

1
2
3
4

Mathematics Framework
Second Field Review Draft
March 2022
Page 1 of 86

5
6
7

Mathematics Framework
Chapter 5 Data Science, TK–12
Second Field Review Draft

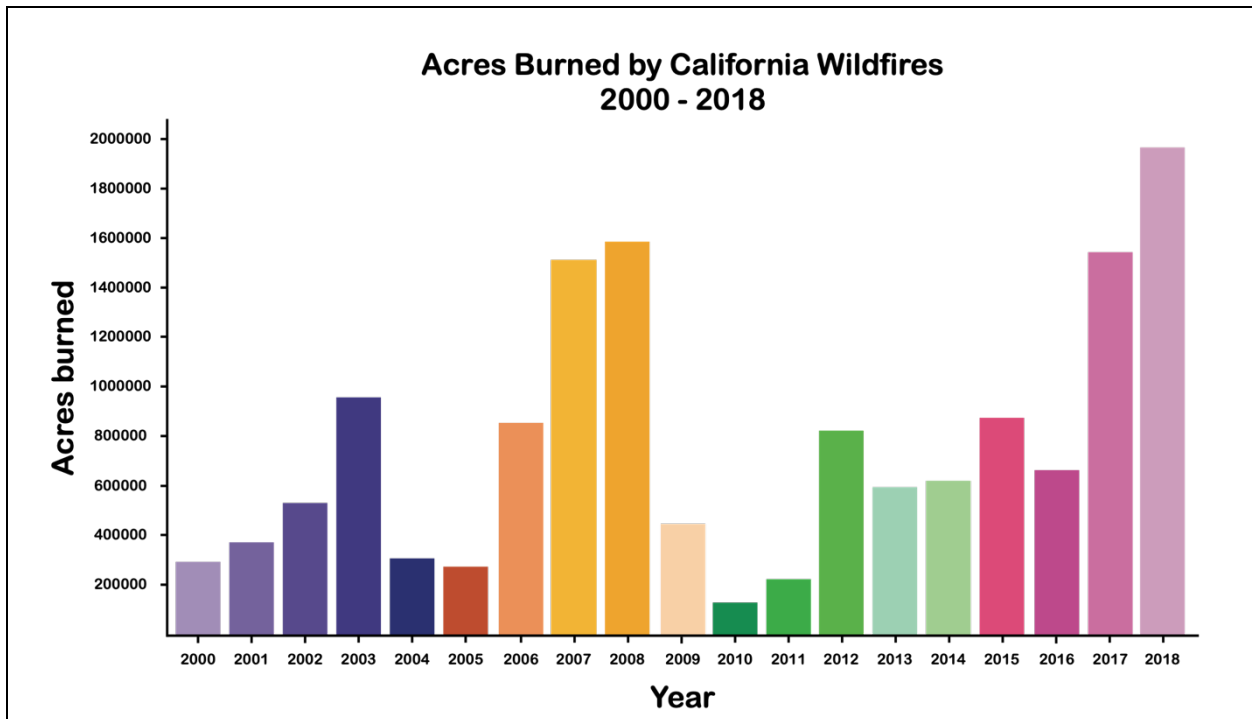
8	Mathematics Framework Chapter 5 Data Science, TK–12	1
9	Introduction	3
10	What is Data Science?	5
11	Big Ideas in Statistics and Data Science	10
12	Driving Investigation and Making Connections	10
13	The Statistical and Data Science Investigation Process	13
14	Data Talks K–12	23
15	Transitioning from Pre–K	27
16	Kindergarten Through Grade Five	27
17	What questions can data help to answer?	28
18	Asking Questions, Collecting and Analyzing Data	31
19	Interpreting and Communicating Results	32
20	Preparing for the Major Data Science Work of Grades Six Through Eight	34
21	Grades Six Through Eight	40
22	Data in the World: Question Asking, Exploration, Interpretation, Decision Making,	
23	Ethics, Technology	41
24	Describing, Displaying, and Comparing Variability (Grades Six Through Seven)	42
25	Sampling to Understand a Population: Randomness, Bias, How Many? (Grades	
26	Seven Through Eight)	45
27	Are They Related? Two Changing Quantities (Grade Eight)	48
28	What Are the Chances? Probability as the Basis for Data-Based Claims	49
29	High School	52
30	Data Science for Equity and Inclusion	53
31	Data for All Students: Living in a World Overloaded with Information	56
32	Advanced high school data science	66
33	Content Learning Outcomes	75
34	Sample Courses	80
35	Conclusion	82
36	Long Descriptions for Chapter 5	83

37 **Note to reader:** The use of the non-binary, singular pronouns *they*, *them*, *their*, *theirs*,
38 *themselves*, and *themselves* in this framework is intentional.

39 **Introduction**

40 The ability to work with and understand data is an essential life skill in a world
41 continuously inundated with data. Data drives students' lives, whether they see them or
42 not; making sense of data, being able to identify data that is misleading, and using data
43 to make decisions are all important for their role as global citizens. It is not only those
44 with professions in data science—almost all occupations now require that employees
45 collect feedback from data and adjust their practice. Stories about the world are
46 illuminated by massive quantities of data, and community members telling and listening
47 to those stories need to be able to make sense of data to understand their health,
48 finances, and news feeds.

49 Figure 5.1. Example of California data for students to explore




50 [Link to long description](#)



51 Source: CalFire, 2020.

52 The numbers are staggering: around 1.7 megabytes of digital data were created and
53 stored *every second for every person on Earth* in 2020, and the vast majority of data
54 goes unanalyzed (Petrov, 2021). Our lives are increasingly subject to data-driven
55 algorithms that determine aspects of our daily experience, including the ads we see,
56 which neighborhoods receive business or public investment, who gets screened more
57 closely at the airport, who receives favorable terms on monetary loans, and which
58 medical procedures are recommended or approved.

59 All California students should graduate from high school with data literacy and have
60 access to options to learn an introduction to data science in their K–12 experience. Data
61 literacy refers to the ability to reason with and about data, to make good decisions
62 based on data, to ask questions of data, and to use statistical reasoning. Data science
63 is an emerging discipline that includes understanding principles of data collection, data
64 manipulation, data analysis, inference, and interpretation and communication. The
65 California Common Core State Standards in Mathematics (CA CCSSM) set out the
66 learning of statistics K–12. Viewing the CA CCSSM through a data science lens can
67 highlight the statistical ideas in the standards and increase their relevance and meaning
68 for students.

69 **Definition: Data** are observations or measurements in context. In a given context, a
70 unit of observation (a member of a population) may have multiple attributes measured
71 or observed; each of these attributes is a **variable**. Often, data are recorded in a table in
72 which rows represent units of observation and columns represent variables. For
73 instance, in the table below, the countries are the units of observation and the variables
74 are Flag, 2020 population, Region, and Highest elevation in meters.

Country	Flag	2020 population	Region	Highest elevation (m)
China		1,440,000,000	Eastern Asia	8848

Country	Flag	2020 population	Region	Highest elevation (m)
India		1,370,000,000	Southern Asia	8586
United States		330,600,000	North America	6190

75 **What is Data Science?**

76 Data Science is the process of uncovering the stories hidden within data. It involves
77 formulating questions, and collecting, cleaning, wrangling, analyzing, and visualizing
78 data (that is often huge and complex) to uncover patterns and trends and communicate
79 them to others. Professional data scientists draw upon mathematics, statistics, and
80 computer science, and think critically about the qualitative features of a data set to find
81 meaning and communicate the results of their inquiries. Data scientists work together to
82 address uncertainty in data while avoiding bias (Finzer, 2013).

83 The terms *statistics* and *data science* both refer to the processes and tools of finding
84 meaning in data, and the distinction is still a matter of discussion. Statistics traditionally
85 uses theoretical tools to build and evaluate proposed mathematical models, using data
86 from a population of interest. *Data science* highlights the expansion in computing and
87 visualization tools that have made many more techniques available for finding meaning
88 in data—often relying on innovative visualizations of complex data that enable major
89 features to be identified and explored further. Because *statistics* has become
90 synonymous in much of TK–12 education with a very limited set of procedures (mean,
91 median, standard deviation, interquartile range, correlation, and linear regression, along
92 with a few data visualizations such as line plots and scatter plots), this framework uses
93 *data science* to emphasize the full statistical and data science investigation process
94 (see below). Students need to experience statistical tools in the process of investigating
95 authentic questions. Statistics has been overlooked in recent years, when districts and

96 schools prioritize some content over others, but this is important to change as it is a
97 critical area of content for students in the twenty-first century.

98 The professional statistics community does not have a limited definition of statistics: The
99 Pre-K–12 Guidelines for Assessment and Instruction in Statistics Education II (GAISE II;
100 Bargagliotti, Franklin, Arnold, Gould, Johnson, Perez, and Spangler, 2020) is a
101 professional report from the American Statistical Association (ASA) setting out
102 guidelines for assessment and instruction Pre-K–12 in statistics and data science, and
103 is an important resource for this area of mathematical science. GAISE II emphasizes
104 the following:

- 105 1. The importance of asking questions throughout the statistical problem-solving
106 process (formulating a statistical investigative question, composing data
107 collection questions, interrogating existing data, analyzing data, and interpreting
108 results), and how this process remains at the forefront of statistical reasoning for
109 all studies involving data
- 110 2. The consideration of different data and variable types, the importance of carefully
111 planning how to collect data or how to consider data to help answer statistical
112 investigative questions, and the process of collecting, cleaning, interrogating, and
113 analyzing the data
- 114 3. The inclusion of multivariate thinking (using 2 or more variables) throughout all
115 Pre-K–12 educational levels
- 116 4. The role of probabilistic thinking – considering the likelihood of an event - in
117 quantifying randomness throughout all levels
- 118 5. The recognition that modern statistical practice is intertwined with technology,
119 and the importance of incorporating technology as feasible
- 120 6. The enhanced importance of clearly and accurately communicating statistical
121 information
- 122 7. The role of assessment at the school level, especially items that measure
123 conceptual understanding and require statistical reasoning involving the
124 statistical problem-solving process. (GAISE II, 2020, 2)

125 Students should be able to draw on the Standards for Mathematical Practices (SMP)
126 through a statistical lens articulated in the ASA’s Statistical Education of Teachers
127 (SET) report. For example, students should reason abstractly and quantitatively by
128 engaging in statistical thinking while considering where data come from (SMP.2), apply
129 statistical models to “include descriptions of the variability present in data (SMP.4), and
130 consider available tools such as calculators, spreadsheets, applets, statistical
131 packages, and graphical displays to help facilitate the statistical problem-solving
132 process (SMP.5). When students participate in the analysis of large datasets, they
133 should be able to decide which questions matter, and identify which ones can be
134 answered with a given dataset (SMP.4). The statistical problem-solving process, shown
135 in Figure 5.3, is used within the process—with students formulating questions, collecting
136 and analyzing data and communicating their results. Further, students should
137 understand some of the ways in which data are frequently misunderstood or misused
138 and should understand the content and implications of their own digital data footprints.
139 Finally, students should be prepared to pursue additional study directed towards fields
140 which include more intensive work with data, such as designing data collection, deciding
141 on statistical measures appropriate to the questions under consideration, or making
142 conclusions and claims based on data.

143 Particular aspects of the CA CCSSM help build the data understanding and skills that
144 high school graduates require. However, the progression—from counting, categorizing,
145 and simple picture graphs, to the complex skills and understanding that older students
146 may develop—requires careful thought and considerably more focus through the K–12
147 curriculum than most students have historically experienced. The study of data
148 continues to expand and broaden. The types of data being collected are vast and the
149 types of techniques used to analyze data adhere to a strong reliance on computational
150 tools. The statistical problem-solving process is important as it provides the foundation
151 for finding meaning in data. Data science and statistics are the science of working with
152 data. The development of statistics and data science mastery articulated in this chapter
153 represents a contemporary lens through which to examine the CA CCSSM.

154 Educators regularly use data at the student and classroom level to try to drive
155 instructional decisions. However, a data-science perspective can help educators create
156 experiences in which their students learn to “read and write the world with mathematics”
157 (Gutstein, 2003). As emphasized throughout this framework, students must experience
158 mathematics as tools for making sense of and impacting their worlds.

159 Educators should be encouraged to bring data science and statistics directly into their
160 classroom in ways that create meaningful student experiences. Students can explore
161 statistics and data science as tools for making sense of and impacting their worlds. The
162 statistical problem-solving process (GAISE II) helps students formulate statistical
163 investigative questions, take in information by collecting primary data or considering
164 secondary data, analyze the data to identify relationships and patterns, and (in many
165 cases) interpret results to answer the question and propose changes to impact the way
166 the world works.

167 Students who are exposed to and have the capacity to understand data concepts at an
168 early age begin to develop data literacy and data sense in parallel with number sense.
169 As students progress through school they should learn different approaches to data
170 analysis, culminating in the investigation of large data sets using appropriate
171 technological tools.

172 As students learn the investigative statistical and data science process they should
173 always consider meaning and context. In the past, some learning of statistics was
174 removed from situational settings, leading students to learn abstract methods. Data
175 science involves developing meaning and communicating about a data-rich situation; it
176 should remain in its context. Teachers can use local data sets that give students
177 opportunities to ask questions that are meaningful to them, that can help their local
178 community, or school, allowing students to experience using mathematics to be an
179 engaged citizen. Statistics and data science is about studying situations—asking
180 questions such as: Who collected the data? How was it collected? What is the unit of
181 analysis? Teachers can ask students to turn and talk to their partners and groups about
182 these questions.

183 In this chapter, we present the progression of data literacy and data science standards
184 and the types of experiences that help build the necessary skills and understandings.
185 Four important principles in the learning of data science are outlined here:

- 186 1. Students should experience working with data from a context that is meaningful
187 to them personally. They should have opportunities to solve problems of value to
188 the students and to their schools and communities.
- 189 2. Students should learn to engage with real data that include multiple variables. At
190 first students can learn to understand two variables with bivariate data, as they
191 progress through the grades they can learn to handle multivariable data and
192 multivariate thinking. Multivariable data often includes 3 or more variables. For
193 example, students could categorize their favorite toys considering their size,
194 fluffiness, and type of animal (3 variables).
- 195 3. Data investigations should be investigative and collaborative, with students
196 working together to learn the data science and statistical investigative process.
- 197 4. Familiarity with technology and modern tools should progress through the
198 grades.

199 As discussed in more detail in Chapter 2, it is more effective for teachers to plan around
200 big ideas than sets of mathematical methods, and to choose rich tasks that elicit big
201 ideas. In this chapter we set out the big ideas of data science that build to the kind of
202 connected understanding needed.

203 **Usage note:** In Latin, the word *data* is the plural of *datum*. However, in English, *data* is
204 now also commonly used with singular verbs and refers to a collection of data points.
205 Thus, “the data shows a correlation...” is more common than “the data show a
206 correlation....” In this chapter we most often use the word *data* in this way—to refer to a
207 collection of data points—and in these contexts it takes singular verbs.

208 Two important sources for contexts in which to explore data science are

- 209 ● The California Next Generation Science Standards (CA NGSS) (California
- 210 Department of Education, 2013a)
- 211 ● The California Environmental Principles and Concepts (EP&Cs) (California
- 212 Department of Education, 2013b)

213 **Big Ideas in Statistics and Data Science**

214 The table below presents the big ideas that will be addressed in each grade level band.

TK–5	6–8	9–12: all students	11–12: advanced data science
<ul style="list-style-type: none"> ● Data for understanding. What questions can we ask? What data do we need to answer it? ● Defining data: What is data, how and where is data collected? ● Representing and interpreting data: What does data look like and what does it mean? 	<ul style="list-style-type: none"> ● Data in the world: exploration, interpretation, decision making, ethics ● Variability: Describing, displaying, and comparing ● Sampling to understand a population: randomness, bias, how many? ● Are they related? Multivariate thinking ● What are the chances? Probability as the basis for data-based claims 	<ul style="list-style-type: none"> ● Interpreting categorical and quantitative data ● Making inferences and justifying conclusions ● From statistics to data science: messy data, computational tools 	<ul style="list-style-type: none"> ● The role of data in the world ● Formulating statistical investigative questions ● Collecting and considering data ● Computational tools, including programming, for analyzing data

215 **Driving Investigation and Making Connections**

216 Since motivating students to care about mathematics is crucial to forming meaningful

217 content connections, this Framework describes instruction that is situated in student

218 investigations, falling in one of three **Drivers of Investigation** (DIs), which provide the

219 “why” of learning mathematics. These Drivers are then paired with **Content**
220 **Connections** (CCs), which provide the “how and what” of mathematics (the CA
221 CCSSM standards) to be learned in an activity. So, the Drivers of Investigation propel
222 the learning of the content framed in the Content Connections.

223 ***Drivers of Investigation (DIs)***

224 The Content Connections should be developed through investigation of questions in
225 authentic contexts; these investigations will naturally fall into one or more of the
226 following Drivers of Investigation. The Drivers of Investigation are meant to serve a
227 purpose similar to that of the Crosscutting Concepts in the CA NGSS, as unifying
228 reasons that both elicit curiosity and provide the motivation for deeply engaging with
229 authentic mathematics. In practical use, teachers can use these to frame questions or
230 activities at the outset for the class period, the week, or longer; or refer to these in the
231 middle of an investigation (perhaps in response to the “Why are we doing this again?”
232 questions that often crop up), or circle back to these at the conclusion of an activity to
233 help students see “why it all matters.” Their purpose is to pique and leverage students’
234 innate wonder about the world, the future of the world, and their role in that future, in
235 order to foster a deeper understanding of the Content Connections and grow into a
236 perspective that mathematics itself is a lively, flexible endeavor by which we can
237 appreciate and understand so much of the inner workings of our world. The Drivers of
238 Investigation are:

- 239 ● Driver of Investigation 1: Making Sense of the World (Understand and Explain)
- 240 ● Driver of Investigation 2: Predicting What Could Happen (Predict)
- 241 ● Driver of Investigation 3: Impacting the Future (Affect)

242 ***Content Connections (CCs)***

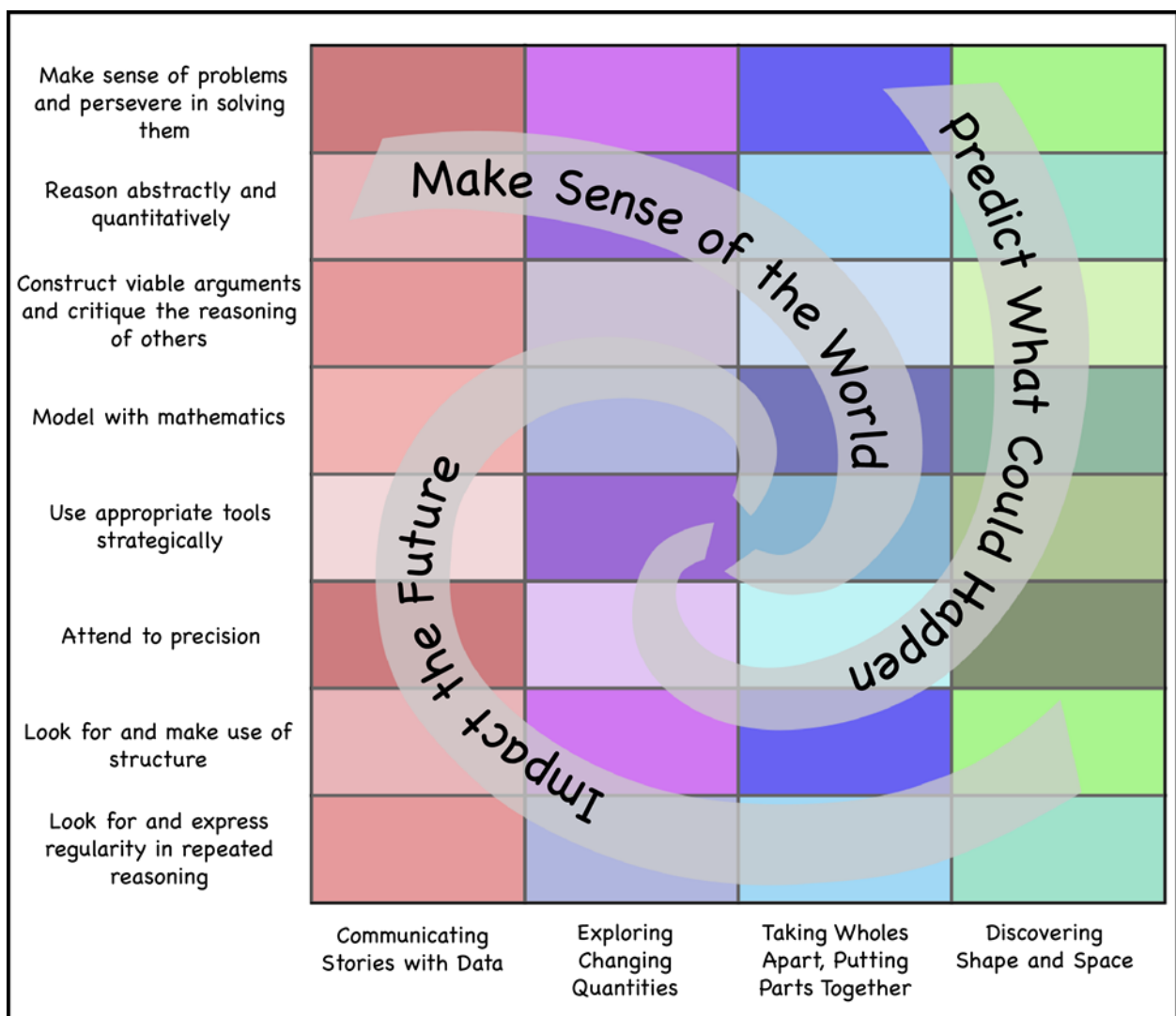
243 The four Content Connections described in the framework organize content and provide
244 mathematical coherence through the grades:

