

R for Everyone — An Interactive, Ethics-Aware Data Science Micro-Curriculum for Middle School

Tene Ariyo

November 21, 2025

Introduction and Motivation

Data has become a native language of civic life and work. People scroll past charts on social media, navigate dashboards, and encounter numerical claims about health, climate, and finance every week. To participate effectively and build strong data literacy, they need guided practice with the full statistical problem-solving process: posing investigable questions, collecting or considering data, analyzing patterns and variability, and interpreting results with appropriate attention to uncertainty and audience.

Why Middle school?

We target middle school because it is the moment when many students begin owning or regularly using personal devices (e.g., iPhones, tablets, laptops), which means they already search, click, and interact with data daily. Meeting them at this stage lets us turn informal scrolling into intentional inquiry asking good questions, collecting or considering data responsibly, and making first plots they can explain. Yet access to formal data science learning opportunities remains uneven, and many middle school curricula still treat statistics as a short end-of-year unit rather than a core problem-solving lens (Bargagliotti et al. 2020).

National reports and coalitions such as Data Science 4 Everyone argue that data literacy is becoming a foundational civic competency, comparable to reading and algebra, and that students should experience tools and tasks that resemble those used by practicing data scientists, not just hand-computed summaries. At the same time, research shows that data literacy outcomes are falling fastest in these middle school grades: between 2019 and 2022, 8th-grade scores on NAEP’s Data Analysis, Statistics, and Probability category declined by about 10 points, roughly a full grade level outpacing drops in other math areas (Drozda 2023).

The Guidelines for Assessment and Instruction in Statistics Education (GAISE II) emphasize a four-step statistical problem-solving process, **formulating questions, collecting or considering data, analyzing data, and interpreting results**, and recommend that

instruction center authentic questions and data sources (Bargagliotti et al. 2020). Complementary resources such as youcubed’s “Data Big Ideas” provide grade-level guidance for how these ideas can unfold in the middle grades, including attention to variability, fairness, and interpretation of visual displays (YouCubed, n.d.).

This capstone is situated at the intersection of these recommendations and the emerging literature on **embedding ethics in mathematics and statistics education**. Instead of treating ethical considerations as stand-alone lessons, recent work suggests integrating short, recurring “micro-prompts” that ask students to reflect on bias, missing data, and who is affected by decisions based on data (Chiodo, Müller, and Shah 2025). My goal is to design a small set of learnr-based lessons that implement this approach in R, at a level appropriate for middle school students and their teachers.

Curriculum Design and Methods

Framing the four-lesson arc

The four-lesson sequence is explicitly organized around the GAISE problem-solving cycle:

1. **Formulate questions**
2. **Collect or consider data**
3. **Analyze data**
4. **Interpret results**

Using the youcubed “Data Big Ideas” for grades 6–7 as a guide, I mapped these phases onto a four-lesson arc that can fit into a short unit:

- **Lesson 1 – Ask Good Questions**
Focus: distinguishing investigable vs. non-investigable questions, identifying variables, and beginning to think about sampling and bias.
- **Lesson 2 – Collect & Consider Data**
Focus: examining a small, structured dataset; understanding rows and columns; recognizing variability; and running first R commands (`nrow()`, `head()`, simple plots).
- **Lesson 3 – Analyze Patterns**
Focus: creating and interpreting graphical displays (dotplots, bar charts, or boxplots), comparing groups, and quantifying simple differences.
- **Lesson 4 – Interpret and Communicate**
Focus: synthesizing findings, articulating limitations, and reflecting on how data and ethics intersect in drawing conclusions.

The aim is to make R feel approachable to middle school students while keeping the tasks anchored in authentic question-posing and interpretation, rather than in syntax for its own sake.

Data and Instructional Artifacts

In this project, “data” appears in several layers:

1. **Contextual data** from reports and coalitions (e.g., NAEP results, Data Science 4 Everyone, GAISE) that motivate the need for early data science experiences.
2. **Instructional data** embedded directly inside the learnr tutorials (e.g., a small reading-time survey in Lesson 2).
3. **Validation data**, which will eventually include teacher feedback collected via a short Google Forms rubric.

Contextual data

I draw on existing reports and position papers primarily as background rather than as objects of original statistical analysis. These documents provide:

- Evidence that data literacy gaps exist and have equity implications.
- Conceptual frameworks for structuring a coherent progression of ideas across grades.
- Recommendations for integrating technology and real-world datasets.

Rather than re-analyzing those large datasets, I use them to justify the design choices in my four-lesson sequence and to anchor the learning goals for each lesson.

