

Introductory Data Science: A Blueprint to Navigate Curricular, Pedagogical, and Computational Challenges

Elijah Meyer

Department of Statistics, Duke University

and

Mine Çetinkaya-Rundel

Department of Statistics, Duke University

April 6, 2023

Abstract

The text of your abstract. 200 or fewer words.

Keywords: Data Science, Curriculum, Pedagogy

1 Introduction

2 The Course

In the following sections, we describe our introductory data science course, offered to predominantly freshman level students at Duke University, titled Introduction to Data Science and Statistical Thinking (STA199). Often, these students are **(interested in / major + minor in / description of typical student demographic)**. Students enrolling in STA199 commonly have little to no statistics, data science, or coding experience. In a typical semester, this course seats roughly 179 students, which is considered to be large by all measures. Both the lack of experience and large class size are identified as two common hurdles by faculty when trying to create and instruct an introductory data science course ([Schwab-McCoy et al. 2021](#), [Kokkelenberg et al. 2008](#)). However, by the end of this course, students are able to use data both the statistical software R and GitHub to create data visualizations, investigate patterns, and model outcomes in a reproducible format.

This course is built on four large-scale learning objectives: learn to explore, visualize and analyze data in a reproducible and shareable manner; gain experience in data wrangling and munging, exploratory data analysis, predictive modeling, and data visualization; work on problems and case studies inspired by and based on real-world questions and data; learn to effectively communicate results through written assignments and project presentation. These objectives are accomplished through interactive lectures and labs that present content, problems, and case studies inspired by and based on real-world questions and data.

When teaching, instructors are committed to the Kaplan Way learning model (Figure 1) “that combines a scientific, evidence-based design philosophy with a straightforward educational approach to learning” ([Schweser 2023](#), pg. 3).

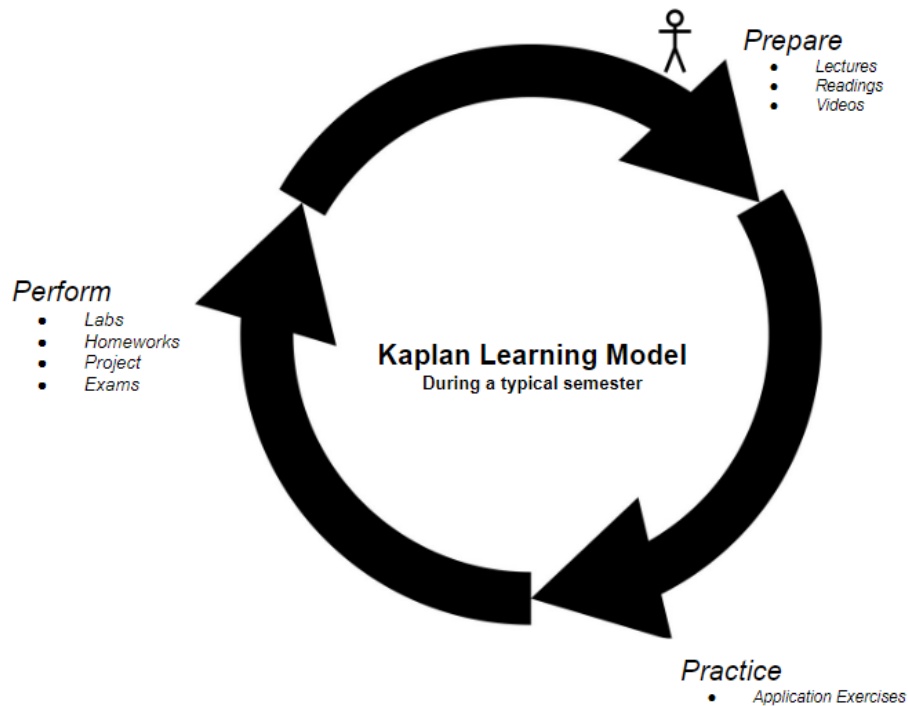


Fig 1: Kaplan Learning Model for STA 199

This model is composed of three phases: prepare, practice, and perform. All three phases are designed to guide the instructor in facilitating an overall quality learning experience for the student. Each of these phases are performed during a typical week within a semester, accomplished through preparation materials, lectures, in-class application exercises (AEs), labs, and assessments.

