DSA 610 Redesign, Lecture 3 Outline

# Lecture Outline: Basic Database Types and the Conceptual Model **Duration**: 50 minutes

#### 1. Introduction to Database Types (5 minutes)

- **Objective:** Provide an overview of the variety of database types and their significance in data management.
- Content:
  - Brief introduction to databases and their role in storing, organizing, and managing data.
  - Overview of relational and non-relational databases.
  - Introduction to the concept of the "conceptual model" in database design.

#### 2. Relational Databases and SQL (10 minutes)

- **Objective:** Explain the basics of relational databases and the SQL query language.
- Content:
  - Relational Databases:
    - Definition: A database structured to recognize relations among stored items of information.
    - Structure: Tables, rows, and columns (schema-based).
    - Key Concepts: Primary keys, foreign keys, relationships (one-to-many, many-tomany).
    - SQL (Structured Query Language):
      - Definition: The standard language for managing and manipulating relational databases.
      - Basic SQL Operations:
        - SELECT: Querying data from tables.
        - INSERT: Adding new records.
        - UPDATE: Modifying existing records.
        - DELETE: Removing records.
      - **Example:** Simple SQL query to select data from a table.

#### 3. NoSQL Databases (10 minutes)

- **Objective:** Introduce NoSQL databases, including their types and use cases.
- Content:
  - **Definition of NoSQL:** 
    - A category of database management systems that do not use the traditional relational database model.
    - Advantages: Flexibility, scalability, handling unstructured data.
  - Types of NoSQL Databases:
    - Document Stores (e.g., MongoDB):
      - Store data in JSON-like documents.
      - Ideal for hierarchical data structures.
    - Key-Value Stores (e.g., Redis):
      - Simple data storage in key-value pairs.
      - High-performance use cases like caching.
    - Column-Family Stores (e.g., Cassandra):
      - Store data in columns rather than rows.

- Suitable for large-scale distributed systems.
- Graph Databases (e.g., Neo4j):
  - Designed to represent data as interconnected nodes and edges.
  - Ideal for relationship-heavy data like social networks.
- Comparison with SQL Databases:
  - Schema flexibility, horizontal scalability, and eventual consistency.

# 4. XML and JSON in Databases (10 minutes)

- Objective: Discuss the role of XML and JSON as data formats and their integration with databases.
- Content:
  - XML (Extensible Markup Language):
    - **Definition:** A markup language for encoding documents in a format that is both human-readable and machine-readable.
    - Usage in Databases:
      - Commonly used in document stores and as a format for data interchange.
      - Can be stored and queried directly in certain database systems.
      - **Example:** Sample XML structure and discussion on storing it in a database.
  - JSON (JavaScript Object Notation):
    - **Definition:** A lightweight data-interchange format that is easy for humans to read and write and easy for machines to parse and generate.
    - Usage in Databases:
      - Widely used in document stores like MongoDB.
      - Increasingly supported in relational databases (e.g., PostgreSQL).
    - **Example:** Sample JSON structure and discussion on its advantages over XML in certain use cases.

# 5. Object-Oriented and Graph Databases (10 minutes)

- **Objective:** Explore object-oriented and graph databases, their features, and use cases.
- Content:
  - **Object-Oriented Databases:** 
    - **Definition:** Databases that integrate object-oriented programming principles, storing data as objects.
    - Key Concepts: Inheritance, encapsulation, polymorphism.
    - **Usage:** Often used in applications where data and logic are tightly integrated.
    - **Example:** Discuss how an object-oriented database stores and manages complex data types.
  - Graph Databases:
    - **Definition:** Databases designed to represent and store data in the form of graphs (nodes and edges).
    - Key Concepts: Nodes (entities), edges (relationships), properties (attributes of nodes/edges).
    - Use Cases: Social networks, recommendation engines, fraud detection.
    - **Example:** Overview of a simple graph structure representing a social network and how queries are executed in graph databases.

- **Objective:** Introduce the concept of the conceptual model in database design and its importance.
- Content:
  - Definition of Conceptual Model:
    - An abstract framework that describes the structure of a database without focusing on physical implementation details.
  - Importance:
    - Helps in visualizing the data structure and relationships before implementation.
    - Guides the logical and physical design of the database.
  - Components of Conceptual Models:
    - Entities: Key objects or concepts within the domain (e.g., Customer, Product).
    - Attributes: Properties of entities (e.g., Customer Name, Product Price).
    - Relationships: Connections between entities (e.g., Customer "purchases" Product).
  - **Tools:** Mention of Entity-Relationship Diagrams (ERDs) as a way to represent conceptual models.



