

Table of Contents

1. Data Science Lifecycle
2. Python Fundamentals
3. Central Tendency
4. Pandas Fundamentals
5. Distributions
6. Data Visualization
7. Seaborn Fundamentals
8. Data Collection
9. HTML
10. CSS
11. JavaScript
12. D3
13. Data Stories: Introduction
14. Data Stories: Central Insights

8/21/23

Data Science Lifecycle

1. Data scientists use techniques to transform data into a visual representation that can be easier to understand by humans
2. Data science platform market is expected to grow by upwards of 20% annually.
3. Data science generally falls under math, statistics, and computer science.
4. The Life Cycle
 - a. Question
 - b. Collect Data
 - c. Wrangle Data
 - d. Analyze Data
 - e. Visualize Information
 - f. Communicate Information
5. All steps in the life cycle are all fluid

9/5/23

Python Fundamentals

1. Datasets
 - a. The collection of data
 - b. Types of datasets
 - i. Lists
 1. Ordered, changeable, duplicates allowed
 - ii. Dictionaries
 1. Ordered, changeable, duplicates not allowed
 - iii. Sets
 1. Unordered, unchangeable*, duplicates not allowed
 - iv. Tuples
 1. Unordered, unchangeable, duplicates allowed
2. Representing datasets with code
 - a. Column-oriented
 - i. Grouping by features
 - b. Row-oriented
 - i. Grouping by a single observation

3. Indexing

- a. Used to access values of a collection type
- b. Python syntax to access values
 - i. List
 - 1. name[index]
 - ii. Dictionary
 - 1. name[key]
 - iii. Set
 - 1. for loop
 - iv. Tuple
 - 1. Name[index]

4. Iteration

- a. Can repeat processes with loops or recursion in Python
- b. Python loop types
 - i. While loop
 - 1. while condition: statements
 - ii. For loop
 - 1. for thing in collection: statements

5. Useful methods

- a. [Dictionaries](#)
 - i. values()
 - ii. items()
 - iii. keys()
- b. [Lists](#)
 - i. len()
 - ii. append()
 - iii. sort()
- c. [Other](#)
 - i. range()
 - ii. print()
 - iii. split()
 - iv. type()
 - v. int()
 - vi. str()

9/5/23

Central Tendency

1. Measures of Central Tendency
 - a. Statistical measures that help describe the behavior of a collection of data points
 - b. Mean
 - i. The average of all the values in a dataset
 - ii. Summation of all the values, divided by the count of values
 - iii. Can be misleading, because outliers skew the result
 - c. Median
 - i. The value in the direct center of a sorted dataset
 - ii. Gives a more proportional representation of data that excludes outliers
 - d. Mode
 - i. The most frequently occurring value in a dataset
 - ii. Most useful in a relatively large sample size
2. The center of a dataset is a good measure of determining the behavior or distribution of a dataset
 - a. Gives examples of a whole dataset, not data points individually.
3. Distribution
 - a. Shows how often data occurs in a dataset
4. Outliers
 - a. Unusually large or small values
 - b. Skew the result of a mean in a dataset
5. Bimodal
 - a. When two values are most common
6. Unimodal
 - a. When one value is the most common value
7. Symmetric distribution
 - a. When the mean and median are the same
8. Skewed distribution
 - a. When the dataset is offset by outliers, causing the mean to be an inaccurate representation of the population

9/15/23

Pandas Fundamentals

1. Pandas
 - a. Pandas is a python library that can make analyzing data easier.
2. Dataframes
 - a. A pandas object that is used to store a dataset
 - b. Information is organized in rows and columns
 - c. Dataframes simplify common operations such as sorting data
3. Series
 - a. A pandas object used to create dataframes
 - b. Seen as a one-dimensional list of data
 - i. Think of it as a single column of a dataframe
4. Indexing into Dataframes
 - a. [df.loc\[\]](#)
 - i. name.loc[row_label, col_label]
 - b. [df.iloc\[\]](#)
 - i. name.iloc[row_index, col_index]
5. Selection
 - a. The process of accessing a subset of a dataframe
 - b. Can select subsets using loc and iloc
6. Filtering
 - a. Selecting values of a dataset where certain conditions are true
 - b. df[condition]
7. Combining Dataframes
 - a. Concatenating
 - i. Naively combines along an axis
 - b. Merge
 - i. Combine through a shared column
 - c. Join
 - i. Combine using shared indices
 - ii. Inner Join
 1. Keep similar pieces
 - iii. Left Outer Join
 1. Keep the left
 - iv. Right Outer Join
 1. Keep the right

- v. Full Outer Join
 - 1. Keep everything

9/19/23

Distributions

1. Distributions
 - a. Graphs that tell us about some characteristic of a population
 - b. Mean and median are important parts of the graphs
 - c. Tells us about the shape and spread of data
2. Normal distribution
 - a. The mean, median, and mode are all the same
 - b. Empirical Rule
 - i. 68% of data is within 1 standard deviation from the mean
 - ii. 95% within 2 standard deviations
 - iii. 99.7% within 3 standard deviations
 - c. Unimodal
 - i. Only one peak
3. Standard deviation
 - a. The average distance between any point and the mean
4. Skewed distribution
 - a. Skew is towards outliers
 - i. Can be seen on graphs by a "tail"
5. Bimodal
 - a. Has two peaks on a graph
6. Uniform distribution
 - a. Each value has the same frequency

10/2/23

Data Visualization

1. Data Visualizations
 - a. A graph or picture that helps humans understand important patterns in a dataset

10/2/23

Seaborn Fundamentals

1. Seaborn
 - a. A python library that can make visualizing data easier
2. Bar Charts
 - a. A graph type that uses bars to depict a value associated with a category
3. Histogram
 - a. A graph that shows the frequency distribution of a variable in a dataset
4. Scatterplots
 - a. A graph that uses points to show the relationship between 2 quantitative variables in a dataset.