- 245 ● Content Connection 1: Communicating Stories with Data

- 246 • Content Connection 2: Exploring Changing Quantities
- 247 • Content Connection 3: Taking Wholes Apart, Putting Parts Together
- 248 • Content Connection 4: Discovering Shape and Space

249 The relationship between content, mathematical practices and the drivers of
 250 investigation is highlighted in Figure 2:

251 Figure 5.2: The Drivers of Investigation



252
 253 [Link to long description](#)

254 Big ideas that drive design of instructional activities will link one or more Content

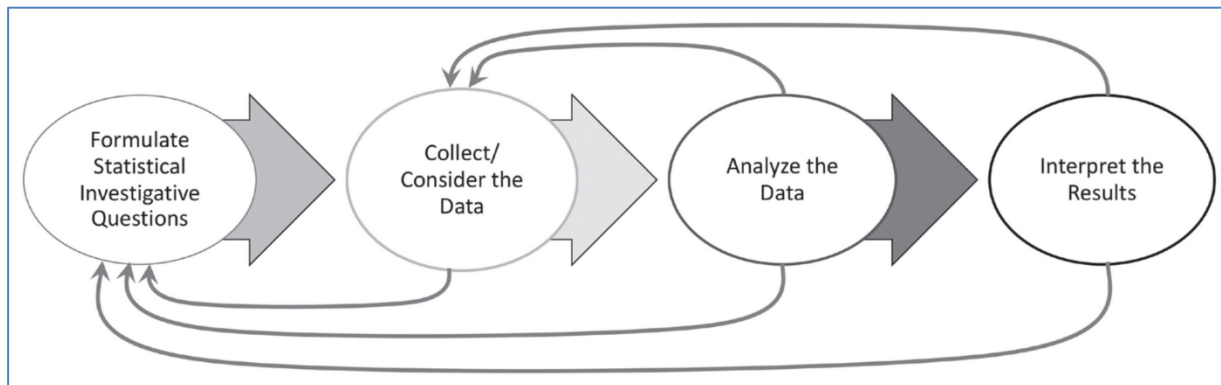
255 Connections and one or more Standards for Mathematical Practice (SMPs) with a
256 Driver of Investigation, so that students can Communicate Stories with Data *in order to*
257 Predict What Could Happen, or Illuminate Changing Quantities *in order to* Impact the
258 Future. The aim of the Drivers of Investigation is to ensure that there is always a reason
259 to care about mathematical work—and that investigations provide opportunities for
260 students to make sense, predict, and/or affect the world.

261 This chapter especially addresses Content Connection 1; investigations will live in all
262 three Drivers of Investigation.

263 **The Statistical and Data Science Investigation Process**

264 Statistical and data science investigation is a four-part process, as outlined in Figure 5.3
265 below.

266 Figure 5.3: The statistical and data science investigation process

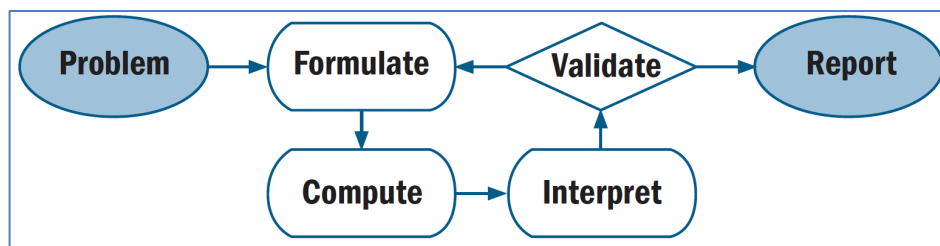


267

268 Source: Bargagliotti et al., 2020.

269 This process has many similarities with the mathematical modeling cycle (graphic
270 below; see Content Connection 2 in Chapter 8 for more discussion of the modeling
271 cycle). Importantly, both include multiple opportunities to revisit earlier steps, in order to
272 revise based on new information or deepening understanding. The tools used in the
273 corresponding “Analyze” (statistical investigation process) or “Compute” (modeling
274 cycle) might differ somewhat, but these two graphics should be understood as
275 describing very similar processes, with slightly different emphases.

276 Figure 5.4: The Mathematical Modeling Cycle



277

278 (1) Asking Questions

279 Formulating questions that anticipate variability should be the beginning of the
280 investigative process. Examples of such questions include:

- 281 ● How fast will my plant grow?
- 282 ● Do plants exposed to more sunlight grow faster?
- 283 ● How does sunlight affect the growth of a plant?

284 These questions contrast with questions that are not investigative and have one
285 answer, such as: How tall is my plant? While questions start the investigative process,
286 students should be encouraged to ask questions throughout the investigative process.
287 (GAISE II, 15).

288 Recent work in the data feminism movement (see, for example, D'Ignazio and Klein,
289 2020) draws attention to the need to understand not just the context of the data, but the
290 motivation behind data collection and to ask questions about who has been included or
291 excluded from data.

292 Survey questions will be important to students' investigations. These are questions
293 designed to elicit data from people in order to address a statistical question, such as the
294 length of time it takes to ride a bus to school.

295 Arnold (2007) notes: "Any question whose investigation requires repeated counting,
296 measuring, or categorizing is one that data helps to answer." Students learn to use data
297 in increasingly sophisticated ways. Early questions are primarily about description,

298 beginning with categorizing and counting, expanding into questions in measurement
299 situations (at first length/distance; later time, area, volume, and rates). Describing
300 relationships between two varying quantities develops as students move through the
301 grades, as do formal quantitative calculations.

302 The Common Online Data Analysis Platform (CODAP) provides a set of databases that
303 will be interesting to school students, such as data on earthquakes, mammals, stars and
304 cities, and an accessible data investigation online tool. Students can be encouraged to
305 ask questions of the data. For example, a data set of mammals may raise the question,
306 “Is the size of mammals related to the length of time they sleep?” (see vignette 1).
307 Students can investigate questions using graphing tools that compare variables,
308 statistical tools, a mapping tool and others.

309 Multivariate thinking comes naturally to humans, and students can develop curiosity
310 about all sorts of data and situations. Young students may ask questions with one
311 variable—such as what is the average age of my class?—but as they get older teachers
312 should encourage bivariate and multivariate (three or more variables) thinking. “Are
313 older students at my school more likely to read more books” would be an example of
314 bivariate data collection. “What are the important factors affecting the growth rate when
315 growing bean plants from seeds?” is an example of a multivariate investigative question,
316 as many variables (e.g., amount of sun, amount of water, temperature, soil nutrients)
317 may affect the plants’ growth rate.

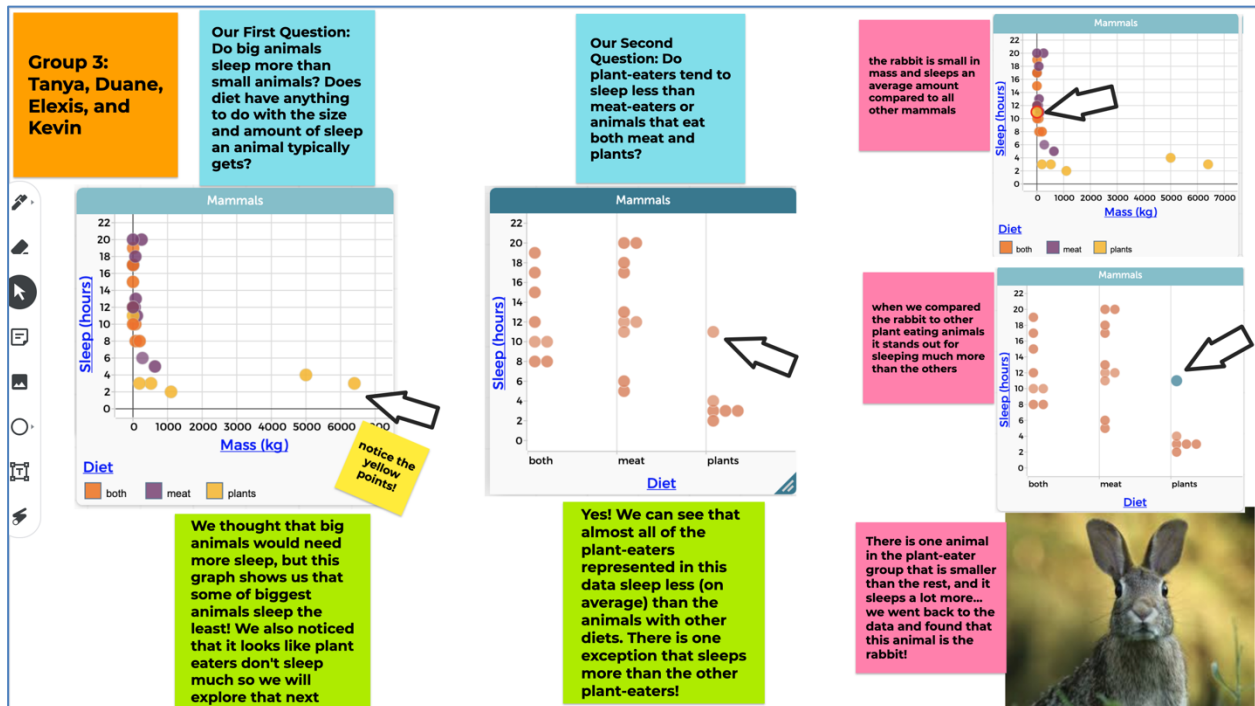
318 ***Vignette 1: CODAP***

319 A group of three students work to explore a CODAP database of 27 mammals
320 (Concord, 2021):

321 The database provides variables such as the height, mass, speed, life-span, and sleep
322 hours of the mammals. The students quickly become curious and ask questions like,
323 “Do bigger animals sleep longer?” They plot the two variables with the graph tool and
324 start to notice a relationship—in the opposite way than the one they thought—it seems
325 the bigger animals sleep less. The students start an animated conversation discussing
326 the reasons this might be, is it because they are more likely to be predators? They then

327 move on to investigate another relationship—who sleeps more, plant or animal eaters?
 328 The students again notice a relationship as well as an outlier (the rabbit) so they wonder
 329 about the rabbit, and look at more rabbit data. The students' investigation of bivariate
 330 data and their relationships is filled with moments of curiosity and excitement, as well as
 331 important learning.

332 Figure 5.5 Student Work from an Investigation of Mammals



333
 334 [Link to long description](#)

335 Source: Concord, 2021.

336 **(2) Collecting and Considering Data**

337 Sometimes students may collect their own data when investigating a question. For
 338 example, they may ask how far do students travel to school? or they may consider two
 339 variables, such as: Are students happier on sunny days? Or they may consider which
 340 plants are most prevalent in their local area. In all of these cases, students could collect
 341 data by observing plants or surveying students.

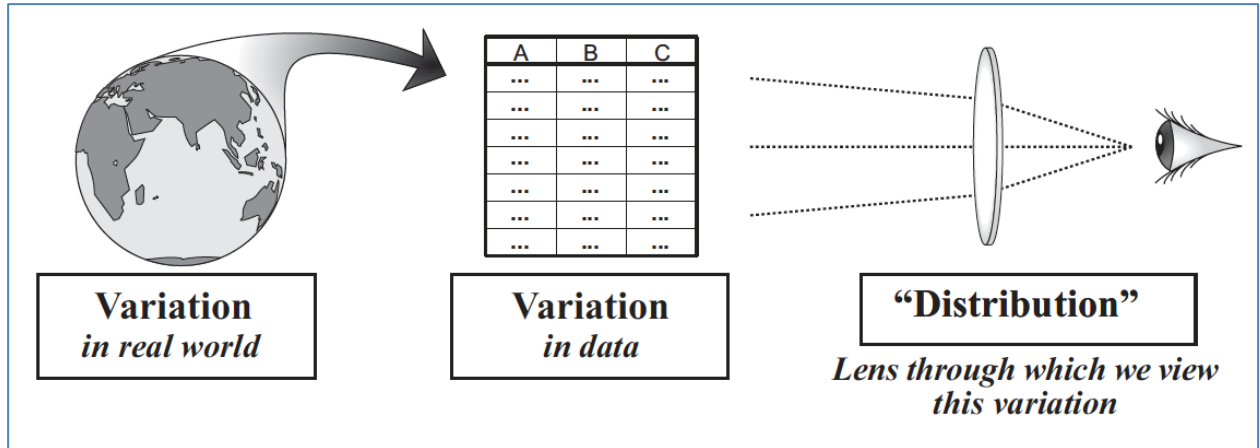
342 Data collection, or consideration of existing data, can be an opportunity for family
343 engagement and support, connecting school with home and community. For example,
344 GAISE II (Bargagliotti et al., 2020, 62–67) includes a vignette of an investigation *Dollar*
345 *Street—Pictures as Data* in which students explore an existing data set documenting
346 families’ living spaces and consumption habits, in order to explore the statistical
347 investigative question, “*How are people’s concepts of family and living spaces similar or*
348 *different across the world?*” This exploration offers many opportunities to challenge
349 cultural misconceptions and “...to realize that our daily lives are often more similar than
350 different.” The California Department of Education provides advice on data that is
351 protected by privacy laws.

352 A key characteristic of data science is asking questions of “big data”—a data set that is
353 complex, messy, and includes many variables. Students can ask questions of big data
354 sets and different students in a class may ask different questions. In a high-school data-
355 science class students can learn to clean data sets—removing any data that is
356 incorrect, corrupted, incorrectly formatted, duplicated, or incorrect in some other way.
357 This is an important part of the work of a data scientist. High school students can also
358 learn to download and upload data, and develop the more sophisticated “data moves”
359 that are important to learn if students are tackling real data sets. Students who take a
360 data science course in high school can also be introduced to programming, in order to
361 interrogate and analyze data.

362 After acquiring or collecting data, students should ask questions about the data—How
363 do the variables differ? How were they collected? Who or what was included in the data
364 collection? This helps students develop an understanding of variability.

365 High school students taking a course in data science may consider more complex
366 conceptions of data science, that are located in the idea of variation, see for example
367 Figure 5.6 from Wild (2006).

368 Figure 5.6



369

370 Details of the understandings that may be developed in a high-school course are
 371 outlined later in this chapter.

372 Sometimes students may first find or be given data, and then ask a question of the
 373 data—reversing the order of 1) and 2).

374 **(3) Analyzing Data and Developing Meaning**

375 In the younger grades, students can analyze and develop meaning from data as they
 376 represent it in different ways, using picture graphs, line graphs, bar graphs and other
 377 forms of data visualization. From sixth grade, students can learn more formal methods
 378 to understand data. The field of statistics has been described as the study of variation,
 379 and students learn about variation when they receive opportunities to consider the
 380 distribution of data. Measures such as mean, median and mode are measures of the
 381 center of a distribution that students learn in middle school. CODAP tools allow students
 382 to see distributions of data and to see, visually, that the spread of a distribution will
 383 impact measures of center. In high school, students will learn about measures of spread
 384 and about regression lines.

385 One of the features of data science is the possibility of predicting outcomes, such as the
 386 cable news programs' predictions of election outcomes. Developing understanding of
 387 what a prediction means, and how to compare predictive strength of one model over
 388 another is not simple and should be developed as a learning trajectory spanning several
 389 grades. Students who specialize in high school can learn about cross-validation

390 techniques. Much of the work of professional data scientists is concerned with
391 quantifying error from predictions.

392 **(4) Interpreting and Communicating Results**

393 Students learn to interpret data in increasingly sophisticated ways. Young students may
394 make statements about their data or create data visualizations to communicate results.
395 They may describe the difference between two groups. Even in the early grades,
396 teachers can have conversations with students about generalizability—how much can
397 we generalize from the data we have collected to broader populations? As students
398 move through the grades they can learn to generalize more formally and to include
399 statements of probability and certainty.

400 A data scientist does not just perform calculation, and an important part of data science
401 is the communication of results. Whereas statistics used to rely on bar charts, pie charts
402 and other familiar representations, data science has created multiple forms of
403 visualizations that represent data. Vignette 2 provides an example.

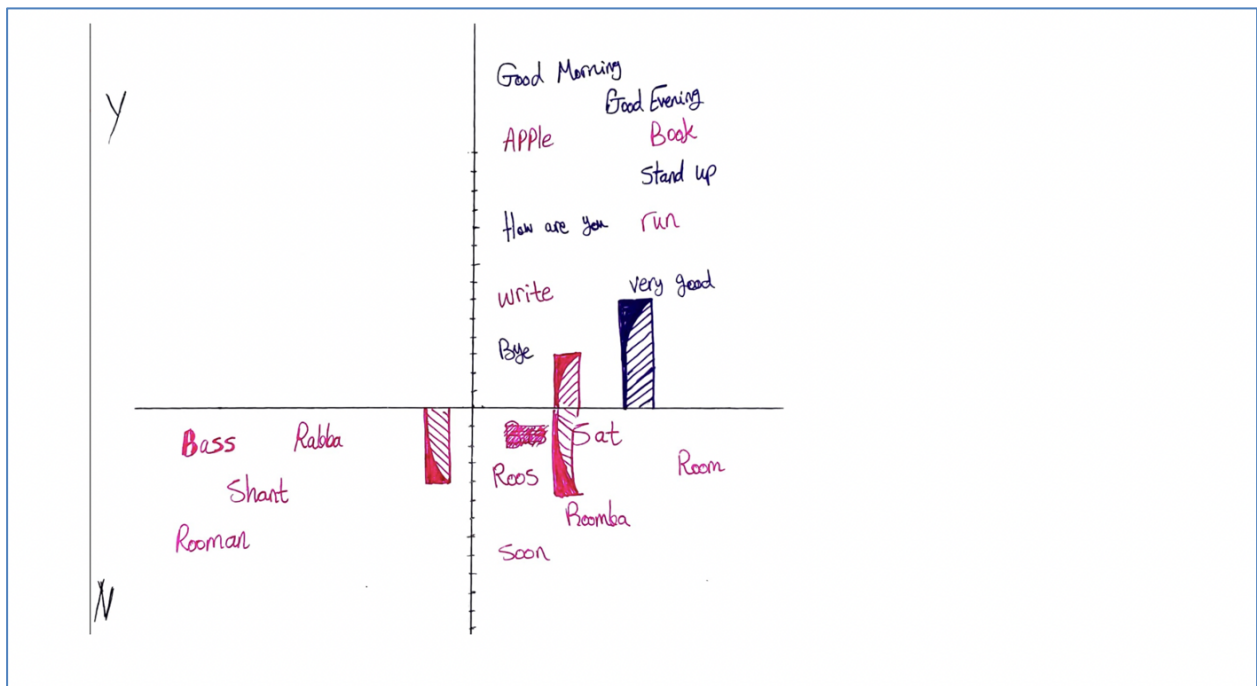
404 Data science is about developing understanding of a situation, it involves holistic
405 thinking, interpretation of meaning, and the communication of complex ideas. An
406 effective data communication draws from writing, and visualizing as well as calculating.

407 ***Vignette 2: Dear Data***

408 Rico shares with his class of students the true story of two designers, who lived on
409 different sides of the Atlantic Ocean—one in London, one in New York. For an entire
410 year the two designers mailed each other a postcard every week, that included data
411 from their lives, that they represented in creative and visual ways. The data
412 representations included multiple variables. For example, some weeks the designers
413 recorded all their moments of indecision, in another they recorded all the times that they
414 laughed. The students looked at some of the data visualizations the designers produced
415 and discussed what they could learn and how they could interpret the different
416 variables.