During the prepare phase, students are completing readings, watching videos, and listening to lectures. These modes of learning ultimately helps students' develop a new foundation of data science concepts. The goal of the prepare phase is to put students in a position where they can build upon what they are learning, and create new knowledge through the practice phase. The practice phase is designed to be an opportunity for students to reinforce the new information gained, as well as uncover new concepts in data science. This is achieved through the use of interactive AEs administered in class where students either work alone,

in groups, or with the instructor in live coding sessions. Finally, students enter the perform phase, showcasing their progress made in the previous two phases. This is typically done through assessments such as homework, labs, exams, and projects. This learning model is a continuous cycle throughout the semester as new topics are introduced.

Topics taught through this cycle fall under two major units: Unit 1 - Exploring data; Unit 2 - Making rigorous conclusions. In Unit 1, students become first introduced to R, R-studio, and Github. During this unit, students start to create data visualizations and learn how to both import and manipulate data to be better suited for modeling. In Unit 2, students extend their investigations with data to include modeling. Specifically, students fit a variety of models (simple linear regression, multiple linear regression, logistic regression), and learn the fundamentals of hypothesis testing and confidence intervals. For a more complete description of the topics taught and data sets used in creating these lessons, please see *A fresh look at introductory data science (?)*.

In this paper, we describe the preparation and implementation process necessary to run STA199, in its entirety. This includes details of a teaching team used to instruct, technology we've chosen to use, as well as our pedagogical choices that go into a typical week of teaching. This comprehensive description is encouraged to be used as a flexible framework on how to create, set up, and implement an introduction to data science course similar to STA199. In this description, we articulate first hand experiences, suggestions, and provide example code and lessons that aim to support newer instructors who are developing or teaching an introductory data science course; instructors with courses that are increasing in size; and instructors who want to implement more technical tools into their curricula and classroom.

3 Teaching Team

We structure a teaching team to help account for the difficulties a large class size can bring. Our teaching team consists of one instructor and multiple teaching assistants (TAs). The responsibilities of any TA is to both support the instructor in charge of the class, and support the students in the classroom. These TAs range from undergraduate to PH.D. level students, and vary in teaching experience. **(Writing on TA selection process).**
Once selected... (writing on TA training).

Once training is complete, TAs are assigned roles that indicate their responsibilities during the semester. These roles include *course organizer*, *head TA*, *lab leader*, and *lab helper*. Often, these roles are given based on the academic level of student, with more academically experienced TAs taking on the roles of course organizer, head TA, and lab leader, where as students with less experience (i.e. undergraduate students) take on the role of lab helper.

Lab sections are held once a week, and are facilitated in person by both a lab leader and lab helper. The responsibilities of a lab helper are supporting both the students and lab leader as they see fit. Examples of this may include setting up the classroom before class, or conducting small group conversations when students have questions about the material. The lab leader is responsible for facilitating the lab. This may involve giving a brief introduction and wrap up of lab content, as well as being able to answer questions and facilitate conversations among students about the lab assessment material. In addition, both must hold two hours of office hours each week and have grading responsibilities assigned throughout the semester.

Head TA responsibilities can generally be categorized into the following categories: Administrative and pedagogical. Administrative responsibilities include the organization and

distribution of TA responsibilities throughout the semester. It is imperative that the head TA and instructor clearly communicate expectations with each other to establish exactly how rules and responsibilities are assigned to TAs. This includes distributing grading assignments and deadlines to both lab helpers and leaders, weekly. The head TA also makes sure that all TAs complete grading within a week and spot check the grading accuracy and quality of all written feedback given. Other administrative duties include reminding other TAs about bi-weekly payroll deadlines and ensure TAs are working their allotted hours per week (and not more). Pedagogically, head TAs are responsible for creating or reviewing answer keys and grading rubrics for homework and lab assessments as the instructor sees fit. Each head TA is also assigned to instruct one lab section during the semester. Before becoming a Head TA, there is additional training that specializes head TAs in their administrative responsibilities.

