DSA 610 Redesign, Lecture 2 Outline

## Lecture Outline: Data Creation Methods and Their Risks/Benefits Duration: 50 minutes

## 1. Introduction to Data Creation Methods (5 minutes)

- **Objective:** Provide an overview of various data creation methods and their significance in data science.
- Content:
  - Definition of data creation and its role in generating information for analysis.
  - Overview of different data creation methods: structured forms, IoT devices, manual data entry, automated data generation, etc.
  - Importance of understanding the context, source, and quality of the data.

## 2. Data Creation via Structured Forms (10 minutes)

- **Objective:** Discuss structured data collection methods such as medical forms and surveys.
- Content:
  - Medical Forms & Surveys:
    - Standardized templates for collecting specific data points (e.g., patient information, health metrics).
    - Importance of standardized data for consistent analysis and interoperability across systems.
  - o Benefits:
    - High level of control over the data being collected.
    - Easier to analyze due to the structured format.
    - Enables comparison across different datasets due to standardization.
  - Risks/Cons:
    - Potential for incomplete or inaccurate data due to human error.
    - Standardization may limit the depth or richness of the data collected.
    - Privacy concerns, especially in sensitive fields like healthcare.
  - **Case Study Example:** Using electronic health records (EHR) to demonstrate the use of structured forms in a real-world scenario.

#### 3. Data Creation via IoT Devices (15 minutes)

- **Objective:** Explore how IoT devices contribute to data creation, along with their risks and benefits.
- Content:
  - Introduction to IoT:
    - Definition: Internet of Things (IoT) refers to interconnected devices that collect and transmit data automatically.
    - Examples: Wearables (e.g., fitness trackers), smart home devices, industrial sensors, healthcare monitors.
  - Benefits:
    - Continuous, real-time data collection.
    - Ability to collect large volumes of data (Big Data).
    - Enables predictive analytics and automation.
  - Risks/Cons:
    - Data privacy concerns, especially with personal or sensitive information.

- Security vulnerabilities: IoT devices can be targets for hacking.
- Data quality issues: Inconsistent or inaccurate data due to device malfunctions.
- **Case Study Example:** Use of IoT in smart cities for traffic management, showing both the benefits and the challenges of using IoT-generated data.

### 4. Automated Data Creation Methods (8 minutes)

- **Objective:** Discuss automated data generation methods such as web scraping and data logs.
- Content:
  - Web Scraping:
    - Collecting data from websites using automated scripts.
    - Common uses: Price tracking, sentiment analysis from social media.
  - Data Logs:
    - Automatically generated records of system events (e.g., server logs, transaction logs).
    - Importance for monitoring, auditing, and security.
  - o Benefits:
    - High efficiency and speed in data collection.
    - Can collect vast amounts of data from diverse sources.
    - Automation reduces the risk of human error.
  - Risks/Cons:
    - Legal and ethical concerns, especially with web scraping (e.g., violating terms of service).
    - Data quality issues if the data is not cleaned and validated.
    - Privacy concerns with sensitive data being logged or scraped.
  - **Example:** Automated web scraping for market analysis, discussing both the advantages and the potential legal risks.

#### 5. Manual Data Entry and Crowdsourced Data Creation (7 minutes)

- **Objective:** Examine the role of manual and crowdsourced data entry in data creation.
- Content:
  - Manual Data Entry:
    - Data input by individuals (e.g., entering survey responses, cataloging information).
    - Common in fields where automation is not feasible or cost-effective.

#### • Crowdsourced Data Creation:

- Data collected from a large group of people, often through platforms like Amazon Mechanical Turk or citizen science projects.
- Examples: Crowdsourced mapping (e.g., OpenStreetMap), public health data collection.
- Benefits:
  - Flexibility in data collection, especially for unique or nuanced datasets.
  - Crowdsourcing can gather diverse perspectives and large volumes of data quickly.
- **Risks/Cons**:
  - Higher risk of human error, leading to inaccurate or inconsistent data.
  - Quality control challenges with crowdsourced data.
  - Potential for bias in manually entered or crowdsourced data.

• **Example:** Crowdsourcing for disaster response data, discussing the pros and cons of relying on public input.

