

# DS-100: Data Speak Louder than Words

## Fall 2025 Syllabus

### Course Description

In this course we will introduce you to three fundamental perspectives for reasoning with data: critical thinking, inferential thinking, and computational thinking. All three of these perspectives are integral to the data-driven research processes that are common in data science, thus allowing you to learn and practice how you can make and test hypotheses, and construct or deconstruct arguments that are rooted in data.

We will first use public data sets (both curated or scraped) focused on socially-relevant themes (e.g., public health, education, and environment) to model and understand real-world phenomena. We will focus on using model summarization, data visualization, and model-based simulations to interpret and communicate our understanding of these real-world phenomena as well as the potential for bringing these derived models to bear on real-world questions and applications (e.g., comparing different policies).

Particular emphasis will be placed on exposing you to and developing your appreciation for the principles underlying data mining and machine learning methods, including regression, classification and clustering, and the statistical concepts of measurement error and prediction. We will teach you critical concepts and skills in computer programming (Python), linear regression, and statistical inference. We will also delve into dilemmas surrounding data analysis such as balancing individual privacy and social utility.

This course uses a learning model where students use large language models (LLMs), commonly referred to as AI or GenAI, as learning partners for concept exploration while developing authentic data science competency through guided practice and individual verification.

### Hub Learning Outcomes

#### Social Inquiry I (SO1)

**Learning Outcome #1:** Students will identify and apply major concepts used in the social sciences to explain individual and collective human behavior including, for example, the workings of social groups, institutions, networks, and the role of the individual in them.

We will employ hands-on analysis of real-world datasets, including curated economic data, data scraped from digital collections, social networks, and more. In this context, the course will expose you to social and legal issues surrounding data analysis, including issues of privacy and data ownership, and will highlight the many ways in which data could be used (or misused).

In this course we will be looking at data from multiple vantage points. For example, by looking at data characterizing COVID-19 infections, hospitalizations, deaths, vaccinations, we will be able to differentiate between phenomena (e.g., correlations) identified at the macro scale (federal and state) versus those identified at the micro scale (cities and communities) and draw conclusions or make statements supported with evidence from data (e.g., impact of socioeconomic background).

We will encourage you to apply what you learn on societally-relevant case studies of your choice (e.g., case studies similar to those presented in <https://www.callingbullshit.org/>) by applying the tools and techniques covered in class to analyze data sets in order to support or debunk hypotheses.

### **Digital/Multimedia Expression (DME)**

**Learning Outcome #1:** Students will be able to craft and deliver responsible, considered, and well-structured arguments using media and modes of expression appropriate to the situation.

We will use real data to understand relationships and patterns while also introducing critical concepts and skills in computer programming and statistical inference. In order to build your arguments, you will use multimodal data analysis and visualization in ways that are appropriate to the task at hand. This will include:

- Generation and interpretation of scatter plots, histograms, bar charts, and box plots
- Making predictions using simple regression
- Characterizing data quality and communicating associated uncertainties
- Establishing confidence in reproducible predictions
- Reaching defensible conclusions about real-world questions

These skills will be taught and evaluated through learning logs and discussion-section/in-class activities, as well as in checkpoints and projects

**Learning Outcome #2:** Students will be able to demonstrate an understanding of the capabilities of various communication technologies and be able to use these technologies ethically and effectively.

As part of DS-100, we will introduce you to multiple forms of data visualization and presentation, including histograms, scatterplots, word clouds, heat maps, infographics, etc. Each one of these forms of communication can be particularly effective (or even misleading) in certain settings. For example, the choice of different scales (e.g., absolute vs relative change) on an axis could over or under-emphasize particular conclusions from the data.

Given the multitude of sources from which the data is collected, you will be exposed to proper ways of handling the data. For example, to preserve the privacy of individuals or communities in a large data set, and be introduced to the use of randomization techniques (blurring the data). As another example, to deal with the scale of data it may be necessary to only consider/analyze a subset of all observations. In that context, we will introduce you to various ways in which selection bias may influence conclusions you may be able to reach with implications on reproducibility.

**Learning Outcome #3:** Students will be able to demonstrate an understanding of the fundamentals of visual communication, such as principles governing design, time-based and interactive media, and the audio-visual representation of qualitative and quantitative data.

We will teach you how to use Python to organize and manipulate data in tables, and to visualize data effectively. Furthermore, you will be able to use computation to help your data tell a story through fundamental principles and methods of data visualization. The data used throughout this course will include longitudinal data (time series over long-time scales), geospatial data (data overlaid on apps), or both. These modalities will offer you different ways to interact with the data. For example, with time series data, you will be able to develop animations to show how phenomena or inferences may evolve over time. As another example, with geospatial data sets, you will be able to develop animations or heat maps that may project different messages/narratives based on the level of aggregation (e.g., achieved by zooming in and out).

