom. To do so, students collect the counts  
495 of crayons or markers in each of the boxes, choosing from different methods of  
496 counting, such as using dots, arrays, or tallies. The counts are reported to the teacher,  
497 who constructs a basic dot plot, graphically representing the class counts, as shown in  
498 figure 5.4. Students can use the counts to first describe the shape of their data by  
499 describing where they see “hills” or where the data form a group or crowd (cluster).  
500 Students can use their counts to find the smallest count and largest count (range) and  
501 to describe the most common crayon box size. Further extensions might include asking

502 which box sizes occur more often or asking students to visually estimate where the  
503 middle of the plot would be. These early activities help expose students to key ideas  
504 (median and frequency). The teacher might ask the students to use the data to predict.  
505 “Imagine if we put all the crayon boxes into a large black bag. If we closed our eyes and  
506 reached in, which crayon box size (2, 4, 8, or 16) do we think we would grab first?”

507 Figure 5.4 A Teacher’s Dot Plot of the Data to Determine the Most Common Crayon  
508 Box Size



509

510 [Long description of figure 5.4](#)

511 The following elementary school snapshot illustrates many of the ideas discussed in this  
512 section. Experiences related to thematic topics (*understanding and describing variability*  
513 *in data; data collection, sampling, and random processes; and comparing distributions*

514 *and identifying associations between variables*) are included in parentheses where  
515 relevant.

516 ***Snapshot: Logan's Early Elementary School Explorations with Data***

517 In first grade, student teams were asked to think of two similar things at school for which  
518 they were not sure which was taller and then to find a way to compare the objects'  
519 heights. A variety of materials was available to use in the comparison (*data collection*).  
520 Logan's team was able to compare the height of the slide in front of the school with the  
521 height of the slide behind school, measuring the height of both using towers of large  
522 DUPLO® bricks. The whole class used their data to discuss how much taller the slide in  
523 front of the school was compared to the one in back. Afterward, Logan wanted to build  
524 DUPLO® towers to measure height and length of lots of things and was disappointed  
525 that the class did not have enough bricks to measure the height of the school (and that  
526 their teacher would not let them climb the school).

527 In another class activity, students recorded the length of each day (sunrise to sunset) by  
528 looking at the weather station in the main office. The class maintained a visible running  
529 tally of the number of school days with less than 11 hours of daylight, 11 to 13 hours,  
530 and more than 13 hours for the entire year (*understanding and describing variability in*  
531 *data*). The class then discussed what the students thought might happen to the number  
532 of hours of daylight in the future and checked the data a month later to see whether  
533 their predictions were correct.

534 The students in Logan's second grade class made their own yardsticks by marking a  
535 blank wooden rod in inches, using only a three-inch by five-inch card to measure the  
536 marks. The class then used their yardsticks extensively to measure objects of interest to  
537 the nearest inch. Later, they added centimeter markings to the other side of the  
538 yardstick and discovered that measuring the same things with smaller units led to larger  
539 number measurements and improved the quality of data (*data collection*).

540 When choosing an activity for measuring time, Logan's group decided to time and  
541 record the amount of time in a week that team members spent reading in school and



542 then to compare those measurements over several weeks. (This activity had the benefit  
543 that team members read much more during those weeks!) Other teams measured time  
544 spent playing outside, listening to announcements, and working at math stations (*data*  
545 *collection*). Teams made line plots of their data and compared the line plots of different  
546 activities to discuss how students typically spend their school time (*understanding and*  
547 *describing variability in data, and comparing distributions*).

548 (*end snapshot*)

### 549 **Grades Three Through Five**

#### 550 **Understanding and Describing Variability in Data and Data Distributions**

551 Data collected from observations and measurement often vary; that is to say, the values  
552 reported are not identical. Variability is a term used to describe how much the values  
553 differ from each other or are consistent. If the lengths of 10 NBA basketball courts are  
554 measured, the values would be expected to be very nearly identical, but if the numbers  
555 of pages in third grade math textbooks from different publishers were counted, these  
556 values would be expected to differ. Foundational work in variability begins in elementary  
557 mathematics when students count, measure, observe, and describe their data. The use  
558 of simple plots to identify patterns is an integral part of preparing students for the  
559 statistical concepts in variability that are covered in grades six through eight.

560 When working with visualizations of data, students not only should consider the most  
561 popular value in a data set (the mode) but also should describe the shape and spread of  
562 data distributions. Identifying the maximum and minimum values of quantitative data  
563 sets can help students appreciate the concept of range as a measure of spread, and  
564 looking for clusters and gaps in a distribution can begin to help them attend to the  
565 shapes of data sets. As students engage in experiences in which they produce their  
566 own data through measurement, teachers should highlight for students the variation that  
567 results. Measuring the same variable on multiple individuals or objects, for example,  
568 results in data that vary, and students should consider the causes or sources that might  
569 have given rise to the variation they have observed, working as they do so to

570 differentiate between variation and error. For example, if students plant a particular  
571 variety of flower seed at multiple locations around the school, then measure the plants'  
572 height and the amount of sunlight each month, they can conduct investigations into the  
573 ways that plant growth and sunlight relate to each other. They should discuss and  
574 describe any patterns in their data and discuss reasons for the variability. Upper  
575 elementary students should have the opportunity to represent their data through plots  
576 that they themselves create. This process helps students notice variation within the data  
577 as well as communicate their thinking in multiple ways.

578 Students in grades three through five refine their measurements of length and time and  
579 expand the set of units they use, adding area and volume measurement to their  
580 repertoires. By the fifth grade, students should understand that data sets can include  
581 both categorical and numerical data. They should recognize that an individual instance  
582 or object can possess attributes that exemplify these different types of variables, and  
583 they should have gained experience measuring, characterizing, and analyzing such  
584 diverse types of data and associating them together. Mathematical and scientific work  
585 that can reveal variability of measured dimensions, mass, and volume present natural  
586 opportunities for students to explore variation in their different measurements  
587 graphically—such as in a common classroom activity in which students compare the  
588 mass or volume of objects to other dimensions.

589 The following snapshot on Logan's fifth grade exploration with data illustrates many of  
590 the ideas discussed in this section. Experiences related to thematic topics  
591 (*understanding and describing variability in data; data collection, sampling, and random*  
592 *processes; and comparing distributions and identifying associations between variables*)  
593 are included in parentheses where relevant.

594 ***Snapshot: An Example of Logan's Fifth Grade Exploration with Data***

595 By fifth grade, Logan and classmates had constructed many line plots and thus often  
596 wondered about quantities that vary on repeated measurement, such as the following:

597       • The cartons of milk from lunch say they each contain 8 fluid ounces, but yours  
598       feels heavier than mine. Does my container have less milk? Are you getting more  
599       milk than me?

600       • The weather site says the average high temperature here is 57°F (degrees  
601       Fahrenheit) in November, but today it got up to 65°F. How can we check whether  
602       this month is near average?

603       To explore the first question, the school donated 20 cartons of milk to the experiment so  
604       students could measure the volume (*data collection*). When they examined the line plot  
605       of their measured volumes, they saw that it had a tightly clustered shape, with a  
606       minimum measurement of 7.8 fluid ounces and a maximum of 8.2 fluid ounces, and that  
607       the most frequent value was 7.9 ounces (*understanding and describing variability in*  
608       *data*). One student in the group thought that some milk probably remained in the  
609       containers, so the group spent a while trying to figure out how they might identify how  
610       much had been left inside. The teams came up with several methods, laying the  
611       groundwork to talk about random measurement.

612       For the second question, the class recorded the daily high temperature for each day of  
613       the month, recorded these temperatures on a line plot which also had marked the  
614       “average” high temperature from the weather site, and used the line plot at the end of  
615       the month to discuss whether their measurement was consistent with the stated  
616       average (without computing an average of the data).

617       Fifth grade does not extend the expected set of data representations, but students do  
618       use line plots in a sophisticated way that sets the stage for understanding the most  
619       common measure of center for a data set—the mean (commonly called the average)—  
620       in sixth grade. Namely, fifth grade students use a line plot to decide how a repeatedly  
621       measured quantity could be redistributed equally (5.MD.2): “Given different  
622       measurements of liquid in identical beakers, find the amount of liquid each beaker  
623       would contain if the total amount in all the beakers were redistributed equally.”

624       (*end snapshot*)

625 Although the data visualizations mastered by fifth grade include only picture graphs, bar  
626 graphs, and line plots, students do not need to be restricted to these. Each of these  
627 represents repeated measurements of a single varying quantity; science curricula, and  
628 many questions of interest in general, require the consideration of relationships between  
629 *two or more different* changing quantities, such as erosion and time (NGSS 4-ESS2-1  
630 Earth's Systems) or length or direction of shadows and time (NGSS 5-ESS1-2 Earth's  
631 Place in the Universe). Reasoning that involves such multiple variables is an important  
632 aspect of modern encounters with data, and students should experience this kind of  
633 reasoning at all levels (SMP.2). These science investigations represent an excellent  
634 opportunity to compare distributions between variables by posing questions such as  
635 "How does the shadow length change between fall and spring?"

636 In recent years, new technological tools and developments have prompted an explosion  
637 in interesting data visualizations, many of which are quite comprehensible to young  
638 students with some exploration. Technology and the power of computing play an  
639 important role in data science and can be incorporated into the kindergarten through  
640 grade five experience. The ability to use technology to collect, organize, represent, and  
641 share data is fundamental to the development of data literacy. California's 2018  
642 Computer Science Standards include computer-based data sorting, categorizing, and  
643 visualizing for students in kindergarten through grade two and for grades three through  
644 five (CS K–2.DA.8, K–2.DA.9, 3–5.DA.8). Working toward these standards is important  
645 preparation for middle and high school use of data software to visualize and interpret  
646 large data sets. Experiences with different types of visualizations will further expand  
647 students' sense-making opportunities and encourage them to think about what they can  
648 understand by looking at data sets in different ways (SMP.4). Newspapers and online  
649 news sources offer specific examples; student-gathered examples help to build buy-in  
650 for a "Can we figure out what this visualization is trying to help us understand?" routine.

### 651 **Data Collection, Sampling, and Random Processes**

652 Remaining alert to student wonderings about their everyday experiences—perhaps in  
653 attendance, weather, or lunch-count data—may generate opportunities for the class to  
654 explore how the collection of data can help answer questions asked by the class or the

655 teacher. In addition to their own observations, students should gain exposure to  
656 designing and using surveys and simple experiments as ways to collect data. By  
657 producing their own data from their classroom or community (“How does age of  
658 students relate to their enjoyment of school? Does time on social media apps increase  
659 with age? How much waste is generated by different companies/our school?”) students  
660 recognize data as having context and deriving from observation and measurement, and  
661 they come to see how mathematics and data are tools to help think about their worlds  
662 (SMP.4).

663 As students seek data to address authentic questions similar to those described above,  
664 they should also encounter opportunities to help determine how data might be produced  
665 and to consider how their choices might impact the data. When conducting classroom  
666 surveys, students can begin to grapple with very basic ideas of fairness, laying a  
667 foundation for sampling. For example, “Is it fair to interview only my friends?” or “Should  
668 I measure only my tallest seedling?” Also, they should consider their own measurement  
669 techniques and how confident they are that all of the students measured the same way  
670 (so that if someone else measured, such as for height or amount of sunlight, they would  
671 get the same results) (SMP.6). Students should also encounter data collected by other  
672 people for a similar purpose.

673 Students at the elementary level often express informal wonderings about probability,  
674 randomness, and uncertainty. For example, “How likely is it that it rains when we have  
675 recess?” or “Can we predict who will come through the door next or what color cube we  
676 will draw out of a bag?” Randomness is a complex idea encompassing uncertainty and  
677 a level of predictability. When (blindly) drawing a cube out of a bag containing three blue  
678 cubes, two red cubes, and one yellow cube, nobody can predict with certainty what will  
679 happen on a single draw. But, over many draws, the person who always predicts a blue  
680 cube will be right about half the time. Activities that demonstrate this concept can be  
681 used to generate data for many of the explorations of the thematic topics above, which  
682 will leave students well-prepared for a more formal treatment of randomness and  
683 probability in middle school. At this point, students should begin to conceive of  
684 probability as a general measure—e.g., not likely, likely, very likely chance that

685 something will happen—and should see it as a basic measure of certainty or  
686 uncertainty.