417 After the discussion, Rico asked his students to collect data over at least a 24-hour
418 period, collecting data on something that interested them, recording at least two
419 variables. When the students came back to class with their data, Rico organized the
420 students into groups and asked them to create data visualizations together, supporting
421 each other to consider ways they would represent different variables. In the discussion
422 Rico paid attention to the language needs of the students, and the ways that the activity
423 aligned with principles of Universal Design for Learning (UDL). Students were excited to
424 make their data visualizations, such as the following:

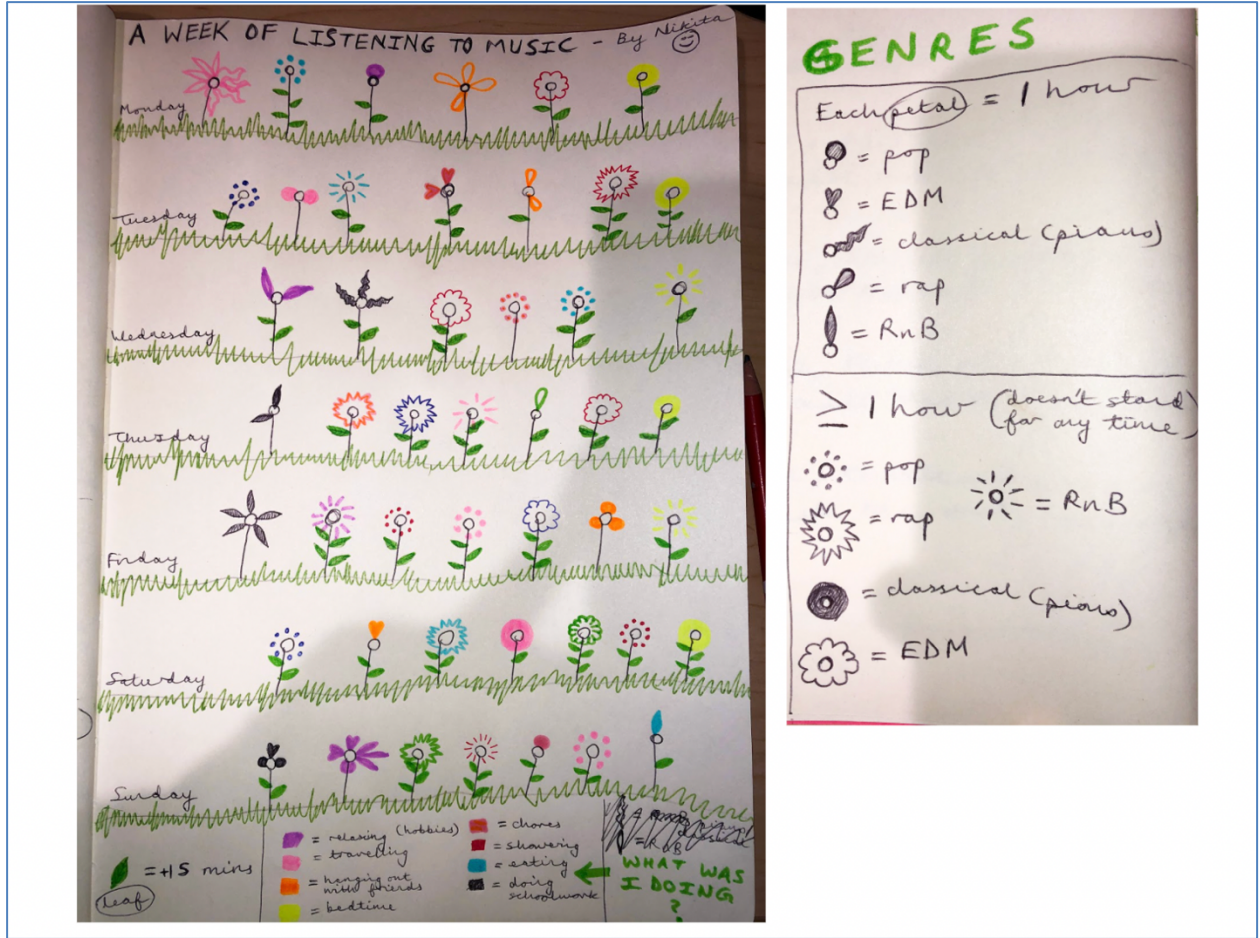
425 Abdu's question: *How many times does Abdu's 6-year-old sister use English and non-*
426 *English words she knows or does not know while pretending to be a teacher?*



427

428 [Link to long description](#)

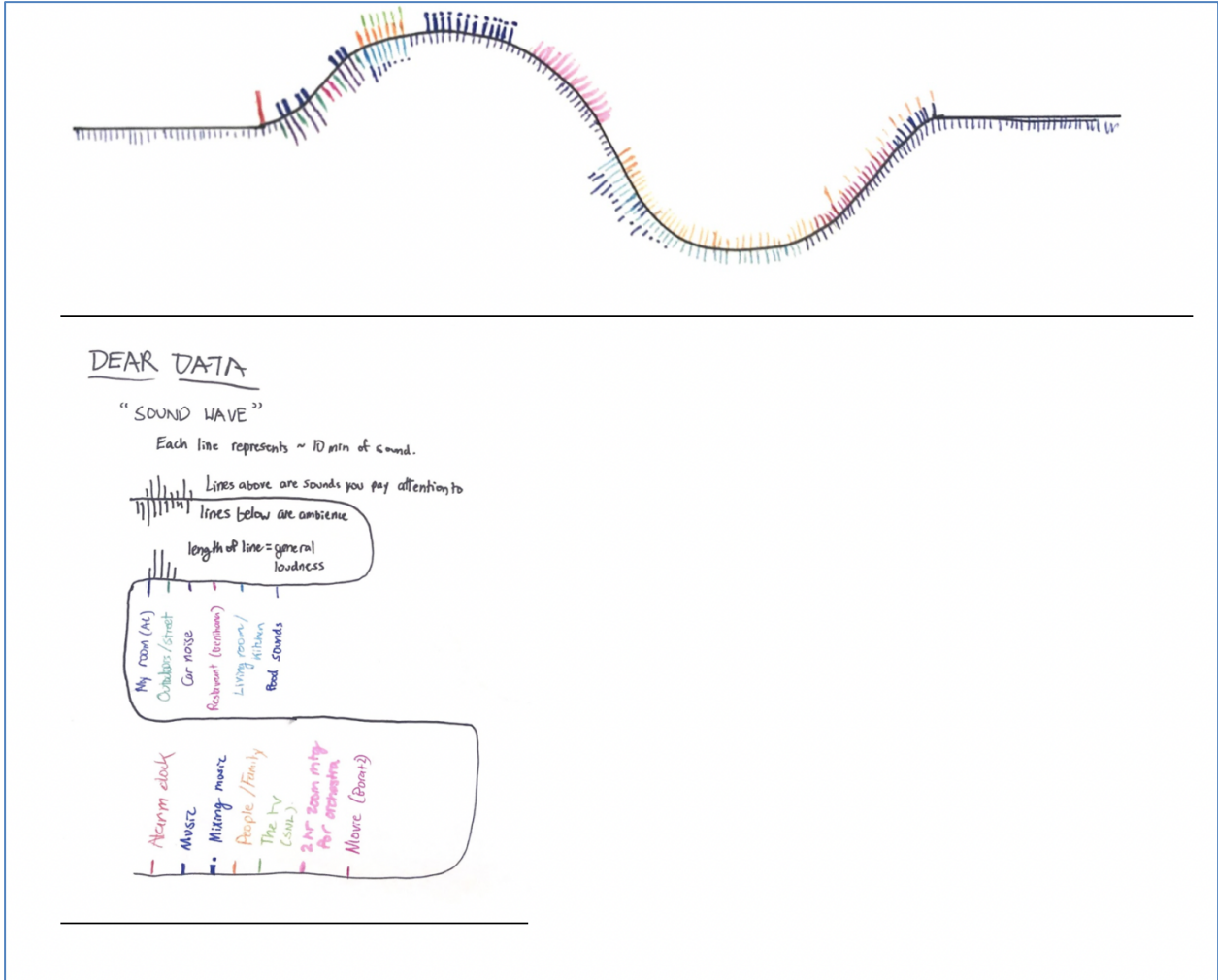
429 Nikita's description: *One week of listening to music, what type of genre it was, and what*
430 *Nikita was doing.*



431

432 [Link to long description](#)

433 Nathan's description: Representation of sound length, level of loudness, and how much
 434 attention was given to it.



435

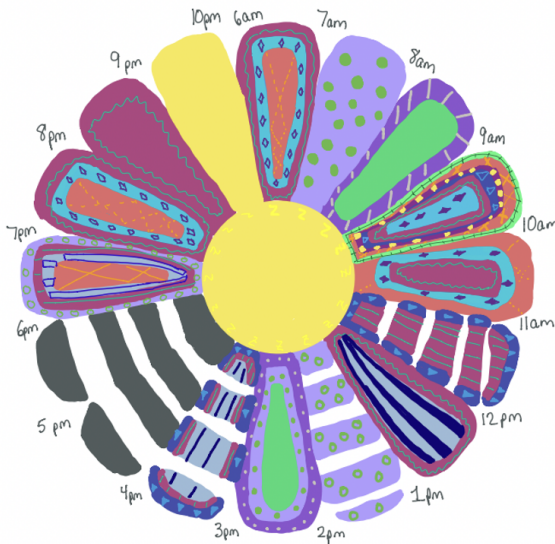
436 [Link to long description](#)

437 Visualization of Kira's Dog's Interactions.

Dear Data:

For a day, between 6am-10pm while I was awake, I recorded my interactions with my dog - Daisy, a golden doodle - and her (sometimes sassy) behaviors. Usually, she is my study buddy for the day.

This data is from Wednesday 11/4/20. Below you will find the key. Something to note, starting from the outer/most layer of a petal and going inward accounts for the order of the actions.



My Actions	# of Occurrences	Daisy's Responses	# of Occurrences
• Physical Pet nn = belly rub ■ = pat	7	• Roll over ↕ = from sitting ↕ = from standing	4
• Show a treat ☼ = cheese ☼ = cookie	2	• Walk away ■ = I changed her □ = distraction	4
• Call name ○ = actual name (course) ⊙ = course name (Daisy, Frankie, etc.)	4	• Come Solid fill = when called	2
• Talk to Daisy x = scolding x' = positive	6	• Paw (beg) ▲ = unprompted △ = Prompted	3
• No interaction	2	• Sleep Z = her bed Z = other	16
• Gave Daisy a shower □ = Paws ■ = full	1	• Muddy # = digging = = rain	1

Not in class = solid White in class = dashed

438

439 [Link to long description](#)

440 The students made their visualizations using Google Jamboards. After they made them,
441 Rico asked the groups to look at the work of other groups and provide feedback to each
442 other on a sticky note. The students were excited to see the ways the different variables
443 related to each other and the ways they could be represented.

444 Data Talks K-12

445 Data talks are short classroom discussions to help students develop data literacy. This
446 pedagogical strategy is similar in structure to a number talk, but instead of numbers
447 students are shown a data visual and asked what interests them. The idea of a data talk
448 was inspired by a New York Times weekly section called, "What's Going On In This
449 Graph?" Students can submit their own ideas to a member of the American Statistical
450 Association, who reveals their thinking on the data in the graphs. In the classroom the
451 teacher can guide the discussions and help students develop important understandings.
452 However, it is important to recognize that teachers do not have to be an expert in the

453 topic of the data visualization—instead teachers can guide and encourage curiosity and
454 question asking. One way to support thinking and speaking like a mathematician is to
455 incorporate writing activities or math journals, which allow students to process learning
456 and continue questioning. These activities help all students gain and exchange
457 information and ideas, and support the California English Language Development
458 Standards’ three communicative modes (collaborative, interpretive, and productive), and
459 allow them to apply knowledge of language to academic tasks using various linguistic
460 resources.

461 If questions cannot be answered by the teacher or students they can be investigated
462 further. Data talks are intended to pique students’ curiosity and encourage question
463 asking, and to help them understand and “read” the data-filled world in which they live.
464 Many of the data visualizations illustrate how multiple variables can be incorporated into
465 one graphic—allowing students to think multivariately.

466 Grades with younger students can use data visualizations with no or few numbers, or
467 smaller numbers, as in Figure 5.7.

468 Figure 5.7

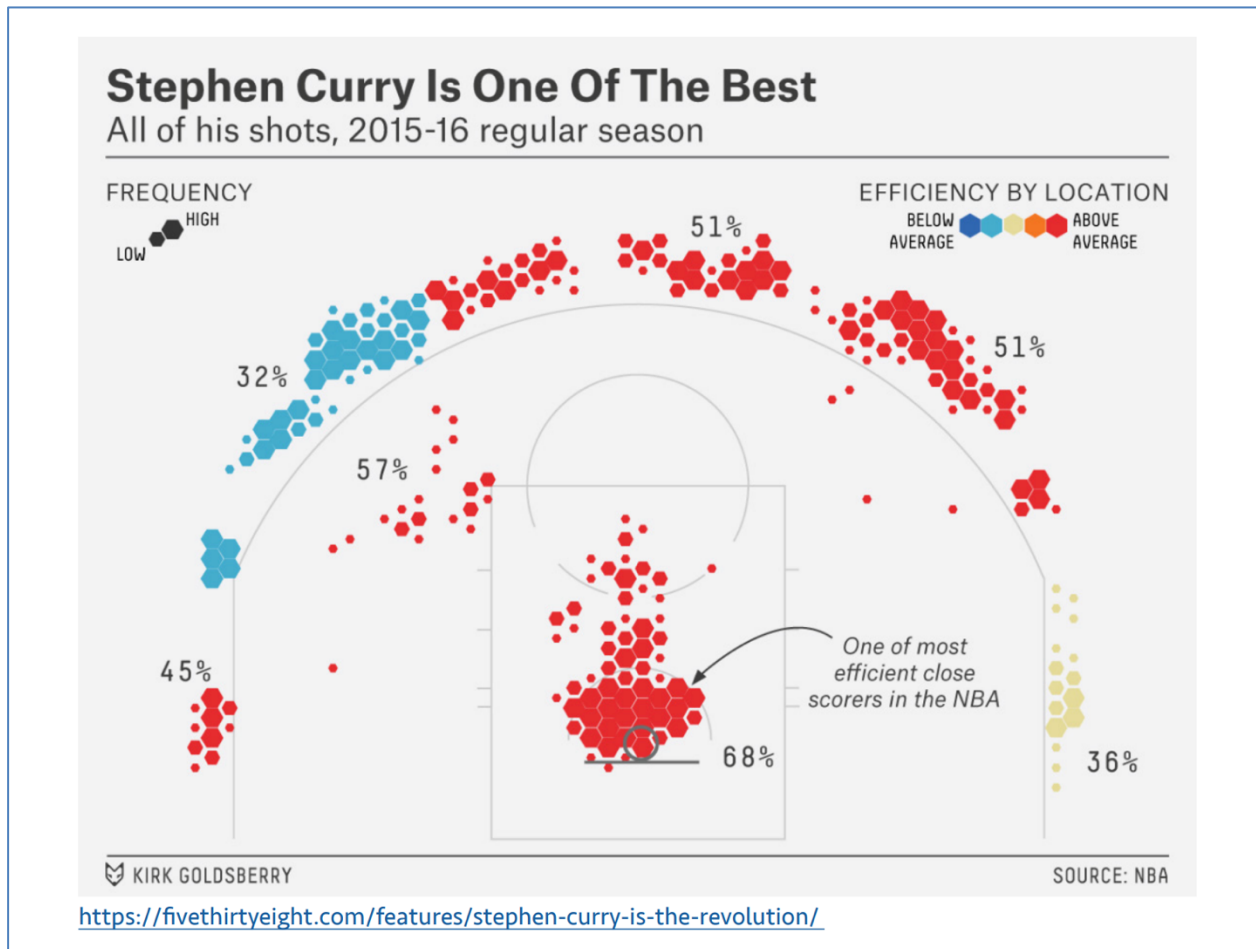


469

470 Source: Youcubed, 2020.

471 Students learn how to calculate percentages in grade 6 but they are usually able to
472 interpret the meaning of percentages, such as those shown in Figure 5.8 below, from
473 earlier grades:

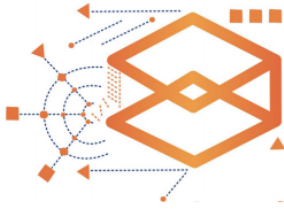
474 Figure 5.8



475
476 Source: Morris, 2015.

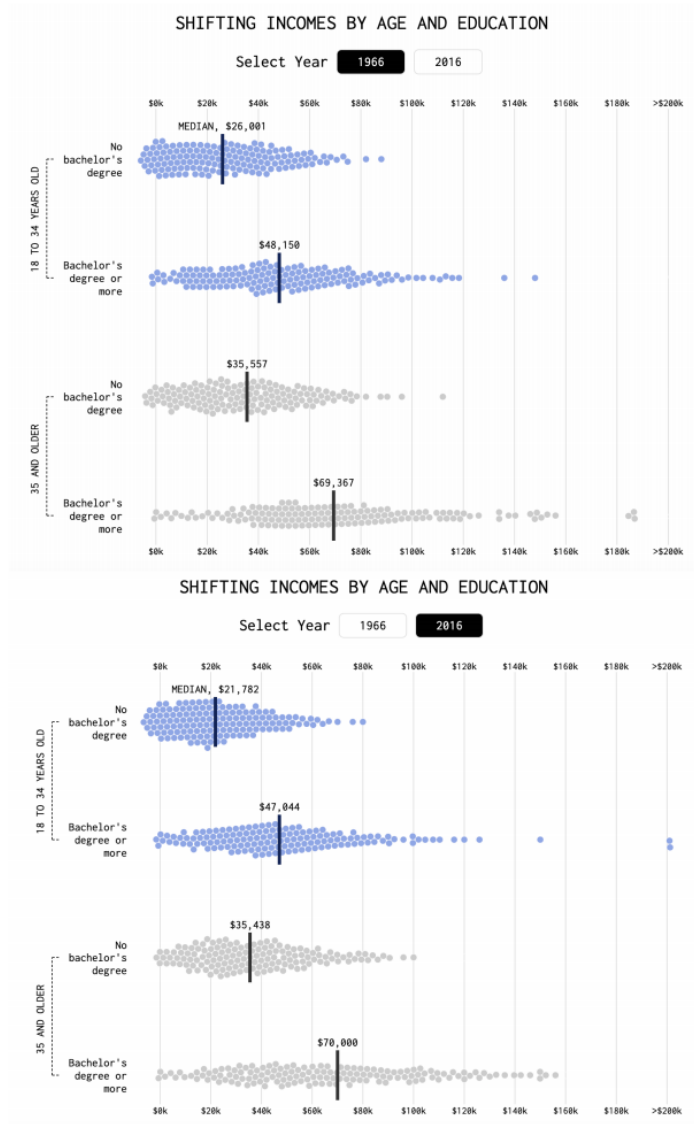
477 In higher grades data visualizations can include more complex data representations like
478 those in Figure 5.9:

479 Figure 5.9



Youcubed Data Talk Shifting Incomes

What do you notice?
What do you wonder?
What is going on in this data visualization?



<https://flowingdata.com/2017/05/02/shifting-incomes-for-young-people/>

480

481 This image establishes a context for a talk around how income shifts based on age and
482 education. The data is from the years 1966 to 2016. The income distribution is shown

483 for people 18–34-years old and those 35 and older, with each group split into “No
484 bachelor’s degree” and “Bachelor’s degree or more” groups. Above the data, questions
485 are posed: What do you notice? What do you wonder? What is going on in this data
486 visualization?

487 **Transitioning from Pre–K**

488 Before kindergarten, children begin to describe their world in language, identifying
489 characteristics of objects, places, people, and events: *The ball is red. My classroom is*
490 *warm. My teacher is in their twenties. Our trip to the park was too short.* Identifying
491 characteristics is the beginning of data, and wondering about characteristics—including
492 countable characteristics—is the beginning of asking questions that data can help to
493 answer. In the California Preschool Learning Foundations, this content is located under
494 the heading of “Algebra and Functions (Classification and Patterning),” in which children
495 “sort and classify objects in their everyday environment,” (by one attribute at around 48
496 months and by more than one attribute at around 60 months of age); and in
497 “Measurement,” in which students compare and order objects directly at around 48
498 months of age and may use an intermediate object to compare at around 60 months of
499 age (Preschool Learning Foundations, Volume 1). These preschool activities directly
500 enable the types of kindergarten through grade five learning trajectory described below.

501 **Kindergarten Through Grade Five**

502 The big ideas of data in these early grades include:

- 503 ● Data for understanding. What questions can we ask? What data do we need to
504 answer it?
- 505 ● Defining data: What is data and how where data collected?
- 506 ● Representing and interpreting data: What does data look like and what does it
507 mean?

508 These ideas are represented in the most important pedagogical and practical process
509 through which data plays a role in making sense of the world, outlined in these four
510 steps, from GAISE II.