The course organizer is expected to work across each section of STA199, instead of working with a single instructor. Their responsibilities include creating rubrics for and working through homework and lab assignments. Additionally the course organizer, along with the instructor, answers real time questions virtually during labs asked by lab leaders. Questions often range from content related to technical questions about GitHub and R. Finally, the course organizer is responsible for handling all requested assignment extensions from students. This includes filing away student exemptions, providing extensions for extreme circumstances, and enforcing the late work policy outlined in the syllabus when necessary.

We typically have one course organizer, one head TA, six lab leaders, and six lab helpers. Although this is our proposed team structure, we want to emphasize that there is great flexibility in coming up with how a teaching team is built and operates for an introductory to data science course. Among any team, we encourage a system designed to alleviate

the grading responsibilities of a large class from a single individual, dispersed into many among the team. When grading, it is suggested that each individual is properly trained on how to grade assessments, and expectations on how to provide feedback are clear. Unclear feedback given has often been a point of contention from students in the past. For each assignment, we recommend having one member of the team grade one problem across the entirety of the class roster. This helps adhere to more consistent grading, and encourages more grading questions to be asked as they arise, earlier in the grading process.

Through our experiences, it has been imperative that everyone within the team is communicating with each other. A team with many different roles poses risk for the instructor to be unaware of how or what decisions are being enacted at the grading and lab levels of the course. Thus, it is recommended that the instructor trains everyone on the teaching team to use a communication system that allows every member to communicate any questions they may have, or decisions they make, to the entire team. In the past, we have used the software *Slack*, with appropriately named channels such as *grading-questions*, where TAs can post grading examples and questions about grading so the instructor can clearly respond with their expectations. Further, it should be noted that the head TA should not be treated as a “bridge of communication” for the instructor to the rest of the teaching team. It is critical that the instructor is in consistent contact with all members of the teaching team in making sure all lab leaders and helpers understand the course content, how to grade, and know what’s expected of them in their assigned role. We recommend holding a weekly meeting with all members of the teaching team to ensure this. When members are unable to come, it is an expectation that they watch a recorded video of the meeting and reach out if they have any questions about what was discussed.

4 Technology

We use R, R-studio, and GitHub to create interactive lessons, and assign pre-created assessments to individual or groups of students. In the following sections, we detail the computing infrastructure needed to do so. This includes details and examples with R, R-studio, and GitHub, from the instructor's perspective, to set up lectures, AEs, labs and homework. First, we justify our decision to use GitHub in STA199 before detailing necessary steps and student information needed so creation can take place.

4.1 Why GitHub

We choose to use GitHub for STA199 because it easily and efficiently allows us to create and administer assessments and interactive lessons to a large class size in individual GitHub repositories. Further, within these repositories, instructors can provide template code and other resources necessary to best facilitate an interactive lesson where students can code together with the instructor during lecture. Additionally, teaching through GitHub provides the ability to re-use what you currently create as a template for subsequent semesters. Thus, this system both rewards and encourages time investment into the creation of assessments and lesson plans. Moreover, exposing students to version control software early in their academic career helps them develop good habits as researchers and emphasize the importance of reproducibility in science. These are just some of the many benefits of using GitHub in an introductory data science course. To take advantage of these benefits as an instructor, we first need students to create their own GitHub account.