687 Interpreting data is a matter of making inferences from the data available. Although  
688 students will encounter quantitative and nuanced techniques for making inferences in  
689 later grades, they should nevertheless encounter opportunities to make claims and infer  
690 conclusions across their kindergarten through grade five years (SMP.3). When they do,  
691 students should learn to wonder whether patterns or trends they notice extend to larger  
692 populations (including considering ways in which a group might not be representative of  
693 the larger population). Additionally, students should learn that good claims draw upon  
694 data as evidence and that they always come hand in hand with a degree of uncertainty.  
695 Modeling the use of appropriate terminology such as “tends to,” “typical,” “usually,” and  
696 “similar” can help lay important groundwork for this concept (Rugin, 2019).

697 Upper elementary students begin to reason more abstractly and work toward using all  
698 four operations to solve problems. The classroom, home, and community present  
699 meaningful opportunities for students to apply tools to measure and describe the world  
700 around them, including collecting data for one or more attributes. For example, students  
701 may be challenged to discover which location (inside or outside the classroom) has the  
702 “best” types of tables for collaborative group work. As a part of this task, students  
703 explore the idea of “best” and discuss features such as size of the tabletop, height of the  
704 table, and shape. The teacher guides the students to consider which attributes of the  
705 table could be measured and then provides a template for student pairs to collect their  
706 observations and measurements—e.g., of width, height, and shape of the table. After  
707 collecting data, students notice that some of the table shapes were hard to measure  
708 because of their unexpected shape—e.g., trapezoids, kidney beans, and circles.

709 The teacher guides a discussion by asking students what was fun, easy, or hard about  
710 the data collection process. While students share, the teacher tracks some of the  
711 challenges and prompts students to brainstorm reasons they encountered the  
712 challenge, as shown in figure 5.5.

713 Figure 5.5 Tracking Challenges and Reasons for Challenges

Challenge	Reasons for the Challenge
We measured the same type of table but found different widths or heights.	Maybe we measured wrong. We measured different units. We didn't round up or down in the same way. Some of us used 1/2s, like 5.5.
We only measured certain tables.	We didn't know what to do for circle or trapezoid tables. The shape was weird (irregular). We ignored certain tables. Tables are packed away in the cafeteria or were being used in the library so we couldn't measure them.

714 The process of collecting and working with data helps students recognize that data  
715 collection is a human endeavor and includes decisions and sometimes errors made by  
716 people. Activities such as the tabletop area exercise help students begin to develop a  
717 list of features important when designing surveys or experiments. Is it fair to choose the  
718 biggest tables in one room and then only the smallest tables to measure in the next  
719 (randomness)? What happens when someone designs a data collection method but it  
720 doesn't work in the real world—for example, the large tables in the cafeteria were stored  
721 away? How many tables (number of cases or sample size) should be measured?  
722 Statistical investigations which expose students to data that are sometimes messy or  
723 which involve a process that can be ambiguous (SMP.1) are crucial to data science.  
724 Students can also come to recognize that they can ask questions about the data that  
725 are given to them—e.g., questions about how or when the data were collected and for  
726 what purposes. As students ask questions, the data required to answer them may not  
727 be accessible or possible to collect firsthand. Data gathered by others (such as other  
728 students in the discussion) can help to answer questions students generate about their  
729 own communities and can open up discussion about randomness and probability.

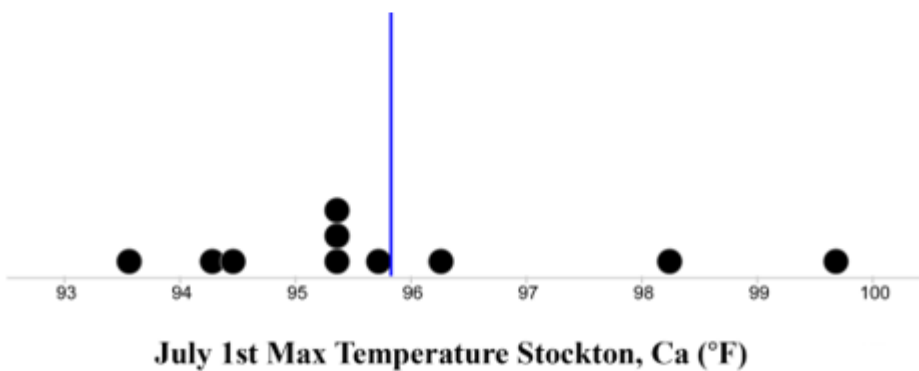
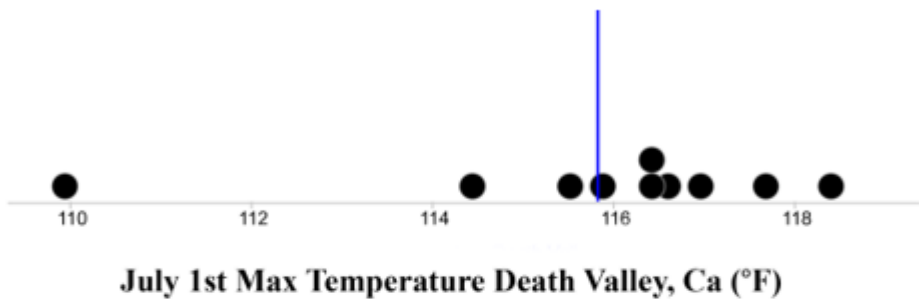
730 **Comparing Distributions and Identifying Associations Between Variables**  
731 In kindergarten through fifth grade, students are invited to ask questions about the world  
732 around them, especially about objects inside the classroom. When students compare  
733 measurements and frequently reported values between groups, they are engaged in

734 tasks that lay the groundwork for more sophisticated comparisons. Invitations for  
735 students to describe their data and make predictions about why the values might vary  
736 are important early opportunities. Sentence starters, tally tables, dot plots, line plots, pie  
737 charts, and bar graphs are all important tools that help students to describe patterns  
738 within data, especially when making comparisons between groups.

739 Notably, many questions that students might wonder about in science and other fields  
740 will not be fully answerable using the tools and mathematical understanding available to  
741 them in kindergarten through grade five. It is important that teachers have resources for  
742 helping students figure out which aspects of questions can be investigated with  
743 currently available tools and that teachers have some understanding of technological  
744 tools which students will encounter later. For example, many students will wonder about  
745 relationships between two different variables: “If I get up earlier, do I feel tired earlier in  
746 the afternoon at school? Do students who skip lunch eat more candy in the afternoon?”  
747 When one of the variables is categorical (like for the skipping lunch question), separate  
748 line plots can be made for each category and the line plots compared. When both  
749 variables are quantitative, students can input data into an online graphing and  
750 visualization tool such as CODAP (The Common Online Data Analysis Platform),  
751 Desmos, or TinkerPlots, and then investigate the relationships by plotting their data on  
752 graphs, observing their distributions, and adding line plots. Another option is that one of  
753 the variables can be made into a categorical variable by defining categories in terms of  
754 the quantitative variable. For instance, waking-up times could be classified into “early”  
755 and “late” (ideally with a student-generated cut-point between early and late) and then  
756 dot plots of “time in the evening when I felt tired” created for each category. Science  
757 investigations can provide opportunities for students to compare mean values between  
758 two locations or experimental conditions, as shown in figure 5.6. As in the figure, the  
759 use of stickers to create line plots with different symbols or graphing programs can help  
760 with making plots quickly and easily. Note that dot plots are not formally introduced until  
761 sixth grade.

762 Figure 5.6 Temperature Plots to Compare Mean Values for Two Cities in California





763

764 [Long description of figure 5.6](#)

765 Source: Generated with CODAP NOAA Plugin

766 As students are invited to ask questions through data and measurement, teachers  
 767 should be mindful of the types of comparisons being generated. Questions such as  
 768 “What time will it be when the next person walks into the classroom?” or “Which book in  
 769 the classroom is the most read?” compare events or objects within a shared space and  
 770 are generally preferred. Questions and data collection tied to personal characteristics  
 771 (“Who is the shortest in our class?”) or that serve as potential markers for economic or  
 772 social status (“What brand of shoes is the most popular in our class?”) usually should be  
 773 avoided.

774 Grades three through five provide opportunities to investigate questions using data that  
 775 should include volume and mass measurement (grams, kilograms, and liters, but not  
 776 compound units such as  $\text{cm}^3$ ) in addition to the length, time, and money contexts from  
 777 earlier grades (3.MD.2). Time measurements are refined to the nearest minute (3.MD.1)  
 778 and length now includes half- and quarter-inches (3.MD.4). Increased ability to report

779 lengths more precisely helps students begin to notice that some data can fall into  
780 specific counts (discrete) while other types of data (measuring length of objects in  
781 millimeters) are continuous. A significant context for data-investigation questions is  
782 classification and analysis of two-dimensional shapes (4.G.2). Incorporating this  
783 geometry standard to help build data understanding can foster the important practice of  
784 analyzing by attributes—one instance of SMP.7 (Look for and make use of structure).

785 In this grade band, students extend the set of units they work with (4.MD.1) and can  
786 generate data about area for more complex shapes. Fifth graders deepen their  
787 understanding of volume to include unit cubes, making this an important context for  
788 data-inquiry questions. For example, a teacher could invite students to build a structure  
789 out of multi-link cubes and then collect data from the class by inquiring into how many  
790 cubes they used, the height and width of their structures, or which colors they used.  
791 Invitations to collect data on multiple variables produce data sets that allow for students  
792 to compare measures across plots, such as comparing the average amount of time it  
793 takes for students to walk to school versus drive to school.

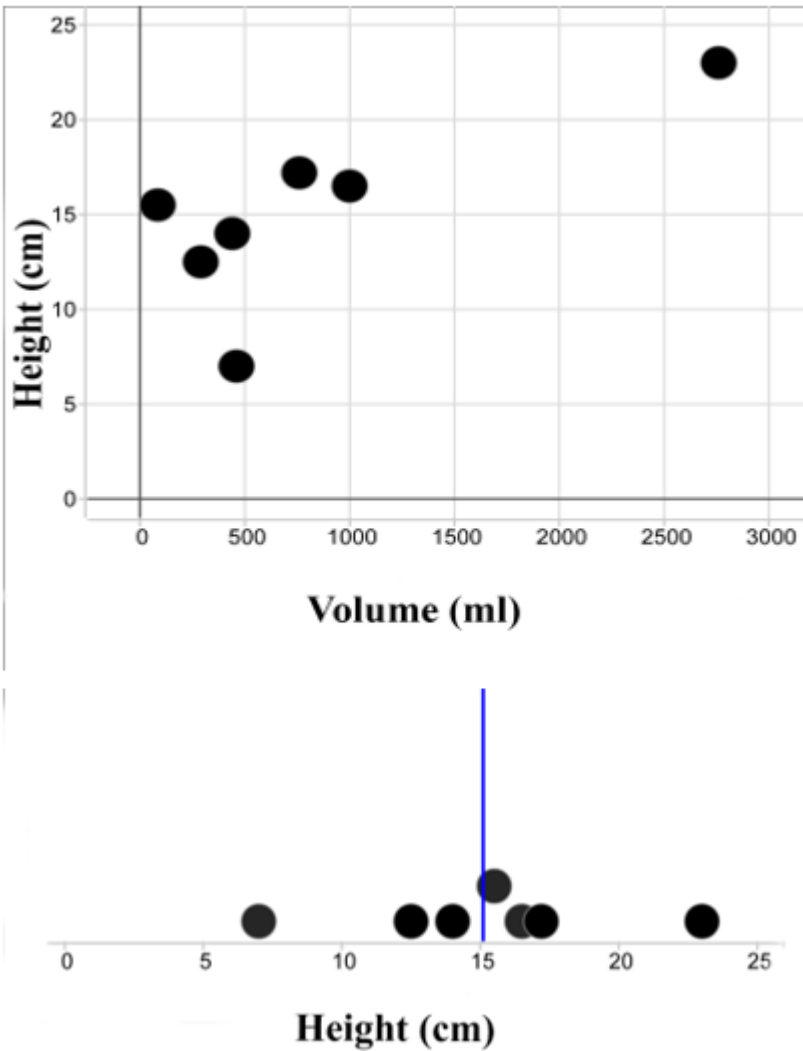
794 The snapshot of Logan’s third and fourth grade explorations with data describes two  
795 possible encounters a student in this grade band might have with the thematic topics  
796 (*understanding and describing variability in data; data collection, sampling, and random*  
797 *processes; and comparing distributions and identifying associations between variables*)  
798 and are highlighted in parentheses where relevant.

799 ***Snapshot: Logan’s Third and Fourth Grade Explorations with Data***

800 In third grade, as mass and volume became characteristics to measure, Logan’s class  
801 used length, height, mass, and volume measurements they had collected to examine  
802 sets of objects (*data collection*). In the science corner, the line plots of the masses  
803 looked quite different from the line plots of the lengths/heights of the objects, as did the  
804 line plots of volume, height, and mass of all objects in the room which hold water  
805 (vases, cups, etc.) (*understanding and describing variability in data and data*  
806 *distributions; see figure 5.7*). Logan’s team had a great disagreement about whether a  
807 taller vase should hold more water than a shorter vase (*comparing distributions and*

808 *identifying associations between variables*); the class eventually decided that it was  
809 usually but not always true that taller vases hold more water.