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

- **Objective:** Summarize key points and address any questions.
- Content:
  - Recap of different database types, their use cases, and the importance of conceptual modeling.
  - Highlight the importance of choosing the right database type based on the specific needs of a project.
  - Encourage students to explore databases hands-on, experimenting with different types and querying methods.
  - Open the floor for questions and discussion.

- Different database types (relational, NoSQL, object-oriented, graph) cater to different data needs and use cases.
- SQL databases are structured and ideal for transactional data, while NoSQL databases offer flexibility for unstructured or semi-structured data.
- XML and JSON are common data interchange formats that integrate with various database systems.
- Conceptual modeling is a crucial step in designing effective databases that align with business requirements.

#### Resources:

W3Schools: <u>https://www.w3schools.com/sql/ https://www.w3schools.com/databases/</u> Databases and NoSQL tutorials: <u>https://www.tutorialspoint.com/database\_tutorial/</u> Database Normalization:

https://www.youtube.com/watch?v=GFQaEYEc8\_8&pp=ygUgRGF0YWJhc2Ugbm9ybWFsaXphdGlvbiBleH BsYWluZWQ%3D

Entity-Relationship Diagrams: <u>https://www.lucidchart.com/pages/er-diagrams</u>

https://www.youtube.com/watch?v=WuZt1X3kztI&pp=ygUkRW50aXR5LVJlbGF0aW9uc2hpcCBkaWFncm FtIHR1dG9yaWFs https://medium.com/@elizabethskachkov/entity-relationship-diagrams-anexplanation-1c478499b77f

UML Class Diagrams:

https://www.youtube.com/watch?v=ao1ESgly2Ws&pp=ygUSVU1MIENsYXNzIERpYWdyYW1z Document Databases: https://www.mongodb.com/university

# Lecture Outline: Data Structures Essential for Data Analysis

Duration: 50 minutes

# 1. Introduction to Data Structures (5 minutes)

- **Objective:** Provide an overview of the importance of data structures in data analysis.
- Content:
  - Define data structures and their role in storing and organizing data for analysis.
  - Brief introduction to the types of data structures commonly used in data analysis: lists, dictionaries, arrays, dataframes, etc.
  - Overview of the differences between basic and complex data structures.

# 2. Lists and Tuples (8 minutes)

- **Objective:** Explain the basic data structures: lists and tuples, and their use cases in data analysis.
- Content:
  - Lists:
    - **Definition:** Ordered, mutable collections of items.
    - **Syntax Example:** my\_list = [1, 2, 3, 4]
    - Use Cases: Storing sequences of data, iterating over elements.
    - **Operations:** Adding, removing, sorting elements.
    - Tuples:
      - **Definition:** Ordered, immutable collections of items.
      - Syntax Example: my\_tuple = (1, 2, 3, 4)
      - Use Cases: Storing fixed data, multiple return values from functions.
      - Comparison with Lists:

- Lists are mutable (can change), while tuples are immutable (cannot change).
- Tuples can be faster and are often used for data that shouldn't be modified.
- Example Use Case:
  - Using a list to store a series of data points, and using a tuple to store constant values (like configuration settings).

#### 3. Dictionaries (10 minutes)

- **Objective:** Introduce dictionaries, their structure, and how they are used in data analysis.
- Content:
  - **Definition:** Collections of key-value pairs, where keys are unique identifiers for their associated values.
  - o Syntax Example: my\_dict = {'key1': 'value1', 'key2': 'value2'}
  - Use Cases:
    - Storing and accessing data by key (e.g., mapping IDs to user data).
    - Fast lookups and efficient management of data that doesn't require order.
  - Operations:
    - Adding, updating, and removing key-value pairs.
    - Iterating over keys, values, or key-value pairs.
    - Handling missing keys with methods like .get() or using default dictionaries.
  - **Example Use Case:** 
    - Using a dictionary to store and retrieve user information by user ID.

#### 4. Arrays (10 minutes)

- **Objective:** Explain the role of arrays in data analysis, particularly when dealing with numerical data.
- Content:
  - **Definition:** Fixed-size, homogeneous collections of elements (all elements are of the same data type).
  - In Python:
    - Often used through libraries like NumPy for numerical computations.
    - Syntax Example: import numpy as np; my\_array = np.array([1, 2, 3, 4])
  - Use Cases:
    - Efficient storage and computation of large datasets, especially in mathematical operations and algorithms.
    - Commonly used for matrix operations, statistical calculations, and machine learning.
  - **Comparison with Lists:** 
    - Arrays are more memory-efficient and provide faster computations for large numerical datasets.
    - Unlike lists, arrays require all elements to be of the same type.
  - Example Use Case:
    - Using an array to perform element-wise mathematical operations or to store and manipulate a dataset for machine learning.

