

# EXPLORATORY DATA ANALYSIS

A Complete Guide (Zero to Hero)



# Table of Contents

## **Chapter 1: Introduction to EDA**

1. Introduction to Exploratory Data Analysis (EDA).
2. Importance of EDA in Data Science.
3. The Purposes and Objectives of EDA.
4. Types of EDA.
5. Tools and Libraries for EDA.
6. Statistical Concepts Used in EDA.
7. Graphs and Charts Used in EDA.
8. Best Practices in EDA.
9. Common Pitfalls and How to Avoid Them.
10. Future Trends in EDA.

## **Chapter 2: Setting Up The Computer For EDA**

- Step 1: Install Python.
- Step 2: Install Visual Studio Code (VS Code).
- Step 3: Install Python Extension for VS Code.
- Step 4: Install Jupyter Extension for VS Code.
- Step 5: Install Jupyter via the Command Line.
- Step 6: Install Essential Python Libraries for EDA.
- Step 8: Configure Jupyter Notebook in VS Code.

## **Chapter 3: Exploring The Libraries**

1. Pandas
2. NumPy
3. Matplotlib
4. Seaborn
5. Plotly
6. SciPy
7. Jupyter Notebooks

# **Exploratory Data Analysis Complete Guide**

Chapter 1: Introduction to EDA

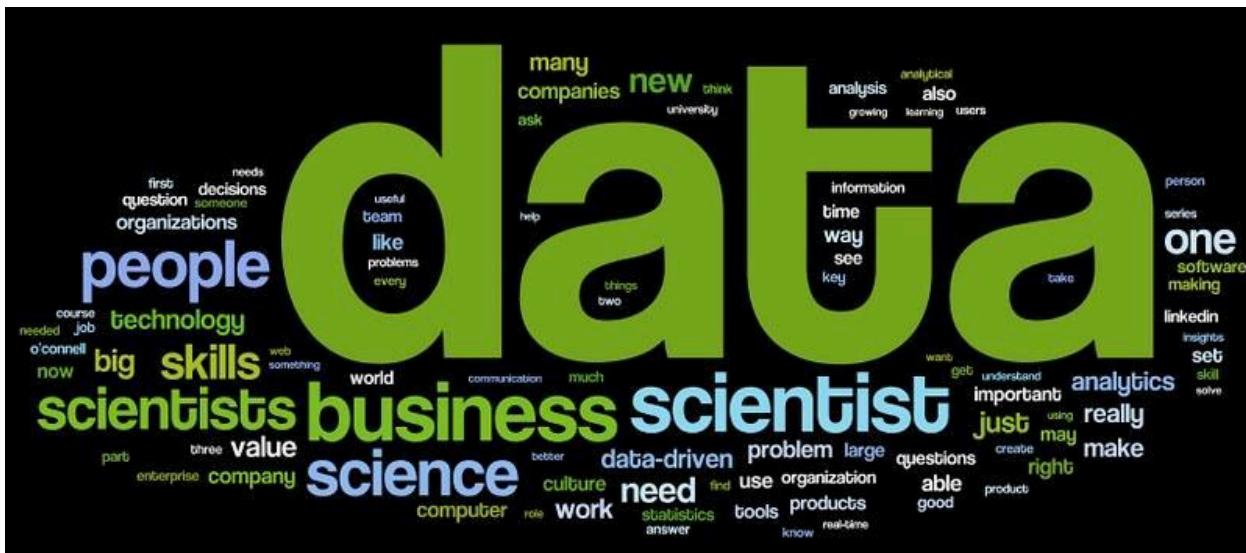
# Exploratory Data Analysis (EDA)

Welcome to the exciting world of Exploratory Data Analysis (EDA)! Imagine you're a detective, and your data is the crime scene. EDA is your magnifying glass, helping you uncover hidden patterns, understand relationships, and ultimately solve the mystery your data holds.

This ebook will equip you with the tools and knowledge to become a data detective. Let's dive into each topic:

## 1. Introduction to Exploratory Data Analysis (EDA).

Exploratory Data Analysis (EDA) is a crucial initial step in data science projects. It involves analyzing and visualizing data to understand its key characteristics, uncovering patterns, and identifying relationships between variables. It refers to studying and exploring record sets to apprehend their predominant traits, discover patterns, locate outliers, and identify relationships between variables. EDA is normally carried out as a preliminary step before undertaking extra formal statistical analyses or modeling. This preliminary analysis ensures that your data is ready for more advanced modeling and analysis.



## 2. Importance of EDA in Data Science.

Just like a detective wouldn't jump straight to conclusions, EDA is crucial in data science. It lays the groundwork for everything that comes after. Through EDA, you can identify potential issues with the data, refine your research questions, and ultimately build more robust models. Think of it as building a solid foundation for your data science project. Without EDA, you might miss important insights or make decisions based on flawed data.

## 3. The Purposes and Objectives of EDA.

The main purposes of EDA are to:

- Summarize the main characteristics of the data.
- Identify patterns and relationships.
- Spot anomalies and outliers.