- 511 1. Ask a question (SMP.1: Make sense of problems and persevere in solving them)
- 512 2. Collect and consider data
- 513 3. Analyze data and develop meaning (SMP.2, SMP.4, SMP.7)
- 514 4. Interpret and communicate results (SMP.4, SMP.5)

515 An important distinction to consider is between *categorical* (non-numerical, or
 516 qualitative) data and *measurement* or *quantitative* data. For instance, consider a set of
 517 colored blocks in the classroom. “Color” is a categorical or qualitative variable that
 518 students could observe about each block. “This block is 15 centimeters long” is a
 519 measurement data point. The standards develop categorical data in kindergarten
 520 through grade three and measurement data beginning in grade two.

521 Figure 5.10: Examples of Categorical and Quantitative Data

Categorical Data	Quantitative (or Measurement) Data
<ul style="list-style-type: none"> • Color (red, green, blue, yellow) of blocks in the class set • Species of trees on the school grounds • Identification of schools in the district as “elementary school,” “middle school,” or “high school.” 	<ul style="list-style-type: none"> • Height (or circumference of trunk, or biomass) of trees on the school grounds • Number of pages (or weight, or height) of books in the classroom • Annual income for households in a census tract

522 **What questions can data help to answer?**

523 All work with data should begin with noticing and wondering: “I notice that...” or “I
 524 wonder what...” or “I wonder how many....” To prompt wonder, teachers can ask: “What
 525 do you notice or wonder about here [in this context], that we could (count/measure/keep
 526 track of) to figure out or explore further?” To establish effective routines, and to support
 527 language development in “I wonder” activities, it can be effective to provide these
 528 examples as sentence starters.

529 As students gain confidence in their ability to speak like mathematicians, statisticians
 530 and/or data scientists, the teacher should encourage students to generate questions
 531 themselves to build their agency in using mathematics to make sense of their worlds. A
 532 weekly whole-class “I wonder” routine—in which students propose questions to

533 investigate by collecting data—would build a powerful practice of observing the world
534 with a data lens, contributing to students’ development of modeling with mathematics
535 (SMP.4).

536 In kindergarten, students compare the number of objects in different categories
537 (K.CC.C.6) to answer “Which has more?” questions (*I wonder whether there are more*
538 *square blocks or more triangular blocks on the desk?*). At first, the teacher suggests or
539 specifies categories; eventually students generate ideas for classification. They also
540 directly compare (as opposed to measuring with a unit or an intermediate) objects with
541 common measurable/countable attributes to see which has more (K.MD.A.2, K.G.B.4) (I
542 wonder which shape has more sides? Which kind of block is heaviest?—using a
543 balance or informal one-in-each-hand comparison, rather than a scale). “I wonder...”
544 questions should explore both of these: two-category, “Which is more?” questions and
545 comparison of objects according to length, height, weight, and countable attributes like
546 number of sides. Student-generated questions provide opportunities to work on
547 precision of language as well—for example, when students are asked to clarify what
548 they mean by “bigger.” Mathematics discussions that are rooted in academic language
549 will help students understand mathematical concepts more deeply as well as discover
550 new ones. As the years progress, students or teachers may reach beyond the
551 classroom to find contexts for their: “I wonder...” questions.

552 In addition to questions that can be answered with a single value, students can start to
553 pose statistical investigative questions that involve multiple variables such as, I wonder
554 if plants grow more with more sunlight? Or I wonder if age affects which color people
555 like?

556 In first grade, measurement of length and time are the contexts to emphasize in
557 generating questions (1.MD.A.2, 1.MD.B.3), along with continued work categorizing and
558 counting objects (1.MD.C.4) and categorizing geometric objects by attributes (1.G.A.1).
559 Second graders should continue to explore questions in length measurement
560 (2.MD.D.9) and time (2.MD.C.7) contexts, and add money contexts (2.MD.C.8). When
561 selecting “I wonder” questions, it is important to avoid situations that serve as markers
562 for economic or social status, e.g., “I wonder who has the most expensive backpack,” “I

563 wonder who is the most popular kid in school,” or, perhaps less obvious, “I wonder who
564 has the newest shoes.” It is similarly important to avoid questions about students’
565 physical attributes, even those that seem innocuous such as height or arm length.
566 Instead, some good questions to wonder about might be “I wonder what time it will be
567 when the next person walks into the classroom,” or “I wonder which book in the
568 classroom is the most read,” comparing events or objects rather than personal
569 characteristics.

570 In third grade, contexts for questions to investigate using data should expand to include
571 volume and mass measurement (grams, kilograms, and liters, but not compound units
572 such as cm^3) in addition to the length, time, and money contexts from earlier grades
573 (3.MD.A.2). Time measurements are refined to the nearest minute (3.MD.A.1) and
574 length now includes half- and quarter-inches (3.MD.B.4). Beginning ideas of area give
575 another possible context, limited here to areas that can be covered by a whole number
576 of unit squares (3.MD.C.5, 3.MD.C.6).

577 In fourth grade, a significant context for data-investigation questions is classification and
578 analysis of two-dimensional shapes (4.G.A.2). Incorporating this geometry standard to
579 help build data understanding can foster the important practice of analyzing by
580 attributes—one instance of SMP.7 (Look for and make use of structure). Fourth-grade
581 students also extend the set of units they work with (4.MD.A.1) and can generate data
582 about area for more complex shapes. Fifth graders deepen their understanding of
583 volume to include unit cubes, making this an important context for data-inquiry
584 questions. A teacher could invite students to build a structure out of multi-link cubes and
585 then collect data from the class by asking, for example, how many cubes they use in
586 each of their different structures they built, or the height and width of their structures,
587 and color of the blocks. Students can collect data on multiple variables.

588 In kindergarten through grade five, “I wonder...” questions come primarily from personal
589 experience. See below for additional examples.

590 **Asking Questions, Collecting and Analyzing Data**

591 Questions invite inquiry. An important part of students' kindergarten through grade five
592 experience should involve coming to recognize that, when they choose and pose
593 questions, they can collect or analyze data to find answers (SMP.4). Some of the most
594 valuable conversations about data occur when students notice patterns in a data set
595 and begin asking questions. Remaining alert for these everyday moments—perhaps in
596 attendance, weather, or lunch-count data—may generate opportunities for discussing
597 statistical investigative questions and exploring how data can help answer them.

598 As students pose authentic questions reflective of those described above, they should
599 also encounter opportunities to help determine how data might be produced to answer
600 them. In addition to producing data directly through their own observations, students
601 should gain exposure to designing and using surveys and simple experiments to
602 generate data. By producing their own data from their classroom or community (*How*
603 *does age of students relate to their enjoyment of school? Does time on social media*
604 *apps increase with age? How much waste is generated by different companies/our*
605 *school?*), students recognize data as having context and deriving from observation and
606 measurement, and they come to see data (and mathematics more broadly) as a tool to
607 help think about their worlds. Data gathered by others (such as those in the data talks)
608 can help to answer questions students generate about their own communities.

609 When choosing data tasks that include categorizing and counting, consider the grade
610 level expectations for counting (up to 10 objects scattered, or up to 20 if arranged in a
611 line, array, or circle, in kindergarten [K.CC.B.5], 120 by the end of first grade
612 [1.NBT.A.1], and up to 1,000 by the end of second grade [2.NBT.A.2]). Such tasks can
613 also be structured to build place value understanding.

614 In kindergarten, once students notice things in a context and wonder about a question,
615 they describe measurable, countable, and observable attributes of objects or situations
616 (K.MD.A.1, K.G.A.1, K.G.B.4), and classify objects and count the number in each
617 category (K.MD.B.3), such as categorizing a set of cubes by color. In this last context,
618 both “this cube is red” and “there are 13 red cubes in the set” are data points. Notably,

619 most work on *number* in kindergarten should be with numbers representing quantities of
620 objects (SMP.2); thus, most numbers encountered in kindergarten are actually data.

621 In first grade, students explore their time and length questions by measuring lengths of
622 objects which are a whole number of units (1.MD.A.2) and telling and writing time in
623 hours and half-hours. Counting and categorization situations should include up to three
624 categories (1.MD.C.4). Second graders measure length to the nearest whole unit
625 (2.MD.D.9), using different standard units (centimeters, meters, inches, feet) (2.MD.A.3)
626 and several tools (2.MD.A.1) and measure time to the nearest five minutes (2.MD.C.7).

627 Students in grades three through five refine their measurements of lengths and time,
628 and expand the set of units they use; and they add area and volume measurement to
629 their repertoires (as described above in “What questions”). By the fifth grade, students
630 should understand that data sets can include different types of variables, such as
631 categorical and quantitative. They should recognize that an individual instance or object
632 can possess attributes that exemplify these different types, and should have gained
633 experience measuring, characterizing and analyzing such diverse types of data and
634 associating them together.

635 An important understanding that students need to develop through kindergarten through
636 grade five is the idea of variability and variables. When students ask questions such as:
637 How high are the plants in the classroom? they are considering one variable: height.
638 When they consider whether older students spend more time on social media apps,
639 they are collecting bivariate data—with two variables—age and time. When they make
640 their own data visualizations, as seen in Vignette 2, they may collect data on multiple
641 variables. Multivariable thinking is important to develop through the grades.

642 **Interpreting and Communicating Results**

643 Sorting objects into two categories and representing these categories by their count
644 (K.MD.B.3) is a first example of students representing data to help make sense of their
645 worlds (SMP.4). First-grade students organize up to three categories and ask and
646 answer questions about the relative sizes of categories and about the total number of
647 data points.

648 Second grade begins an expanded focus on data representation, introducing line plots
649 (whole number units only; 2.MD.D.9), picture graphs, and bar graphs. These graphs can
650 be used to answer put-together, take-apart, and compare questions (2.MD.D.10). In
651 third grade, *scaled* picture and bar graphs are added as a tool for visualizing “how many
652 more” questions (3.MD.B.3), and line plots may have half-unit and quarter-unit markings
653 as appropriate (3.MD.B.4). In fourth grade, line plots may display additional fractional
654 units (to eighth-units), and be used to answer additional questions about differences—
655 between maximum and minimum measurement, for example.

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

663 While the data visualizations mastered by fifth grade only include picture graphs, bar
664 graphs, and line plots, students do not need to be restricted to these. Each of these
665 represents repeated measurements of a *single* varying quantity; science curricula in
666 particular, and many questions of interest in general, require the consideration of
667 relationships between *two or more different* changing quantities, such as erosion and
668 time (NGSS 4-ESS2-1 Earth’s Systems) or length or direction of shadows and time
669 (NGSS 5-ESS1-2 Earth's Place in the Universe). Such reasoning involving multiple
670 variables is an important aspect of modern encounters with data, and students should
671 experience it at all levels. Although the scatter plot, a crucial data representation tool for
672 two varying quantities, is not mastered until eighth grade (8.SP.1), it must be explored
673 informally much earlier for students to be able to meet the eighth-grade expectations.
674 For example, students can plot quantities changing over time (e.g., height of a plant,
675 length of the day, high temperature for the day, temperature of a glass of water every
676 minute for an hour), with time on the horizontal axis and the changing quantity on the

677 vertical. Once such a plot is created, it is an excellent context for a “notice and wonder”
678 discussion.

679 In recent years, new technological tools and developments in data science have
680 prompted an explosion in interesting data visualizations, many of which are quite
681 comprehensible to young students with some exploration. Experiences with different
682 visualizations will further expand students’ sense-making opportunities and encourage
683 them to think about what they can understand looking at data sets in different ways.
684 Newspapers and online news sources offer specific examples; student-gathered
685 examples help to build buy-in for a “can we figure out what this visualization is trying to
686 help us to understand?” routine.

687 Interpreting data is a matter of making inferences from the data available. While
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 both to wonder whether patterns or trends they notice in data
692 extend beyond the particular group that generated the data, and to be skeptical about
693 such extensions to larger populations (including considering ways in which the group
694 might not be representative of the larger population). Additionally, students should learn
695 that good claims draw upon data as evidence and that they always come hand in hand
696 with a degree of uncertainty. Modeling the use of appropriate terminology such as
697 “tends to,” “typical,” “usually,” and “similar” can help lay important groundwork for this
698 concept (Rugin, 2019).

699 **Preparing for the Major Data Science Work of Grades Six Through** 700 **Eight**

701 *Understanding Variability:* Variability is everywhere, and understanding variability is the
702 core of developing data sense. While outcomes for variability and distributions are not in
703 the standards until middle school, it is essential that kindergarten through grade five
704 students encounter many experiences with variation, including counting, measuring,
705 and observing quantities and characteristics that vary in order to be prepared for the

706 first big idea in the grades six through eight section below. In particular, their encounters
707 with data representations should highlight important ideas that set the stage for more
708 involved work with distributions.

709 When working with visualizations of data, students should consider not only the most
710 popular value in a dataset (the mode) but also describe the shape and spread of data
711 distributions. Identifying the maximum and minimum values of quantitative datasets can
712 help students appreciate the concept of range as a measure, and looking for clusters
713 and gaps in a distribution can begin to help them attend to its shape. As they engage in
714 experiences where they produce their own data through measurement, teachers should
715 highlight for students the variation that results. Measuring the same variable on multiple
716 individuals or objects, for example, results in data that vary, and students should
717 consider the causes or sources that might have given rise to the variation they have
718 observed, working as they do so to differentiate between variation and error. For
719 example, if students plant a particular variety of flower seed at multiple locations around
720 the school, then measure the plants' height and the amount of sunlight each month,
721 they can conduct investigations into the ways plant growth and sunlight relate to each
722 other. They should discuss and describe any patterns in their bivariate data, and
723 discuss reasons for the variability. Finally, they should consider their own measurement
724 techniques, and how confident they are that they all measured the same way (so that if
725 someone else measured, they would get the same height or sunlight).

726 *Randomness, probability, and uncertainty:* Randomness is a complex idea
727 encompassing uncertainty *and* a level of predictability. When (blindly) drawing a cube
728 out of a bag containing three blue cubes, two red cubes, and one yellow cube, nobody
729 can predict with certainty what will happen on a single draw. But, over many draws, the
730 person who always predicts a blue cube will be right about half the time. Activities that
731 demonstrate this can be used to generate data for many of the explorations of the big
732 ideas above, which will leave students well-prepared for a more formal treatment of
733 randomness and probability in middle school. At this point, students should begin to
734 conceive of probability as a measure of the chance that something will happen, seeing it
735 as a basic measure of certainty or uncertainty.

736 *Technology:* California’s 2018 Computer Science Standards include computer-based
737 data sorting, categorizing, and visualizing for students in kindergarten through grade
738 two and grades three through five (CS K–2.DA.8, CS K–2.DA.9, CS 3–5.DA.8). These
739 standards are important preparation for middle and high school use of data software to
740 visualize and interpret large data sets.

741 Finally, it is worth noting that (as in science and other fields) many questions that
742 students might wonder about will not be fully answerable using tools designed for
743 kindergarten through grade five. It is important that teachers have resources for helping
744 students figure out which aspects of questions can be investigated with currently
745 available tools, and have some understanding of data science tools which students will
746 encounter later. For example, many will wonder about relationships between two
747 different variables: *If I get up earlier, do I feel tired earlier in the afternoon at school? Do*
748 *students who skip lunch eat more candy in the afternoon?* When one of the variables is
749 categorical (like the skipping lunch question), separate line plots can be made for each
750 category and the line plots compared. When both variables are quantitative, students
751 could input data into CODAP and 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.

757 ***Vignette 3: Logan from Kindergarten through Grade Five***

758 A small sampling of Logan’s data science experiences in kindergarten through grade
759 five is described below. This is not intended to capture *all* of their data science
760 experience, only to indicate a development towards powerful uses of data to understand
761 their world. In each grade, Logan generated questions and gathered data (steps 1–3 in
762 the process described at the beginning of this kindergarten through grade five section)
763 and represented and interpreted data (steps 4–5).

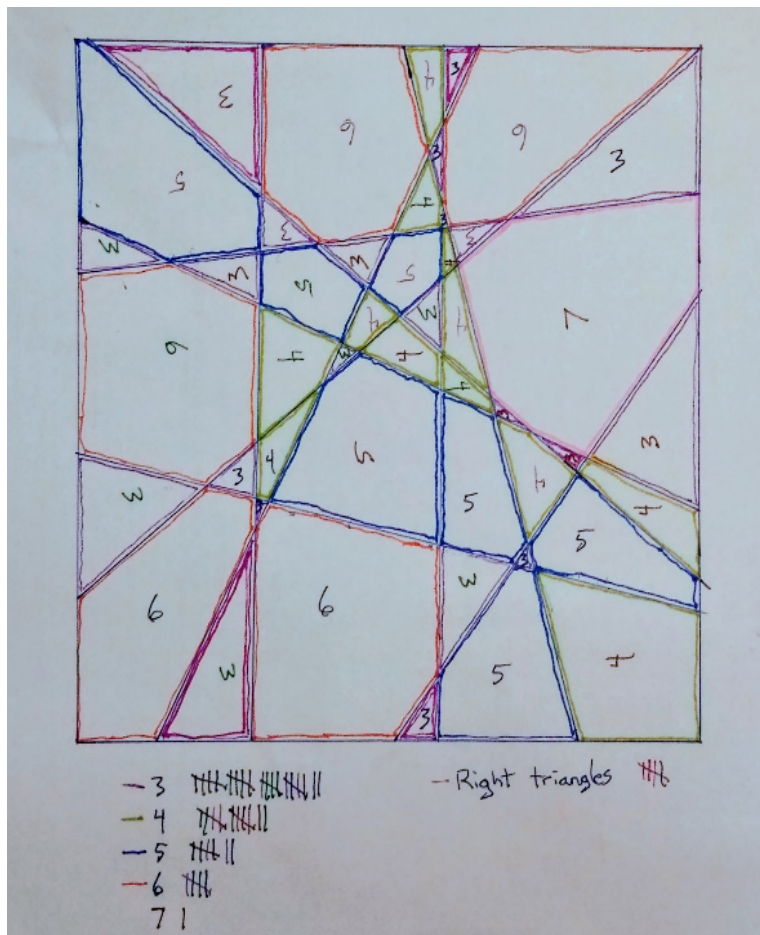
764 Logan was described as an active child and gravitated towards the kindergartners who
765 ran and climbed. Logan’s teacher asked lots of questions of students about what they
766 noticed in the classroom and around school, both inside and out. These ranged from
767 specific “how many in each category” questions (how many classroom doors of each
768 color are there?), to direct comparison questions (which slide is taller, the blue one or
769 the green one? How do you know?), to, eventually, *types* of questions: What are some
770 things at school whose size we could compare? Recording students’ observations and
771 category counts allowed all students to pose and answer “relative size” questions.

772 In first grade, student teams were asked to think of two similar things at school, such
773 that they weren’t sure which was taller, and then to find a way to compare their heights.
774 A variety of materials was available to use in the comparison. Logan’s team was able to
775 compare the height of the slide in front of the school with the height of the slide behind
776 school, measuring the height of both using towers of large DUPLO® bricks. The whole
777 class used their data to discuss how much taller the slide in front was. After this, Logan
778 wanted to build DUPLO® towers to measure height and length of lots of things, and was
779 disappointed that the class didn’t have enough bricks to measure the height of the
780 school (and that their teacher wouldn’t let them climb the school). As a class, students
781 checked the length of the day (sunrise to sunset) each day, and maintained a running
782 visible tally of the number of school days with less than 11 hours of daylight, 11 to 13
783 hours, and more than 13 hours, for the entire year. As a class, they discussed what they
784 thought might happen to the number of hours of daylight in the future and checked the
785 data a month later to see whether their predictions were correct.

786 Logan’s second grade class marked their own yard sticks (marking a wooden blank in
787 inches using only a three-inch by five-inch card), and then used them extensively to
788 measure objects of interest to the nearest inch. Later, they added centimeter markings
789 to the other side of the yardstick, and discovered that measuring the same things with
790 smaller units led to larger number measurements. When choosing an activity to time,
791 Logan’s group decided to time and record the amount of time in a week that team
792 members spent reading in school, and compare those measurements over several
793 weeks (this had the benefit that team members read much more during those weeks!).

794 Other teams measured time spent playing outside, listening to announcements, and
795 working at math stations. Teams made line plots of their data, and compared the line
796 plots of different activities to discuss how students typically spend their school time.

797 As mass and volume became available in third grade as characteristics to measure,
798 Logan's class used length/height, mass, and volume measurements to examine
799 collections of objects. The line plots of the masses and the lengths/heights of the
800 objects in the science corner looked quite different from each other; similarly for the line
801 plots of volume, height, and mass of all objects in the room which hold water (vases,
802 cups, etc.). Logan's team had a great disagreement about whether a taller vase should
803 hold more water than a shorter vase; the class eventually decided that this was usually
804 but not always true.