Instructional data: the Lesson 2 reading survey

Lesson 2 introduces students to small, manageable datasets that they can explore within R. For example, a simplified reading survey dataset records how many minutes each student read the previous evening and which device they typically use for reading or homework.

In the learnr tutorial, this dataset is created behind the scenes so that students can immediately focus on interpreting rows, columns, and variability. A simplified version of the dataset creation appears below:

```
set.seed(123)

survey <- tibble(
  student_id = 1:12,
```

```

grade      = rep(7, 12),
minutes_read = c(20, 35, 0, 15, 40, 30, 10, 25, 5, 50, 15, 45),
device      = c(
  "phone", "tablet", "phone", "phone",
  "tablet", "other", "other", "phone",
  "tablet", "other", "phone", "tablet"
)
)
survey

```

```

## # A tibble: 12 x 4
##   student_id grade minutes_read device
##       <int> <dbl>         <dbl> <chr>
## 1         1     7           20 phone
## 2         2     7           35 tablet
## 3         3     7            0 phone
## 4         4     7           15 phone
## 5         5     7           40 tablet
## 6         6     7           30 other
## 7         7     7           10 other
## 8         8     7           25 phone
## 9         9     7            5 tablet
## 10        10     7           50 other
## 11        11     7           15 phone
## 12        12     7           45 tablet

```

Students are then guided through basic “first look” operations:

```
nrow(survey)
```

```
## [1] 12
```

```
min(survey$minutes_read)
```

```
## [1] 0
```

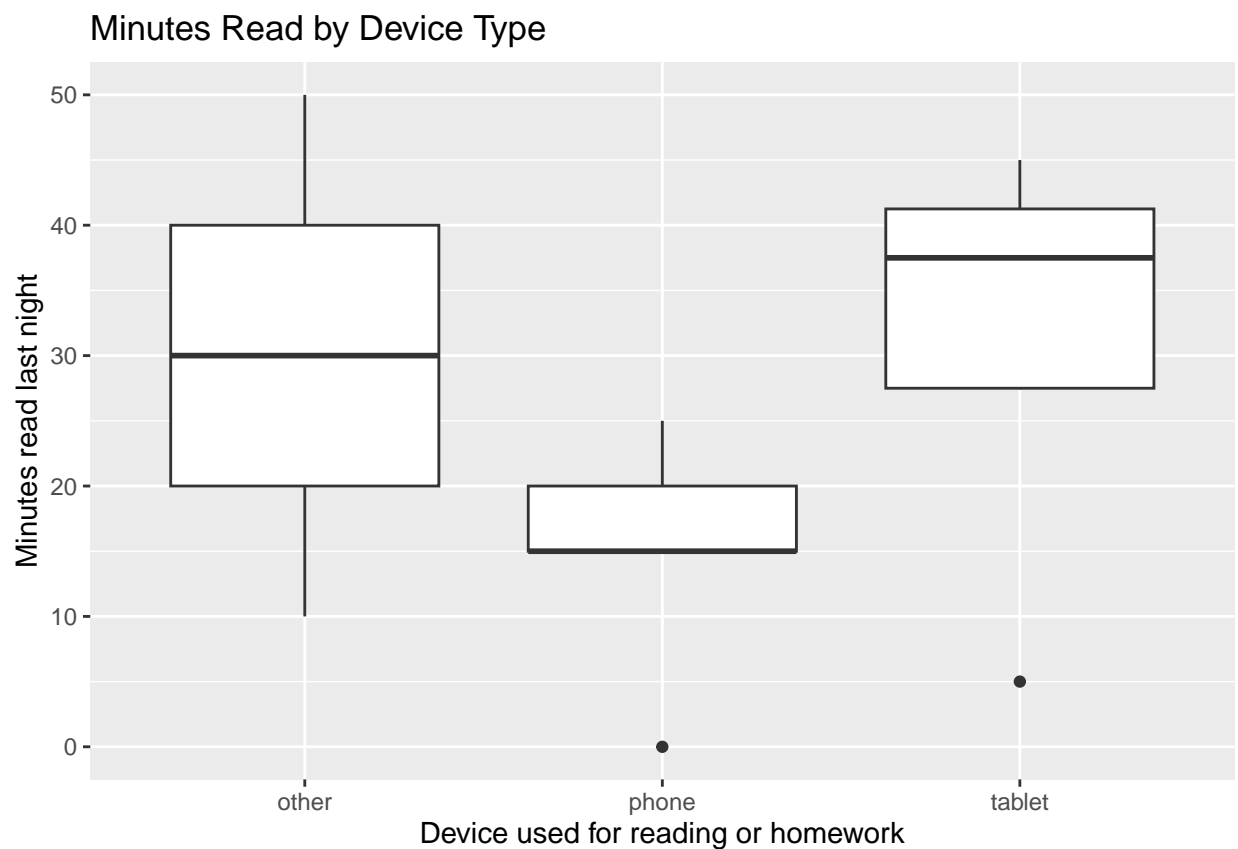
```
max(survey$minutes_read)
```

```
## [1] 50
```

These operations are not just technical; they prepare students to reason about **variability** (“Why do the reading minutes differ across students?”) and to confront **ethical questions** about who is represented in the data and who is missing.