## 6. Ethical Considerations and Data Privacy (5 minutes)

- **Objective:** Discuss the ethical implications and privacy concerns related to data creation methods.
- Content:
  - Ethical Issues:
    - Informed consent, especially in medical data collection.
    - Transparency in data collection processes.
    - Avoiding bias in data collection and ensuring representativeness.
  - Data Privacy:
    - Importance of data anonymization and encryption, particularly with sensitive data.
    - Compliance with data protection regulations (e.g., GDPR, HIPAA).
  - **Real-World Examples:** Discuss cases where data privacy was breached and the consequences.

## 7. Conclusion & Q&A (5 minutes)

- **Objective:** Summarize key points and address any questions.
- Content:
  - Recap of different data creation methods and their respective benefits and risks.
  - Emphasis on the importance of understanding the source and quality of data.
  - Encourage critical thinking about ethical and privacy issues in data creation.
  - Open the floor for questions and discussion.

#### Key Takeaways

- Different data creation methods serve various purposes, each with unique benefits and risks.
- Understanding the context and quality of data is crucial for meaningful analysis.
- Ethical considerations and privacy concerns are paramount in the data creation process.

#### Resources:

Strauss, A., & Corbin, J. Basics of Qualitative Research: Techniques and Procedures for Developing Grounded Theory

Tabachnick, B. G., & Fidell, L. S. Using Multivariate Statistics Creswell, J. W., & Plano Clark, V. L. Designing and Conducting Mixed Methods Research

# Lecture Outline: Data Collection Methods, Data Discovery, and Experimental Design Duration: 50 minutes

#### 1. Introduction to Data Collection and Discovery (5 minutes)

- **Objective:** Provide an overview of the importance of data collection and discovery in the context of data analysis.
- Content:
  - Introduction to the role of data collection in data analysis.
  - Brief overview of data discovery as the process of identifying and understanding available data.

• The connection between data collection methods, data discovery, and experimental design.

## 2. Data Collection Methods (15 minutes)

- **Objective:** Discuss various data collection methods, their applications, and considerations.
- Content:
  - Primary Data Collection:
    - Surveys and Questionnaires:
      - Designing effective surveys, types of questions (open vs. closed-ended).
      - Benefits: Control over data collection, specific to the research question.
      - Challenges: Response bias, designing unbiased questions.
    - Interviews and Focus Groups:
      - Gathering qualitative insights through direct interaction.
      - Benefits: In-depth understanding, flexibility in responses.
      - Challenges: Time-consuming, potential interviewer bias.
    - Observational Studies:
      - Collecting data through direct observation without interference.
      - Benefits: Real-world context, natural behavior.
      - Challenges: Observer bias, lack of control over variables.
    - Experiments:
      - Collecting data under controlled conditions to test hypotheses.
      - Benefits: High internal validity, ability to establish causality.
      - Challenges: Limited external validity, ethical considerations.
  - Secondary Data Collection:
    - Existing Databases:
      - Using previously collected data (e.g., government databases, company records).
      - Benefits: Cost-effective, large datasets.
      - Challenges: Limited control over data quality, relevance to current research.
    - Web Scraping:
      - Collecting data from websites using automated scripts.
      - Benefits: Access to real-time data, large-scale data collection.
      - Challenges: Legal and ethical considerations, data quality issues.
  - **Case Study Example:** Compare and contrast primary and secondary data collection using a real-world example (e.g., health studies).

## 3. Data Discovery (10 minutes)

- **Objective:** Explore the process of data discovery, including finding and assessing relevant data for analysis.
- Content:
  - Definition of Data Discovery:
    - The process of identifying, gathering, and understanding data relevant to a specific analysis.
  - Data Sources and Repositories:
    - Public data repositories (e.g., UCI Machine Learning Repository, Kaggle).
    - Private datasets within organizations.
    - How to assess the relevance and quality of discovered data.

- Data Profiling:
  - Analyzing the structure, content, and quality of discovered data.
  - Techniques for assessing data quality (e.g., checking for missing values, outliers).
- Ethical Considerations:
  - Ensuring data privacy and compliance with regulations (e.g., GDPR).
  - Transparency in data usage and sourcing.
- **Practical Example:** Demonstration of data discovery and profiling using a sample dataset (e.g., public health data).

