

# **INTEGRATING ADVANCED DATA ANALYSIS TECHNIQUES AND COMPUTATIONAL SKILLS INTO EARLY COLLEGE CURRICULA**

JAMIS J. PERRETT, PH.D.

DEPARTMENT OF STATISTICS

BRIGHAM YOUNG UNIVERSITY

# OUTLINE

- FOCUS ON STUDENT PREPARATION
- INTRO STATS EXAMPLE
- DATA SCIENCE ECOSYSTEM

FOCUS ON STUDENT  
PREPARATION

# PREPARING STUDENT OF TODAY FOR THE CHALLENGES OF TOMORROW

## Thought Questions

- Even if it is relevant today, will it still be relevant when students graduate?
- What will students remember 5 years from now?
- In which ways will students use statistics in their future careers?
- What are employers looking for?

Employers have indicated a need for data science students to have the following preparation:

- Python – coding languages
- Cloud Tools
  - Amazon Web Services (AWS)
  - Snowflake
  - DataBricks
- SQL – querying data
- Dashboarding – visualizing data

# WHERE DO STUDENTS GET THEIR TECHNOLOGY SKILLS?



Many undergrads have no formal programming skills



CS & IS departments offer dedicated computing courses



How does the computing get implemented in the statistics course?

# SOME IDEAS:

- For service courses, talk to the other departments, regularly. Collaborate on the course development
- For Intro courses,
  - Decide to use coding or point & click, depending on the need.
- For higher-level courses, use appropriate tools and stay current. Technology changes daily (DBMS list at 379 as of 1/5/2025, db-engines.com/en/ranking)
- Teach students how to use generative AI to aid in their learning process
  - Generative AI:
    - Maybe the best research assistant you will ever have
    - Can do the work of 25 programmers
  - Understand prompt engineering

# SOME IDEAS:

- IMPLEMENT PROGRAMMING IN STATISTICAL METHODS COURSES
  - USE R/PYTHON FOR COMPUTATIONS
  - INTRODUCE OTHER TOOLS LIKE RSTUDIO, VS CODE, ANACONDA, GITHUB, TABLEAU AS A MEANS OF FACILITATING THE COURSE
- USE MORE ADVANCED TECHNOLOGY IN INTRO COURSES
  - RSHINY APP OR WEBAPP TO FACILITATE/AUTOMATE CALCULATIONS
- OFFER “DATA SCIENCE” TOOLS COURSES THAT INTRODUCE MORE ADVANCED TOOLS
  - CLOUD, DATABASES, SQL, TABLEAU, ETC.
- TEACH STUDENTS HOW TO USE GENERATIVE AI TO AID IN THEIR LEARNING PROCESS
  - GENERATIVE AI: MAYBE THE BEST RESEARCH ASSISTANT YOU WILL EVER HAVE
  - UNDERSTAND PROMPT ENGINEERING

# INTRO STATS

BYU STAT 121 for various majors, no prerequisites,  
sometimes only statistics class students will take

AN INSIGHTFUL EXAMPLE

# WE START WITH REAL DATA

- BEGINNING OF SEMESTER SURVEY OF 34 QUESTIONS:
- COLLEGE (CATEGORICAL)
- SEMESTERS (INTEGER)
- IN A RELATIONSHIP (Y/N)
- MAC OR PC (BINARY)
- DIET (CATEGORICAL)
- HEIGHT AND MOTHER'S HEIGHT (CONTINUOUS NUMERIC)
- STREAMING SUBSCRIPTIONS (CATEGORICAL)

When the data is about you, you are more interested in solving the problem

# R-SHINY APP, POINT & CLICK TO AUTOMATE THE FORMULAS

The screenshot shows the Stat 121 Analysis Tool shiny app. The left sidebar lists categories: Exploratory Data Analysis, Normal Probability Calculator, Central Limit Theorem, Analysis for Means, Analysis For Proportions, and Regression. The main area is titled "Exploratory Data Analysis" and contains three sections: 1) Dataset Selection, 2) Select Variable, and 3) Graphical Exploratory Data Analysis. Section 1) Dataset Selection is active, showing options to use a preexisting dataset (selected) or upload one, and a dropdown menu for "Select dataset" containing "401K Returns". It also includes a description of the dataset and a sample size of 40. A checkbox for "Display Dataset" is present. Section 2) Select Variable shows a dropdown menu for selecting a variable to explore, with "Year" currently selected. Section 3) Graphical Exploratory Data Analysis is the final section.

## Measures of Spread

- Standard deviation ( $s$ ): “average distance from the mean”

$$s = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n - 1}}$$

- Inter-quartile Range = IQR =  $Q_3 - Q_1$  = spread of the middle 50% of the data.
- Range = Max - Min

## TOOLS CAN REINFORCE THE TOPICS COVERED IN THE COURSE (OUTLINE)

- Exploratory Data Analysis
- Normal Probability Calculator
- Central Limit Theorem
- Analysis for Means:
  - One Mean
  - Two Means
  - ANOVA
- Analysis for Proportions:
  - One Proportion
  - Two Proportions
  - Chi-Square
- Regression:
  - Simple Linear Regression
  - Multiple Linear Regression

Stat 121 Analysis Tool ≡

Exploratory Data Analysis

Normal Probability Calculator

Central Limit Theorem

Analysis for Means <

Analysis For Proportions <

Regression <

## Exploratory Data Analysis

### 1) Dataset Selection

**Data Selection**

Use Preexisting Dataset  
 Upload Your Own Dataset

**Select dataset:**

401K Returns

Description: Data on the annual return of a 401K portfolio.