4.2 Setting up a GitHub account

To participate in interactive lessons and assessments, students must set up a GitHub account. This is done on the first day of class, and often, students are given time during class to sign up. Following tips from “Happy Git with R” (?), we suggest students do the following when creating their name:

- Incorporate their actual name
- Reuse their username from other contexts if you can
- Pick a username they will be comfortable revealing to a future boss
- Be as unique as possible in as few characters as possible. Shorter is better than longer
- Avoid words with special meaning in programming (i.e., NA)

Once students create their account, we suggest getting this information from them in a survey. This normally can be done through your learning management system. It is critical to reiterate to students that spelling and capitalization matter when they provide their GitHub username. We suggest asking the question as follows:

(should we include the questions in a highlighted text box of some sort?)

What is your GitHub username?

Answer this question by ONLY typing your GitHub name and nothing else.

Make sure you triple-check the spelling so I don’t add random strangers to our course organization on GitHub

Once this information is collected, it must be exported and stored as a `.csv` file to be read into R later. We recommend having the following structure of data collected from the

survey.

last_name	first_name	github_username
-----------	------------	-----------------

Fig 2: Example Data Collection

If you, as the instructor, do not have a GitHub account, you will need to create one as well. This student information will be used to enroll students in your created GitHub organization for your course.

4.3 GitHub organization

GitHub organizations are shared accounts where instructors and students can collaborate across many projects at once. Using your GitHub account, you can create a new GitHub organization by clicking on your profile icon in the upper right hand corner, clicking *Settings*, *Access*, *New Organization*. It's suggested to name this organization the name of your class and the current semester you teaching in (i.e., sta199-s23). Once your organization is created, you can use packages within R and R-studio to invite students to enroll. Once enrolled, we can create and efficiently distribute assessments and live coding activities for them to access on their own laptop device.

5 Why R & R-Studio

R is a statistical programming language for computing and modeling while R-studio is an integrated development environment for R (?). We choose these resources because they are freely available to download and are in high demand in different data science job op-

portunities (cite?). In addition, this programming language has a comprehensive library of packages that allows students to create data visualizations and seamlessly analyze data, aligning with our learning objectives for the course. Further, this software is compatible with online versions that students can access without having to go through a local install-ment if they do not wish or are unable. There exists online versions that can be set up for students, with capabilities of pre-populating their software with the appropriate packages needed for the semester. Addressing these hurdles for students early is advantageous any instructor who is working with a large amount of students that are inexperienced with coding.

5.1 Setting Up Students R & R-Studio

In an introductory course, it is recommended to minimize student frustration and distraction through the use of pre-packaged computational infrastructure (?). Per this recommendation, STA199 has students use R and R-studio through a *Duke Container*. Duke containers provides instructors at Duke university the opportunity to facilitate the use of different software, such as R and R-studio, through an online container instead of needing students to locally download both programs. Additionally, instructors have the ability to manage and install packages students will need for the semester, helping provide a neat and well organized starting experience with a new statistical language. **(insert discussion on how this is done).**

(Transition to R-studio cloud discussion? What alternatives can instructors use to imitate Duke containers if this is something they do not have access to at their university.)

5.2 Using R & R-Studio to Manage GitHub organization

To invite students into your GitHub organization, we will use a myriad of functions from the `gh_class` package. (insert discuss on what the `gh-class` package is).

To start, we suggest creating a separate GitHub repository that will contain all code used to invite students to your organization and to create individual student repositories. Within a document in this new repository, we must first define our GitHub organization as a character string to be used in later function. We suggest calling this `org`.

```
org <- "sta-199-s23"
```

The following code below invites students to your GitHub organization using your defined `org` and the GitHub usernames collected from your students in a previous survey (see section Setting up a GitHub account). For the example below, we chose to name the student survey data as “roster.csv”.

```
# invite students -----

roster <- read_csv("roster.csv")

org_invite(
  org = org,
  user = roster$github_name
)
```

Once the invites are sent, students will need to accept the invitation to the organization.

Streamlining your course through R, R-studio, and GitHub greatly alleviates common

challenges that arise when working with a large class size, such as activity and assessment distribution. This includes the creation and distribution of in-class AEs, lectures, labs, and homework assignments.