## 4. Introduction to Experimental Design (15 minutes)

- **Objective:** Introduce the principles of experimental design, focusing on the relationship between experimental design and data collection.
- Content:
  - **Definition and Purpose:** 
    - Designing experiments to test hypotheses and understand causality.
    - Importance of experimental design in ensuring valid and reliable results.
  - Key Concepts:
    - Variables: Independent, dependent, and confounding variables.
    - **Control Groups:** Establishing a baseline for comparison.
    - **Randomization:** Reducing bias by randomly assigning participants or conditions.
    - **Replication:** Ensuring that results can be reproduced.
  - Types of Experimental Designs:
    - Between-Subjects Design: Different groups are exposed to different conditions.
    - Within-Subjects Design: The same group is exposed to all conditions.
    - Factorial Design: Studying the effects of two or more variables simultaneously.
  - Challenges in Experimental Design:
    - Ensuring internal and external validity.
    - Managing ethical considerations, especially in human studies.
  - **Example:** Walkthrough of designing a simple experiment, such as testing the effect of different teaching methods on student performance.

#### 5. Integrating Data Collection, Discovery, and Experimental Design (5 minutes)

- **Objective:** Show how data collection, discovery, and experimental design are interconnected in the context of data analysis.
- Content:
  - The role of data discovery in identifying data for both observational and experimental studies.
  - How experimental design informs the data collection process, ensuring that collected data is relevant and reliable.
  - The iterative nature of these processes: Data discovery can inform experimental design, and experimental outcomes can guide further data collection.

#### 6. Conclusion & Q&A (5 minutes)

- **Objective:** Summarize the key points and open the floor for questions.
- Content:
  - Recap of data collection methods, data discovery, and experimental design.
  - Emphasize the importance of careful planning and ethical considerations in data analysis.

- Encourage students to think critically about how data is collected and analyzed in their own projects.
- Open the floor for questions and discussion.

#### Key Takeaways

- Data collection methods vary widely and should be chosen based on the research question and context.
- Data discovery is a crucial step in identifying relevant and high-quality data for analysis.
- Experimental design is fundamental to testing hypotheses and establishing causality, with careful consideration of variables and control.

## **Resources**:

Fowler, F. J. Survey Research Methods Trochim, W. M. K., & Donnelly, J. P. The Research Methods Knowledge Base Salmons, J. Doing Qualitative Research Online Shamoo, A. E., & Resnik, D. B. Responsible Conduct of Research

## Lecture Outline: Importing Data and Creating Basic Graphs in Python Duration: 50 minutes

## 1. Introduction to Data Importing and Visualization in Python (5 minutes)

- **Objective:** Provide an overview of the importance of data importing and visualization in data analysis.
- Content:
  - Brief introduction to pandas for data manipulation and the importance of data importing.
  - Overview of Python's graphing capabilities using different libraries.
  - Outline of what will be covered: importing various file types into a pandas DataFrame, and creating basic graphs using different packages.

## 2. Importing Data into a Pandas DataFrame (20 minutes)

- **Objective:** Teach students how to import different file types (CSV, Excel, JSON) into a pandas DataFrame.
- Content:
  - Importing CSV Files:
    - Syntax: pd.read\_csv('file\_path.csv')
    - **Example:** Importing a CSV file into a DataFrame.
    - Discuss options like specifying the delimiter, handling missing values, and setting an index column.
  - Importing Excel Files:
    - **Syntax:** pd.read\_excel('file\_path.xlsx', sheet\_name='Sheet1')
    - **Example:** Importing data from an Excel file, specifying the sheet name.
    - Handling multiple sheets, working with Excel-specific options.
  - Importing JSON Files:
    - **Syntax:** pd.read\_json('file\_path.json')
    - **Example:** Importing a JSON file into a DataFrame.

- Discussing the structure of JSON data (nested data) and how to handle it in pandas.
- Hands-on Practice:
  - Quick demonstration of importing different file types using sample data.
  - Encourage students to try importing a dataset of their choice on their own systems.

### 3. Creating Basic Graphs Using Matplotlib (10 minutes)

- **Objective:** Introduce Matplotlib for creating basic graphs.
- Content:
  - Introduction to Matplotlib:
    - Overview of Matplotlib as a versatile plotting library.
    - **Syntax:** Importing Matplotlib: import matplotlib.pyplot as plt
  - Creating Basic Plots:
    - Line Plot: plt.plot(x, y)
    - Bar Chart: plt.bar(x, y)
    - Histogram: plt.hist(data)
  - Customizing Plots:
    - Adding titles, labels, legends, and customizing colors.

## Syntax Example:

plt.title('Sample Title') plt.xlabel('X-axis Label') plt.ylabel('Y-axis Label')

٠

#### • Practice Example:

 Create a line plot and a bar chart using sample data, and demonstrate basic customization options.