In lesson 3, a simple plot of reading minutes by device type serves as a springboard for both statistical and ethical interpretation:

```
ggplot(survey, aes(x = device, y = minutes_read)) +  
  geom_boxplot() +  
  labs(  
    title = "Minutes Read by Device Type",  
    x = "Device used for reading or homework",  
    y = "Minutes read last night"  
  )
```



Planned teacher validation data

To support iterative refinement, I also designed a short **Teacher Validation Rubric** delivered through Google Forms. The rubric asks teachers (or prospective teacher partners) to rate each lesson on clarity, grade-level appropriateness, alignment with the GAISE framework, and usefulness of the embedded ethics prompts. Open-ended items invite sugges-

tions for improvement and concerns about classroom implementation. *Link to the form:* <https://forms.gle/5UdDv4Q3WCDaYhu37>

At the time of this draft, I have not yet collected classroom teacher responses. Instead, I have relied on ongoing feedback from my capstone advisor, who has helped refine lesson goals, wording, and pacing. In future work, responses to the rubric will function as implementation data, complementing the design analysis in this write-up.

Writing and iterating on prompts

I developed the lesson prompts in several passes:

1. **Initial drafting from frameworks**

I began by translating the GAISE steps and youcubed big ideas into concrete middle school tasks. For Lesson 1, this meant building a set of example questions that varied along dimensions such as specificity, measurability, and scope.

2. **Embedding ethics micro-prompts**

Using the *Teaching Resources for Embedding Ethics in Mathematics* as a guide, I aimed to include **short, recurring touches** on ethics rather than one large “ethics day”. This meant adding brief reflective questions in context, such as asking students who might be left out of a survey or how device access could shape the data.

3. **Feedback and refinement**

Throughout the process, I shared lesson drafts with my capstone advisor, who provided feedback on pacing, clarity, and the cognitive load imposed by R syntax. Teacher feedback via the validation rubric is still an ongoing process, and I plan to incorporate those perspectives in future iterations of both the lessons and this write-up.

Analysis of Lessons 1 and 2

In this section, I analyze Lessons 1 and 2 as designed artifacts, focusing on their alignment with the GAISE framework and the ethics-embedding goals.

Lesson 1: Formulating questions and noticing bias

Lesson 1 asks students to distinguish between questions that can be investigated with data and those that cannot. For instance, students might be given a list such as:

- “What’s the best candy”
- “Which candy is the most popular in our class”

- “Is chocolate good”

They then choose the answer they think is **investigable with data**. They are then asked to rewrite this vague question, “Do people read a lot?”, to make it investigable. This directly supports the GAISE emphasis on formulating clear, measurable questions.

An embedded ethics micro-prompt invites students to reflect on sampling bias:

If you only survey friends sitting near you, what sampling bias might that introduce, and how could you reduce it? Write 2–3 sentences.

This prompt encourages students to notice how their own social networks and habits might limit the data they collect. Rather than asking for a single correct answer, the prompt invites argumentation: students must articulate a potential bias (e.g., only surveying people with similar interests or schedules) and propose an improvement (e.g., random sampling, making sure different groups are represented). This shift—from finding “the right answer” to **arguing a position using ideas about bias and representation**—is directly inspired by the ethics resource’s recommendations for short, reflective tasks.

Screenshots from Lesson 1 can illustrate these elements in context:

Warm-up: What makes a statistical question?

A good statistical question expects **different answers** and defines what we will measure. Some questions are subjective (“best”), but we can still study opinions statistically if we measure them (e.g., “What percent of students say chocolate is their favorite?”). For example, “Which snack is most popular in our class?” is investigable because we can count votes; “What’s the best snack?” is ambiguous unless we define how we’ll measure “best.”

Quiz

Which is *investigable with data*?

- ☐ What’s the best candy?
- ☐ Which candy is most popular in our class?
- ☐ Is chocolate good?

Submit Answer

- *Figure 1.*

Ethics Micro-Prompt (2–3 sentences)

Consider the question: “Which snack is most popular in our class?”

