**Title:**
**First Task Summary Report – Big Data Analytics Intern**

**Date:**
5/11/24

---

**Objective:**
Perform exploratory data analysis (EDA) on the provided dataset to understand its structure, identify missing or inconsistent values, and prepare it for further analysis.

**Dataset Description:**
The dataset contains details of AI-generated art, including artwork ID, artist information, style, creation date, and popularity scores.

---

**1. Tasks Completed**

**Environment Setup:**

- Configured Jupyter Notebook and Python IDE.

- Verified access to data repositories.

- Ensured connectivity to required data sources.

**Exploratory Data Analysis:**

- **Dataset Structure and Data Types:**

    o The dataset contains **12 columns** and **10,000 rows**.

    o Data types include object, float, and datetime formats.

    o Key columns:

        ▪ **Artwork_ID, Artist_Name, Art_Style, Creation_Date, Medium, Popularity_Score.**

- **Summary Statistics:**

    o **Categorical Data:**

        ▪ Most frequent artist: *MidJourney* (appears 694 times).

        ▪ Most common style: *Cubism* (1,050 artworks).

        ▪ Popular tools: *DeepDream* (used 2,032 times).

    o **Numerical Data:**

        ▪ Average Popularity Score: **2,508.19**

        ▪ Score Range: **50.85 to 4,999.62**

- **Missing or Inconsistent Values:**

    o No missing values identified in any columns.

**Data Cleaning and Transformation:**

- Converted the **Creation_Date** column from string to datetime format for consistency.

- Ensured all columns are clean and appropriately typed.

**2. Key Findings**

1. **Dataset Observations:**

   o Comprehensive dataset without missing values or inconsistencies.

   o Clear trends in art styles, tools used, and popularity scores.

2. **Insights Gained:**

   o *Cubism* is the most prominent art style.

   o South America is the most frequently represented region.

   o Popular tools like *DeepDream* dominate AI-generated art trends.

**3.Learnings**

**Learnings:**

- Improved understanding of EDA workflows, including type conversions and summary report generation.

- Strengthened Python and Pandas skills for large-scale data analysis.

**4.PDF file link**

https://drive.google.com/file/d/15FBF0m-EUMpREJdbjvsjq8n8UBGFTvP0/view?usp=drive_link

**Title: Task 2 -Setting Up Web Scraping Tools**

**Date:**
6/11/24-7/11/24

---

**1. Task Overview**

Objective:
To install, configure, and demonstrate the functionality of web scraping tools BeautifulSoup and Scrapy, along with documenting the process and dependencies.

Key Deliverables:

1. Three separate test scripts for each tool (BeautifulSoup and Scrapy).

2. A setup guide with step-by-step instructions and screenshots of command-line output

**2. Tasks Completed**

**Library Installation:**

1. Installed necessary libraries via pip:

> **pip install beautifulsoup4**

2. Verified installations:

> **pip list**

**Tool Testing:**

1. BeautifulSoup:

- Wrote a simple script to fetch and parse static HTML from a sample website:

import requests

from bs4 import BeautifulSoup

# URL of the webpage to scrape

url = 'http://quotes.toscrape.com/'

# Send GET request to fetch the webpage

response = requests.get(url)

# Check if the request was successful

if response.status_code == 200:

    print("Successfully fetched the webpage!")

    # Parse the HTML content of the page

```python
    soup = BeautifulSoup(response.text, 'html.parser')

    # Find all quote elements on the page

    quotes = soup.find_all('div', class_='quote')

    # Loop through each quote and extract the text, author, and tags

    for quote in quotes:

        # Extract the text of the quote

        text = quote.find('span', class_='text').text

        # Extract the author of the quote

        author = quote.find('small', class_='author').text

        # Extract the tags associated with the quote

        tags = [tag.text for tag in quote.find_all('a', class_='tag')]

        # Print the extracted information

        print(f"Quote: {text}")

        print(f"Author: {author}")

        print(f"Tags: {', '.join(tags)}")

        print('-' * 80)

else:

    print(f"Failed to fetch the webpage. Status code: {response.status_code}")
```

**Execution Command:**

```
Python quotes.py
```

```
C:\Windows\System32\cmd.e   X   +   ∨                                                                    —   □

Microsoft Windows [Version 10.0.22631.4391]
(c) Microsoft Corporation. All rights reserved.

C:\Users\PURVI\Downloads>python quotes.py
Successfully fetched the webpage!
Quote: "The world as we have created it is a process of our thinking. It cannot be changed without changing our thinking."
Author: Albert Einstein
Tags: change, deep-thoughts, thinking, world
------------------------------------------------------------------------------------
Quote: "It is our choices, Harry, that show what we truly are, far more than our abilities."
Author: J.K. Rowling
Tags: abilities, choices
------------------------------------------------------------------------------------
Quote: "There are only two ways to live your life. One is as though nothing is a miracle. The other is as though everything is a miracle."
Author: Albert Einstein
Tags: inspirational, life, live, miracle, miracles
------------------------------------------------------------------------------------
Quote: "The person, be it gentleman or lady, who has not pleasure in a good novel, must be intolerably stupid."
Author: Jane Austen
Tags: aliteracy, books, classic, humor
------------------------------------------------------------------------------------
Quote: "Imperfection is beauty, madness is genius and it's better to be absolutely ridiculous than absolutely boring."
Author: Marilyn Monroe
Tags: be-yourself, inspirational
------------------------------------------------------------------------------------
Quote: "Try not to become a man of success. Rather become a man of value."
Author: Albert Einstein
Tags: adulthood, success, value
------------------------------------------------------------------------------------
Quote: "It is better to be hated for what you are than to be loved for what you are not."
Author: André Gide
Tags: life, love
------------------------------------------------------------------------------------
Quote: "I have not failed. I've just found 10,000 ways that won't work."
Author: Thomas A. Edison
Tags: edison, failure, inspirational, paraphrased
------------------------------------------------------------------------------------
Quote: "A woman is like a tea bag; you never know how strong it is until it's in hot water."
Author: Eleanor Roosevelt
Tags: misattributed-eleanor-roosevelt
------------------------------------------------------------------------------------
```

**2. Scrapy:**

- Created a basic Scrapy project and crawler:

To get started, first install Scrapy:

```
pip install scrapy
```

Then create a Scrapy project:

```
scrapy startproject quotes_scraper
```

Edited the spider to scrape titles from a static website:

```python
import scrapy

class QuotesSpider(scrapy.Spider):

    name = "quotes"  # Name of the spider

    start_urls = [

        'http://quotes.toscrape.com/page/1/',  # Starting URL

    ]

    def parse(self, response):

        # Loop through each quote block

        for quote in response.css('div.quote'):

            # Extract text, author, and tags

            yield {

                'text': quote.css('span.text::text').get(),

                'author': quote.css('small.author::text').get(),

                'tags': quote.css('div.tags a.tag::text').getall(),

            }

        # Follow the "Next" page link if it exists

        next_page = response.css('li.next a::attr(href)').get()

        if next_page is not None:

            yield response.follow(next_page, self.parse)
```

- **Execution Command:**

```
scrapy crawl quotes
```

**4. Key Findings**

- **BeautifulSoup:** Ideal for static websites or small-scale scraping tasks.

- **Scrapy:** Best suited for complex projects with multi-page scraping requirements.