#### 5. DataFrames (12 minutes)

• **Objective:** Explore the DataFrame structure, its features, and why it is essential in data analysis.

- Content:
  - **Definition:** A two-dimensional, size-mutable, and potentially heterogeneous tabular data structure with labeled axes (rows and columns).
  - In Python:
    - Primarily used through the pandas library.
    - Syntax Example: import pandas as pd; my\_dataframe = pd.DataFrame({'Column1': [1, 2], 'Column2': [3, 4]})
  - Use Cases:
    - Storing, manipulating, and analyzing large datasets in a tabular format.
    - Handling missing data, merging, grouping, and reshaping data.
    - Widely used for data cleaning, preprocessing, and exploratory data analysis (EDA).
  - **Operations:** 
    - Accessing data using labels (df['Column']), slicing rows and columns, and performing aggregate operations.
    - Merging and joining multiple DataFrames.
    - Applying functions across rows or columns with .apply().
  - Comparison with Arrays and Dictionaries:
    - DataFrames provide more functionality and flexibility for data manipulation compared to arrays and dictionaries.
    - Like arrays, DataFrames are optimized for numerical operations but also support complex data manipulation.
  - Example Use Case:
    - Loading a CSV file into a DataFrame, cleaning the data, and performing exploratory analysis.

#### 6. Comparison and Choosing the Right Data Structure (5 minutes)

- **Objective:** Compare the different data structures and provide guidance on when to use each.
  - Content:
    - Lists vs. Arrays:
      - Use lists for simple collections of data that may include different data types.
      - Use arrays when working with numerical data and need efficient computation.
    - Dictionaries vs. DataFrames:
      - Use dictionaries for key-value pair storage, especially when quick lookups by key are needed.
      - Use DataFrames when working with structured data in rows and columns, especially for data analysis and manipulation.
    - Choosing the Right Structure:
      - Consider the nature of the data (size, type, operations needed).
      - Performance considerations (memory usage, speed).
      - Flexibility vs. specificity: DataFrames are versatile, while arrays offer optimized performance for numerical data.

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

- **Objective:** Summarize the key points and address any remaining questions.
- Content:
  - Recap of the main data structures: lists, tuples, dictionaries, arrays, and DataFrames.

- Highlight the importance of choosing the right data structure for efficient and effective data analysis.
- Encourage students to practice working with these structures through exercises and projects.
- Open the floor for questions and discussion.

#### Key Takeaways

- Different data structures (lists, tuples, dictionaries, arrays, DataFrames) are essential tools in data analysis, each with unique characteristics and use cases.
- Understanding the differences and similarities between these structures is crucial for efficient data manipulation and analysis.
- The choice of data structure depends on the type of data, the operations needed, and the specific analysis context.

#### **Resources**:

Python Tutorials: <u>https://docs.python.org/3/tutorial/</u> <u>https://www.w3schools.com/python/</u> NumPy's official quickstart guide: <u>https://numpy.org/doc/stable/user/absolute\_beginners.html</u> Pandas 10-minute guide: <u>https://pandas.pydata.org/docs/user\_guide/10min.html</u>

# Lecture Outline: Data Types in Python and Appropriate Graphs for Data Visualization **Duration:** 50 minutes

#### **1.** Introduction to Data Types in Python (5 minutes)

- **Objective:** Provide an overview of Python data types and their significance in data analysis.
- Content:
  - Brief introduction to data types as a fundamental concept in programming.
  - Overview of Python's basic data types and their role in data analysis.
  - Explanation of why understanding data types is important for choosing the right visualizations.

#### 2. Basic Python Data Types (15 minutes)

- **Objective:** Explain the basic data types in Python and their characteristics.
- Content:
  - Numeric Types:
    - Integers: Whole numbers (e.g., int\_var = 5).
    - **Floats:** Numbers with decimals (e.g., float\_var = 3.14).
    - Use Cases: Storing quantitative data, performing mathematical operations.
  - Strings:
    - Definition: Sequence of characters (e.g., str\_var = "Hello").
    - Use Cases: Storing categorical data, text processing, labels in graphs.
  - Booleans:
    - **Definition:** True/False values (e.g., bool\_var = True).
    - Use Cases: Conditional operations, binary categorizations (e.g., yes/no, success/failure).
  - Lists:
    - Definition: Ordered, mutable collections of items (e.g., list\_var = [1, 2, 3]).
    - Use Cases: Storing sequences of data that can be of mixed types.
  - Tuples:

- Definition: Ordered, immutable collections of items (e.g., tuple\_var = (1, 2, 3)).
- Use Cases: Storing fixed data that shouldn't be modified.
- Dictionaries:
  - Definition: Collections of key-value pairs (e.g., dict\_var = {'key1': 'value1'}).
  - Use Cases: Storing and accessing data by key, often used for mapping.
- Example Use Case:
  - Discuss how to choose the correct data type based on the nature of the data (e.g., using floats for continuous numerical data).