#### 4. Creating Graphs Using Seaborn (8 minutes)

- **Objective:** Introduce Seaborn for creating more advanced, aesthetically pleasing plots.
- Content:

#### • Introduction to Seaborn:

- Overview of Seaborn as a higher-level interface for Matplotlib, focused on statistical plots.
- **Syntax:** Importing Seaborn: import seaborn as sns
- Creating Basic Plots:
  - Scatter Plot: sns.scatterplot(x='column\_name', y='column\_name', data=df)
  - Box Plot: sns.boxplot(x='column\_name', y='column\_name', data=df)
  - Pair Plot: sns.pairplot(df)
- Customization and Aesthetics:
  - Using Seaborn's built-in themes and color palettes.
  - Customizing plots with additional options, such as adding hue for categorical variables.
- Practice Example:
  - Create a scatter plot and a box plot using sample data, and demonstrate customization options.

## 5. Creating Graphs Using Plotly (10 minutes)

- **Objective:** Introduce Plotly for creating interactive plots.
- Content:
  - Introduction to Plotly:
    - Overview of Plotly as a library for creating interactive, web-based plots.
    - Syntax: Importing Plotly: import plotly.express as px
  - Creating Basic Interactive Plots:
    - Scatter Plot: px.scatter(df, x='column\_name', y='column\_name')
    - Bar Chart: px.bar(df, x='column\_name', y='column\_name')
    - Line Chart: px.line(df, x='column\_name', y='column\_name')
  - Interactivity Features:
    - Demonstrating hover, zoom, and export options in Plotly plots.
    - Discussing how these features can enhance data exploration and presentation.
  - Practice Example:
    - Create an interactive scatter plot and a bar chart using sample data, showing how to interact with the plots.

## 6. Conclusion & Q&A (5 minutes)

- **Objective:** Summarize key points and address any questions.
- Content:
  - Recap of how to import different file types into pandas and the basics of creating graphs using Matplotlib, Seaborn, and Plotly.
  - Highlight the strengths of each graphing library and when to use them.
  - Encourage students to explore further by creating plots with their own datasets.
  - Open the floor for questions and discussion.

## Key Takeaways

- Importing data from different file types into pandas is fundamental for data analysis in Python.
- Matplotlib, Seaborn, and Plotly each offer unique advantages for creating visualizations, from basic static plots to advanced interactive graphs.
- Understanding the basics of these tools is essential for effectively analyzing and presenting data.

## Resources:

Python Graph Gallery <u>https://python-graph-gallery.com/</u> Plotting Graphs in Python <u>https://www.geeksforgeeks.org/graph-plotting-in-python-set-1/</u> Plotly Graphing Gallery <u>https://plotly.com/python/</u> Python Graphs <u>https://www.tutorialspoint.com/python\_data\_structure/python\_graphs.htm</u>

Importing from Sheets requires quite a bit extra work, so I'm including a general example here if you need it.

**To import a Google Sheets file into Python**, you can use the `gspread` library along with `oauth2client` for authentication. Here's a step-by-step guide to help you get started:

```
    **Install the required libraries**:
    ``bash
    pip install gspread oauth2client
```

- 2. \*\*Set up Google Cloud Console\*\*:
  - Log in to Google Cloud Console.
  - Create a new project or select an existing one.
  - Enable the Google Sheets API and Google Drive API.
  - Navigate to "APIs and Services" > "Credentials" to set up a service account.
  - Create a service account and download the JSON key file.
- 3. \*\*Share your Google Sheet\*\*:
  - Open your Google Sheet.
  - Click the "Share" button and add the service account email address found in the JSON key file.
  - Grant Editor permissions.

4. \*\*Authenticate and import the Google Sheets file\*\*:

```
```python
import gspread
from oauth2client.service account import ServiceAccountCredentials
```

# Define the scope and credentials

```
scope = ['https://www.googleapis.com/auth/spreadsheets', 'https://www.googleapis.com/auth/drive']
credentials = ServiceAccountCredentials.from_json_keyfile_name('path/to/your/credentials.json',
scope)
```

```
# Authorize and open the Google Sheets file
client = gspread.authorize(credentials)
sheet = client.open('Your_Sheet_Name').sheet1
```

```
# Get all values from the sheet
data = sheet.get_all_values()
```

```
# Print the first few rows
for row in data[:5]:
    print(row)
...
```

Replace `'path/to/your/credentials.json'` with the path to your JSON key file and `'Your\_Sheet\_Name'` with the name of your Google Sheet.

This should help you get started with importing Google Sheets files into Python.