810 Figure 5.7 Logan’s Vase Measurement Data Visualized in CODAP



811

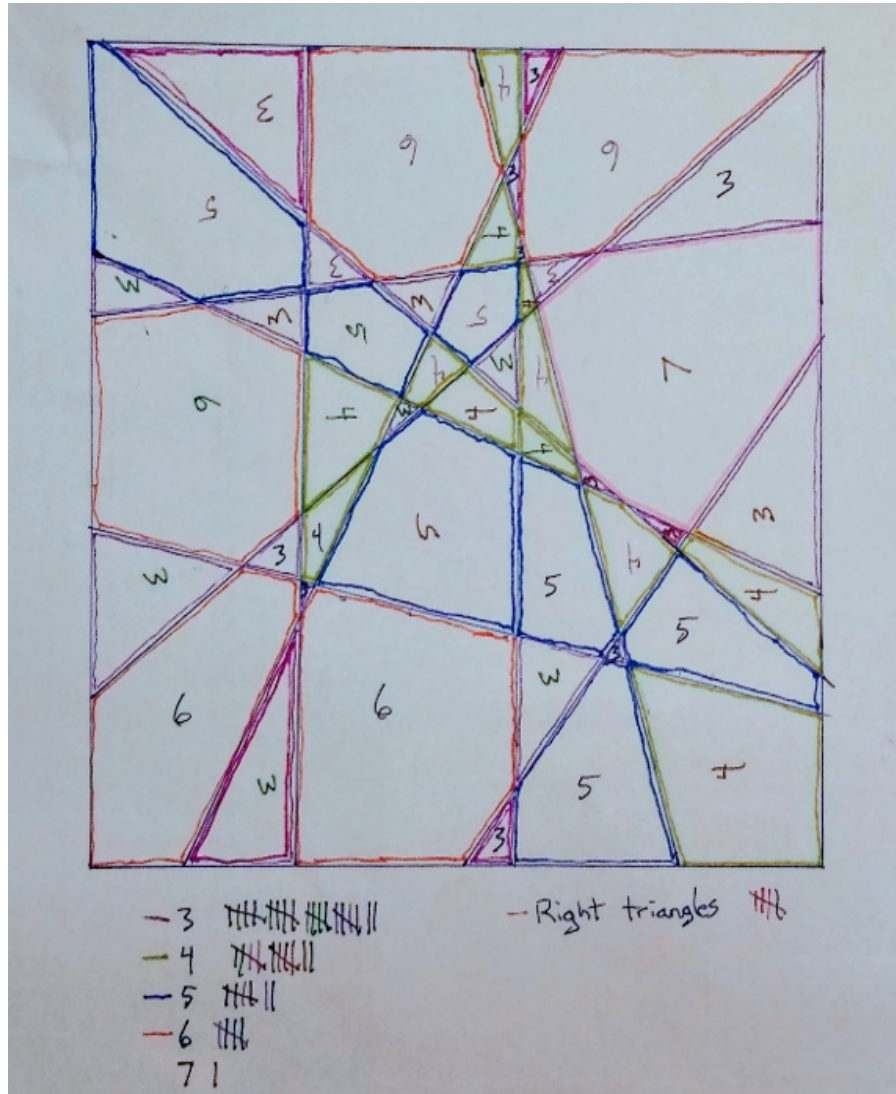
812 [Long description of figure 5.7](#)

813 Source: Collected student data in CODAP

814 One of Logan’s favorite activities in fourth grade was one that combined data work with  
815 classifying shapes by attributes: creating collaborative art pieces. For this activity, each  
816 team had a 1/2-meter by 1/2-meter square on the board, and each student in the team  
817 drew in two edge-to-edge straight lines of their choice, using their meter sticks. Then

818 one student in class chose a shape to try to find in the drawings, and each team  
819 outlined each new instance of that shape they found and described how they knew it  
820 was a triangle, rectangle, right triangle, quadrilateral, etc.; this process was repeated for  
821 several other shapes. Team members collected data by making an individual card to  
822 represent each piece of artwork (as shown in figure 5.8), using the card to represent the  
823 different variables they measured for each piece (how many triangles, how many  
824 instances of each color, how clean or messy each line was, etc.). When they had made  
825 a full set of cards, they sorted them in various ways, then made a table to compare the  
826 tallies for the different pieces (*understanding and describing variability in data and data*  
827 *distributions*), discussing the different features of the art and the process of creating it  
828 that might help explain the variations in their data.

829 Figure 5.8 Using Data to Classify Shapes



830

831 [Long description of figure 5.8](#)

832 (end snapshot)

### 833 **Grades Six Through Eight**

834 As in earlier grades, students in grades six through eight can understand their world via  
 835 a process that begins with wondering questions. This grade span is also the beginning  
 836 of when students experience the mathematical modeling cycle (Pelesko, 2015) and  
 837 investigations in science (NGSS Lead States, 2013). In middle school, students develop  
 838 a formal understanding of several key ideas in statistics, including describing  
 839 distributions and variability in data and random processes. Students begin informally

840 comparing and identifying associations in eighth grade in preparation for work in high  
841 school that develops a more formal understanding of linear models and statistical tests.

842 At the middle school level, students should encounter data sets that are small (a few to  
843 a dozen to a few hundred data points) and, when possible, can encounter larger data  
844 sets that contain thousands of data points. Many statistical concepts are more intuitive  
845 and accessible when illustrated with large data sets.

846 The following sections provide examples of how the same concepts can be illustrated  
847 with “little data” versus “big data.” Working with bigger data sets requires the use of  
848 computational tools, some of which may require programming skills. California has  
849 adopted K–12 computer science standards which can be consulted to determine what  
850 level programming is appropriate for middle school and high school (California State  
851 Board of Education, 2022). Working with complex data sets provides students with  
852 opportunities to engage in multivariate explorations, conduct simulations, and quickly  
853 create plots to reveal patterns—all nearly impossible to do by hand. Knowing when and  
854 how to leverage the power of computational tools is a crucial skill in data science.  
855 Students and teachers will need additional support with selecting and using these tools.

### 856 ***Understanding and Describing Variability in Data and Data Distributions***

857 Sixth grade students build on earlier experiences by distinguishing between statistical  
858 questions that can be investigated using data that varies (e.g., analysis of social media  
859 usage by age of students) versus questions for which there are no variations in (correct)  
860 responses (How many days are there in January?) (6.SP.1). When considering a  
861 statistical question, they understand that the variation in numerical data has a  
862 distribution which can be described by its center (first the median, then the mean); by its  
863 variability (also called spread, which is described both qualitatively and via a numerical  
864 measure—either interquartile range [IQR], range, or mean absolute deviation); and by  
865 an overall shape (including descriptors such as symmetric, skewed left or right, peak,  
866 gap, and outlier) (6.SP.2, 6.SP.3). As students explore data sets, they can produce  
867 visual representations of the distributions of their data; they can look at the shape of  
868 distributions that have different measures of center and spread and can develop visual

869 understandings of the shape of distributions. In sixth grade, visual representations of  
870 distributions include box plots and histograms, adding to the line plots (called dot plots  
871 from grade six onward) from earlier grades (6.SP.4). In addition, students learn to report  
872 and interpret measures of center and variability, and descriptions of distributions, in the  
873 context in which the data arose (6.SP.5.d).

874 Students should have experiences, beginning in sixth grade, deciding which measure of  
875 center is a more useful descriptor of a typical value for data sets with different shapes.  
876 Because the mean is sensitive to extreme values, the median is often a more useful  
877 measure for skewed distributions; in this case, the interquartile range is a useful  
878 measure of variability. For some distributions—with multiple clusters, for example—  
879 students may decide that neither median nor mean is a useful measure and might  
880 decide that a single number cannot reasonably represent a typical value (6.SP.5).

881 The following snapshot illustrates several themes discussed in this section, with  
882 thematic topics included in parentheses where relevant.

883 ***Snapshot: Óscar’s Visual Proof for Finding a Mean***

884 Óscar did not enjoy learning about mean, median, and mode. He often confused the  
885 different measures and felt they had little meaning. His parent contacted Maria, his  
886 teacher, to let her know that Óscar had been expressing frustration about the meaning  
887 of the terms since his last assessment. Óscar was not alone; Maria knew that many of  
888 the students were still struggling with the meanings of these measures of average.  
889 Based on results from an electronic, anonymous “exit ticket” survey used as formative  
890 assessment, Maria approached the students with the idea to build physical models so  
891 they could experience the averages in visual and physical ways, encouraging important  
892 brain connections.

893 Maria gave her students cubes and asked them to make six different towers of cubes  
894 that represented the numbers 1, 6, 3, 2, 4, and 2. She asked them how they might  
895 construct a physical proof to show the mean of the numbers. Some of the students were  
896 able to calculate the answer; however, she kept pushing them to build a visual proof

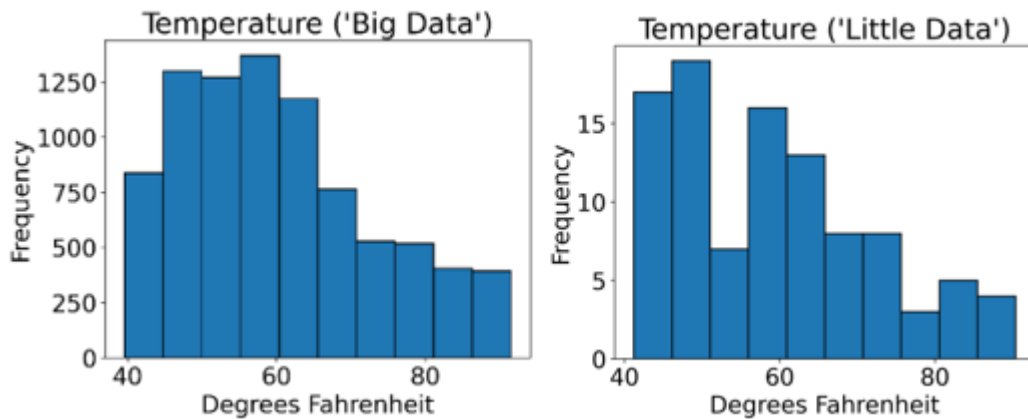
897 while remaining open to multiple means of representation. This strategy, based on  
898 specific UDL guidelines, allowed Maria to ensure scaffolds and supports would exist to  
899 help highlight the patterns of language and draw on background knowledge to express  
900 what students know in ways that are authentic and meaningful. Óscar and his group  
901 members came up with the idea of moving the cubes from tower to tower to show that  
902 they could make six towers that were all the same height. They just needed to average  
903 out all of the blocks (*understanding and describing variability in data and data*  
904 *distributions*). Óscar and his group excitedly explained to the class how they had made  
905 a physical proof of finding the mean of the blocks (*understanding and describing*  
906 *variability in data and data distributions*). They shared the calculation with the class and  
907 compared it to the method they used of moving the blocks. After her students had  
908 discussed finding mean, Maria asked them to make a visual proof for the median and  
909 the mode.

910 (*end snapshot*)

911 A key characteristic of data science is asking questions of “big data”—a data set that  
912 has many cases and variables and needs to be manipulated and analyzed  
913 computationally. One benefit of working with bigger data sets is that discerning the  
914 shape of a distribution is often easier. Both histograms in figure 5.9 show distributions of  
915 air temperature measurements taken in Sacramento, CA, in 2013. The table on the left  
916 shows 8,700 measurements taken hourly from January 1 to December 31, 2013, while  
917 the table on the right contains 100 measurements randomly sampled from the larger  
918 data set on the left. The “big data” distribution on the left shows a fairly smooth  
919 distribution with a rightward skew and modal temperatures of approximately 55 degrees.  
920 The “little data” distribution on the right also shows a rightward skew but contains peaks  
921 and valleys, with modal temperatures of approximately 45 degrees.

922 Figure 5.9 Comparing Distributions for Large and Small Data Sets





923

924 [Long description of figure 5.9](#)

925 Source: Data from the National Oceanic and Atmospheric Administration, National  
 926 Centers for Environmental Information

927 The following snapshot describes a classroom scenario in which students investigate  
 928 hurricane data from multiple years and use a range of data displays to understand the  
 929 science of hurricanes and to generate additional questions. Thematic topics are shown  
 930 in parentheses within the body of the snapshot where relevant.