Bias (quick definition): a systematic error in how we **sample** or **measure** that pushes results in one direction.

If you only survey friends sitting near you, what sampling bias might that introduce, and how could you reduce it? Write 2–3 sentences.



- *Figure 2.*

Lesson 2: Considering data and variability

Lesson 2 moves from questions to data. Students are given the **survey** dataset shown earlier and are asked to interpret it:

- What does each **row** represent?
- What does each **column** represent?
- Which variables are numeric, and which are categorical?

Students then use R code to gain a first sense of the data:

```
min(survey$minutes_read)
```

```
## [1] 0
```

```
max(survey$minutes_read)
```

```
## [1] 50
```

In narrative form, the tutorial explains that even if all students answer the same question, their responses naturally differ. This **variability** is a central big idea:

Some students read 0 minutes, some read 15–25 minutes, and some read 40–50 minutes. This natural difference across students is called variability.

An interactive exercise asks students to compute or interpret the range of reading minutes, reinforcing the idea that data spread is meaningful rather than just noise.

Ethically, Lesson 2 invites students to discuss questions such as:

- How might the sample or the way we asked the question introduce bias?
- How could we make the data collection more fair?

These short prompts are designed to surface issues of access, equity, and representation without requiring a full separate ethics lecture.

Cont'd: Lessons 3 and 4

The roles of Lessons 3 and 4 in the sequence are:

Lesson 3 – Analyze Patterns

Lesson 3 will deepen students' ability to **analyze patterns** in data using graphical displays and simple numerical summaries. Building on the Lesson 2 survey, students might:

- Create boxplots or bar charts of reading minutes.
- Compare distributions across groups (e.g., by device type or by another categorical variable).
- Begin to articulate claims such as “On average, group A tends to read more than group B,” supported by plots and summary statistics.

Ethics prompts in Lesson 3 will continue to foreground representation. For example, students might consider how an unbalanced sample (e.g., very few students in one category) affects the reliability of their comparisons and who might be overlooked when decisions are based on such data.

A link to the live learnr tutorial will be inserted here once Lesson 3 is fully deployed:

Lesson 4 – Interpret and Communicate

Lesson 4 is intended to culminate the unit by asking students to **interpret results and communicate conclusions** in written or oral form. Possible tasks include:

- Writing a short “data story” about reading habits in their class, citing plots and summary statistics.

- Explaining limitations of their data (sample size, sampling method, missing groups).
- Reflecting explicitly on how ethical considerations—such as sampling bias or unequal device access—should shape the way they talk about their findings.

Here, the ethics emphasis is on **responsible communication**: students are encouraged to recognize that data summaries can influence decisions and perceptions, and that it is important to acknowledge uncertainty and limitations.

Ethical Challenges and Responses

The ethics component of this project is documented in detail below. I summarize the main challenges I encountered and the design choices I made to address them, drawing on the *Teaching Resources for Embedding Ethics in Mathematics* and related literature.

1. Data Privacy and Consent

A major ethical consideration in this project was ensuring that any data used, especially if it involved students, teachers, or minors, respected privacy and consent standards. Because my *R for Everyone* curriculum is designed for middle school classrooms, I needed to make sure no identifiable or sensitive information from students would ever be stored, shared, or analyzed.

How I addressed it:

- Used only simulated or publicly available datasets rather than collecting student-level data.
- Designed the learnr lessons to process data locally in RStudio or RStudio Cloud on the user’s device, without exporting responses externally.
- Removed any automatic data-logging features from quizzes or free-response questions to avoid saving identifiable information.

2. Informed Participation

Since this project may later involve classroom pilots, I also considered what it would mean for teachers and students to participate voluntarily and knowingly.

How I addressed it:

- Created a teacher validation rubric for expert review rather than using live student data in this phase.
- Included a short consent statement in the teacher feedback form clarifying that participation is voluntary, anonymous, and for educational research purposes only.

3. Algorithmic Bias and Representation

Teaching data science to younger audiences raises the risk of unintentionally reinforcing biased or incomplete views of data. For instance, datasets that overrepresent certain regions or demographics could lead students to draw skewed conclusions.

How I addressed it:

- Added discussion prompts encouraging students to question data sources, collection methods, and who might be excluded.
- Included a short “ethics micro-prompt” in each lesson, for example: “*Whose data is missing here, and how might that affect the story this graph tells?*”

4. Accessibility and Digital Equity

Because not all schools or students have the same access to technology, the project needed to be designed with varying levels of connectivity and device availability in mind.