805

806 One of Logan’s favorite activities in fourth grade was one that combined data work with
807 classifying shapes by attributes: collaborative art pieces: Each team had a 1/2-meter by
808 1/2-meter square on the board, and each student in the team drew in two edge-to-edge
809 straight lines of their choice, using their meter sticks. Then one student in class chose a
810 shape to try to find in the drawings, and each team outlined each new instance of that
811 shape they found and described how they knew it was a triangle, rectangle, right
812 triangle, quadrilateral, etc.; this was repeated for several other shapes. They made an
813 individual card to represent each piece of artwork, using the card to represent the
814 artwork the different variables they measured for each piece (how many triangles, how
815 many instances of each color, how clean or messy each line was, etc.). When they had
816 made a full set of cards, they sorted them in various ways, then made a table to
817 compare the tallies for the different pieces, discussing how the different features of the
818 art and the process of creating it that might help explain the variations in their data.

819 By fifth grade, Logan and their classmates had constructed many line plots, and thus
820 often wondered about quantities that vary on repeated measurement: The cartons of
821 milk from lunch say they each contain 8 fluid ounces, but yours feels heavier than mine;
822 am I getting less or are you getting extra? The weather site says the average high
823 temperature here is 57°F (degrees Fahrenheit) in November, but today it got up to 65°F.
824 How can we check whether this month is near average? To explore the first, the school
825 donated 20 cartons of milk to the experiment. When they examined the dot plot of their
826 measured volumes, they saw that it had a tightly clustered shape, with a minimum
827 measurement of 7.8 fluid ounces and a maximum of 8.2 fluid ounces, and that the most
828 frequent value was 7.9 ounces. One student in the group thought that some milk
829 probably remained in the containers, so the group spent a while trying to figure out how
830 they might identify how much had been left inside (teams came up with several
831 methods). For the second, the class recorded the daily high all month, recorded them
832 on a line plot which also had marked the “average” high temperature from the weather
833 site, and used the line plot at the end of the month to discuss whether it was consistent
834 with the stated average (without computing an average of the data).

835 Students like Logan, with a rich variety of experiences using data to explore contexts
836 and questions of interest, will be well-prepared to use mathematics, and data in
837 particular, to make sense of their ever-expanding worlds. Data sets will get larger, and
838 contexts wider, in ensuing years.

839 **Grades Six Through Eight**

840 Middle school includes a big expansion in important ideas. The big ideas of data
841 science include

- 842 ● Data in the world: exploration, interpretation, decision making, ethics
- 843 ● Variability: Describing, displaying, and comparing
- 844 ● Sampling to understand a population: randomness, bias, how many?
- 845 ● Are they related? Multivariate thinking
- 846 ● What are the chances? Probability as the basis for data-based claims

847 As in earlier grades, students experience data science as a tool to help understand their
848 worlds via a process that begins with wondering questions. This is also the beginning of
849 the mathematical modeling cycle (Pelesko, 2015) and the statistical and data science
850 exploration process, and of investigations in science (NGSS Lead States, 2013).

851 The GAISE II Statistical and Data Science Exploration Process:

- 852 1. Curiosity and question asking
- 853 2. Collect and consider data
- 854 3. Analyze data and develop meaning
- 855 4. Interpret and communicate results

856 This process, beginning with noticing and wondering, often gets lost in the details of
857 step 3, which contains the different statistical methods that have been developed for the
858 analysis of data. It is crucial to keep all work with data tied to authentic questions.

859 Prediction is a key activity that builds student ownership of the process and conclusions;
860 it also builds a habit of asking “does this make sense?” (SMP.1) by comparing results
861 with expectations.

862 **Data in the World: Question Asking, Exploration, Interpretation,**
863 **Decision Making, Ethics, Technology**

864 What functions does data science play in the modern world?

- 865 ● *Question Asking and Exploration:* Data science and statistical exploration starts
866 with questions that are posed by students. When students are invited to wonder
867 about situations, and when they are given interesting datasets they will become
868 curious and can ask questions of data that they can explore and investigate.
869 Data exploration includes understanding the context and situation, data should
870 never be abstracted from their context. Students can look for hidden patterns and
871 associations. Any patterns or associations discovered can lead to new
872 conjectures or questions to investigate further. In eighth grade, students can
873 begin this process with datasets that include multiple variables, such as those
874 examples provided in the vignette above. As mentioned in the chapter
875 introduction, vast quantities of data are collected every day, and only a small
876 fraction are analyzed.
- 877 ● *Interpretation:* Every encounter with data should revisit the context from which
878 the data originated, interpreting results of data analysis in that context. This
879 includes answering any questions that began the encounter and reporting any
880 other associations or patterns that were discovered.
- 881 ● *Decision making:* Commonly, data is used to inform decisions following the
882 question/data/represent/interpret process. Often, however, data is used to justify
883 and explain a decision, even if data didn't play a meaningful role in the decision.
884 There is a high risk for abuse, however, when including data collections that
885 support the predetermined decision and leaving out those that do not.
- 886 ● *Ethics:* Modern, ubiquitous data collection raises a host of ethical questions, both
887 about how and what data is gathered and stored, who is included or excluded in
888 the data, and how that data is used and presented. Middle-school students need
889 to understand their own online data footprint (for example, how companies
890 aggregate information about individuals to create detailed profiles) and should

891 confront scenarios in which they must make decisions in hypothetical situations
892 involving data exposure, consent for data collection, etc.

- 893 • *Technology*: California’s 2018 Computer Science Standards expect students to
894 make use of computers for data organization and visualization (CSS 6–8.DA.8).
895 More importantly, given the amount of data collected and stored today, real-world
896 datasets are incomprehensible without such computer assistance. Students
897 should use modern data software extensively, especially for organizing and
898 displaying features of data set.

899 **Describing, Displaying, and Comparing Variability (Grades Six** 900 **Through Seven)**

901 Sixth-grade students build on earlier experiences by distinguishing between statistical
902 questions, which can be investigated using data that varies (analysis of social media
903 usage by age of students), and questions without variations in (correct) responses (How
904 many days are there in January?) (6.SP.A.1). When considering a statistical question,
905 they understand that the variation in numerical data has a distribution which can be
906 described by its center (first the median, then the mean), its variability (also called
907 spread—described both qualitatively and via a numerical measure, either inter-quartile
908 range [IQR], range, or mean absolute deviation), and an overall shape (including
909 descriptors such as symmetric, skewed left or right, peak, gap, and outlier) (6.SP.A.2,
910 6.SP.A.3). As students explore datasets, they can produce visual representations of the
911 distribution of their data; they can look at the shape of distributions that have different
912 measures of center and spread, and develop visual understandings of the shape of
913 distributions.

914 Students should have experiences, beginning in sixth grade, deciding which measure of
915 center is a more useful descriptor of a typical value for data sets with different shapes.
916 Because the mean is sensitive to extreme values, the median is often a more useful
917 measure for skewed distributions; in this case, the inter quartile range is a useful
918 measure of variability. For some distributions—with multiple clusters, for example—

919 students may decide that neither median nor mean is a useful measure, and might
920 decide that a single number cannot reasonably represent a typical value (6.SP.B.5).

921 Two tasks that reinforce the notion of these standard measures and replace rote
922 disconnected calculation with conceptual thought are:

923 1. Students form a “name count line” creating a human graph to depict how many
924 letters are in their first name. (All the students with five-letter names stand in a
925 line, all those with four letters form a similar line to one side, and those with six
926 letters form a line to the other, etc.) Then the teacher instructs one student from
927 each end of the human graph to sit down. After repeating this multiple times, only
928 one or two student(s) are still standing. If one, that student represents the median
929 name length. If two, the median name length is halfway between the name
930 lengths of the standing students.

931 2. Students are invited to explore the CODAP dataset of four elephant seals
932 (Concord, 2021):

933 The dataset includes data on the paths taken by the seals – visible on a mapping
934 tool, the distance they swim, their latitude and longitude, the depth and
935 temperature of the water and more.

936 In groups students are invited to explore the data and form investigative
937 questions. Students start by plotting different variables with the graph tool, to
938 consider the shapes of distributions. They choose to display the mean and
939 median and consider how the measures of center relate to the visual distribution
940 of the data. They form questions they are curious about: Do certain seals prefer
941 deeper water? Does the distance seals swim relate to the temperature of the
942 water? As students explore these questions, they plot two variables on a graph
943 and consider the slope of the relationship, they even add a third variable which is
944 shown through color coding. Students learn to be comfortable investigating data,
945 making use of measures to learn about their data.

946 Visual representations of distributions include box plots and histograms in sixth grade,
947 adding to the line plots (called dot plots from grade six onward) from earlier grades
948 (6.SP.B.4). In addition, students learn to report and interpret measures of center and
949 variability, and descriptions of distributions, in the context in which the data arose
950 (6.SP.B.5). Seventh- and eighth-grade standards do not include additional
951 representations of single-variable data sets, but these students should continue to
952 create visual displays of such distributions.

953 In seventh grade, comparisons between two populations with similar variables is a
954 context in which students describe and create visual displays of data. They can plot
955 data and draw from different statistical methods such as creating box plots and dot plots
956 to informally assess the degree of overlap of two populations, and students should be
957 able to describe the difference between the two centers in terms of the measure of
958 variability they use for the distributions.

959 ***Vignette 4***

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

969 Maria gave her students cubes and asked them to make 6 different towers of cubes that
970 represented the numbers 1, 6, 3, 2, 4 and 2. She asked them how they might construct
971 a physical proof to show the mean of the numbers. Some of the students were able to
972 calculate the answer; however, she kept pushing them to build a visual proof while
973 remaining open to multiple means of representation. This strategy, based on specific
974 UDL guidelines, allowed Maria to ensure scaffolds and supports would exist to help

975 highlight the patterns of language, and draw on background knowledge to express what
976 they know in ways that are authentic and meaningful. Óscar and his group members
977 came up with the idea of moving the cubes from tower to tower to show that they could
978 make six towers that were all the same height. They just needed to average out all of
979 the blocks. Óscar and his group excitedly explained to the class how they had made a
980 physical proof of finding the mean of the blocks. They shared the calculation with the
981 class and compared it to the method they used of moving the blocks. After her students
982 had discussed finding mean, Maria asked them to make a visual proof for the median
983 and the mode.

984 **Sampling to Understand a Population: Randomness, Bias, How** 985 **Many? (Grades Seven Through Eight)**

986 Prior to seventh grade, students' work with data has focused exclusively on using data
987 to understand, describe, and compare the particular collection of objects or situations
988 that were observed or measured. For example, to calculate the median highest
989 temperature on school days in September, students would record the highest
990 temperature on *each* school day.

991 Seventh grade includes the first introduction to *sampling*, the process of collecting data
992 from a subset of a population in an attempt to understand or describe the whole
993 population. This represents a big jump in sophistication from earlier work. Early
994 experiences with sampling should first describe the measured variables for the sample
995 (favorite lunch, number of minutes looking at screens, recorded for all students in the
996 sample for one week), followed by team and class discussions about whether the
997 description extends. For instance, if all students who come in to play basketball before
998 school are asked to track their screen usage for the week, the class should discuss
999 whether they expect the average of 862 minutes to be close to the average for everyone
1000 at school—and if not all teenagers in this age range—then perhaps close to the average
1001 for some smaller, definable group of students. Many similar discussions, with some
1002 obviously non-representative samples, help students understand the idea of a *random*
1003 *sample*.

1004 If researchers decide to gather data from 40 members of the population, then their
1005 collection of 40 members is *random* if it is chosen in such a way that every possible
1006 subset of size 40 has an equal chance of being selected. It is important for students to
1007 have multiple experiences selecting samples from known populations in ways that are
1008 random (for instance, drawing numbered ping-pong balls from an opaque bag or
1009 drawing student names on identical slips of paper from a hat) *and* in ways that are not
1010 random (for instance, asking survey questions only of the students who sit near you in
1011 class). The goal is an understanding that random sampling tends to produce samples
1012 that are *representative* of the population—that is, their distribution of the quantities
1013 under consideration are close to the distribution for the population as a whole
1014 (7.SP.A.1)—and a sense for the variability when using samples to make inferences and
1015 estimates for a population (7.SP.A.2).

1016 Non-random sampling (such as attempting to understand the school as a whole by
1017 collecting data only from one’s friends, or by asking about eating habits at the gym after
1018 school, produce *biased* conclusions, even when the bias in the sample selection might
1019 not be obviously linked to the quantity being measured in the measurement or
1020 observation. *Bias* does not here refer to temperament or outlook (prejudice), which is
1021 one meaning of the word; instead; it means a *systematic error*.

1022 Once teachers implement ways to ensure random sampling becomes a tool for student
1023 learning, the pool of questions empower their inquiry expands greatly: “I wonder how
1024 long on average it takes students from different grades to get from home to school?”
1025 “How do students who live in different areas spend time at the weekends?” “How much
1026 food is wasted in the lunchroom every month?” These kinds of questions could form a
1027 data exploration where students consider their sample, which variables can be defined
1028 and collected, and engage in the four-part exploration process.

1029 Sampling is introduced in the seventh-grade standards and does not appear again until
1030 high school, but much of the eighth-grade work with *bivariate* (two variable) data will
1031 make use of sampling, so it is important to continue activities that help understand
1032 *random sampling* through eighth grade as well. Students often believe that arbitrary
1033 sampling schemes (first 10 students I meet or every tenth student alphabetically) are

1034 random; they need to understand the difference between these schemes and choosing
1035 *by chance* so that every possible sample has an equal likelihood of being selected.

1036 **Vignette 5**

1037 Understanding the ways Rosa's seventh-grade students have responded to the
1038 probability activities offered through her instruction has influence the next steps in her
1039 planning. Overall, Rosa has not been satisfied with student understanding of random
1040 sampling. She decides to give students a more visual and physical experience of the
1041 concept. Her plan calls for six paper bags filled with differently colored cubes. The sum
1042 of cubes and the color distribution of the cubes in the bags reflect the following:

1043 Bag One, 15 total: 15 blue

1044 Bag Two, 12 total: 11 blue and 1 red

1045 Bag Three, 20 total: 15 blue, 4 yellow, 1 red

1046 Bag Four, 10 total: 5 red and 5 yellow

1047 Bag Five, 12 total: 5 blue, 4 red, 3 yellow

1048 Bag Six, 20 total: 8 blue, 8 red, 4 yellow

1049 Rosa explained the task: Students would determine the contents of each bag through
1050 sampling. She chose not to tell them how many times to sample but she did tell them to
1051 sample from the bags by selecting one cube at a time and then putting it back into the
1052 bag. Rosa also asked students to determine the chance of drawing a blue cube from
1053 each bag.

1054 Students engaged in the activity, brainstorming methods for collecting and recording
1055 their information. When each group of students felt satisfied with their determinations of
1056 the number of cubes and color distributions of the contents of each bag, she asked
1057 them to choose which bag belonged to which card showing the contents of each bag. In
1058 setting up the lesson, Rosa filled the bags differently and made sure to have two
1059 different bags where the probability of drawing a blue cube would be one and another

1060 would be zero. After the activity and class discussion, Rosa was happy to hear her
1061 students later, talking about situations where the probability was one or zero as well as
1062 everything in between. Her students recognized the number of times they sampled
1063 usually led to better predictions about the contents of the bags. They also realized
1064 sampling without replacement would have shown them the exact contents of the bag.
1065 The class engaged in a rich conversation about sampling with and without replacement,
1066 recognizing that it would be unproductive to draw all the cubes if there were a million.

1067 **Are They Related? Two Changing Quantities (Grade Eight)**

1068 Prior to grade seven, students work with a single collection of data measuring a single
1069 variable. In grade seven, they compare the same variable measured across two
1070 populations, either by actually measuring the whole populations or obtaining estimates
1071 for the distributions via sampling.

1072 In eighth grade, the focus is *bivariate data*: Two quantities or categorical variables
1073 measured or observed across a population, or across a sample drawn from a population
1074 (8.SP.A.1). This work has important connections with linear equations and modeling.

1075 The *scatter plot* as a visual representation of *quantitative* bivariate data is one of the
1076 most important ideas introduced here. A survey of students collecting both time and
1077 distance for traveling from home to school might reveal *clusters*, *outliers*, and any of
1078 various types of *association* (positive, negative, linear, non-linear). Students should
1079 describe such patterns in a scatter plot and interpret them in the context of the data
1080 (8.SP.A.1).

1081 Students can explore large, relevant datasets—such as earthquake data from
1082 California—and explore bivariate relationships between the location of earthquakes in
1083 the database and the magnitude of the earthquakes. They can plot the data using
1084 graphing tools and consider associations, data distributions, and relationships.

1085 If students vary the weight added to a simple cart and measure the distance it travels
1086 when released at the top of a ramp, then plot the results on a distance (vertical axis) vs
1087 added weight (horizontal axis), they will likely see a relationship. This association

1088 between the two variables can then be *modeled* by a line if the association appears
1089 roughly linear (line-shaped). In eighth grade, students choose a line to fit the data by
1090 visual approximation on the scatter plot, and compare and argue for whose line fits
1091 “best” (8.SP.A.2). They then interpret the meaning of the slope and intercept of their
1092 chosen model line, and use the line to make predictions for one variable when the other
1093 variable is specified (8.SP.A.3).

1094 Finally, eighth-grade students use two-way frequency tables as tools to see
1095 associations in bivariate *categorical* data (8.SP.A.4). For instance, they might survey
1096 their class members’ favorite color and favorite genre of books, then input the data into
1097 a spreadsheet, organize the data and calculate relative frequencies in rows to explore
1098 possible relationships between the two variables.

1099 **What Are the Chances? Probability as the Basis for Data-Based**
1100 **Claims**

1101 Randomly selecting from a population and measuring a characteristic (in which variation
1102 is expected across the population) is a *chance process*: It may result in different results
1103 and its outcomes follow some *distribution*.

1104 Probability expresses the chance of an outcome as a number between 0 and 1
1105 (7.SP.C.5). Probability is combined with statistics in the grade seven standards;
1106 statistics and probability are historically linked because statistical claims and estimates
1107 are based on the mathematical field of probability. Models that draw from data science
1108 and offer predictions of events, such as voting in elections, draw from probabilistic
1109 reasoning.

1110 Students sometimes struggle to see clear connections between probability and
1111 statistics, especially when their experiences focus on procedures and calculation rather
1112 than exploration, context, and interpretation. There is much work with probability that
1113 does not support statistical reasoning (e.g., calculating theoretical probabilities for the
1114 sum of two dice without using those theoretical probabilities to decide whether a given

1115 pair of dice are likely fair), and middle-school probability experiences should be carefully
1116 designed to support reasoning with interesting and meaningful data.

1117 In seventh grade, students gather data to estimate the probability of outcomes by
1118 observing their long-run relative frequency; that is, they compute *experimental*
1119 *probability*. Consider repeating the same experiment 150 times: draw a marble from a
1120 bag with marbles in it, record its color, then put the marble back in the bag. If we get a
1121 blue marble 32 times, our estimate for the probability of getting blue on any particular
1122 draw is $32/150$ (7.SP.C.6, 7.SP.C.7.B).

1123 Compare the marble experiment just described to another, placing the following marbles
1124 in a bag (all identical except for color): 16 blue marbles, 31 red marbles, 16 green
1125 marbles, and 12 white marbles (75 total marbles). If you blindly pull a marble from the
1126 bag, what is the probability that you will get a blue marble? If you perform this 150 times
1127 (putting the marble back each time), about how many times do you expect to get a blue
1128 marble? After calculating this expectation, students might construct an algorithm or
1129 pseudo-code to run the simulation 150 or 1500 or 15,000 times to compare with their
1130 theoretical expectations (CSS 6–8.AP.10).

1131 Note the difference between the questions in the previous two paragraphs: In the first,
1132 students use long-run relative frequency to estimate probability; in the second, students
1133 build a (*theoretical*) probability model and use it to estimate long-run frequency
1134 (7.SP.C.7). If a marble experiment is then performed and relative frequencies of
1135 outcomes do not seem close to predictions from the probability model, then students
1136 need to be able to discuss possible sources of discrepancy (7.SP.C.7): Perhaps the
1137 green marbles have a different texture and tend to be drawn more frequently than
1138 predicted. Maybe somebody changed the mix of marbles in the bag. Or perhaps not
1139 enough draws were performed to see the relative frequencies approach the probability
1140 model.

1141 Finally, seventh-grade students find probabilities of compound events (events which are
1142 made up of several simple events; for example, drawing two marbles from the bag of 75
1143 described above and getting one white and one blue marble) (7.SP.C.8).

1144 The specific calculations above are not central to the data science progression, but
1145 recognition that some events (repeat the draw five times, get all blue; or repeat the draw
1146 five times, obtain WBWWB in that order) are *much* less likely than others (repeat the
1147 draw five times, get three white and two blue) is key to understanding claims made from
1148 data.