In all of the above learning outcomes, we note that some of your work products will be in the form of multimedia reports, in which data visualization is coupled with narratives or video clips. For example, in a report on deforestation due to climate change, you may add audio or video clips to demonstrate change over time. You may also include your own narration to supplement and/or add texture to the graphs, heatmaps, etc.

### **Research and Information Literacy (RIL) Learning Outcomes**

We will teach you critical concepts and skills in computer programming and statistical inference, in conjunction with hands-on analysis of real-world datasets, including economic data, document collections, geographical data, and social networks. In discussion sections, you will work in small teams, working under the supervision of the teaching fellow to frame a question or test a hypothesis using a set of potential data sources. The key phases of that process are the exploration and identification of relevant data sets, the formulation and reformulation of the questions based on the identified data, the development of a set of data processing/analytics steps leading to an answer, and the interpretation and/or validation of the answer. To a large extent, going through these phases mirrors the six steps of the data science research process.

## Books & Tools

- **Core Text:** *Inferential Thinking* by Ani Adhikari and John DeNero, with contributions by David Wagner and Henry Milner: [inferentialthinking.com](http://inferentialthinking.com)
- **Optional Resource:** *Calling Bullshit* ([callingbullshit.org](http://callingbullshit.org)) for case studies.
- **Programming Environment:** Python with industry-standard libraries (numpy, pandas). We will help set up your environment so everyone can access the same tools.

## Course Platforms

- **Blackboard Ultra** – Primary hub for announcements, weekly schedule, policies, and resources.
- **Gradescope** – Submit and receive feedback on learning logs, in-class work, checkpoints, and projects; also for regrade requests.
- **Piazza** – Questions, discussion, and project info. Use Piazza instead of emailing the teaching team.
- **TerrierGPT** – Our GenAI exploration platform ([terriergpt.bu.edu](http://terriergpt.bu.edu)) with access to OpenAI, Anthropic, Amazon, and open-source models.
  - Provided credits should be sufficient. If not, budget up to **\$60** for a commercial subscription. If this would be a hardship, contact the instructor—alternatives will be arranged. *No student is penalized for inability to pay.*

## Assignments & Grading

### Grade Distribution

Category	Weight
GenAI Exploration Portfolio (weekly learning logs)	15%
Understanding Demonstrations & In-Class Activities	20%
Course Checkpoints (3, in-person, closed-book)	25%
Project (Proposal 5%, Milestone 5%, Final 30%)	40%

### How Grading Works

- **Completion-based (Check / Check+ / Missing):** Used for learning logs and in-class activities. Lowest ~10% dropped automatically. Completing 90% = B+; sustained Check+ work can raise to A range.

- **Rubric / Points:** Used for checkpoints and projects. Checkpoints are closed-book, in person, no GenAI allowed. Projects graded on clarity, accuracy, reproducibility, ethics, and communication.

### Timing & Late Work

- Logs: due before Tuesday class; no late work (auto-drop covers exceptions).
- In-class activities: must be done in class; no make-ups (auto-drop applies).
- Checkpoints: scheduled in advance with 1 week study materials. No GenAI tools. Make-ups only with documented, university-approved reasons.
- Projects: Proposal and Final Deliverable may be up to 48h late (−10%); not accepted after. Milestone is completion-based and must be on time.

## How to Succeed

1. **Engage actively with learning logs.** Use GenAI tools to explore ideas, then reflect in your own words.
2. **Engage consistently in class.** Our model depends on in-class verification and discussion. Attendance isn't graded separately, but engagement is.
3. **Come prepared for Friday sections.** You'll demonstrate your understanding individually—practice concepts beforehand.
4. **Use office hours.** We welcome you for questions about concepts, projects, or GenAI strategies.
5. **Ask questions on Piazza.** It's the fastest way to get help from peers and staff.
6. **Build community.** Use ERC tutoring, writing support, and peer study groups (outside individual assessments).

## Integrity & GenAI Policy

We follow BU's [Academic Conduct Code](#) and the [CDS GAIA Policy](#), adapted for this course. Discussing ideas is encouraged; copying is not. Always cite collaborators, data, libraries, and GenAI use.

### GenAI Zones

- **Green (Encouraged):** Learning logs, concept exploration, debugging help, project brainstorming.

- **Yellow (Allowed with care):** Project execution (analysis must be student-driven), writing drafts, study groups.
- **Red (Prohibited):** Checkpoints, in-class demonstrations, or any assessment explicitly marked “individual verification.”

Violations—especially GenAI use in the Red Zone—will be treated as academic misconduct.

**Tone of conduct:** Disagreement is welcome; disrespect is not. We aim to model the community we want for our University and industry.

## Support & Accessibility

- **Disability accommodations:** Contact the Office for Disability Services (617-353-3658, [access@bu.edu](mailto:access@bu.edu)). Share your letter privately with the instructor.
- **Academic & well-being support:** ERC tutoring, writing support, and BU Student Wellbeing resources (see Blackboard).

## Regrades

Submit requests via Gradescope within **7 days** of score release, identifying the specific error. Scores may go up, down, or remain unchanged; decisions are final.