How I addressed it:

- Structured the learnr lessons so they can be run offline in RStudio or RStudio Cloud using minimal computing resources.
- Provided printable summary worksheets and screenshots for teachers with limited device access.
- Planned for open-source release (e.g., via shinyapps.io) so that all materials remain free to use and adapt.

5. Pedagogical Transparency

Finally, I wanted to ensure that the project models ethical data literacy itself: students should learn not only *how* to analyze data but *why* ethical reflection matters.

How I addressed it:

- Embedded reflection checkpoints after each coding task, asking students to consider reliability, limitations, and implications of their analyses.
- Used plain-language explanations to make ethics approachable rather than abstract or punitive.
- Documented all datasets, their origins, and transformations in reproducible R Markdown cells.

Summary

Throughout this capstone, ethical considerations guided both the design (how data are collected and represented) and the delivery (how learners interact with R safely and thoughtfully). The overarching goal was not just to avoid harm, but to teach ethical awareness as a core part of data science practice.

Conclusion and Future Work

This capstone project demonstrates that it is possible to create **digestible, ethics-aware data science lessons for middle school students** using R and learnr. By grounding the four-lesson sequence in the GAISE statistical problem-solving framework and the youcubed data big ideas, and by incorporating recurring ethics micro-prompts inspired by recent work on embedding ethics in mathematics, the project offers one concrete model for integrating tool use, statistical thinking, and ethical reflection.

The analysis in this write-up focuses on design choices and alignment with established frameworks rather than on classroom outcome data, reflecting the current stage of the project.

Future work will involve:

- Collecting teacher feedback through the validation rubric and possibly through semi-structured interviews.
- Iterating on lesson content, pacing, and ethics prompts in response to that feedback.
- Exploring opportunities to pilot the lessons with local middle schools or after-school programs and to study how students engage with both the statistical and ethical dimensions.

Ultimately, the project aims not only to provide a workable set of materials, but also to contribute to the broader conversation about what it means to teach data science **for everyone**, in ways that are technically accessible, intellectually honest, and ethically attentive.

Attached below are the links to the Lessons

- **Lesson 1 tutorial link:** <https://tene.shinyapps.io/Lesson1/>
- **Lesson 2 tutorial link:** <https://tene.shinyapps.io/Lesson2/>
- **Lesson 3 tutorial link:** <https://tene.shinyapps.io/Lesson3/>
- **Lesson 4 tutorial link:** <https://tene.shinyapps.io/Lesson4/>

Bargagliotti, Anna, Christine Franklin, Pip Arnold, Rob Gould, Sheri Johnson, Leticia Perez, and Denise A. Spangler. 2020. *Pre-k–12 Guidelines for Assessment and Instruction in Statistics Education II (GAISE II): A Framework for Statistics and Data Science Education*. Alexandria, VA: American Statistical Association. https://www.amstat.org/docs/default-source/amstat-documents/gaiseiiprek-12_full.pdf.

- Chiodo, Maurice, Dennis Müller, and Rehan Shah. 2025. “Teaching Resources for Embedding Ethics in Mathematics: Exercises, Projects, and Handouts.” arXiv. <https://doi.org/10.48550/arXiv.2310.08467>.
- Drozda, Zarek. 2023. “Data Science Is Vital to Student Success. So Why Are Outcomes Going Down?” Data Science 4 Everyone. https://d2c47c0f-cda1-420e-ae77-f1279b4ac7ea.usrfiles.com/ugd/d2c47c_067273b2bef041c7ae08cab6c7a3be8c.pdf.
- Education, Stanford Graduate School of. 2021. “How to Teach Data Science in k–12 Schools? Stanford-Led Team Launches ‘Big Ideas.’” October 13, 2021. <https://ed.stanford.edu/news/how-teach-data-science-k-12-schools-stanford-led-team-launches-big-ideas>.
- FiveThirtyEight. 2023. “FiveThirtyEight Data Repository.” 2023. <https://github.com/fivethirtyeight/data>.
- qmd-lab. 2025. “Closeread: A Quarto Format for Scrollytelling.” 2025. <https://closeread.dev/guide/>.
- Schloerke, Barret, Garrick Aden-Buie, RStudio, and PBC Posit Software. 2025. *Learnr: Interactive Tutorials for r*. <https://rstudio.github.io/learnr/>.
- Xie, Yihui. 2024. *Bookdown: Authoring Books and Technical Documents with r Markdown*. <https://bookdown.org/yihui/bookdown/>.
- YouCubed. n.d. “Data Science k–10 Big Ideas.” <https://www.youcubed.org/data-big-ideas/>.