# Stat 470/670: Exploratory Data Analysis

Meeting time: Tuesdays and Thursdays, 12:45-2pmMeeting location: WH 111Website: jfukuyama.github.io/teaching/stat670jfukuyama.github.io/teaching/stat670Instructor: Prof. Julia Fukuyamajfukuyama at iu dot eduOffice hours: Tuesdays 2:30-3:30pm, Wednesdays 1-2pm Swain East 225yyu3 at iu dot edu

## **Course Overview**

Office hours: TBA

Graphical and modeling techniques for exploring data, with an emphasis on visualization, interpretation, and clear communication of findings. Use of modern software tools for data manipulation and visualization. Connections to traditional statistical methods.

### Textbooks

We will be drawing heavily on Cleveland's *Visualizing Data* and Hadley Wickham's *ggplot2: Elegant Graphics for Data Analysis*. Both of these are available online through the IU library.

Also useful will be *R* for Data Science by Wickham and Grolemund, available online at http://r4ds.had.co.nz.

Readings and notes for topics not covered in the textbooks will be posted to the course website and to canvas.

### **Class Structure**

Classes will be a combination of lecture and tool demonstration. It will generally be helpful for you to have an R session open to follow along with the code. Slides or notes, with R code, will be posted to the class website before each lecture.

We will also have a regular set of in-class exercises where you will try out some of the sorts of analyses we have recently covered in class, think through why we would want to perform such an analysis, and think through what the implications of the analysis might be.

### Assessment

Grades will be assigned based on:

In-class exercises and presentations, worth 30% of the grade. By the second week of class, you should identify a dataset. A couple of times over the course of the semester, we will take some time in class for you to (1) apply some of the techniques we have talked about to your dataset, and (2) think about why someone would want to do that analysis and what

the implications of that analysis would be. You will be required to upload your analysis and a response to the question of "why would you do this analysis" to canvas. You will also be required to present your results at least one time over the course of the semeste.

- Mini project, worth 30% of the grade. This will involve more substantial data analysis and a writeup.
- Final project, 40% of the grade, on a dataset and question of your choice.

There will be no final exam; the last responsibility for the course will be the report for the final project due on the last day of class.

All the assignments will be graded on how well the material is presented in addition to accuracy. This means there should be no extraneous material, plots should be readable, and text and figures should be formatted nicely.

## Topics

There are two categories of topics: *what* to do and *how* to do it. In the *what* to do category, we will cover:

- Univariate data: measures of center and spread, transformations, visualization.
- Bivariate data: Simple regression, curve fitting,
- Trivariate/Hypervariate data: Multiple regression, model selection, principal components.
- Binary responses: Logistic regression, residuals.
- Categorical data: Contingency tables, correspondence analysis.
- Distance data: Multi-dimensional scaling, non-linear dimensionality reduction.
- Graph data: Descriptive statistics, spectral methods, visualization.
- Dangers of EDA and remedies: Multiple comparisons, data splitting, cross validation.
- Other topics according to time and interest.

In the *how* to do it category, we will cover

- ggplot2 for plotting.
- tidy-verse methods for data wrangling.

By the end of the course, you should feel comfortable using R to visualize and model many kinds of data. Given a dateset, you should be able to visualize the data, generate hypotheses about the relationships among the variables, investigate those hypotheses, and communicate your results.

## **Course Policies**

### Late Policy

For the mini project, late work will be penalized at 10% per 24 hours. Final projects cannot be turned in late. Special accommodations for the mini project may be granted if you ask very early.

### **Academic Integrity**

You are expected to abide by the guidelines of the IU Code of Student Rights, Responsibilities, and Conduct (http://studentcode.iu.edu/responsibilities/academic-misconduct.html) regarding cheating and plagiarism. Any ideas or materials taken from another source must be fully acknowledged and cited.

### **Disability Accommodation**

Please contact me if you require assistance or academic accommodations for a disability. You should establish your eligibility for disability support services through the Office of Disability Services for Students in Wells Library W302, 812-855-7578.

### Incomplete

Incomplete grades will be assigned in accordance with the incomplete policy described in the graduate school bulletin https://bulletins.iu.edu/iu/gradschool/2024-2025/policies/grading.shtml.