1149 In fact, most statistical claims depend on a comparison of a (theoretical and
1150 hypothetical) probability model with observed data, as in 7.SP.C.7. To prepare middle-
1151 school students for future data science work, teachers should offer experiences that
1152 develop an awareness that more data tends to produce relative frequencies closer to
1153 actual probabilities.

1154 Invite students to explore rich datasets, such as the distribution of births in the US—and
1155 consider questions of probability that they can explore, like the chance that two people
1156 share the same birthday. This is a question that could be explored theoretically or
1157 experimentally.

1158 ***Vignette 6***

1159 Quincey started middle school without a lot of interest in math class. Quincey had
1160 always been interested in how the world works, and science and social studies were
1161 their favorite classes. Quincey had not had much experience with math class content
1162 connecting to their areas of interest.

1163 Quincey’s sixth-grade math teacher, Leonora, saw the value of tapping into student
1164 interest to ensure math content reflected their real-world experiences. Leonora knew
1165 that the data science standards in sixth grade would give Quincey an opportunity to use
1166 real data to understand that they could question the data and make connections
1167 between mathematics and life. One strategy Leonora decided to use was an activity to
1168 explore the “shape” of data: The context is hurricanes in the Atlantic Ocean using real
1169 data collected from five years of hurricanes spread over four decades. Quincey showed
1170 real interest and engaged in the lesson’s opening discussion of 2017 hurricane data
1171 displayed on a line plot. Quincey and the class were really interested in the number of
1172 hurricanes that were in category 0—tropical storms.

1173 Next, students worked in groups where they studied hurricane category data for the
1174 years 1977, 1987, 1997, and 2007. Each decade’s data was presented in different
1175 ways: bar graph, line plot, tables and sentences. Quincey enjoyed the analysis and was
1176 taken with the different ways of displaying data as well as the changes in the spread of
1177 data.

1178 Quincey asked important questions about the science of hurricanes. *How do they*
1179 *develop?* he wondered. *What makes them get larger? What is the difference between a*
1180 *category 3 storm and a category 5 storm?* At the close of the lesson, Leonora was
1181 convinced that students understood that different visual displays of data can make it
1182 easier to see the shape of data. The shape of the data on the displays helped students
1183 see how a situation might be changing over time. The class reflected that the changes
1184 were easier to see in line plots and histograms versus the data being shared in writing
1185 or in a table of values. Quincey decided to further investigate the number of category 4
1186 and 5 hurricanes over the past 100 years and how these storms become stronger, and
1187 they set out to gather more data and ask questions of the data. Others in the class
1188 decided to investigate why the number of category 4 and 5 storms are increasing.

1189 **High School**

1190 This outline for data science in high school organizes two sections of guidance:
1191 (1) experiences and expertise in data science common for all high school students, and
1192 (2) experiences and expertise for a high school pathway with a data science focus
1193 (expanding on the pathway outline in Chapter 8).

1194 Modern life depends on computer technology. Devices like laptops, phones, so-called
1195 “smart” appliances, medical records systems, exercise trackers, GPS-location
1196 recording, or payment methods, computers facilitate most transactions. Every
1197 interaction with a computer generates data about that interaction (which is collected and
1198 saved)—but, to most people, very little of this data is analyzed and interpreted.

1199 Even as computers have led to the collection of vast amounts of data, computational
1200 tools (including computer hardware capabilities and advances in algorithms) have
1201 dramatically altered the available methods for making use of and communicating

1202 interpretations of data. In fact, meaningful analysis of large or multivariate data sets is
1203 impossible without computer tools.

1204 For many questions about which students might wonder, existing data sources might
1205 provide the necessary information. Designing data collection to obtain exactly the
1206 desired data for answering a specific question (the classical statistical experiment
1207 approach, still the main approach in kindergarten through grade eight) is expanded to
1208 include techniques for analyzing multivariate data, critical questioning skills to
1209 interrogate pre-existing data's suitability for the investigation, and ways to access and
1210 acquire data through the internet. This understanding extraction uses two processes:
1211 (1) data description methods, both visual and numerical, to investigate conjectures and
1212 discover patterns; and (2) model-building to test conjectures, make predictions of future
1213 observations, and evaluate the predictive success. These huge, many-variable existing
1214 data sets are not collected in order to answer a particular question, do not typically
1215 represent random samples, and are often missing data or are otherwise “messy.”

1216 Data science should be understood as a broad term encompassing many tools relevant
1217 to learning from data. These include tools of traditional statistics classes, but also
1218 include computational and programming tools to address the massive size and
1219 complexity of many of today’s data sets, and disciplinary knowledge of the field
1220 generating the data. Thus, data science is an inherently interdisciplinary field that uses
1221 scientific and statistical methods and processes to derive understanding, insight, and
1222 predictive ability from (often unstructured) data (Dhar, 2013).

1223 **Data Science for Equity and Inclusion**

1224 Educators can offer social and emotional support to students by designing engaging
1225 lessons that allow students to connect in meaningful ways with content. Traditional
1226 mathematics lessons that have taught the subject as a set of procedures to follow have
1227 resulted in widespread disengagement as students see no relevance for their lives. This
1228 is particularly harmful for students of color and for girls—who receive additional harmful
1229 messages that mathematics is not for them. The data science field provides
1230 opportunities for equitable practice, with multiple opportunities for students to pursue

1231 answers to wonderings and to accept the reality that all students can excel in data
1232 science fields.

1233 Studies by Walton and colleagues (2015) show that many students, particularly girls
1234 and students of color, do not feel that they belong in certain disciplines. These feelings
1235 often due to a history of negative and off-putting messages (Chestnut et al., 2018).
1236 Other studies have shown that different topics and teaching approaches can lead to
1237 feelings of belonging or not belonging (Boaler, 2019; Boaler, Cordero, and Dieckmann,
1238 2019). Data science holds promise for teachers seeking to create climates of belonging
1239 for students, inviting them to investigate real data that is likely relevant to their lives.
1240 This meaningful engagement can create opportunities for students to develop self-
1241 confidence and self-efficacy. When teachers empower students to become data
1242 investigators who ask questions and respond to issues with data, they can take an
1243 active role in their development of self-motivation and goal-setting strategies. Important
1244 principles in the teaching of data science, that will offer the greatest chance for social,
1245 emotional, and academic development, include the following:

1246 ● *Mindset and Belonging Messages*

1247 Teachers should remind students that data science fields welcome all people.
1248 Informed by successful interventions in mindset and belonging strategies,
1249 teachers can remind students that struggle represents an important part of
1250 learning; all students struggle at times, and that successful students respond to
1251 times of difficulty using the strategies developed and practiced over time. Share
1252 with students examples of successful people inside data science, that highlight
1253 gender and racial diversity.

1254 ● *Use Real Data*

1255 Data science represents an opportunity for students to question real sets of data,
1256 developing social awareness, and investment in the solutions they discover.
1257 When working with secondary data sets (data obtained from others, rather than
1258 collected by students), teachers should choose meaningful content selected to
1259 create a connection with their learning and secure opportunities to hear the
1260 perspective of others, which will help them develop empathy. When teachers use
1261 local data sets they can also help students feel like they are important members

1262 of their community—as they explore questions and find answers to local
1263 problems that they can help with real data. Identifying problems and finding
1264 solutions will help students develop skills to make responsible decisions.

1265 Some teachers worry that they cannot provide culturally sustaining connections
1266 for their classes because they lack expertise in the cultures of all their students,
1267 but real data sets from different communities provide opportunities for students to
1268 bring their own knowledge and expertise to data rich problems. There should
1269 also be times when students are invited to collect data from their own community
1270 and build their own data sets. Students can pose questions that are important to
1271 them, including those with cultural meaning, collecting data from their own lives
1272 and communities. As Paris (2012) describes students will be fostering and
1273 sustaining “linguistic, literate, and cultural pluralism.” The act of collecting data
1274 provides an important learning opportunity for students to understand decisions
1275 that need to be made around the collection and organization of data as well as
1276 how to deal with uncertainty in their data. Students will be the ones with
1277 important expertise in these investigations.

1278 ● *Focus on Collaboration and Communication*

1279 Data science is a field in which people collaborate, connecting ideas to solve
1280 difficult problems with data. Meaningful collaborations typically reflect
1281 perspectives from diverse groups of students who come together to work
1282 effectively with different ideas being valued and developed. Creating
1283 opportunities for this kind of group work makes open problems accessible to
1284 students in an environment where differences thrive, where it is safe to share
1285 their ideas, and where students have the tools to work respectfully to reach
1286 solutions. Group work is usually much more effective when implemented in
1287 meaningful ways. For example, students may start their work in structured and
1288 unstructured conversations where each group member shares their thoughts.
1289 Collaborative classrooms founded in engaged listening and the capacity to
1290 articulate verbally as they build on each other’s ideas, are places where students
1291 feel valued and where they develop Important relationship skills of
1292 communication, social engagement and teamwork.

1293 **Data for All Students: Living in a World Overloaded with Information**

1294 With data serving as the bases of large-scale decisions and predictions, all California
1295 high school graduates need data acumen, as described in this chapter’s introduction:
1296 skills in interpreting and visualizing data, making and critiquing data-based arguments,
1297 and some facility with data software. It is crucial for students to develop the ability to
1298 identify types of questions that are subject to exploration through data; just as crucial is
1299 their understanding of some misuses of data and of one’s own online data footprint. As
1300 in earlier grades, teachers should provide opportunities for students to generate and
1301 investigate their own “I wonder” questions in given contexts. All statistics standards are
1302 identified as *modeling* standards, reflecting the origin of all work with data in authentic
1303 questions about the world.

1304 “I wonder” just initiates the learning, however. Students must ultimately formulate
1305 statistical investigative questions, pose data collection questions, interrogate existing
1306 data, pose analysis questions, analyze data, and formulate, interpret, and communicate
1307 findings. Thus, questioning is a central practice throughout the statistical problem-
1308 solving process. GAISE II summarizes this process:

1309 The statistical problem-solving process typically starts with a statistical
1310 investigative question, followed by a study designed to collect data that aligns
1311 with answering the question. Analysis of the data is also guided by questioning.
1312 Constant questioning and interrogation of the data throughout the statistical
1313 problem-solving process can lead to the posing of new statistical investigative
1314 questions.

1315 Often when considering secondary data, the data need to first be interrogated—
1316 how were measurements made, what type of data were selected, what is the
1317 meaning of the data, and what was the study design to collect the data. Once a
1318 better understanding of the data has been gained, then one can judge whether
1319 the data set is appropriate for exploring the original statistical investigative
1320 question or one can pose statistical investigative questions that can be explored
1321 with the secondary data set. (Bargagliotti et al., 2020)

1322 The process of *interrogating* secondary data (data collected by anyone other than the
1323 person doing the analysis) is crucial. Public datasets are not collected specifically to
1324 answer students' questions. So before using such datasets, students will need to
1325 evaluate the appropriateness for their purposes: How do the measures, methods, and
1326 scope of the dataset match the statistical investigative question(s) that interest us? For
1327 what purpose(s) were the data collected?

1328 Using data to answer authentic questions about the world is a powerful antidote to the
1329 famous student retort, "when will I ever need to know this?" Chapter 8 of this framework
1330 encourages the use of data science contexts as a way to frame many of students'
1331 explorations to develop content and practice standards in all domains. In this section,
1332 we propose an outline for data science understanding, considering the understandings
1333 high school students should develop.

1334 Students enter high school with significant, relevant experiences that high-school
1335 teachers can use to enhance their work. The CA CCSSM introduction to High School
1336 Probability and Statistics summarizes work in prior grades:

1337 Data are gathered, displayed, summarized, examined, and interpreted to
1338 discover patterns and deviations from patterns. Quantitative data can be
1339 described in terms of key characteristics: measures of shape, center, and spread
1340 [variability]. The shape of a data distribution might be described as symmetric,
1341 skewed, flat, or bell shaped, and it might be summarized by a statistic measuring
1342 center (such as mean or median) and a statistic measuring spread (such as
1343 standard deviation or interquartile range). Different distributions can be compared
1344 numerically using these statistics or compared visually using plots. Knowledge of
1345 center and spread are not enough to describe a distribution. Which statistics to
1346 compare, which plots to use, and what the results of a comparison might mean,
1347 depend on the question to be investigated and the real-life actions to be taken.

1348 The big ideas of data science for all students in high school are identified in the
1349 statistics cluster headings in the standards, with an additional big idea discussed here,
1350 in response to the changing approaches to data described above. The first two are

1351 described in more detail, with additional examples, in the Draft High School Progression
1352 on Statistics and Probability (Daro, 2019).

- 1353 • Interpreting Categorical and Quantitative Data
- 1354 • Making Inferences and Justifying Conclusions
- 1355 • From Statistics to Data Science

1356 ***Understanding the Role of Data in the World***

1357 Students should develop an understanding of what qualifies as data and the many types
1358 of data that exist. They also learn how data is generated and collected, and the
1359 existence of extremely large amounts of data created by our digital lives. Students
1360 consider their own privacy and data footprint. Throughout data-based investigations,
1361 students discuss the ethics and consequences of collecting and using big data, and the
1362 ways data is collected, including the bias that may be present in the data collection or
1363 selection process. Students evaluate and critique data-based claims and arguments; in
1364 particular, they distinguish correlation and causation. Students understand that all data
1365 and data-based arguments have several sources of bias and are able to identify them.
1366 They understand the importance of communicating with data and making data-based
1367 arguments.

1368 ***Interpreting Categorical and Quantitative Data***

1369 High-school students continue their work with the representations of data introduced in
1370 kindergarten through grade eight. However, the major work interpreting data in high
1371 school relies on the use of functions as models of associations in two-variable
1372 quantitative data.

1373 Building on kindergarten through grade eight experiences, high-school students
1374 continue to visualize and represent *single-variable* data with dot plots, histograms, and
1375 box plots; use measures of center and spread to describe such distributions; and
1376 compare distributions from different populations or samples using these representations
1377 and statistics (S-ID.1–3).

1378 When available data includes two (or more) measurements or characteristics for each
1379 observation, students' tools for representing and interpreting relationships between
1380 pairs of variables depend on the nature of each variable.

- 1381 ● If both are categorical, two-way frequency tables give an important summary that
1382 reveals relationships when interpreted in the context of the data.
- 1383 ● If one is categorical and the other quantitative, students can treat each category
1384 as a separate population and compare the quantitative data for the different
1385 categories as in the single-variable paragraph above.
- 1386 ● If both variables are quantitative, the scatter plot is the standard visual
1387 representation.

1388 Building on earlier experiences with multivariate, investigative questions and data,
1389 students should also move to examining large (many variable) data sets, make multi-
1390 variable conjectures (or pose multi-variable statistical questions), and create data
1391 visualizations that could potentially refute those conjectures. Many technologies make it
1392 easy to display three-variable relationships and students should practice examining “off-
1393 the-shelf” visualizations with more variables.

1394 ***Modeling association in numerical data using functions***

1395 Once a scatter plot is created, an association between the two variables may become
1396 visually identifiable. Fitting a function to the data is the creation of a mathematical model
1397 for the association. This begins in eighth grade with visual fitting of a linear model. While
1398 the type of function that is used most frequently is a line (a linear function), students
1399 also need experiences with plotting associations that are clearly non-linear, as well as
1400 experiment with fitting other types of functions (quadratic, exponential).

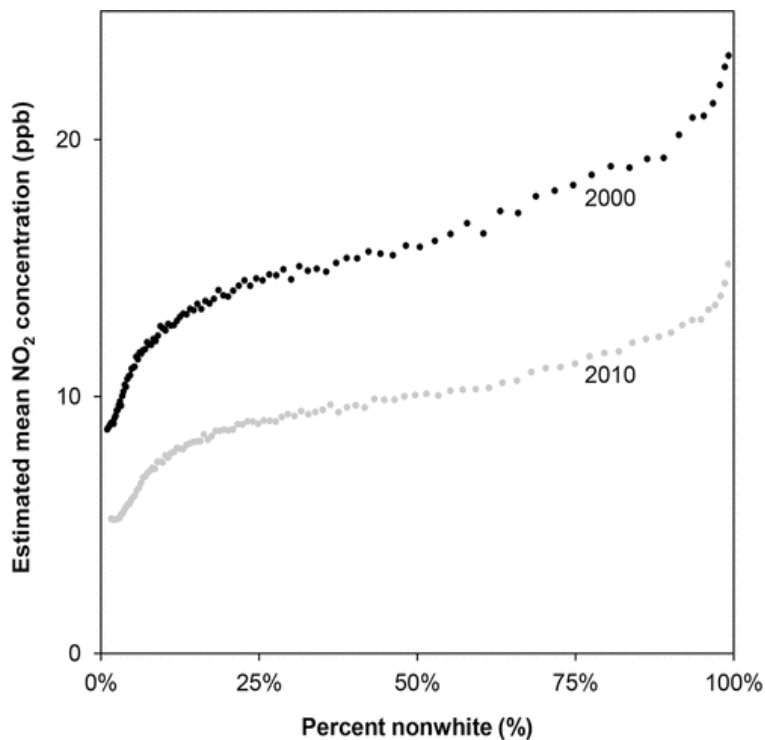
1401 Any standard data software (including spreadsheets, Desmos, Geogebra, CODAP) will
1402 fit lines, quadratic functions, and exponential functions to given data. The specific
1403 standard technique for identifying a line (or quadratic or exponential function) of best fit
1404 (least-squares regression) is *not* an expectation; but students should have experiences
1405 fitting lines and some other functions visually (by adjusting parameters on appropriate

1406 function types in graphing software) and using appropriate software tools which perform
1407 the regression behind the scenes.

1408 Most importantly, functions that model associations must be used to solve problems
1409 (e.g., prediction of the value of one variable given other[s]) (S-ID.6.a), and must be
1410 interpreted in the context of the data (S-ID.7).

1411 Important examples of modeling association in numerical data arise in many contexts in
1412 science, history, physical education, and social studies. Appendix C of the *History–*
1413 *Social Science Framework* (CDE, 2017) outlines expectations that students develop
1414 *Chronological and Spatial Thinking*, including analyzing change over time; both *time*
1415 and *space* provide opportunities for finding meaningful quantities that vary together. For
1416 example, students might wonder whether pollution exposure is related to wealth, and
1417 either find zip-code level data on both air pollution and income, or find existing research
1418 like the graph here, and work to understand and explain it. (The graph in Figure 5.11
1419 represents both change over time and change in space.)

1420 Figure 5.11



1421

1422 Source: Clark, 2017.

1423 ***Making Inferences and Justifying Conclusions***

1424 Making conclusions and generalizations about a population from a sample (S-IC.1) is
1425 the goal of *inferential* statistics, as opposed to *descriptive* statistics. Students work with
1426 random samples beginning in seventh grade, their first experience trying to understand
1427 a population without gathering data about all of its members. This strand of high-school
1428 data work is the foundation for most meaningful use of statistics for making decisions.

1429 Students must decide whether a result observed through data is consistent with a
1430 mathematical model of the process that generates the data (S-IC.2). For instance, if a
1431 student estimates that 30 percent of the students at the school grow food at home, the
1432 estimates offer a mathematical model that gives them an idea of what proportion to
1433 expect in a sample. If they then survey five, randomly-chosen students, and all say they
1434 grow food at home, then the student should be able to reason as follows: *If* 30 percent
1435 of students grow food at home, then the chances of five randomly-chosen students all
1436 being among those 30 percent of students is $(.3)^5 = .00243 = .342$ percent, or less than
1437 half of one percent. Thus, the student might doubt—that is, they might *reject*—the 30-
1438 percent hypothesis. Students should have many experiences of simple situations like
1439 this to understand how decisions based on data rely on probability, and are not
1440 *guaranteed* to produce correct answers to the original question.

1441 Students should work with data originating in four different methods of data production,
1442 including at least some student-generated questions and student-gathered data. These
1443 methods are (1) *census* data; that is, data that contain measurements on every member
1444 of the target population (such as the database of crimes occurring in a given city in a
1445 given time frame, or rain gauge data for a given location, which captures all precipitation
1446 at that location—census data is first encountered in early elementary grades); (2)
1447 surveys administered to random samples (to estimate population values, or *parameters*,
1448 for the surveyed quantities); (3) randomized experiments (to compare treatments and
1449 demonstrate cause); and (4) observational studies (to study characteristics or quantities
1450 when random selection or assignment is not possible) (S-IC.3). The Draft High School

1451 Progression on Statistics and Probability (Daro, 2019) contains detailed examples
1452 describing the expectations in the standards.