931 **Snapshot: Quincey’s Investigation of Hurricane Data**

932 The sixth-grade math teacher, Leonora, decided to have her students explore the  
 933 “shape” of some weather data. The context is hurricanes in the Atlantic Ocean and uses  
 934 real data collected from 5 years of hurricanes at successive 10-year intervals. One  
 935 student, Quincey, showed real interest and engaged in the lesson’s opening discussion  
 936 of 2017 hurricane data displayed on a line plot (*understanding and describing variability*  
 937 *in data and data distributions*). Quincey and the class were really interested in the  
 938 number of hurricanes that were in the tropical storm category.

939 Next, students worked in groups to study hurricane category data for the years 1977,  
 940 1987, 1997, and 2007 (*data collection, sampling, and random processes*). Each  
 941 decade’s data were presented in different ways: bar graphs, line plots, tables, and

942 sentences. Quincey enjoyed the analysis and was taken with the different ways of  
943 displaying data as well as the changes in the spread of data from decade to decade.

944 Quincey asked important questions about the science of hurricanes. “How do they  
945 develop? What makes them get larger? What is the difference between a category 3  
946 storm and a category 5 storm?” At the close of the lesson, Leonora was convinced that  
947 students understood that different visual displays of data can make it easier to  
948 recognize how a situation might be changing over time (*comparing distributions and*  
949 *identifying associations between variables*). The class reflected that the changes were  
950 easier to see in line plots and histograms than through the data being shared in writing  
951 or in a table of values. Quincey decided to further investigate the number of category 4  
952 and 5 hurricanes over the past 100 years and how these storms become stronger, and  
953 Quincey set out to gather more data and ask questions of the data. Others in the class  
954 decided to investigate why the number of category 4 and 5 storms are increasing.

955 (*end snapshot*)

## 956 ***Data Collection, Sampling, and Random Processes***

### 957 **Sampling**

958 Prior to seventh grade, students’ work with data focused exclusively on using data to  
959 understand, describe, and compare the particular collection of objects or situations that  
960 have been collected by observations, experiments, or measurements.

961 Seventh grade includes the first introduction to sampling, the process of collecting data  
962 from a subset of a population in an attempt to understand or describe the whole  
963 population. This focus represents a big jump in sophistication from earlier work. As an  
964 example, suppose all students who come in to play basketball before school are asked  
965 to track their screen usage for the week. The class analyzes the data and determines  
966 that those sampled spent an average of 862 minutes on the screen. Small- and whole-  
967 group discussions invite students to consider whether this sample can extend to (i.e., is  
968 representative of) the entire student population at the school or to all students who are  
969 the same age. Although the 862 minutes may be the typical screen time for the defined

970 group of students where the data were collected, it may not extend to everyone.  
971 Explorations and discussions, including considering some obviously nonrepresentative  
972 samples, can help students understand the idea of a random sample.

973 It is important for students to have multiple experiences selecting samples from known  
974 populations in ways that are random (for instance, drawing numbered ping-pong balls  
975 from an opaque bag or drawing student names on identical slips of paper from a hat)  
976 and in ways that are not random (for instance, asking survey questions only of the  
977 students who sit near you in class). The goal is for students to develop an  
978 understanding that random sampling tends to produce samples that are representative  
979 of the population—that is, their distribution of the quantities under consideration are  
980 close to the distribution for the population as a whole (7.SP.1)—and for students to have  
981 a sense of the variability when using samples to make inferences and estimates for a  
982 population (7.SP.2). Many computational tools enable students to quickly draw samples  
983 from data sets using a variety of methods (e.g., randomly, selecting every fifth record or  
984 the first 100 records).

985 Although sampling is not explicitly named again in the standards until high school,  
986 eighth grade students may benefit from additional opportunities to deepen their  
987 understanding of sampling as they work with bivariate data. Random sampling can  
988 become a tool to engage in data explorations of interest to students—e.g., “I wonder  
989 how long on average it takes students from different grades to get from home to  
990 school?” or “How much food is wasted in the lunchroom every month?”

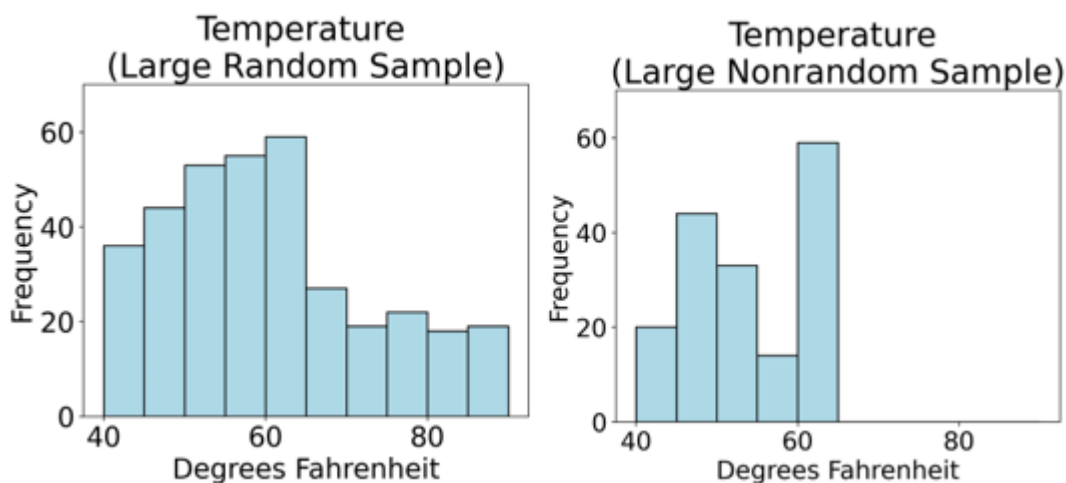
991 Nonrandom sampling (such as attempting to understand the school as a whole by  
992 collecting data only from one’s friends, or by asking about eating habits at the gym after  
993 school) produces biased conclusions, even when the bias in the sample selection might  
994 not be obviously linked to the quantity being measured in the measurement or  
995 observation. Bias in the statistical setting does not refer to temperament or outlook  
996 (prejudice), which is one meaning of the word; instead, it means a systematic error.

997 Students often believe that arbitrary sampling schemes (first 10 students I meet or every  
998 tenth student alphabetically) are random; they need to understand the difference

999 between these schemes and choosing by chance so that every possible sample has an  
1000 equal likelihood of being selected.

1001 Figure 5.10 has two samples drawn from the Sacramento temperature data first shown  
1002 in figure 5.9 above. The table on the left shows a random sample of 365 measurements  
1003 from the original data set, while the table on the right shows a nonrandom sample of  
1004 365 measurements: the first measurement of the day for every day in 2013. The shape  
1005 of the distribution on the left is much closer to the shape of the original distribution that  
1006 contained 8,700 measurements (figure 5.9).

1007 Figure 5.10 Comparing Random and Nonrandom Samples



1008

1009 [Long description of figure 5.10](#)

1010 Source: Data from the National Oceanic and Atmospheric Administration, National  
1011 Centers for Environmental Information

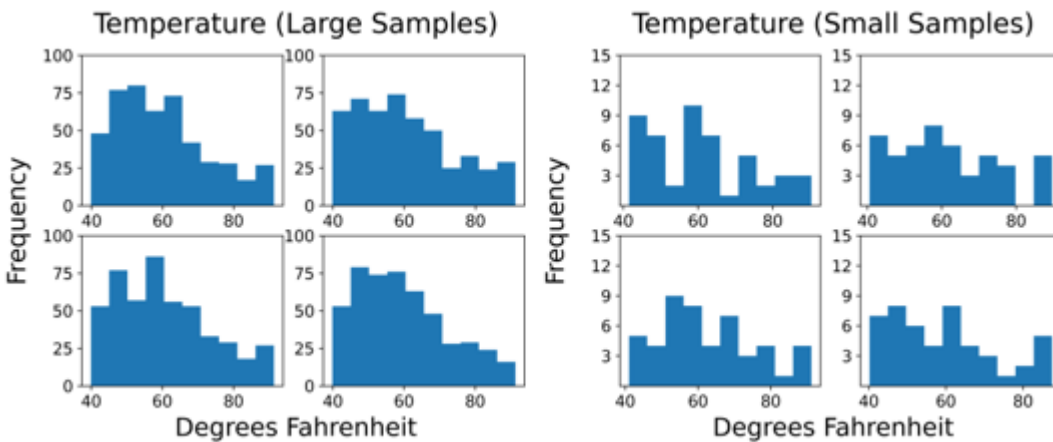
### 1012 **Probability and Random Processes**

1013 Randomly selecting from a population and measuring a characteristic (in which variation  
1014 is expected across the population) is a chance process: It may result in different results  
1015 and its outcomes follow some distribution.

1016 Although students can generate samples through non-computational means,  
1017 computational tools can enable students to quickly and easily draw samples from a data  
1018 set and visualize or summarize each sample in order to compare and contrast the

1019 results of different sampling methods. The tables on the left in Figure 5.11 show four  
1020 different random samples of 500 data points from the same Sacramento temperature  
1021 data shown in figure 5.9, and tables on the right show four different random but much  
1022 smaller samples of 50 data points from this data set. Students should notice that the  
1023 shapes of the small samples are much more variable compared to the shapes of the  
1024 large samples, even though all samples were generated randomly. A small sample,  
1025 even if random, is less likely to be representative of the population than is a large  
1026 random sample.

1027 Figure 5.11 Comparing Distributions for Large and Small Random Samples



1028

1029 [Long description of figure 5.11](#)

1030 *Source:* Author (data from the National Oceanic and Atmospheric Administration,  
1031 National Centers for Environmental Information)

1032 Probability expresses the chance of an outcome as a number between 0 and 1  
1033 (7.SP.5). Probability is combined with statistics in the grade seven standards.

1034 Statistics and probability are historically linked because statistical claims and estimates  
1035 are based on the mathematical field of probability. Using sampled data to predict  
1036 events, such as elections outcomes, is based on probabilistic reasoning.

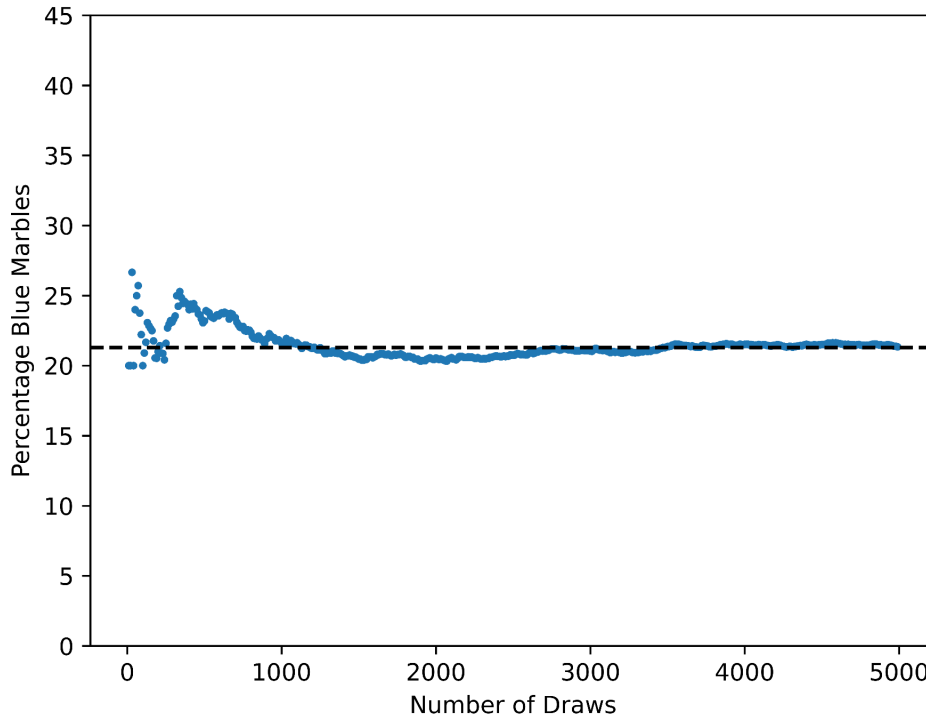
1037 Students sometimes struggle to see clear connections between probability and  
1038 statistics, especially when their experiences focus on procedures and calculations

1039 rather than exploration, context, and interpretation. Statistics produces estimates for  
1040 parameters in probabilistic models. There is much work with probability that does not  
1041 support statistical reasoning and may not be applicable to a setting of interest (e.g.,  
1042 calculating theoretical probabilities for the sum of two dice without using those  
1043 theoretical probabilities to decide whether a given pair of dice are likely fair), and middle  
1044 school probability experiences should be carefully designed to support reasoning with  
1045 interesting and meaningful data.