In the following sections, we discuss pedagogy used within our course, strategies on how to implement such pedagogy in the classroom, and detail the creation and distribution of AEs and lab assessments using R, R-studio, and GitHub during a typical week in the semester (Figure 2).

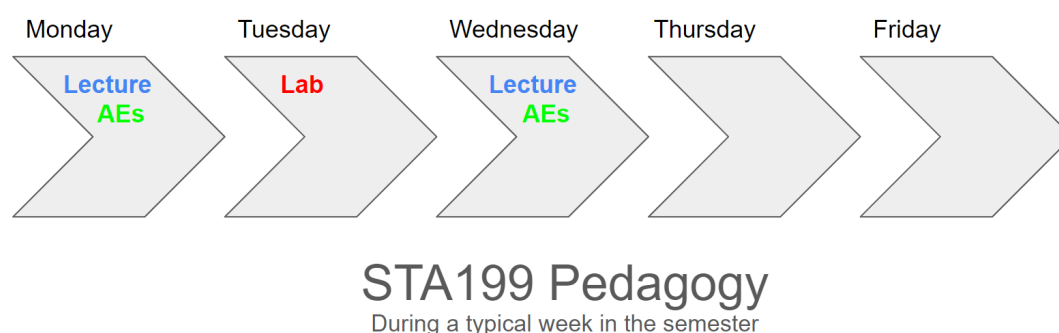


Fig 2: Model of pedagogical choices used in STA199 for a typical week

6 Pedagogy

In STA199, we have chosen a combination of teaching methods, interactive activities, and learning assessments to help prepare introductory data science students the tools they need to be successful outside of university or in future coursework. Our pedagogy includes a combination of lectures, facilitating in-class AEs, and running a lab, to provide students an opportunity to perform what they’ve learned.

6.1 Lecture & Application Exercises

During a typical week in the semester, class is held twice a week. Class is a combination of lecture and interactive AEs. The amount of time dedicated to lecture and AEs can and should vary by instructor. Typically, we allot the first 15 minutes of class for lecture. During this time, students are introduced to new content that will be re-addressed in the AE. To help maximize engagement of such a large class size, we allot more class time to individual and team-based AEs where students are asked to code on their own, together, or along with the instructor. This is followed by a roughly 5-10 minute allotment for a wrap up lecture, summarizing the highlights of AE.

Application exercises are structured interactive lessons that allow students to learn through doing in class. Students are expected to bring their laptops to class to participate in AEs. This expectation is made clear prior to the first day of teaching. The purpose of these exercises are to give students the opportunity to practice apply the statistical concepts and code introduced in the prepare videos, readings, and anything introduced during lecture.

When first starting to create an AE, we suggest streamlining the process through a GitHub template repository cloned into R-studio as a project. We choose to use *Quarto*, within R-studio, to create AEs. This choice is intentional, to provide students the opportunity to practice writing code and answering questions in a reproducible format supported by R and R-studio. When designing questions for an AE, we suggest creating a mix of coding and concept questions (e.g., fill in the code blanks, short response questions) that encourage students to follow along with instructor demonstrations, and also provides students an opportunity to answer questions on their own. This format has largely been accepted and appreciated by students: “I really enjoy the AE’s and that you take the time to walk us through the code and answer questions; I like how the ae gives us a chance to practice

the skills on our own after class as well.” Additionally, we suggest that questions are scaffold, meaning that if students fall behind or type incorrect code when trying to initially follow along, that they have access to correct code that grants them the ability to continue engaging for the remaining AE. This is especially important at the beginning of the semester to minimize students’ frustration around a new coding language. This may include having answers to questions in a separate Quarto document that students have access to, or within the same document as the AE. Moreover, we suggest clearly labeling where students are expected to type out answers in text or code throughout the exercise to further streamline their involvement and continue to minimize frustration.