1453 Teaching with surveys and experiments must include a link between the random
1454 selection or assignment and the ability to reason probabilistically to make claims. With a
1455 survey, the random sampling allows generalizing to a population. With an experiment,
1456 the random assignment allows causal conclusions but not generalization to a broader
1457 population—unless the sample in the experiment was randomly selected from some
1458 larger population. For example, medical studies (experiments) must use willing
1459 volunteers and thus are not random samples of the overall population; this makes it
1460 much harder to draw broadly-applicable conclusions.

1461 In a college statistics course or the data science course outlined below, students will
1462 learn ways to quantify the comparisons between gathered data and hypothesized
1463 population parameters (margins of error and p -values). Making sense of these,
1464 however, requires an understanding of the role of randomization in the data gathering.

1465 When using a sample mean or proportion to estimate a population mean or proportion,
1466 students use simulation models to estimate a margin of error, instead of formulaic
1467 calculations. Briefly, the process is to use data simulation software to draw many
1468 random samples from a hypothetical population, and to see how often a result is
1469 obtained that is as extreme as the sample mean or proportion. Doing this process for
1470 hypothetical populations with many different mean or proportion parameters helps
1471 students see that there is a range of population parameters that often (more than 5
1472 percent of the time) produce simulated sample means or proportions that are as
1473 extreme as (or more extreme than) the actual sample mean or proportion. This range of
1474 population parameters is the (simulation-based) confidence interval, given as (sample
1475 mean or proportion \pm margin of error). Note the probabilistic argument here: *If* the
1476 population mean or proportion were outside of the confidence interval, *then* sample
1477 means or proportions as extreme as we obtained in our random sample would be rare.
1478 So, we expect that the true population mean or proportion is within the confidence
1479 interval (*but cannot be certain* that it is!).

1480 A similar process is used to evaluate confidence in a randomized experiment, in which
1481 subjects are randomly assigned to two or more treatment groups. (Treatment could
1482 mean medical treatment, or assignment of different tasks, or being shown different
1483 motivational videos, etc.) Some quantity is then measured for each subject, and the
1484 investigator then has to decide from the results whether a treatment, say treatment A,
1485 produced any effect on the measured quantity. Simply having a different mean for each
1486 treatment group is not enough, as we expect variation in the measurement and thus
1487 between groups. In this case, all of the treatment groups are pooled into a population
1488 and then re-sampled (randomly) many times, to see how often the re-sampled mean or
1489 proportion is at least as extreme as the actual treatment A group difference. If such
1490 differences are rare, the experiment is taken as evidence that treatment A caused a
1491 change in the measured quantity.

1492 ***From Statistics to Data Science***

1493 For questions about community, society, or natural systems beyond students'
1494 immediate experience—and thus beyond their ability to gather data directly—existing
1495 data can often be identified from online sources. Since students frequently encounter
1496 claims made from such large data sets, it is crucial that all students have experiences in
1497 which they explore the ways in which such claims are made. A major difference
1498 between the classical statistical approach in kindergarten through grade eight and the
1499 “big data” of the growing field of data science is the richness and complexity of available
1500 data sets, even more so than their sheer size.

1501 Many sophisticated approaches to working with rich, complex data sets are left to a
1502 data-science course in the data science pathway; but *all* high-school students should
1503 exercise and refine their understanding of data exploration, causal inference, and
1504 statistical reasoning using large, real world data sets. As students work with these data
1505 sets, they can draw upon the data science understandings they have developed in their
1506 kindergarten through grade eight mathematics lessons. Instruction should emphasize
1507 opportunities for questioning and interpreting, rather than technical procedures.

1508 Data exploration begins with a search for available data about a context of interest. The
1509 data set is then examined for hidden patterns and associations (usually via visual
1510 representations). Any patterns or associations discovered can lead to new hypotheses
1511 or questions to investigate further. Students began this process in eighth grade, and
1512 continue in high school with experiences in which they examine data sets with multiple
1513 variables measured for each member of the sample. They plot pairs of variables to
1514 decide which ones might show associations. Important discussions for students to
1515 engage in when working with existing data sets include

- 1516 ● Prior to exploring: Do we *expect* any of these variables to be associated? Why?
- 1517 ● Might the association we see just be a result of the way in which the data was
1518 collected, rather than truly reflective of the population? What features of the data
1519 collection might make conclusions suspect, and what features might give
1520 confidence? Note that a large sample size is not enough to have confidence in
1521 conclusions.
- 1522 ● Can we think of possible explanations for the association(s) we see? Can we
1523 think of ways we could decide which explanations might be accurate?

1524 After data exploration identifies some association(s) of interest, the stage of model
1525 building follows. Technical methods are reserved for the specialized data science
1526 course below, but all students need to explore questions such as:

- 1527 ● Could we use some variables to predict others? This is a hugely important use of
1528 data, since some factors are easier to measure or observe than others. In
1529 medicine and many other fields, this often takes the form of trying to predict
1530 future outcomes using presently-available information.
- 1531 ● If we could only know measure one variable to try to predict a variable of interest,
1532 which one would we pick? Why? What if we could measure two? Which second
1533 variable gives us the *newest* information for prediction?

1534 Most importantly, high school students (like kindergarten through grade eight students)
1535 must experience data science as a set of tools for making sense of their worlds in ways
1536 that matter to them.

1537 ***Vignette 7: Data on environmental threats to health***

1538 In this example (Lieberman and Brown, 2020), students compared CalEnviroScreen
1539 data related to four environmental topics that are known to affect human health:
1540 (1) water (using data on groundwater threats, impaired water, and drinking water);
1541 (2) toxic chemicals (using data on pesticides, cleanups, and toxic releases); (3) air
1542 pollution (using data on the ozone, particulate matter [PM 2.5], diesel, and traffic); and
1543 (4) waste (using data on hazardous waste and solid waste). They compared these
1544 results against environmental impacts using data for asthma, low birth weight, and
1545 cardiovascular disease (California Health Standards 1.13.P; 2.3.P; 3.3.P, 3.4.P).

1546 In preparation for their analysis and reporting, the teacher reviewed California’s
1547 Environmental Principles and Concepts (EP&Cs) with students by asking them to
1548 identify one that is directly related to their environmental health problem. Based on their
1549 data analysis, students identified environmental health and environmental justice
1550 concerns related to water pollution in the local community and observed that they
1551 differentially affected various parts of the community. Their conclusion was that the key
1552 factors in the differential environmental health impacts were related to “Environmental
1553 Principle V: Decisions affecting resources and natural systems are based on a wide
1554 range of considerations and decision-making processes.”

1555 Depending on the focus of their individual environmental health study, students are
1556 encouraged to choose two variables to analyze such as the impact of water quality on
1557 low birth weight, or the impact of toxic chemicals on the incidence of cardiovascular
1558 disease. or the impact of air quality on asthma. After collecting the data for these
1559 variables, the students will use technology to create a scatter plot of the data, fit a
1560 function to the data, and create a symbolic representation for the function. Students will
1561 be able to connect the parameters of the symbolic representation to the context of the
1562 data. After a class discussion about the comparison of different variables, students
1563 should be guided to focus on the combinations of variables that make the most sense
1564 for their investigations.

1565 Following their research and analysis, student teams reported back to the class,
1566 summarizing their quantitative comparisons using charts to depict the results about
1567 water, toxic chemicals, air pollution, and waste (Health 4.1.P). In their presentation, they
1568 used graphs to compare the environmental effects they discovered with the
1569 environmental health impacts they analyzed (ELA SL.9-12.1; ELA SL.9-12.2; ELA SL.9-
1570 12.4; ELA SL.9-12.5; Health 5.3.P).

1571 Several of the teams mentioned that they observed a pattern that relates to the socio-
1572 economic conditions in the communities they compared. Some of the students
1573 mentioned that they see these issues as directly related to EP&C V, because the places
1574 where waste, toxic chemicals, and manufacturing facilities are located depend on a
1575 variety of political, economic, and social factors. The teacher explained that differential
1576 environmental health impacts on communities with varied socio-economic conditions is
1577 a major health topic identified as “environmental justice,” a term that came into use in
1578 the 1980s when residents of an African-American community in North Carolina
1579 protested the siting of a landfill to store soil contaminated with polychlorinated-biphenyls
1580 (PCBs). These residents knew the health hazards associated with this toxin and
1581 responded by demanding that their health and well-being be protected by the
1582 government. The landfill proposal went forward, but the protests spurred the federal
1583 government to study the issue. The findings show that many of the nation’s landfill sites
1584 are located in communities of color. The environmental justice movement has grown to
1585 focus on a more equitable distribution of environmental benefits and burdens. Since
1586 many of the students expressed a strong interest in this topic, they requested a guest
1587 speaker from a community-based health organization to provide additional information
1588 and answer students’ questions about environmental justice (Health 8.1.P.; 8.2.P).

1589 **Advanced high school data science**

1590 The traditional sequence of high school courses—Algebra, Geometry, Algebra 2—was
1591 standardized in the United States following the “Committee of Ten” reports in the 1890s.
1592 The course sequence—which was primarily designed to give students a foundation for
1593 calculus—has seen little change since the Space Race in the 1960s. With the rapid

1594 expansion of information available to all in the form of data, far more students pursue
1595 statistics classes than calculus, and may be better served by a data science course as a
1596 culminating high school mathematical science experience. In addition to the importance
1597 of the data science content—to twenty-first century jobs and to a wide range of college
1598 majors—many students are more engaged by open-ended explorations of important
1599 data sets, drawing upon important mathematical principles and tools, than by many
1600 traditional courses organized around mathematical techniques. This framework provides
1601 design principles and content outcomes for such a course.

1602 California high schools offer upper-level data science courses in two ways. In the first
1603 model, students have a common experience in grades nine and ten, with pathways
1604 branching at grade eleven. Some districts have designed and are offering eleventh
1605 grade data science courses as an option for this third year of high school mathematics;
1606 in this case, the ninth and tenth grade courses need to be designed to include the
1607 important high school algebra and geometry, or integrated mathematics, or MIC 1 and 2
1608 standards. The second model is a data science course as a fourth-year course,
1609 following a coherent three-year pathway that builds the “for all students” data science
1610 understanding outlined in the previous section. In either case data science can be open
1611 to all students and not require any advancement in middle school – thus opening STEM
1612 pathways for many more students and combatting the persistent inequities in STEM.
1613 The design principles and content outcomes below are flexible enough to be
1614 implemented in either model, with appropriate adjustments for students’ prior
1615 experiences. If students are intending to pursue STEM majors in college (including data
1616 science), it is advised that they take courses that, at minimum, allow them to enter
1617 college having completed the prerequisites for calculus. Chapter 8 and Appendix A give
1618 further detail on the different pathways offering options to schools and students.

1619 ***Design Principles***

1620 These principles provide guidelines for design of curricular materials and classroom
1621 instruction for a data science course, in order to support a coherent and engaging
1622 experience for students. These principles should be used by developers to build
1623 curricular materials that are true to the vision of the course, as well as by educators

1624 reviewing materials and developing a repertoire of pedagogical strategies for use in
1625 teaching the course. Many students and teachers already engage in these behaviors; in
1626 these cases, these design principles will be seen as reinforcing and supportive. The
1627 spirit of this framework recognizes that, at some levels, everyone is a learner, and
1628 everyone is growing an understanding of mathematics, each other, and the world we
1629 share.

1630 Design Principle: Active Learning

1631 The course provides regular opportunities for students to actively engage in data
1632 explorations using a variety of different instructional strategies (e.g., hands-on and
1633 technology-based activities, small group collaborative work, facilitated student
1634 discourse, interactive lectures).

1635 *Students Will*

- 1636 • Be active and engaged participants in discussion, in working on data
1637 explorations with classmates, and in making decisions about the direction
1638 of instruction based on their work.
- 1639 • Actively support one another's learning.
- 1640 • Discuss results of their explorations with the instructor and/or classmates
1641 in class.

1642 *Teachers Will*

- 1643 • Provide low-floor, high-ceiling activities and explorations that all students
1644 can access and that extend to high levels. Such activities should provide
1645 meaningful opportunities for exploration and co-creation of mathematical
1646 understanding and data literacy.
- 1647 • Provide interesting and sometimes local data sets and invite students to
1648 ask questions of the data. Encourage different students to pose and
1649 investigate different questions, and to come together to discuss findings.
- 1650 • Facilitate students' active learning of data science through a variety of
1651 instructional strategies, including inquiry, problem solving, critical thinking,
1652 and reflection.

- 1653 • Create a safe, student-driven classroom environment in which all students
1654 feel a sense of belonging to the class and the discipline, are encouraged
1655 to take risks and embrace mistakes, and are able to make decisions about
1656 the direction for instruction through the results of their exploration of data
1657 science. Students’ ideas are at the center of the conversation.

1658 Design Principle: Growth Mindset

1659 Courses support students in developing the tenacity, persistence, and perseverance
1660 necessary for learning data science, for using mathematics and statistics to tackle
1661 authentic problems, and for being successful in post-high school endeavors.

1662 *Students Will*

- 1663 • Make sense of data explorations by drawing on and making connections
1664 with their prior understanding and ideas.
- 1665 • Persevere in solving problems and realize that it is acceptable to say, “I
1666 don’t know what to do next,” but that it is not acceptable to give up
- 1667 • Seek help from different sources to move forward in their investigations.
- 1668 • Compassionately help one another by sharing strategies and solution
1669 paths rather than simply giving answers.
- 1670 • Reflect on mistakes and misconceptions to improve their mathematical
1671 understanding and data literacy.
- 1672 • Understand that struggle is valuable for brain growth and times of struggle
1673 should be valued.
- 1674 • Develop/strengthen a growth mindset to continue to apply in mathematics,
1675 data science, and other areas of their post-high-school life.

1676 *Teachers Will*

- 1677 • Provide information about and model the importance of having a growth
1678 mindset.
- 1679 • Value mistakes and times of struggle.
- 1680 • Facilitate discussions on the value of mistakes, misconceptions, and
1681 struggles.

- 1682 • Give students time to struggle with tasks and ask questions that scaffold
1683 students' thinking without stepping in to do the work for them.
- 1684 • Praise students for their efforts in making sense of mathematical ideas
1685 and for their perseverance in reasoning through problems and in
1686 overcoming setbacks and challenges in the course.
- 1687 • Provide students with low-stakes opportunities to fail and learn from
1688 failure.
- 1689 • Provide regular opportunities for students to self-monitor, evaluate, and
1690 reflect on their learning, both individually and with their peers.

1691 Design Principle: Problem Solving

1692 Courses provide opportunities for students to make sense of problems and persist in
1693 solving them.

1694 *Students Will*

- 1695 • Apply intuition, life experience, and previously learned strategies to solve
1696 unfamiliar problems.
- 1697 • Explore and use multiple solution methods.
- 1698 • Share and discuss different solution pathways and methods.
- 1699 • Be willing to make and learn from mistakes in the problem-solving
1700 process.
- 1701 • Use tools and representations, as needed, to support their thinking and
1702 problem solving.
- 1703 • Develop and justify their own strategies to approach new problems.

1704 *Teachers Will*

- 1705 • Present tasks that require students to find or develop a solution method.
- 1706 • Provide data sets that allow for multiple strategies and solution methods,
1707 including transfer of previously developed skills and strategies to new
1708 contexts.
- 1709 • Provide opportunities to share and discuss different solution methods.
- 1710 • Model the problem-solving process using various strategies.

- 1711 • Encourage and support students to explore and use a variety of
1712 approaches and strategies to make sense of and solve problems.

1713 Design Principle: Authenticity

1714 Courses present data science as a subject and learning that allows us to model and
1715 solve problems that arise in the community.

1716 *Students Will*

- 1717 • Recognize specific ways in which mathematics and data are used in
1718 everyday decision making.
- 1719 • Recognize problems that arise in the real world that can be solved with
1720 data science.
- 1721 • Contribute meaningful questions that can be answered using data
1722 science.
- 1723 • Experience the decision making involved in collecting, cleaning, analyzing,
1724 and visualizing data.

1725 *Teachers Will*

- 1726 • Provide opportunities to ask questions of data sets that are relevant to
1727 students, both in class and on assessments.
- 1728 • Provide opportunities for students to pose questions that can be answered
1729 using data science methods and tools, and answer them.
- 1730 • Provide students with real data to explore and work with, including doing
1731 some of the data cleaning that is often required.

1732 Design Principle: Context and Interdisciplinary Connections

1733 Courses present data science in context and connects data science to various
1734 disciplines and everyday experiences.

1735 *Students Will*

- 1736 • Contribute personal experiences, where appropriate, that connect to
1737 classroom experiences.

- 1738
- Actively seek connections between classroom experiences and the world
- 1739 outside of class.
- 1740
- Examine the ways that data is collected in their day-to-day lives, and
- 1741 consider the ethics and consequences of collecting and using data to
- 1742 make decisions.

1743 *Teachers Will*

- 1744
- Provide opportunities for students to share their personal backgrounds
- 1745 and interests, including cultural values, and help make the connection
- 1746 between what is important in students' lives and future aspirations, and
- 1747 what they are learning in data science.
- 1748
- Provide real and interesting data sets, including some that are local to
- 1749 students.
- 1750
- Invite students into data explorations that illustrate authentic applications.
- 1751
- Provide data explorations that include applications from a variety of
- 1752 academic disciplines, programs of study, and careers, and which are
- 1753 culturally sustaining.

1754 Design Principle: Communication.

1755 The course develops students' ability to communicate their data explorations and

1756 findings in varied ways including with words, data visualizations and numbers.

1757 *Students Will*

- 1758
- Present and explain ideas, reasoning, and representations to one another
- 1759 in pair, small-group, and whole-class discourse using discipline-specific
- 1760 terminology, language constructs, and symbols.
- 1761
- Seek to understand the approaches used by peers by asking clarifying
- 1762 questions, trying out others' strategies, and describing the approaches
- 1763 used by others.
- 1764
- Listen carefully to and critique the reasoning of peers using data to
- 1765 support or counterexamples to refute arguments.
- 1766
- Develop the skills to justify mathematical reasoning with clarity and
- 1767 precision.

- 1768
- Practice constructing data-based arguments with specific audiences in
- 1769
- 1770
- Consider matters of accessibility in designing and executing their
- 1771
- 1772
- Consider the pros and cons of various types of data visualizations and
- 1773
- how they fit the communicative situation.

1774 *Teachers Will*

- 1775
- Introduce concepts in a way that connects students' experiences to course
- 1776
- content and that bridges from informal contextual descriptions to formal
- 1777
- definitions.
- 1778
- Clarify the use of data science terminology and symbols, especially those
- 1779
- used in different contexts or different disciplines.
- 1780
- Engage students in purposeful sharing of ideas in data science,
- 1781
- reasoning, and approaches using varied representations.
- 1782
- Support students in developing active listening skills and in asking
- 1783
- clarifying questions to their peers in a respectful manner that deepen
- 1784
- understanding.
- 1785
- Facilitate discourse by positioning students as authors of ideas who
- 1786
- explain and defend their approaches.
- 1787
- Provide regular opportunities for students to communicate about data
- 1788
- science with a variety of data visualizations.
- 1789
- Scaffold instruction to support students in developing the required reading
- 1790
- and writing skills.

1791 Design Principle: Technology

1792 Courses introduce students to current data science technologies, including

1793 programming, and prepare them to learn and use new ones.

1794 *Students Will*

- 1795
- Use technology to visualize and understand important data science
- 1796
- concepts and as a tool in problem solving.

- 1797
- Understand the necessity of digital tools in cleaning and analyzing large data sets and are able to select appropriate tools for different situations.
- 1798
- 1799
- Develop experience in learning new tools which allows them to try out emerging data science tools in the future.
- 1800
- 1801
- Understand that the use of tools or technology does not replace the need for an understanding of reasonableness of results or how the results apply to a given context.
- 1802
- 1803

1804 *Teachers Will*

- Introduce students to various digital data science tools and support them in understanding the best uses for each.
- 1805
- 1806
- Facilitate student learning of technological platforms through exploration, as this will aid in transferring the knowledge to future platforms.
- 1807
- 1808
- Not be experts in the use of every platform but willing to experiment along with students' questions and model good practices for seeking answers to such questions
- 1809
- 1810
- 1811

1812 Design Principle: Assessment

1813 Courses use project-based assessments to evaluate student progress.

1814 *Students Will*

- Assemble a collection of their work which includes both their mathematical work and reflections on their learning process and their evolving understanding of the field of data science.
- 1815
- 1816
- 1817
- At the end of the course, have a portfolio of data science work that showcases their knowledge of data science as well as their software skills. This portfolio might be shared with a potential employer or educational institution.
- 1818
- 1819
- 1820
- 1821

1822 *Teachers Will*

- Provide students with projects through which they are exposed to new content and demonstrate their ability to use this new content to solve problems. These will include products that demonstrate student learning both for the teacher, and to be included in the students' portfolios.
- 1823
- 1824
- 1825
- 1826

- 1827
- Evaluate student progress throughout the course by considering students' evolving portfolios as well as their reflections on their learning.
- 1828
- 1829
- In the final project of the course, allow students freedom to decide the topic and methods of their data exploration, so that they can bring together the various skills they will have developed over the course, and allow the teacher to assess their progress.
- 1830
- 1831
- 1832
-
-

1833 **Content Learning Outcomes**

1834 This section presents the mathematical content outcomes expected from a high school
1835 data science course. These will be motivated by realistic examples and projects which
1836 will help students develop their basic data science skills as well as a larger
1837 understanding of their contexts and of the importance of data in their lives.