1046 In seventh grade, students gather data to estimate the probability of outcomes by  
1047 observing their long-run relative frequency; that is, they compute experimental  
1048 probability. Consider repeating this experiment 150 times: Draw a marble from a bag  
1049 with marbles in it, record its color, then put the marble back in the bag. If you get a blue  
1050 marble 32 times, your estimate for the probability of getting a blue marble on any  
1051 particular draw is  $32/150$  (7.SP.6, 7.SP.7.b). This is really an estimate of the fraction of  
1052 blue marbles in the bag.

1053 Compare the marble experiment just described to another, placing the following marbles  
1054 in a bag (all identical except for color): 16 blue marbles, 31 red marbles, 16 green  
1055 marbles, and 12 white marbles (75 total marbles). If you blindly pull a marble from the  
1056 bag, what is the probability that you will get a blue marble? If you repeat this 150 times  
1057 (putting the marble back each time), about how many times do you expect to get a blue  
1058 marble? Students can create a theoretical probability model and calculate the expected  
1059 frequency by multiplying the probability of drawing a blue marble in one draw by the  
1060 number of draws ( $16/75 \times 150 = 32$ ). After calculating this expectation that students will  
1061 draw blue marbles 32 out of 150 times, or 21.3 percent of the time, students might try to  
1062 verify their theoretical probability model and create a simulation that can pull a marble  
1063 from the bag 150 or 1500 or 15,000 times. Students can then compare the simulated  
1064 frequency of drawing a blue marble with their theoretical expectations (CSS 6–8.AP.10).  
1065 Figure 5.12 shows the results from such a simulation. The long-run proportion of blue  
1066 marbles reaches an asymptote at the theoretical probability of 21.3 percent.

1067 Figure 5.12 Results From a Simulation Containing 5,000 Trials of the Marble  
1068 Experiment



1069

1070 Note the difference between the two different marble experiments described in the  
1071 previous two paragraphs. In the first, students repeat an experiment many times and  
1072 use the long-run frequency of drawing a blue marble to estimate the theoretical  
1073 probability of drawing a blue marble from the bag. In the second, students build a  
1074 (theoretical) probability model and use it to estimate the long-run frequency of drawing a  
1075 blue marble from the bag (7.SP.7). If the relative frequencies of experimental outcomes  
1076 do not seem close to predictions from the probability model, then students need to be  
1077 able to discuss possible sources of discrepancy (7.SP.7): Perhaps the green marbles  
1078 have a different texture and tend to be drawn more frequently than predicted. Maybe  
1079 somebody changed the mix of marbles in the bag. Or perhaps not enough draws were  
1080 performed to see the relative frequencies approach the probability model.

1081 Finally, seventh grade students find probabilities of compound events (events which are  
1082 made up of several simple events)—for example, drawing two marbles from the bag of  
1083 75 described above and getting one white (W) and one blue (B) marble (7.SP.8).

1084 The recognition that some events (repeat the draw five times, get all blue; or repeat the  
1085 draw five times, obtain WBWWB in that order) are much less likely than others (repeat  
1086 the draw five times, get three white and two blue) is key to understanding claims made  
1087 from statistics.

1088 In fact, most statistical claims depend on a comparison of a (theoretical and  
1089 hypothetical) probability model with observed data, as in 7.SP.7. To prepare middle  
1090 school students for future statistical work, teachers should offer experiences that  
1091 develop an awareness that more data tend to produce relative frequencies closer to  
1092 actual probabilities. Computational tools can support activities that use larger data sets  
1093 and the creation of simulations that enable students to compare experimental and  
1094 theoretical probabilities.

1095 The following snapshot describes how a student in this grade band might encounter  
1096 ideas from this section. Major themes are highlighted in parentheses where relevant.

1097 ***Snapshot: Rosa's Students Experience Random Sampling***

1098 Understanding the ways Rosa's seventh grade students have responded to the  
1099 probability activities offered through her instruction has influenced the next steps in her  
1100 planning. Overall, Rosa has not been satisfied with student understanding of random  
1101 sampling. She decides to give students a more visual and physical experience of the  
1102 concept. Her plan calls for six paper bags filled with differently colored cubes. The sum  
1103 of cubes and the color distribution of the cubes in the bags are as follows:

1104 Bag One, 15 total: 15 blue

1105 Bag Two, 12 total: 11 blue, 1 red

1106 Bag Three, 20 total: 15 blue, 4 yellow, 1 red



1107 Bag Four, 10 total: 5 red, 5 yellow

1108 Bag Five, 12 total: 5 blue, 4 red, 3 yellow

1109 Bag Six, 20 total: 8 blue, 8 red, 4 yellow

1110 Rosa explains the task by telling students they will determine the contents of each bag  
1111 through sampling. She chooses not to tell them how many times to sample but she does  
1112 tell them to sample from the bags by selecting one cube at a time and then putting it  
1113 back into the bag. Rosa also asks students to determine the chance of drawing a blue  
1114 cube from each bag (*data collection, sampling, and random processes*).

1115 Students engage in the activity, brainstorming methods for collecting and recording their  
1116 information. When each group of students feels satisfied with their determinations of the  
1117 number of cubes and color distributions of the contents of each bag, she asks them to  
1118 choose which bag belongs to which card showing the contents of each bag  
1119 (*understanding and describing variability in data and data distributions*). In setting up the  
1120 lesson, Rosa filled the bags differently and made sure to have a bag for which the  
1121 probability of drawing a blue cube would be 1 and another for which it would be 0. After  
1122 the activity and class discussion, Rosa is happy to hear her students later talking about  
1123 situations in which the probability is 1 or 0 and other situations representing everything  
1124 in between. Her students recognized the number of times they sampled usually led to  
1125 better predictions about the contents of the bags. They also realized that sampling  
1126 without replacement would have shown them the exact contents of the bag. The class  
1127 engaged in a rich conversation about sampling with and without replacement,  
1128 recognizing that it would be unproductive to draw all the cubes if there were a million.

1129 (*end snapshot*)

### 1130 ***Comparing Distributions and Identifying Associations Between Variables***

1131 Prior to grade seven, students typically work with a single collection of data that  
1132 measures a single variable. In grade seven, they compare the same variable measured  
1133 across two populations, either by actually measuring the whole populations or obtaining

1134 estimates for the population distributions via sampling. They can plot data and draw  
1135 from different statistical methods such as creating box plots and dot plots to informally  
1136 assess the degree of overlap of two populations.

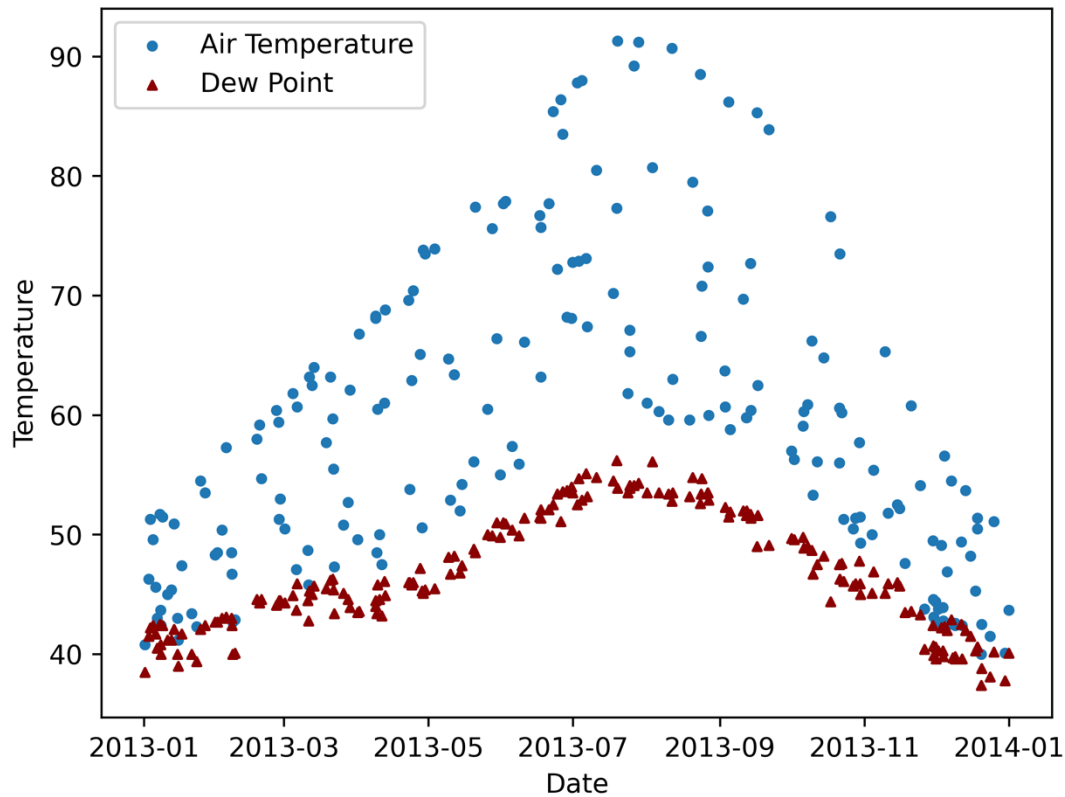
1137 In eighth grade, the focus is on formally describing bivariate data: two quantities or  
1138 categorical variables measured or observed across a population or across a sample  
1139 drawn from a population (8.SP.1). Eighth grade students use two-way frequency tables  
1140 as tools to see associations in bivariate categorical data (8.SP.4). This work has  
1141 important connections to linear equations and modeling.

1142 The scatter plot as a visual representation of quantitative bivariate data is one of the  
1143 most important ideas introduced in eighth grade. A survey of students collecting data on  
1144 both time and distance for traveling from home to school might reveal clusters, outliers,  
1145 and any of various types of association (positive, negative, linear, nonlinear). Students  
1146 should describe such patterns in a scatter plot and interpret the patterns in the context  
1147 of the data (8.SP.1). Once a scatter plot is created, an association between the two  
1148 variables may become visually identifiable. Fitting a function to the data is the creation  
1149 of a mathematical model for the association.

1150 In eighth grade, students choose a line to fit the data by visual approximation on the  
1151 scatter plot, and they compare and argue for whose line fits “best” (8.SP.2). They then  
1152 interpret the meaning of the slope and intercept of their chosen model line and use the  
1153 line to make predictions for one variable when the other variable is specified (8.SP.3).  
1154 Finally, eighth grade students use two-way frequency tables as tools to see  
1155 associations in bivariate categorical data (8.SP.4).

1156 Although the type of function that is used most frequently is a line (a linear function),  
1157 students also need experiences plotting associations that are clearly nonlinear, as in  
1158 figure 5.13, and fitting other types of functions (quadratic, exponential) to the plot.

1159 Figure 5.13 Scatterplot Showing a Nonlinear Association Between Air Temperature and  
1160 Dew Point in Sacramento in 2013



1161

1162 Source: Data from the National Oceanic and Atmospheric Administration, National  
 1163 Centers for Environmental Information

1164 Any standard data software (including spreadsheets, Desmos, Geogebra, CODAP) will  
 1165 fit lines, quadratic functions, and exponential functions to given data. Students are not  
 1166 expected to know the specific standard technique for identifying a line (or quadratic or  
 1167 exponential function) of best fit (least squares regression), but students should have  
 1168 experiences fitting lines and some other functions visually (by adjusting parameters on  
 1169 appropriate function types in graphing software) and using appropriate software tools  
 1170 which perform the regression calculations.

### 1171 **High School**

1172 Students' prior work learning about describing and comparing distributions and random  
 1173 sampling comes together in high school. High school students continue to visualize and

1174 represent univariate data with dot plots, histograms, and box plots; use measures of  
1175 center and spread to describe such distributions (S-ID.4); and compare distributions  
1176 from different populations or samples using these representations and statistics (S-  
1177 ID.1–3). A major difference between students’ data experiences in kindergarten through  
1178 grade eight and what is explored in high school is the richness and complexity of  
1179 available data sets, even more so than their sheer size. As high school students work  
1180 with these data sets, they can draw upon the statistical understandings they have  
1181 developed in their kindergarten through grade eight mathematics lessons. Instruction  
1182 should emphasize opportunities for questioning and interpreting alongside technical  
1183 procedures.