An example on an AE used to help teach **(insert concept)** can be found here: **(insert example AE)**

We highly recommended designing lessons with built in time to have students work on their own. The amount of time can largely depend on the question being asked, and the teaching style of the instructor. We have found that multiple 3-5 minute blocks for students to answer questions without the instructor’s guidance works well. This has been especially effective if the code is then created together as a class, built off student responses, to provide immediate feedback to students before moving on to the next set of material or questions.

Once created, we can distribute this AEs to students within the class GitHub organization by using the function `org_create_assignment` from the `gh_class` package. Example code can be seen below.

```
assignment <- "AE-5" #name of AE repository you want to distribute

student_name <- org_members(org = org, include_admins = FALSE) #name of students to c
```

```
org_create_assignment(  
    org = org,  
    repo = paste0(assignment, "-", student_name),  
    user = student_name,  
    source_repo = paste0(org, "/", assignment)  
)
```

It should be noted that, with such a large class size, **(time out error discussion?)**.

(Code to create “second batch of AEs” if needed?)

Hands on interactive AEs are the backbone of learning in STA199. We have found that providing students with the opportunity to code in real-time during class keeps them more engaged and eager to continue learning about the material, versus simply lecturing for a 75-minute period. This strategy has received positive feedback from students with varying degrees of coding and statistics experience.

As highlighted in previous sections this system, streamlined through GitHub, encourages time investment into the creation of AEs. Investing time into AE creation provides a strong environment for current students to learn, and sets a foundation for AEs to be retooled instead of recreated in subsequent semesters. Future renditions of this course can be modified from previous GitHub organizations efficiently, saving instructors valuable time and energy before and during the semester.

AEs are designed to be a low stakes assessment where students can earn full credit by completing in class. Typically, we make AEs worth completion points that total to be 5% of a student’s end of semester grade. For students that do not show up to class, we

make AEs due three days from when the AE was assigned. At the end of the semester, if student's have completed 80% of AEs, they earn a 100% for their grade. Students turn in AEs by pushing their answers within the AE document up to their GitHub repository. There have been mixed responses from both students and instructors on assigning a grade to AEs. Disadvantages include not currently having an efficient way to implement a hard deadline that students must adhere to, without removing push access from their AE GitHub repository all-together. Further, if an instructor does not finish an AE in class, students who did not attend class can be easily confused on what's expected of them to complete. We suggest tailoring the decision of grading of AEs to your individual course as you see fit.

6.2 Labs

Labs for STA199 meet once a week for 75 minutes, and are facilitated by lab leaders and lab helpers. The purpose of labs are to allow students to apply concepts found in prepare material, lecture, and AEs, to various data analysis scenarios. Lab section sizes are kept to around 30 students so each have more of an opportunity to converse with each other, the lab leader, and lab helper, when working through the lab assessment.

Much like AEs, the instructor or head TA creates and clones a lab GitHub repository to all students that contains any information and any other intangibles (like a data set) by changing the `assignment` object in the above R code (see section Lecture & Application Exercises) to be the name of original lab GitHub repository. In the lab GitHub repository, we choose to create a starter document with places for students to write code under specific question numbers. The actual lab assessment that this document corresponds to can be found on the class website that is referenced during lab time.

Roughly one-third of the way into the semester, students are assigned to groups to complete

a class project. We suggest strategically assigning groups based on a collection of the following information

- Declared or Intended Major
- Year in School
- Suggested times they work on school assignments

From our experience, groups who have significantly different years in school across students have more friction, and tend to work less well together than students that have more similar years in school. Additionally, we highly suggest pairing up students that share common interests using their intended or declared major in school. The class project assigned is an open ended research project where the group collectively decides on a project topic. We have found that students can feel disengaged or left out of the group if they have different interests than the others, and their interests are not reflected when picking a project topic. When groups are assigned, each group during a lab are tasked to come up with a team name. This team name will be the name used to create their team repositories for the remaining labs. We add this information as follows to the roster document.

team_name	last_name	first_name	github_username
-----------	-----------	------------	-----------------

From this point forward in the semester, we choose to have all labs be completed as group work. Introducing group work during labs and through a class project can help students learn from different perspectives, practice their communication skills, and improve their problem solving skills in the context of statistics and data science.