1838 ***Understanding the Role of Data in The World***

1839 Students demonstrate an understanding of what qualifies as data and the different types
1840 of data that exist. They also understand how data is generated and collected, and the
1841 existence of extremely large amounts of data created by their digital lives. Students
1842 consider their own privacy and data footprint. Throughout the course, students discuss
1843 the ethics and consequences of collecting and using big data, and the ways data is
1844 collected, including the bias that may be present in the data collection or selection
1845 process. Students evaluate and critique data-based claims and arguments, in particular,
1846 they distinguish correlation and causation. Students understand that all data and data-
1847 based arguments have several sources of bias and are able to identify them. They
1848 understand the importance of communicating with data and making data-based
1849 arguments. They use multiple different types of data visuals both for analysis and in
1850 order to share their thinking with others. The standards listed below come from
1851 CACSSM Domains: Statistics (S), Functions (F), Number (N) and Vector and Matrices
1852 (VM). They also draw from the California Environmental Principles and Concepts, and
1853 from Computer Science (CS) standards.

- 1854 • Represent data represented by real numbers using dot, box and histograms (S-
1855 ID.1)
- 1856 • Summarize categorical data for two categories in two-way frequency tables (S-
1857 ID.5)
- 1858 • Interpret relative frequencies in the context of the data (S-ID.5)
- 1859 • Recognize possible associations and trends in the data (S-ID.5)
- 1860 • Distinguish between correlation and causation (S-ID.9)
- 1861 • Evaluate the purpose of and differences between sample surveys, experiments
1862 and observational studies and how randomization effects each (S-IC.3)

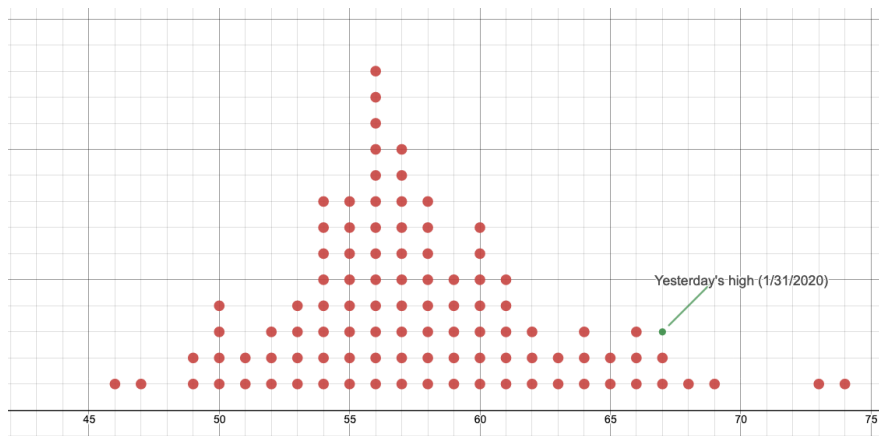
1863 ***Asking Statistical Investigative Questions***

1864 Students are able to identify the types of questions that are subject to exploration
1865 through data as well as formulate their own. They are able to perform exploratory data
1866 analyses to draw preliminary conclusions to explore further. They can do this in a
1867 variety of platforms. Students can look at the data available and identify questions that it
1868 can answer as well as determine what data might be collected in order to answer a
1869 question. Students consider how they might use some of the data they have access to,
1870 in order to predict other variables for which it might be harder to collect data directly.

1871 ***Unraveling the Story That Data Is Telling***

- 1872 • When working with numerical data students can describe a distribution using its
1873 shape, center, and spread (S-ID.3). They are able to make predictions based on
1874 these characteristics, as well as compare distributions to one another. Students
1875 are also able to compare two numerical variables to each other using scatter
1876 plots and can use their understanding of functions (linear, polynomial,
1877 exponential) to fit their data to a curve (using appropriate technological tools) and
1878 use this model to make predictions (S-ID.6). Students are also able to work with
1879 categorical variables in frequency tables as well as use numerical and
1880 categorical variables together in order to answer questions about the data(S-
1881 ID.6). Analyze the shape of data distributions and compare data distribution

1882 using measures of center (mean, median) and spread (interquartile range (IQR),
 1883 standard deviation) of different data sets (S-ID.1,2,3).
 1884 • Interpret differences in shape, center and spread including the effects of outliers
 1885 (S-ID.3)
 1886 • Use mean and standard deviation to fit to a normal distribution and to estimate
 1887 population percentages Know that this procedure is not appropriate for all data
 1888 sets (S-ID.4, SMP.3,5). For example, for data sets that appear to be bell-shaped,
 1889 they use the mean and standard deviation to specify an approximating normal
 1890 distribution and to approximate population percentages in specified ranges (S-
 1891 ID.4). Students might, for example, obtain high temperatures on a specific date
 1892 over the past 100 years from a nearby weather station, create a dot plot, visually
 1893 check for a bell-shaped distribution, and use an approximating normal distribution
 1894 to make a case for whether or not the temperature was consistent with historical
 1895 trends (CA Environmental Principles and Concepts II.a., II.d., V.a., and V.b.).



1896 • Use technological tools (calculators, spreadsheets, and programs) to estimate
 1897 areas under the normal curve (SMP.5, S-ID.4)
 1898 • Represent two variable data on a scatter plot and describe how the variables are
 1899 related (S-ID.6)
 1900 • Fit a linear function on scatter plots where the data suggests a linear fit (S-
 1901 ID.6,7,8)
 1902 • Fit a function to the data to solve problems in context of the data (S-ID.6a)
 1903 • Determine the fit of a function by plotting and analyzing residuals ((S-ID.6b)
 1904

- 1905 • In a linear model interpret slope as a rate of change and the intercept as the
- 1906 constant term in the context of the data (S-ID.7)
- 1907 • Use technology to compute and interpret the correlation coefficient of a linear fit
- 1908 (S-ID.8)
- 1909 • Estimate a line of best fit for a single linear regression (S-ID.6)
- 1910 • Determine and interpret the strength of correlation to determine the best fit. (S-
- 1911 ID.8)
- 1912 • Understand independent and dependent events and the and that two
- 1913 independent events have a probability of occurring together that is a product of
- 1914 their individual probability of occurring (S-CP.2,4)
- 1915 • Conditional probability (S-CP.3,4,5,6)
- 1916 • Construct and interpret two-way frequency tables of data when two categories
- 1917 are associated with each object being classified. Use the two-way table as a
- 1918 sample space to decide if events are independent and to approximate conditional
- 1919 probabilities. (S-CP.4)
- 1920 • Recognize and explain the concepts of conditional probability and independence
- 1921 in everyday language and situations. (S-CP.5)
- 1922 • Calculate expected values and use them to solve problems. (S-MD.2,3,4,5)
- 1923 • Calculate the expected value of a random variable and interpret it as the mean of
- 1924 the probability distribution (S-MD.2)
- 1925 • Use probability to evaluate outcomes of decisions (S-MD.5,6,7)
- 1926 • Recognize situations in which one quantity changes at a constant rate per unit
- 1927 interval relative to another (F-LE.1)
- 1928 • Recognize situations in which a quantity grows or decays by a constant percent
- 1929 rate per unit interval relative to another (F-LE.1)

1930 ***Grappling with Variability and Uncertainty***

1931 Students understand variability is inherent to data and are able to identify multiple

1932 sources of it. They practice collecting and organizing data about their own lives and

1933 communities as well as working with large, real-world, publicly available data sets.

1934 Students consider sampling practices and how they affect the data that is collected.

1935 They can use probability to make decisions and understand the uncertainty that comes
1936 along with predictions.

- 1937 • Know that statistics is a process for making inferences about population
1938 parameters based on random samples of the population (S-IC.1)
- 1939 • Determine if a model from a data generating process or simulation is accurate
1940 (S-IC.2)
- 1941 • Make inferences and justify conclusions from sample surveys, experiments and
1942 observational studies (S-IC.3,4,5,6)
- 1943 • Use data from a sample to estimate population mean or proportion and develop a
1944 margin of error through simulation models (S-IC.4)
- 1945 • Use simulations to decide if differences between parameters are significant (S-
1946 IC.5)
- 1947 • Evaluate reports based on data (S-IC.6)

1948 ***Transforming Data with Technology***

1949 Students understand that data is not always collected/shared/received ready to be
1950 analyzed and it sometimes requires work to prepare it. They can use different digital
1951 tools to clean and transform the data (e.g., merge data, deal with incomplete data,
1952 normalize data). They are familiar with the basics of programming as needed, and are
1953 comfortable editing code or finding the appropriate tools to transform the data in ways
1954 useful to their own data analysis. Students can combine their knowledge of probability
1955 and programming to construct simulations of probabilistic events, and they understand
1956 the basic idea behind machine learning as well as its power and shortcomings.

- 1957 • Perform operations on matrices and use matrices in applications (N-VM)
- 1958 • Use matrices to represent and manipulate data (N-VM.6)
- 1959 • Cleaning names, categories, and strings
- 1960 • Simulation using experimental data
- 1961 • Translate between different bit representations of real-world phenomena, such as
1962 characters, numbers, and images (CS.9-12.DA.8)

- 1963 • Evaluate the tradeoffs in how data elements are organized and where data is
- 1964 stored (CS.9-12.DA.9)
- 1965 • Create clearly named variables that represent different data types and perform
- 1966 operations on their values (CS.6-8.AP.11)
- 1967 • Collect data using computational tools and transform the data to make it more
- 1968 useful and reliable (CS.6-8.DA.8)
- 1969 • Refine computational models to better represent the relationships among
- 1970 different elements of data collected from a phenomenon or process (CS 9-12.
- 1971 DA.11)
- 1972 • Evaluate the ability of models and simulations to test and support the refinement
- 1973 of hypotheses (CS 9-12. S. D.A. 9)
- 1974 • Use data analysis tools and techniques to identify patterns in data representing
- 1975 complex systems (CS 9-12. S. D.A.8)

1976 **Sample Courses**

- 1977 Effective Data Science courses consider how to help students with the following:
- 1978 • Understand how data are used by professionals to address real-world problems.
 - 1979 • Understand that data are used in all facets of modern life.
 - 1980 • Understand how data support science to identify and tackle real-world problems
 - 1981 in our communities.
 - 1982 • Learn about statistical variability.
 - 1983 • Create and analyze statistical graphics to identify patterns in data and to connect
 - 1984 these patterns back to the real world.
 - 1985 • Understand that by treating photos, words, numbers, and sounds as data, we
 - 1986 can gain insight into the real world.
 - 1987 • Learn to analyze data, including: posing questions that can be answered by
 - 1988 considering relations among variables in a data set, using collected data to
 - 1989 generate hypotheses for future data collection, critically evaluating shortcomings
 - 1990 and strengths in the data and the data collection process, and informally
 - 1991 evaluating hypotheses using data at hand.

- 1992 • Learn initial programming, and use programs in the development and analysis of
- 1993 statistical models.
- 1994 • Learn about data ethics, including consideration of where data comes from, who
- 1995 is collecting it, and how it is used.
- 1996 • Refine computational models to better represent the relationships among
- 1997 different elements of data collected and analyzed.
- 1998 • Design algorithms to solve computational problems, using and adapting existing
- 1999 algorithms and creating new ones.

2000 Another sample course begins with a consideration of the meaning of data, the
 2001 importance of communicating data visually, investigating community issues, cleaning
 2002 data, exploratory data analysis, ethical issues around data, creating data dashboards,
 2003 linear and nonlinear regression models, statistics, probability, and forecasting. The
 2004 course is designed to engage students actively and to be flexible enough for teachers to
 2005 include local issues of importance to their communities. While addressing concepts of
 2006 data analysis with rigor, the access and dependence upon current, local, and publicly
 2007 accessible data is a key feature. One goal of the course is that it be open to all students,
 2008 regardless of prior mathematics achievement, all lessons will be “low-floor and high-
 2009 ceiling”—designed so that everyone can access them and they extend to high levels.

2010 Some schools have created a Data Science elective course for students who have
 2011 completed algebra and geometry or integrated 1 and 2. Courses may teach
 2012 distributions, linear regression, matrices, probability, programming, machine learning,
 2013 and statistical inference through investigation-based activities. Course activities may
 2014 include analyzing the link between poverty and obesity or creating and examining
 2015 algorithms that decide, for example, music recommendations, or airport screening.

2016 Districts can design their course to meet A–G course requirements for Mathematics.
 2017 Three data science courses already exist in California, all satisfying A-G requirements.
 2018 These courses can be taken as an alternative to, or in addition to, algebra 2 (see high
 2019 school chapter). The courses also provide lessons that can be used in other courses –

2020 such as integrated 1 and 2, and in algebra and geometry. If students are intending to
2021 pursue STEM majors in college (including data science), it is advised that they take
2022 courses that, at minimum, prepare them to take a calculus course by their freshman
2023 year of college. Chapter 8 and Appendix A give further detail on the different pathways
2024 offering options to schools and students.

2025 An additional example of school-created course for students in grade nine is one
2026 focused on software design and data science. It teaches algebraic, geometric, and
2027 statistical concepts through contexts like video-game design. This course can be an
2028 example of a modernized integrated pathway, teaching the traditional sequence through
2029 modern mediums and applications. The course can also be designed to meet A–G
2030 elective credit requirements.

2031 The different examples of courses and high-school approaches above use different
2032 software and tools, which seems appropriate as data science does not require any
2033 particular software package, it is more important that students learn to ask good
2034 questions and apply an effective tool to help them answer them. Exposure to some
2035 software is essential for those wishing to pursue a full-time career in data science, and
2036 comfort with such programs is increasingly valuable for many other professions that
2037 involve basic data analysis.

2038 In total, over 70 individual high schools and 15 districts offered a data science
2039 mathematics or elective course in California during the 2019–2020 school year that
2040 counted for A–G credit (University of California data). That compares to just 34 high
2041 schools and 6 districts two years before in 2017–2018. This rapid increase in course
2042 offerings is likely an indication of both high interest in and importance of data science
2043 content throughout the curriculum.

2044 **Conclusion**

2045 Life in a data-rich world requires California schools prepare all students to examine
2046 claims justified with data, to understand the probabilistic underpinning of drawing
2047 conclusions from samples, and to see data as a tool to answer many questions of

2048 interest. Developing these abilities requires that students generate questions and work
 2049 with data beginning in kindergarten (or before), and have experiences of increasing
 2050 depth and complexity throughout their school careers. Students who wish to focus extra
 2051 attention on data science should have an opportunity to pursue advanced courses late
 2052 in their high school careers.

2053 Above all, students at all levels should have experiences that build their mathematical
 2054 toolkit for making sense of their worlds.

2055 Long Descriptions for Chapter 5

2056 Figure 5.1. Example of California data for students to explore

2057 Acres Burned by California Wildfires: Starting in 2001, the National Interagency Fire
 2058 Center began keeping more accurate records on the total fire acreage burned in each
 2059 state. Data for California includes the number of fires reported each year from 2000 to
 2060 2018, and the acreage burned.

Year	Fires	Acres
2000	7,622	295,026
2001	9,458	329,126
2002	8,328	969,890
2003	9,116	1,020,460
2004	8,415	264,988
2005	7,162	222,538
2006	8,202	736,022
2007	9,093	1,520,362
2008	6,255	1,593,690
2009	9,159	422,147
2010	6,554	109,529
2011	7,989	168,545
2012	7,950	869,599
2013	9,907	601,635
2014	7,865	625,540
2015	8,745	893,362
2016	6,986	669,534
2017	9,560	1,548,429
2018	8,527	1,975,086

2061 [Return to image.](#)

2062 Figure 5.2: The Drivers of Investigation

2063 Image long description: Three Drivers of Investigation (DIs) provide the “why” of
2064 learning mathematics: Make Sense of the World (Understand and Explain); Predict
2065 What Could Happen (Predict); Impact the Future (Affect). The DIs overlay and pair with
2066 four categories of Content Connections (CCs), which provide the “how and what”
2067 mathematics (CA-CCSSM) is to be learned in an activity: Communicating Stories with
2068 Data; Exploring Changing Quantities; Taking Wholes Apart, Putting Parts Together;
2069 Discovering Shape and Space. The DIs work with the Standards for Mathematical
2070 Practice to propel the learning of the ideas and actions framed in the CCs in ways that
2071 are coherent, focused, and rigorous. The Standards for Mathematical Practice are:
2072 Make sense of problems and persevere in solving them; Reason abstractly and
2073 quantitatively; Construct viable arguments and critique the reasoning of others; Model
2074 with mathematics; Use appropriate tools strategically; Attend to precision; Look for and
2075 make use of structure; Look for and express regularity in repeated reasoning.
2076 [Return to image.](#)

2077 Figure 5.5: Student Work from an Investigation of Mammals

2078 Sample student data collection board with sticky notes, charts, and pictures. One sticky
2079 note has the student names. Another sticky note shares one of the students' question:
2080 Do big animals sleep more than small animals? The students answer this question in
2081 another sticky note, saying that their graph shows that some of the biggest animals
2082 sleep the least. In another sticky note, students ask a second question: do plant eaters
2083 sleep less than meat eaters? They answer this in another sticky note that says: Yes! We
2084 can see that almost all plant eaters represented sleep less on average. Three more
2085 sticky notes share other data findings about the rabbit— One says: the rabbit is small in
2086 mass and sleeps an average amount compared to other mammals. Another says that
2087 when comparing the rabbit to other plant eaters it stands out for sleeping much more.
2088 And a third says: There is one animal smaller than the rest in the plant eaters that
2089 sleeps a lot more - the animal is a rabbit. The students also share graphs that show the
2090 different relationships they investigated (sleep v. diet, mass v. sleep and diet v. sleep).
2091 [Return to image.](#)

2092 Abdu's Question

2093 Abdu's question *How many times does Abdu's 6-year-old sister use English and non-*
2094 *English words she knows or does not know while pretending to be a teacher?* is
2095 represented in a data visualization. Data is plotted on an axis that includes the following
2096 "real words" she knows and used while pretending: good morning; good evening; apple;
2097 book; stand up; how are you; run; bye; very good; and the following non-words she
2098 knows: roos; sat; room; roomba; soon. On the negative X axis, the words bass, rabba,
2099 shant, and rooman are plotted.

2100 [Return to image.](#)

2101 Nikita's Description

2102 Nikita's description of *One week of listening to music, its genre, and what she was doing*
2103 is represented in a data visualization. Days of the week include flowers where each
2104 petal of the flower constitutes one hour of listening. Petals are separated by genre (pop,
2105 EDM, classical, rap, and R and B). Days of the week include flowers *without* petals
2106 whose design reflects when music was listened to for less than one hour. Each leaf on
2107 the stem of the flower indicates 15 minutes of listening. Flowers are colored according
2108 Nikita's actions while listening: lavender is relaxing; pink is traveling; orange is hanging
2109 out with friends; yellow is bedtime; rose is chores; fuchsia is showering; paisley is
2110 eating; navy is doing schoolwork.

2111 [Return to image.](#)

2112 Nathan's Description

2113 Nathan's data visualization that indicates *representation of sound length, level of*
2114 *loudness, and how much attention was given to it*, shows his sound wave. Along the
2115 wave, each line represents 10 minutes of sound. Lines extending above the wave line
2116 represents sounds that warrant attention; lines extending below the represent ambient
2117 sounds. The length of the line indicates the general loudness of the sound. The colors
2118 code the sounds, which include: my room (AC), outdoors/street, car noise, restaurant,
2119 living room/kitchen, food sounds, alarm clock, music, mixing music, people/family, the
2120 TV (SNL), 2-hour Zoom meeting for orchestra, and movie (Borat 2).

2121 [Return to image.](#)

2122 Visualization of Kira's Dog's Interactions

2123 Kira’s dog interactions visualization begins, “Dear Data: For a day, between 6am –
2124 10pm while I was awake, I recorded my interactions with my dog—Daisy, a
2125 goldendoodle—and her (sometimes sassy) behaviors. Usually, she is my study buddy
2126 for the day. This data is from Wednesday 11/4/20. Below you will find the key.
2127 Something to note, starting from the outermost layer of a petal and going inward
2128 accounts for the order of the actions.”

2129 The data is presented in a flower design. Each petal represents an hour of the
2130 timeframe, and each petal is designed according to the table below.

2131 [Return to image.](#)

California Department of Education, March 2022