1184 Data exploration begins with a search for available data about a context of interest. The  
1185 data set is then examined for hidden patterns and associations. At the high school level,  
1186 visualization of data can illustrate unexpected structure. Any patterns or associations  
1187 discovered can lead to new hypotheses or questions to investigate further. Students  
1188 began this process in eighth grade and continue in high school with experiences in  
1189 which they examine data sets with multiple variables that are measured for each  
1190 member of the sample. They plot pairs of variables to decide which ones might show  
1191 associations. Important discussions for students to engage in when working with  
1192 existing data sets include the following:

- 1193 ● Prior to exploring: Do you expect any of these variables to be associated? Why?
- 1194 ● Might the association you see just be a result of the way in which the data were  
1195 collected rather than truly reflective of the population? What features of the data  
1196 collection might make conclusions suspect, and what features might give  
1197 confidence? Note that a large sample size is not enough to have confidence in  
1198 conclusions.
- 1199 ● Can you think of possible explanations for the association(s) you see? Can you  
1200 think of ways you could decide which explanations might be accurate?

1201 After data exploration identifies some association(s) of interest, the stage of model  
1202 building follows. Technical methods are reserved for the specialized statistics or data  
1203 science course, which is described below, but all students need to explore questions  
1204 such as the following:

- 1205 • Could you use some variables to predict others? Doing so is a hugely important  
1206 use of data because some factors are easier to measure or observe than others.  
1207 Medicine and many other fields often require using presently available  
1208 information to try to predict future outcomes.

1209 Most importantly, high school students (like kindergarten through grade eight students)  
1210 must experience statistics as a set of tools for making sense of their worlds in ways that  
1211 matter to them.

### 1212 ***Bringing It All Together: Introduction to Inferential Statistics***

1213 High school students begin learning about inferential statistics, which aim to generalize  
1214 from a sample and draw conclusions about a population (S-IC.1). Students' work with  
1215 inferential statistics is foundational to using data to make decisions. Students must  
1216 decide whether a result observed through data is consistent with a mathematical model  
1217 of the process that generates the data (S-IC.2). For instance, students are asked to  
1218 engage in a thought experiment and consider how many households use gardens as a  
1219 source of food. If a student hypothesizes that 30 percent of the students at the school  
1220 grow food at home, the estimate offers a mathematical model that gives them an idea of  
1221 what proportion to expect in a sample. If they then survey five randomly chosen  
1222 students, and all say they grow food at home, then the student should be able to reason  
1223 as follows: If 30 percent of students grow food at home, then the chances of five  
1224 randomly chosen students all being among those 30 percent of students is  $(.3)^5$   
1225  $= .00243 = .243$  percent, or less than a quarter of 1 percent. Thus, the student might  
1226 doubt—that is, they might reject—the 30-percent hypothesis. Students should have  
1227 many experiences of simple situations like this to understand how decisions based on  
1228 data rely on probability and are not guaranteed to produce correct answers to the  
1229 original question.

1230 Students should work with data that originate from four different methods of data  
1231 production, including at least some student-generated questions and student-gathered  
1232 data. These methods are (1) generating census data, which are data that contain  
1233 measurements on every member of the target population (such as the database of  
1234 crimes occurring in a given city in a given time frame, or rain gauge data for a given  
1235 location, which captures all precipitation at that location—census data is first  
1236 encountered in early elementary grades); (2) administering surveys to random samples  
1237 (to estimate population values, or parameters, for the surveyed quantities); (3)  
1238 conducting randomized experiments (to compare treatments and demonstrate cause);  
1239 and (4) conducting observational studies (to study characteristics or quantities when  
1240 random selection or assignment is not possible) (S-IC.3). The High School Progression  
1241 on Statistics and Probability (Common Core Standards Writing Team, 2022) contains  
1242 detailed examples describing the expectations in the standards.

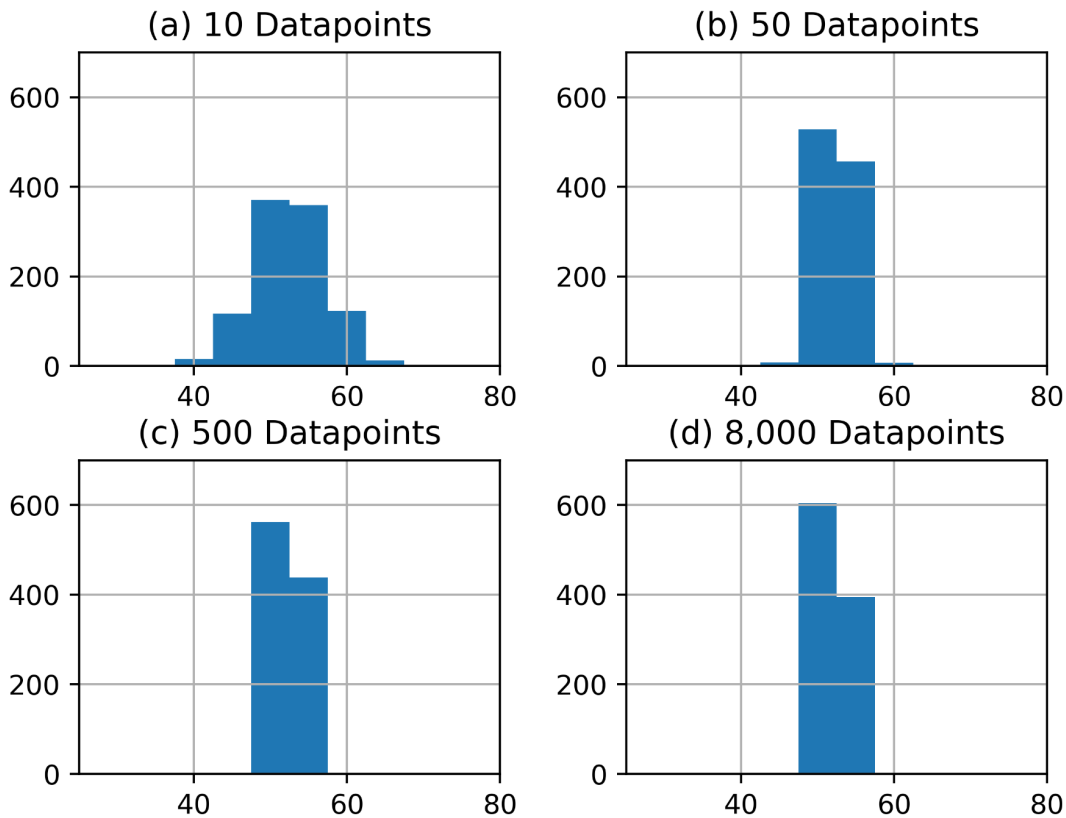
1243 Teaching with surveys and experiments must include a link between the random  
1244 selection or assignment and the ability to reason probabilistically to make claims. With a  
1245 survey, the random sampling allows generalizing to a population. With an experiment,  
1246 the random assignment allows causal conclusions but not generalization to a broader  
1247 population—unless the sample in the experiment was randomly selected from some  
1248 larger population.

1249 When using a sample mean or proportion to estimate a population mean or proportion,  
1250 students use simulation models to estimate a margin of error, instead of using formulaic  
1251 calculations. Briefly, the process is to use data simulation software to draw many  
1252 random samples from a hypothetical population and to see how often a result is  
1253 obtained that is as extreme as the sample mean or proportion. Doing this process for  
1254 hypothetical populations with many different mean or proportion parameters helps  
1255 students see that there is a range of population parameters that often (more than 5  
1256 percent of the time) produce simulated sample means or proportions that are as  
1257 extreme as (or more extreme than) the actual sample mean or proportion. This range of  
1258 population parameters is the (simulation-based) confidence interval, given as a sample  
1259 mean or proportion  $\pm$  margin of error. Note the probabilistic argument here: If the

1260 population mean or proportion were outside of the confidence interval, then sample  
1261 means or proportions as extreme as were obtained in the random sample would be  
1262 rare. So, the true population mean or proportion is expected to be within the confidence  
1263 interval (but one cannot be certain that it is!).

1264 Figure 5.14 shows the results from a simulation that drew 1,000 random samples of  
1265 different sizes from 10,000 normally distributed values with a mean of 50 and a  
1266 standard deviation of 15. Figure 5.14a shows the distribution of sample means obtained  
1267 from 1,000 random samples of 10 values. Although most sample means were close to  
1268 50, the sample means ranged from 38 to 67. The distribution of sample means  
1269 becomes narrower as the samples get larger. Drawing 1,000 random samples of 50  
1270 values yields sample means between 43 and 62. Drawing 1,000 random samples of  
1271 500 values yields an even narrower range of sample means, from 48 to 57. And finally,  
1272 drawing 1,000 random samples of 8,000 individuals only yields sample means that are  
1273 very close to 50. What this exercise shows is that larger random samples are able to  
1274 generate more precise estimates of population parameters, in this case the mean.

1275 Figure 5.14 The Relationship Between Sample Size and the Shape of the Sampling  
1276 Distribution



1277

1278 [Long description of figure 5.14](#)

1279 Source: Data from the National Oceanic and Atmospheric Administration, National  
 1280 Centers for Environmental Information

1281 A similar process is used to evaluate confidence in a randomized experiment in which  
 1282 subjects are randomly assigned to two or more treatment groups. (Treatment could  
 1283 mean medical treatment, or assignment of different tasks, or being shown different  
 1284 motivational videos, and so on.) Some quantity is then measured for each subject, and  
 1285 the investigator then has to decide from the results whether a treatment, say treatment  
 1286 A, produced any effect on the measured quantity. Simply having a different mean for  
 1287 each treatment group is not enough, as variation is expected in the measurement and  
 1288 thus between groups. In this case, all of the treatment groups are pooled into a  
 1289 population and then re-sampled (randomly) many times to see how often the re-  
 1290 sampled mean or proportion is at least as extreme as the actual treatment A group



1291 difference. If such differences are rare, the experiment is taken as evidence that  
1292 treatment A caused a change in the measured quantity.

1293 In the following snapshot, students investigate real-world environmental data and health  
1294 impacts. The snapshot demonstrates how the sequence provided them opportunities to  
1295 use all three thematic topics (shown in parentheses within the snapshot).

1296 ***Snapshot: Data on Environmental Threats to Health***

1297 In this example (Lieberman and Brown, 2020), students compared CalEnviroScreen  
1298 data related to four environmental topics that are known to affect human health:

1299 (1) water (using data on groundwater threats, impaired water, and drinking water);  
1300 (2) toxic chemicals (using data on pesticides, cleanups, and toxic releases); (3) air  
1301 pollution (using data on the ozone, particulate matter [PM 2.5], diesel, and traffic); and  
1302 (4) waste (using data on hazardous waste and solid waste). They compared these  
1303 results against environmental impacts, using data for asthma, low birth weight, and  
1304 cardiovascular disease (California Health Education Standards 1.13.P; 2.3.P; 3.3.P,  
1305 3.4.P).

1306 In preparation for the students' analysis and reporting, the teacher reviewed California's  
1307 EP&Cs with students by asking them to identify one that is directly related to their  
1308 environmental health problem. Based on their data analysis, students identified  
1309 environmental health and environmental justice concerns related to water pollution in  
1310 the local community and observed that they differentially affected various parts of the  
1311 community (*understanding and describing variability in data and data distributions*).  
1312 Their conclusion was that the key factors in the differential environmental health impacts  
1313 were related to "Environmental Principle V: Decisions affecting resources and natural  
1314 systems are based on a wide range of considerations and decision-making processes."

1315 Depending on the focus of their individual environmental health study, students are  
1316 encouraged to choose two variables to analyze, such as the impact of water quality on  
1317 low birth weight, or the impact of toxic chemicals on the incidence of cardiovascular  
1318 disease, or the impact of air quality on asthma (*data collection, sampling, and random*

1319 *processes*). After collecting the data for these variables, the students use technology to  
1320 create a scatter plot of the data, fit a function to the data, and create a symbolic  
1321 representation for the function (*understanding and describing variability in data and data*  
1322 *distributions*). Students are able to connect the parameters of the symbolic  
1323 representation to the context of the data. After a class discussion about the comparison  
1324 of different variables, students should be guided to focus on the combinations of  
1325 variables that make the most sense for their investigations (*comparing distributions and*  
1326 *identifying associations between variables*).

1327 Following their research and analysis, student teams report back to the class,  
1328 summarizing their quantitative comparisons using charts to depict the results about  
1329 water, toxic chemicals, air pollution, and waste (Health 4.1.P). In their presentation, they  
1330 use graphs to compare the environmental effects they discovered from the  
1331 environmental health impacts they analyzed (English Language Arts SL.9-12.1; SL.9-  
1332 12.2; SL.9-12.4; SL.9-12.5; Health 5.3.P).