We can use these new data and change the code found in Lecture & Application Exercises

section to create repositories for each team.

```
teams <- read_csv("roster.csv")

org_create_assignment(

  org = org, # name of your organization

  repo = paste0(assignment, "-", teams$team), #name of the repo you are cloning and t

  user = teams$github_username, #adding individual students to team repo

  team = teams$team_name, #team (or group) they are added to

  source_repo = paste0(org, "/", assignment) #naming repo

)
```

The new expectation is that once groups are formed, one lab assessment will be turned in for each group instead of each individual. This further helps lessen the grading responsibilities to those that are assigned it.

After groups are formed, we give a set of recommendations for groups to help promote a successful and positive group dynamic. This includes:

- Establish a clear line of communication with all members
- Share ideas. Let your voice be heard.
- Teach each other.
- Do not approach group work as a bunch of individual assignments.

To further ensure positive group dynamic, we initiate three peer review surveys. These peer review surveys provide insight into each group's dynamic and may inform the teaching team of issues that may need to be addressed. Questions within each survey range from

having only having instructor only visibility, to being shared with all team members to promote productive conversations. An example question includes: *Estimate the percentage of the total amount of work/effort done by each member, including yourself. Be sure your percentages sum to 100%!* We suggest adding additional questions as deemed necessary to help best understand how everyone is working together within a group.

A new instructor of data science, or one with an increasing class size should think critically if and how they want to implement group work in their classroom. In our experience, merge conflicts, the number of groups, and group formation have been the most difficult aspects of facilitating group work. Students, especially those newer to coding, are extremely hesitant to create merge conflicts. Despite a lab dedicated to the creation and fixing of merge conflicts, students often express frustration and feel as though they are doing something wrong when merge conflicts occur. Some groups choose to collaborate using other forms of technology (e.g., Google documents) before committing and pushing their finished work onto GitHub. We highly advise against this, and try to incentivise students by emphasizing the practical importance of learning how to work through merge conflicts. Secondly, the sheer number of groups created from 179 students creates an additional time investment for all members of the teaching team. This includes making sure all merge conflicts can be fixed accordingly, and helping facilitate an appropriate working group dynamic among all groups. Finally, there have been mixed strategies on how to form groups that ultimately have yielded similar results. Typically groups are assigned by the instructor using the aforementioned questions above, while others have allowed students to select their owns groups. Both strategies have resulted in both positive and fractured group dynamics. **(address literature on group work + group formation?)**. We highly suggest that you consider additional measures and adjust group formation as you

sit fit for your course.

There are advantages and disadvantages of having lab leaders and helpers facilitate the lab. Advantages includes having students learning through additional perspectives, while also providing a potentially more inviting atmosphere for questions to be asked about the material. However, disadvantages may include inconsistencies from what is shared in lab vs lecture. It is critical that lab leaders and helpers are trained to both understand and explain concepts consistently to what is being taught during lecture and through the AEs. Additionally we recommend setting the expectation that lab leaders and helpers become familiar with the AE content before facilitating the labs so they can refer students back to resource they are familiar with.

7 Discussion

8 Conclusion

9 Resources

References

Kokkelenberg, E. C., Dillon, M. & Christy, S. M. (2008), ‘The effects of class size on student grades at a public university’, *Economics of Education Review* **27**(2), 221–233.

URL: <https://www.sciencedirect.com/science/article/pii/S0272775707000271>

Schwab-McCoy, A., Baker, C. M. & Gasper, R. E. (2021), ‘Data science in 2020: Computing, curricula, and challenges for the next 10 years’, *Journal of Statistics and Data Science*

Education **29**(sup1), S40–S50.

URL: <https://doi.org/10.1080/10691898.2020.1851159>

Schweser, K. (2023), ‘Our philosophy’.

URL: <https://www.schweser.com/about-kaplan/philosophy>