Sample size: 40

Display Dataset

### 2) Select Variable

Please select the variable you wish to explore:

Year

### 3) Graphical Exploratory Data Analysis

Which plot would you like to draw?

Stat 121 Analysis Tool

Exploratory Data Analysis

Normal Probability Calculator

Central Limit Theorem

Analysis for Means

Analysis For Proportions

Regression

» Simple Linear Regression

» Multi Linear Regression

## Simple Linear Regression

1) Dataset Selection

Data Selection

Use Preexisting Dataset

Upload Your Own Dataset

Select Dataset

ACT Test Takers

Description: Data on 1000 random ACT test takers.

Sample size: 1000

Display Dataset

Select This Dataset

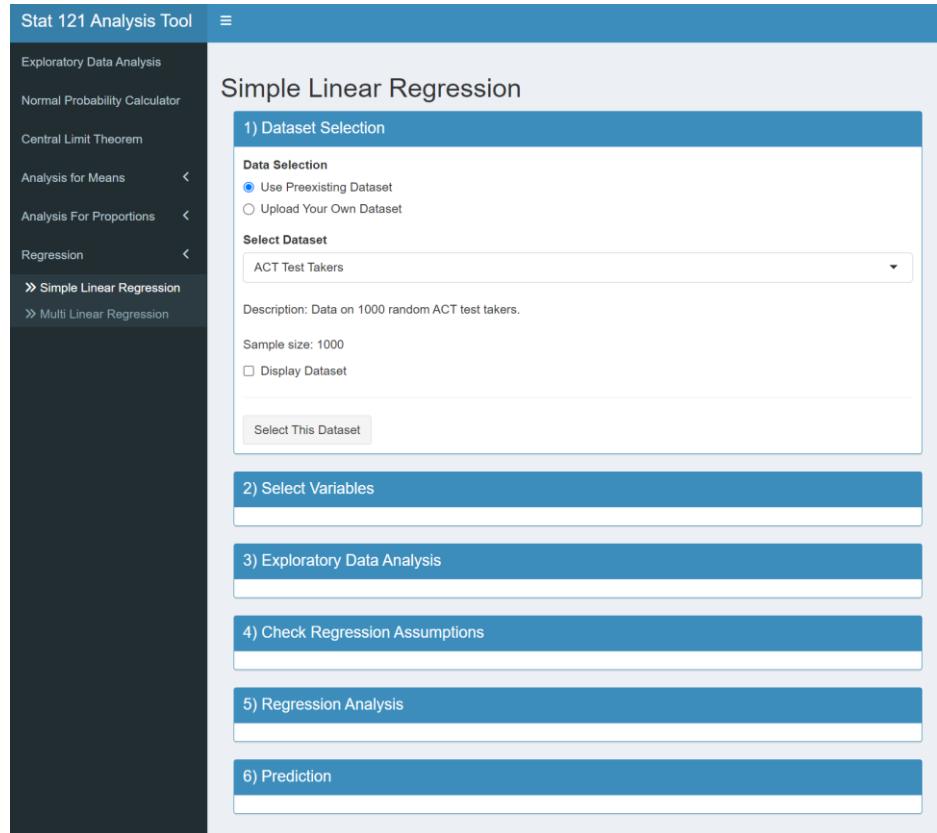
2) Select Variables

3) Exploratory Data Analysis

4) Check Regression Assumptions

5) Regression Analysis

6) Prediction



# TOOLS CAN REINFORCE PROPER STEPS FOR ANALYSIS

1. Data Selection
2. Select Variables
3. Exploratory Data Analysis
4. Check Regression Assumptions
5. Regression Analysis
6. Prediction

# ASA COLLEGE GAISE RECOMMENDATIONS

1. Teach statistics and data science as iterative processes of gleaning insights from data to inform evidence-based decisions.
2. Emphasize effective written and oral communication of results from data, with attention to the scope and limitations of conclusions.
3. Focus on conceptual understanding rather than algebraic manipulation and formulas.
4. Integrate real data with a context and purpose throughout the course. Select data that are meaningful and engaging to the students.
5. Encourage multivariable thinking.
6. Incorporate software/apps to explore concepts and work with data.
7. Emphasize responsible and ethical conduct in the collection and use of data and in their analysis.
8. Employ evidence-based pedagogies that actively engage students in the learning process.
9. Use a variety of formative and summative assessments to improve teaching and learning.
10. Implement a course design that uses inclusive strategies to foster a sense of belonging.

# DATA SCIENCE ECOSYSTEMS

Students planning to be data scientists,  
Prerequisites: coding and intro stats,  
Does not feed another class

AN INSIGHTFUL EXAMPLE

# BASIC OVERVIEW OF DIFFERENT DATA SCIENCE TOOLS

- Linux (SSH): Access remote data and work with permissions
- Containers: Using virtual environments
- **Databases and SQL: Understand how data are stored and how to access them**
- **AWS: Understand cloud environments**
- Dashboards: Interactive data visualizations on demand
- R & Python: Know how to code
- GitHub: Learn version control

# PREPARING STUDENT OF TODAY FOR THE CHALLENGES OF TOMORROW

Repeat Slide

## Thought Questions

- Even if it is relevant today, will it still be relevant when students graduate?
- What will students remember 5 years from now?
- In which ways will students use statistics in their future careers?
- What are employers looking for?

# THANK YOU