# 3. Appropriate Graphs for Individual Variables (10 minutes)

- **Objective:** Discuss appropriate graphs for visualizing individual variables based on their data types.
- Content:
  - Numeric Data:
    - Histograms:
      - Ideal for visualizing the distribution of a single numeric variable.
      - Example: Plotting the distribution of heights in a population.
    - Box Plots:
      - Useful for showing the spread, quartiles, and outliers of a numeric variable.
      - Example: Visualizing the distribution of exam scores.
    - Line Plots:
      - Suitable for time series data or showing trends over continuous intervals.
      - Example: Plotting monthly sales over a year.
  - Categorical Data:
    - Bar Charts:
      - Best for comparing the frequency of different categories.
      - Example: Visualizing the number of students in each grade category.
    - Pie Charts:
      - Used to show the proportion of categories within a whole.
      - Example: Showing the market share of different companies.
      - Count Plots (Seaborn):
        - A variation of bar charts to display counts of categorical variables.
        - Example: Counting occurrences of different job roles in a dataset.
  - Boolean Data:
    - Bar Charts:
      - Effective for binary comparisons.
      - Example: Comparing the number of pass/fail results.

# 4. Visualizing Combinations of Variables (10 minutes)

- **Objective:** Introduce graphs for visualizing relationships between multiple variables.
- Content:
  - Numeric vs. Numeric:
    - Scatter Plots:
      - Ideal for visualizing the relationship between two numeric variables.
      - Example: Plotting the relationship between study hours and test scores.
    - Heatmaps:

- Used for visualizing correlation matrices or showing the intensity of relationships.
- Example: Visualizing the correlation between different financial metrics.
- Categorical vs. Numeric:
  - Box Plots:
    - Useful for comparing the distribution of a numeric variable across different categories.
    - Example: Comparing salaries across different job roles.
  - Violin Plots:
    - Similar to box plots but also show the density of the data.
    - Example: Visualizing the distribution of customer satisfaction scores across regions.
  - Bar Charts (Grouped/Stacked):
    - Show comparisons of categorical data with associated numeric values.
    - Example: Comparing average scores of different departments in a company.
- Categorical vs. Categorical:
  - Stacked Bar Charts:
    - Ideal for showing the relationship between two categorical variables.
    - Example: Comparing gender distribution across different job roles.
    - Mosaic Plots:
      - Useful for showing proportions of two categorical variables.
      - Example: Visualizing the proportion of product categories sold across different regions.

# 5. Selecting Appropriate Graphs for Complex Data Analysis (10 minutes)

- **Objective:** Discuss strategies for choosing the right visualization for complex data analysis involving multiple variables.
- Content:
  - Combining Multiple Variables:
    - Pair Plots (Seaborn):
      - Useful for exploring relationships between multiple numeric variables.
      - Example: Pair plot for visualizing relationships between age, income, and spending score.
    - Facet Grids:
      - Used to plot multiple graphs across different subsets of data.
      - Example: Visualizing income distribution across different age groups and genders.
    - Parallel Coordinates Plots:
      - Good for visualizing multidimensional data and spotting trends.
      - Example: Comparing performance metrics across different departments.
  - Choosing the Right Graph:
    - Consider the nature of the data (e.g., type, size, complexity).
    - Think about the story you want to tell and what the data reveals.
    - Ensure clarity and avoid clutter by choosing the most straightforward visualization.

- **Objective:** Summarize key points and address any remaining questions.
- Content:
  - Recap the importance of understanding Python data types and selecting the appropriate graph for different types of data.
  - Highlight the importance of matching the data structure with the correct visualization technique for effective communication of insights.
  - Encourage students to experiment with different graphs using their datasets.
  - Open the floor for questions and discussion.

#### Key Takeaways

- Different data types in Python (numeric, categorical, boolean) require different approaches to visualization.
- Understanding the characteristics of your data helps in selecting the most appropriate graph to convey insights effectively.
- Combining variables and selecting advanced visualizations can uncover deeper relationships and trends in the data.

#### **Resources**:

Data to Viz: <u>https://www.data-to-viz.com/</u> The Data Visualization Catalog: <u>https://datavizcatalogue.com/</u> Chart Chooser: <u>https://www.highcharts.com/chartchooser/</u>