1333 Several of the teams mention that they observed a pattern that relates to the socio-  
1334 economic conditions in the communities they compared. Some of the students mention  
1335 that they see these issues as directly related to EP&C V because the places where  
1336 waste, toxic chemicals, and manufacturing facilities are located depend on a variety of  
1337 political, economic, and social factors. The teacher explains that differential  
1338 environmental health impacts on communities with varied socioeconomic conditions is a  
1339 major health topic identified as “environmental justice,” a term that came into use in the  
1340 1980s when residents of an African American community in North Carolina protested  
1341 the siting of a landfill to store soil contaminated with polychlorinated biphenyls (PCBs).  
1342 These residents knew the health hazards associated with this toxin and responded by  
1343 demanding that their health and well-being be protected by the government. The landfill  
1344 proposal went forward, but the protests spurred the federal government to study the  
1345 issue. The findings show that many of the nation’s landfill sites are located in  
1346 communities of color. The environmental justice movement has grown to focus on a  
1347 more equitable distribution of environmental benefits and burdens. Since many of the  
1348 students expressed a strong interest in this topic, they request a guest speaker from a

1349 community-based health organization to provide additional information and answer  
1350 students' questions about environmental justice (Health 8.1.P.; 8.2.P).

1351 *(end snapshot)*

### 1352 ***Guidance for High School Data Science Courses***

1353 While all high school students should exercise and refine their understanding of data  
1354 exploration, causal inference, and statistical reasoning using large, real-world data sets,  
1355 many sophisticated approaches to working with rich, complex data sets can be left to an  
1356 advanced statistics or data science course. (Chapter 8 and Appendix A provide further  
1357 detail on the different mathematics pathways available to schools and students.)

1358 With the rapid expansion of information available to all in the form of data, students may  
1359 be interested in a data science course as a culminating high school mathematical  
1360 science experience. In addition to recognizing the importance of the data science  
1361 content—to 21st-century jobs and to a wide range of college majors—many students  
1362 are more engaged with math classes that are taught in a spirit of open-ended  
1363 exploration, drawing upon important mathematical principles and tools rather than more  
1364 traditional teaching methods focused solely on procedures without motivation or  
1365 context.

1366 Effective data science courses consider how to help students with the following:

- 1367 ● Understand how data are used by professionals to address real-world problems
- 1368 ● Understand that data are used in all facets of modern life
- 1369 ● Understand how data support science to identify and tackle real-world problems  
1370 in communities
- 1371 ● Learn about statistical variability
- 1372 ● Use appropriate tools and techniques to make sense of large data sets
- 1373 ● Create and analyze statistical graphics to identify patterns in data and to connect  
1374 these patterns back to the real world
- 1375 ● Understand that by treating photos, words, numbers, and sounds as data, you  
1376 can gain insight into the real world

- 1377 ● Learn to analyze data, including posing questions that can be answered by  
1378 considering relations among variables in a data set, using collected data to  
1379 generate hypotheses for future data collection, critically evaluating shortcomings  
1380 and strengths in the data and the data-collection process, and informally  
1381 evaluating hypotheses using data at hand
- 1382 ● Learn basic programming and use computer programs in the development and  
1383 analysis of statistical models
- 1384 ● Learn about data ethics, including consideration of where data come from, who is  
1385 collecting the data, and how the data are used
- 1386 ● Refine computational models to better represent the relationships among  
1387 different elements of data that are collected and analyzed
- 1388 ● Design algorithms to solve computational problems by using and adapting  
1389 existing algorithms and creating new ones

1390 When designing or planning for a data science course, there are many different  
1391 sequences and approaches. A course might actively engage students by exploring the  
1392 meaning of data, the importance of communicating data visually, the role of cleaning  
1393 data, exploratory data analysis, ethical issues around data, creating data dashboards,  
1394 linear and nonlinear regression models, statistics, probability, and forecasting.

1395 In addition to the mathematics content described above, exposure to some software  
1396 and other technology tools is essential for those wishing to pursue a career in data  
1397 science, and facility with such tools is increasingly valuable for a variety of professionals  
1398 whose work involves basic data analysis. However, data science does not require any  
1399 particular software package, and different data science courses may use different  
1400 software and technology tools depending on the specifics of the course and school  
1401 context. More important than the technology students use is that they learn to ask good  
1402 questions and apply effective mathematical tools to help them answer their questions.

## 1403 **Equitable and Inclusive Instruction**

1404 Educators can offer social and emotional support to students by designing engaging  
1405 lessons that allow students to connect in meaningful ways with content. Traditional  
1406 mathematics lessons that have taught the subject as a set of procedures to follow have  
1407 resulted in widespread disengagement as students see no relevance for their lives. The  
1408 data science field provides multiple opportunities for students to pursue answers to real-  
1409 world problems and their own wonderings and to see that they can excel in quantitative  
1410 fields.

1411 Important principles underlying the teaching of data science that will offer the greatest  
1412 chance for social, emotional, and academic development include the following:

- 1413 ● *Convey Mindset and Belonging Messages*

1414 Informed by successful interventions in mindset and belonging strategies,  
1415 teachers can remind students that struggle represents an important part of  
1416 learning; all students struggle at times, and successful students respond to times  
1417 of difficulty by using strategies they have developed and practiced over time.  
1418 Teachers can share with students examples of successful people within the field  
1419 that highlight gender and racial diversity.

- 1420 ● *Use Real Data*

1421 Modern computing provides an opportunity for students to question real sets of  
1422 data, developing social awareness and investment in the solutions they discover.  
1423 When working with secondary data sets (data obtained from others, rather than  
1424 collected by students), teachers should choose meaningful content selected to  
1425 create a connection with their learning and secure opportunities to hear the  
1426 perspective of others. When teachers use local data sets, they can also help  
1427 students feel like they are important members of their community—as they use  
1428 real data to explore questions and find answers to local problems that they can  
1429 help to address. Identifying problems and finding solutions will help students  
1430 develop skills to make responsible decisions. Some teachers worry that they  
1431 cannot provide culturally sustaining connections for their classes because they

1432 lack expertise in the cultures of all their students, but real data sets from different  
1433 communities provide opportunities for students to bring their own knowledge and  
1434 expertise to data-rich problems. There should also be times when students are  
1435 invited to collect data from their own community and build their own data sets.  
1436 Students can pose questions that are important to them, including those with  
1437 cultural meaning, collecting data from their own lives and communities. As Paris  
1438 (2012) describes, students will thereby be fostering and sustaining “linguistic,  
1439 literate, and cultural pluralism.” The act of collecting data provides an important  
1440 learning opportunity for students to understand decisions that need to be made  
1441 around the collection and organization of data and to understand how to deal  
1442 with uncertainty in their data. Students will be the ones with important expertise  
1443 in these investigations.

1444 ● *Focus on Collaboration and Communication*

1445 Meaningful collaborations typically reflect perspectives from diverse groups of  
1446 students who come together to work effectively with different ideas being valued  
1447 and developed. Creating opportunities for this kind of group work makes an  
1448 environment where differences thrive and where students have the tools to work  
1449 respectfully to reach solutions. For example, students may start their work in  
1450 structured and unstructured conversations in which each group member shares  
1451 their thoughts. Collaborative classrooms founded in engaged listening and the  
1452 capacity to articulate verbally as students build on each other’s ideas are places  
1453 where students feel valued and where they develop important relationship skills  
1454 of communication, social engagement, and teamwork.

1455 **Connecting to the Drivers of Investigation and Content**

1456 **Connections**

1457 This chapter highlights three thematic topics central to data science that exist within the  
1458 CA CCSSM K–12 progressions. Data investigations that leverage the SMPs provide  
1459 students with opportunities to recognize that statistics is a problem-solving process that  
1460 connects the mathematics they are doing in class to their lived experiences and to other

1461 content areas. The chapter also describes how investigations can provide students with  
1462 needed opportunities to collect data within their classrooms and to consider the  
1463 implications of their decisions (such as decisions about sampling). The grade band  
1464 descriptions included examples of how students can use plots to develop an  
1465 understanding of variability in data and the use of plots as a communication tool to  
1466 describe that variability. Additionally, as students pose their wonderings about the world  
1467 around them or try to predict future events using data, there may be opportunities to  
1468 explore probability and identify associations between variables. A focus on these  
1469 thematic topics helps ensure that all students continue to grow in their abilities to work  
1470 with data and lay the groundwork for pathways toward data science.

1471 To help teachers who have been working with the CA CCSSM standards and  
1472 progressions, this chapter focuses on the progressions of statistical and data science  
1473 ideas across grade levels. Additionally, as described in chapter 1, this framework  
1474 encourages teachers to design instruction by using the big ideas of mathematics at  
1475 each grade level as focal points for student investigations. Investigations are guided by  
1476 the three Drivers of Investigation (DIs) which provide the “why” of learning mathematics,  
1477 eight SMPs which provide the “how” of learning mathematics, and the four Content  
1478 Connections (CCs) which provide the “what” of learning mathematics.

1479 To address the mathematical concepts as discussed in this chapter within such an  
1480 approach, educators may want to begin with the kinds of questions they are asking to  
1481 build mathematical understandings. The aim of the DIs is to ensure that there is always  
1482 a reason to care about mathematical work—and that investigations provide  
1483 opportunities for students to make sense, predict, and/or affect the world. Just as the  
1484 three thematic topics progress across the K–12 band, early DI questions are primarily  
1485 about description and begin with categorizing and counting, expanding into questions  
1486 that present measurement situations (initially length/distance; later time, area, volume,  
1487 and rates). As students progress through the grade bands, they begin to investigate  
1488 relationships between two or more varying quantities and perform formal quantitative  
1489 calculations to describe future events or probable outcomes. Figure 5.15 illustrates  
1490 three different questions associated with the three different DIs, each of which enables

1491 elementary students to learn about and use data in ways described earlier in this  
1492 chapter.

1493 Figure 5.15 Using Drivers of Investigation to Frame Questions and Investigations About  
1494 Weather

<b>Driver of Investigation 1: Making Sense of the World (Understand and Explain)</b>	<b>Driver of Investigation 2: Predicting What Could Happen (Predict)</b>	<b>Driver of Investigation 3: Impacting the Future (Affect)</b>
<i>“What are different ways to describe our weather?”</i>	<i>“What will the weather be like tomorrow?”</i>	<i>“How can we use weather data to make recommendations to visitors for packing or outdoor activities?”</i>

1495 When the DIs are coupled with data from relevant contexts, students may be more likely  
1496 to authentically engage in statistical problem-solving. Connecting DIs to contexts  
1497 relevant to students’ backgrounds and interests may also increase the inclusion of and  
1498 participation among girls and students from racial and ethnic groups historically  
1499 underrepresented in STEM fields.

1500 When designing learning sequences, the DI will link CCs and one or more SMPs  
1501 together. While CC1 (Reasoning with Data) is explicitly tied to data, earlier portions of  
1502 this chapter note many places throughout the K–12 experience where data explorations  
1503 might arise and thus support CCs 2, 3, and 4. Figure 5.16 uses the weather situation  
1504 from figure 5.15 to show how investigations using weather data might support each of  
1505 the four CCs. Within each row, the bolded information links back to the primary thematic  
1506 topic covered in this chapter.

1507 Figure 5.16 Using Weather Data to Support Content Connections



Content Connection	Supporting Investigation
<p><b>Content Connection 1: Reasoning with Data</b></p>	<p>Students have the opportunity to count, measure, and classify attributes such as temperature, time, precipitation, or cloud cover. These values can be expressed and explored graphically and then interpreted and shared with peers or the community.  <b>Thematic Topic: Understanding and describing variability in data and data distributions</b></p>
<p><b>Content Connection 2: Exploring Changing Quantities</b></p>	<p>The collection of weather data can provide a rich context for students to explore their numeracy and develop and apply operations (e.g., <i>What is the difference between the highest and lowest daily temperature for Tuesday?</i>); express values in terms of ratio or percent (e.g., <i>What percentage of days were cloudy for the month of March?</i>); and find ways to express patterns between quantities mathematically via multiple representations or express climate patterns through linear equations (e.g., <i>The temperature increase over time can be described by <math>1.2C + 22</math></i>).  <b>Thematic Topic: Data collection</b></p>
<p><b>Content Connection 3: Taking Wholes Apart, Putting Parts Together</b></p>	<p>Describing the weather requires students to decompose the investigation question into attributes that could and should be collected to address the question driving the investigation. Students might explore data on the average amount of different types of precipitation (e.g., rain versus snow) in each calendar month and then combine those data to develop a more detailed understanding of total annual precipitation.  <b>Thematic Topic: Understanding and describing variability in data and data distributions</b></p>
<p><b>Content Connection 4: Discovering Shape and Space</b></p>	<p>While there may be variation in daily temperatures, overall the shape of the plot for temperatures throughout the day or for particular seasons is quite predictable. Length of day or seasonal weather patterns are inversely related between the northern and southern hemispheres. Students might explore 2-D versus 3-D graphical representations and the role of geospatial data in weather forecasting.  <b>Thematic Topic: Comparing distributions and identifying associations between variables</b></p>

1508 **Conclusion**

1509 As readers consider the three subsequent chapters of the framework, they will see

1510 ideas similar to the ones discussed in this chapter, organized to help them learn about

1511 and begin to use the big ideas approach. While the transition between standards  
1512 domains and progressions discussed in this chapter and this new approach will not be  
1513 straightforward for classroom teachers, both emphasize the central idea that students at  
1514 all levels should have experiences that build their mathematical toolkits for making  
1515 sense of their worlds.