10/13/23

Data Collection

1. Techniques
 - a. Observe a sample
 - b. Survey a sample
 - c. Experiment on a sample
 - d. Use data that somebody else has responsibly collected
2. Sourcing Digital Datasets: API Requests
 - a. The act of using HTTP requests in order to access datasets collected and maintained by other people
 - b. Common HTTP requests:
 - i. GET
 1. Requests for information
 2. Only retrieves data
 3. Does not modify data
 - ii. POST
 1. Modify the underlying data
 2. Create new resources
 - iii. PUT
 1. Modify the underlying data
 2. Update existing resources

- iv. DELETE
 - 1. Remove existing resources
- 3. Sourcing Digital Datasets: Web Scraping
 - a. The act of extracting data from websites using the structure of its HTML
 - b. Scraping and crawling exists in legal gray zones
 - i. The TOS determines the legality of web scraping

10/30/23

HTML

- 1. Hypertext Markup Language
 - a. Used to display content on a webpage
 - b. Look for angled brackets <>!
- 2. General Page Structure
 - a. Two major sections
 - i. Head
 - 1. Contains important metadata
 - ii. Body
 - 1. All content that is seen on a page
- 3. Tag Structure
 - a. HTML is made up of tags
 - b. Each tag does something different
 - c. Most have an opening and closing tag
 - d. Example:
 - i. `<h1>Content</h1>`
 - 1. Gives a large heading
- 4. Tag Attributes
 - a. Some tags need more information in order to work
 - i. To do this, you need to use attributes.
 - 1. Example:
 - a. ``
- 5. Important Metadata Tags and Attributes
 - a. Tags:
 - i. `<title>?</title>`
 - ii. `<meta name = "?" content = "?">`
 - iii. `<link rel = "?" href = "?">`

- b. Attributes
 - i. alt = "description"
 - ii. lang = "?"
- 6. Accessibility
 - a. We want to make sure that our websites are accessible to as many people as possible
 - i. Use [these practices](#)
 - b. Considerations
 - i. Low bandwidth users
 - ii. Visually impaired users
 - iii. Low English proficiency users

11/7/23

CSS

- 1. Cascading Style Sheets
 - a. Used to style the content on a web page
 - b. Look for curly braces {}!
- 2. General Structure
 - a. Two major sections
 - i. Selector
 - 1. Targeted HTML tag
 - a. General
 - b. Class
 - c. ID
 - ii. Property
 - 1. Style to be applied
- 3. Class Selectors
 - a. Used to select a subset of the HTML tags used
 - b. Has more priority than the generic HTML tag selector
 - c. Start selector with a period (.) in order to use
- 4. ID Selectors
 - a. Used to style a single HTML tag used
 - b. Has the most priority of all selectors
 - c. Start selector with a hashtag (#) in order to use
- 5. The Box Model

- a. Every HTML Tag makes a box
 - b. Boxes can be styled with CSS to change the default layout of every webpage
- 6. Accessibility
 - a. We want to make sure that our websites are accessible to as many people as possible
 - i. Use [these practices](#)
 - b. Considerations
 - i. Low bandwidth users
 - ii. Visually impaired users
 - iii. Low English proficiency users

11/13/23

JavaScript

- 1. JS
 - a. JavaScript is the programming language of the web
 - b. Used to give websites behavior

11/21/23

D3

- 1. [D3](#)
 - a. A JavaScript library that is used to create beautiful and interactive data visualizations

12/4/23

Data Stories: Introduction

- 1. What purpose does an introduction serve?
 - a. Allows the audience to know the question being asked
 - b. Provides context and important background information to the topic
 - c. Hooks the reader, tailored towards audience
 - d. Outlines how the question will be answered

- e. Defines the tone and flow of the narrative
- 2. Characteristics of a strong introduction
 - a. Includes at least one piece of background information
 - b. Includes a hook (appealing and engaging UI/UX, intriguing content)
 - c. Concise and short (try to keep around a paragraph)
 - d. Presents the question or topic, but intentionally keeps out data and answer to force user engagement
 - e. Builds credibility (certifications, data sources, past projects, reference personal anecdotes)

12/8/23

Data Stories: Central Insights

- 1. Importance of a clear central insight
 - a. Gives story direction
 - b. Helps the writing process because it lets you know where to go
 - c. Maintains focus
 - d. Explicitly defines takeaways to let the reader know what to do with the information
 - i. Helps readers grasp the point of the story
 - ii. Tells readers what they action they should take regarding the information
 - iii. Ensures that readers are concluding what you want them to conclude
 - e. Without a central insight, there is no point to writing a data story
- 2. Characteristics of a great central insight
 - a. Clear and comprehensible
 - i. Think about grammar
 - ii. Engaging
 - b. Actionable
 - i. Should be relevant to the audience you want to present to
 - 1. Relevance builds throughout the story
 - c. Comprehensive
- 3. What to avoid in a central insight
 - a. Bad grammar (not understandable / not clear)
 - b. Don't include the raw analysis

- i. Simplify the analysis in a way that makes the next steps easy to understand
- c. Too long / not engaging
- d. Don't be vague
 - i. Make sure you list set takeaways
- e. Doesn't relate to question or analysis (off-topic)
 - i. Don't show bias
 - ii. Make sure you are honest!
- f. Don't be too general
 - i. The insight should be actionable by most people, but being too general makes it not specific enough for there to be an effective actionable insight.