1516 Life in a data-rich world requires that California schools prepare all students to examine  
1517 claims justified with data, to understand the probabilistic underpinning of drawing  
1518 conclusions from samples, and to see the use of data as a tool for answering many  
1519 questions of interest. Developing these abilities requires that students generate  
1520 questions and work with data beginning in kindergarten (or before) and have  
1521 experiences of increasing depth and complexity throughout their school careers. As  
1522 discussed in chapter 8, students who wish to focus extra attention on data science  
1523 should have an opportunity to pursue advanced courses late in their high school  
1524 careers.

1525 **Additional Resources**

1526 As educators consider how to create rich lessons that integrate data and statistical  
1527 investigations into their classrooms, the following resources, which influenced this  
1528 document, may be helpful (listed in alphabetical order):

1529 **American Statistical Association (ASA):** The following is an excerpt from the  
1530 website’s K–12 page: “The American Statistical Association is dedicated to and involved  
1531 in enhancing statistics education at all levels, including providing resources for K–12  
1532 teachers and teacher educators. Here, you will find information about classroom  
1533 resources, publications in statistics education, guidelines and reports, workshops and  
1534 webinars for teachers, and student competitions.”

1535 <https://www.amstat.org/education/k-12-educators>

1536 **GAISE II:** The Pre-K–12 Guidelines for Assessment and Instruction in Statistics  
1537 Education II (GAISE II) (Bargagliotti et al., 2020) is a professional report from the ASA  
1538 setting out guidelines for assessment and instruction in PreK–12 in statistics and data  
1539 science. The GAISE II is an important resource for this area of mathematical science. It  
1540 includes guidance and examples for skills and concepts at the elementary, middle, and  
1541 high school levels.

1542 [https://www.amstat.org/asa/education/Guidelines-for-Assessment-and-Instruction-in-](https://www.amstat.org/asa/education/Guidelines-for-Assessment-and-Instruction-in-Statistics-Education-Reports.aspx)  
1543 [Statistics-Education-Reports.aspx](https://www.amstat.org/asa/education/Guidelines-for-Assessment-and-Instruction-in-Statistics-Education-Reports.aspx)

1544 **National Council of Teachers of Mathematics (NCTM):** The NCTM published two  
1545 books in its Essential Understanding series devoted to statistics: “Developing Essential  
1546 Understanding of Statistics: Grades 6–8” and a second volume for grades nine through  
1547 twelve.

1548 [https://www.nctm.org/Store/Products/Developing-Essential-Understanding-of-Statistics-](https://www.nctm.org/Store/Products/Developing-Essential-Understanding-of-Statistics-for-Teaching-Mathematics-in-Grades-6-8-(Download)/)  
1549 [for-Teaching-Mathematics-in-Grades-6-8-\(Download\)/](https://www.nctm.org/Store/Products/Developing-Essential-Understanding-of-Statistics-for-Teaching-Mathematics-in-Grades-6-8-(Download)/)

1550 [https://www.nctm.org/Store/Products/Developing-Essential-Understanding-of-Statistics-](https://www.nctm.org/Store/Products/Developing-Essential-Understanding-of-Statistics-for-Teaching-Mathematics-in-Grades-9-12-(Download)/)  
1551 [for-Teaching-Mathematics-in-Grades-9-12-\(Download\)/](https://www.nctm.org/Store/Products/Developing-Essential-Understanding-of-Statistics-for-Teaching-Mathematics-in-Grades-9-12-(Download)/)

1552 **Statistical Education of Teachers (SET):** The following is an excerpt from the  
1553 introduction: “The *SET* report outlines the content and conceptual understanding  
1554 teachers need to know to assist their students develop statistical reasoning skills. *SET*  
1555 is intended for everyone involved in the statistical education of teachers, both the initial  
1556 preparation of prospective teachers and the professional development of practicing  
1557 teachers.”

1558 [https://www.statisticsteacher.org/statistics-teacher-](https://www.statisticsteacher.org/statistics-teacher-publications/#:~:text=of%20practicing%20teachers.-,Download,-Focus%20on%20Statistics)  
1559 [publications/#:~:text=of%20practicing%20teachers.-,Download,-](https://www.statisticsteacher.org/statistics-teacher-publications/#:~:text=of%20practicing%20teachers.-,Download,-Focus%20on%20Statistics)  
1560 [Focus%20on%20Statistics](https://www.statisticsteacher.org/statistics-teacher-publications/#:~:text=of%20practicing%20teachers.-,Download,-Focus%20on%20Statistics)

1561 **Statistics Teacher:** The following is an excerpt from the publication’s website: “In 2016,  
1562 the ASA/NCTM Joint Committee decided the *Statistics Teacher Network* newsletter  
1563 should evolve to the *Statistics Teacher* online journal. The goal of *Statistics Teacher*  
1564 remains to inform and support K–12 teachers. The new publication will continue to  
1565 provide articles about successful classroom practice and announcements of important  
1566 professional development opportunities. It will also more seamlessly integrate peer-  
1567 reviewed lesson plans. Each issue will have dedicated technology and assessment  
1568 columns.”

1569 <https://www.statisticsteacher.org/about-statistics-teacher/>

## 1570 **Long Descriptions of Graphics for Chapter 5**

### 1571 **Figure 5.4. A Teacher’s Dot Plot of the Data to Determine the Most** 1572 **Common Crayon Box Size**

1573 Figure 5.4 is a dot plot displaying information about the number of crayons found inside  
1574 crayon boxes drawn from a large black bag. The vertical axis, which ranges from 0 to  
1575 12, represents the number of boxes drawn of each size. The horizontal axis, which  
1576 ranges from 2 to 16, represents the number of crayons per box.

1577 The dot plot shows:

- 1578 • 2 data points of the 2-crayon box
- 1579 • 6 data points of the 4-crayon box
- 1580 • 13 data points of the 8-crayon box
- 1581 • 6 data points of the 16-crayon box

1582 These data suggest that boxes with eight crayons, with 13 data points, were drawn  
 1583 most often.

1584 [Return to figure 5.4 graphic](#)

1585 **Figure 5.6. Temperature Plots to Compare Mean Values for Two Cities**  
 1586 **in California**

1587 The first graphic shows July 1 Max Temperature for Death Valley, California (degrees  
 1588 Fahrenheit) on a number line. The values shown are 115.88, 109.94, 116.60, 117.68,  
 1589 118.40, 114.44, 116.42, 116.96, and 116.42. Mean result is 115.86 (indicated with a  
 1590 vertical blue line).

1591 The second graphic shows July 1 Max Temperature for Stockton, California (degrees  
 1592 Fahrenheit) on a number line. Values shown are 94.28, 94.46, 93.56, 95.36, 99.68,  
 1593 98.24, 95.72, 96.26, 95.36, and 95.36. Mean result is 95.83 (indicated with a vertical  
 1594 blue line).

1595 [Return to figure 5.6 graphic](#)

1596 **Figure 5.7. Logan’s Vase Measurement Data Visualized in CODAP**

1597 The first figure shows height (in cm) and volume (in ml) on a graph. Data values are as  
 1598 follows:

Height	Volume
23	2760
17.2	760

Height	Volume
16.5	1000
14	440
12.5	290
7	460
15.5	85

1599 The second figure shows height data (in cm) on a number line. Data values are the  
 1600 same as in the table above. Mean value is 15.1 (indicated with a vertical blue line).

1601 [Return to figure 5.7 graphic](#)

1602 **Figure 5.8. Using Data to Classify Shapes**

1603 The figure shows an example of student group work described in the text. It depicts a  
 1604 rectangle with 10 different colored lines running across it, originating and ending at  
 1605 different points around the edges of the rectangle. The crossing lines create polygonal  
 1606 shapes (two-dimensional shapes formed with straight lines) that have different numbers  
 1607 of sides. Each polygonal shape within the rectangle is labeled with the number of sides  
 1608 it has. Underneath the figure are tally mark counts of how many shapes have three  
 1609 sides, four sides, five sides, six sides, seven sides, and right triangles.

1610 [Return to figure 5.8 graphic](#)

1611 **Figure 5.9 Comparing Distributions for Large and Small Data Sets**

1612 Left graph shows Temperature ('Big Data'). The data values are as follows:

Degrees Fahrenheit	Frequency
40–44	850
45–49	1300

Degrees Fahrenheit	Frequency
50–54	1280
55–59	1430
60–64	1210
65–69	780
70–74	540
75–79	520
80–84	400
85–90	390

1613 Right graph shows Temperature ('Little Data'). The data values are as follows:

Degrees Fahrenheit	Frequency
40–44	17
45–49	19
50–54	7
55–59	16
60–64	13
65–69	8
70–74	8
75–79	3
80–84	5
85–90	4

1614 [Return to figure 5.9 graphic](#)

1615 **Figure 5.10. Comparing Random and Nonrandom Samples**

1616 Left graph shows Temperature (Large Random Sample). Data values are as follows:

Degrees Fahrenheit	Frequency
40–44	38
45–49	45
50–54	54
55–59	57
60–64	60
65–69	30
70–74	20
75–79	22
80–84	19
85–90	20

1617 Right graph shows Temperature (Large Nonrandom Sample). Data values are as

1618 follows.

Degrees Fahrenheit	Frequency
40–44	20
45–49	45
50–54	33
55–59	16
60–64	60
65–69	0
70–74	0



Degrees Fahrenheit	Frequency
75–79	0
80–84	0
85–90	0

1619 [Return to figure 5.10 graphic](#)

1620 **Figure 5.11. Comparing Distributions for Large and Small Random**  
 1621 **Samples**

1622 Left side of figure has four small tables for Temperature (Large Samples). They show  
 1623 the following data:

Degrees Fahrenheit	Frequency (table 1)	Frequency (table 2)	Frequency (table 3)	Frequency (table 4)
40–44	49	63	54	52
45–49	79	72	76	78
50–54	80	63	56	75
55–59	65	76	90	76
60–64	75	59	58	62
65–69	45	47	55	58
70–74	32	27	36	28
75–79	30	33	28	26
80–84	20	27	21	25
85–90	25	33	26	20

1624 Right side of figure has four small tables for Temperature (Small Samples). They show  
 1625 the following data:

Degrees Fahrenheit	Frequency (table 1)	Frequency (table 2)	Frequency (table 3)	Frequency (table 4)
40–44	9	7	5	7
45–49	7	5	4	8
50–54	2	6	10	6
55–59	11	9	8	4
60–64	7	6	4	9
65–69	1	3	7	4
70–74	5	5	3	3
75–79	2	4	4	1
80–84	3	0	1	2
85–90	3	5	4	6

1626 [Return to figure 5.11 graphic](#)

1627 **Figure 5.14. The Relationship Between Sample Size and the Shape of**  
 1628 **the Sampling Distribution**

1629 Figure includes four tables showing the distribution of sample means obtained from  
 1630 1,000 random samples of different sizes.

1631 Figure 5.14(a) 10 Datapoints shows the following data:

X-axis	Y-axis
28–32	0
33–37	0
38–42	15
43–47	120

X-axis	Y-axis
48–52	370
53–57	360
58–62	120
63–67	15
68–72	0
73–77	0

1632 Figure 5.14(b) 50 Datapoints shows the following data:

X-axis	Y-axis
28–32	0
33–37	0
38–42	0
43–47	10
48–52	520
53–57	460
58–62	10
63–67	0
68–72	0
73–77	0

1633 Figure 5.15(c) 500 Datapoints shows the following data:

X-axis	Y-axis
28–32	0
33–37	0

X-axis	Y-axis
38–42	0
43–47	0
48–52	570
53–57	430
58–62	0
63–67	0
68–72	0
73–77	0

1634 Figure 5.15(d) 8,000 Datapoints shows the following data:

X-axis	Y-axis
28–32	0
33–37	0
38–42	0
43–47	0
48–52	605
53–57	395
58–62	0
63–67	0
68–72	0
73–77	0

1635 [Return to figure 5.14 graphic](#)