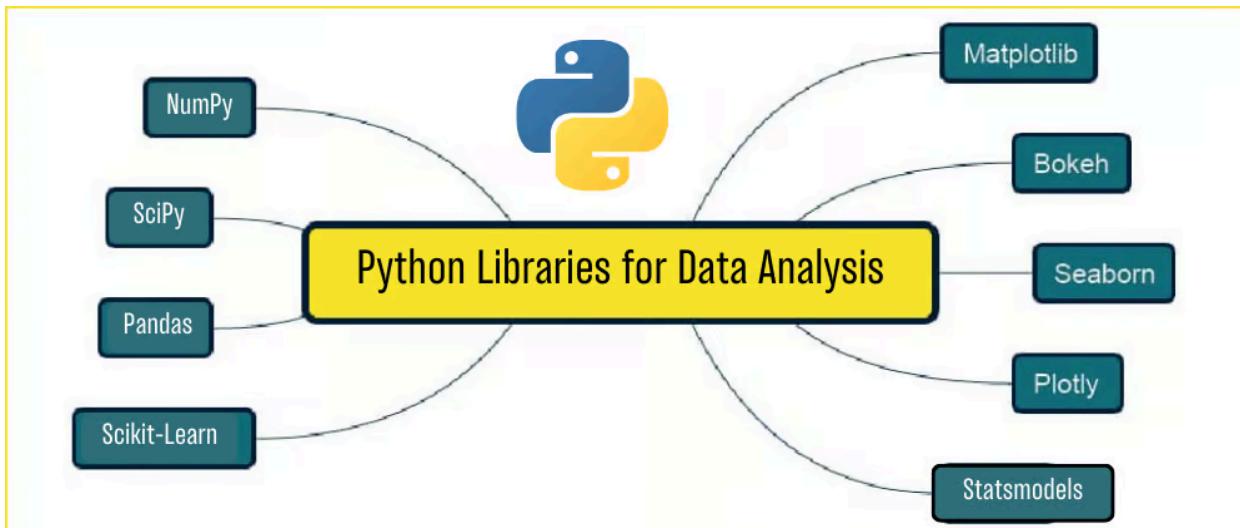
- Test initial hypotheses.
- Determine the data's structure and cleanliness.

The objectives include ensuring data quality, forming a clear understanding of the data, and preparing it for further analysis or modeling.

## 4. Types of EDA.

EDA can be divided into different types based on the techniques used:

- **Descriptive EDA:**  
Summarizes the data with statistics and visualizations.
- **Inferential EDA:**  
Makes inferences about the population based on sample data.
- **Univariate EDA:**  
Focuses on one variable at a time.
- **Bivariate EDA:**  
Examines relationships between two variables.
- **Multivariate EDA:**  
Examines relationships between two variables.



## 5. Tools and Libraries for EDA.

Several tools and libraries make EDA easier and more efficient:

- **Python Libraries:**  
Pandas, NumPy, Matplotlib, Seaborn, Plotly.
- **R Libraries:**  
ggplot2, dplyr, tidyr.
- **Software:**  
Jupyter Notebooks, RStudio, Tableau.

These tools help you manipulate data, perform statistical analysis, and create visualizations.



## 6. Statistical Concepts Used in EDA.

EDA relies on several key statistical concepts:

- **Mean, Median, Mode:**  
Measures of central tendency.
- **Variance and Standard Deviation:**  
Measures of spread.
- **Correlation and Covariance:**  
Measures of relationships between variables.
- **Probability Distributions:**  
Understanding the data's distribution.
- **Hypothesis Testing:**  
Making inferences about the data.

These concepts are your building blocks for interpreting your data's story.



## 7. Graphs and Charts Used in EDA.

A picture is worth a thousand words, and that's especially true in EDA. Some common types of charts include:

- **Histograms:**  
Show the distribution of a single variable.
- **Box Plots:**  
Highlight the distribution and outliers.
- **Scatter Plots:**  
Show relationships between two variables.
- **Bar Charts:**  
Compare different categories.
- **Heatmaps:**  
Visualize correlation matrices or other data patterns.

Remember, these visualizations are meant to help you understand your data, not just impress people with fancy charts.

## 8. Best Practices in EDA.

To get the most out of EDA, follow these best practices:

- **Start Simple:**  
Begin with cleaning your data, basic statistics, and plots.
- **Be Curious:**  
Explore different angles and ask questions.
- **Document Your Steps:**  
Keep track of what you do and why.
- **Clean Your Data:**  
Identify and handle missing values and outliers.

- **Use Visualizations:**  
They often reveal insights that numbers alone can't.

## 9. Common Pitfalls and How to Avoid Them.

Watch out for these common EDA mistakes:

- **Ignoring Data Quality:**  
Always check for and handle missing or erroneous data.
- **Overfitting to Noise:**  
Be cautious not to draw conclusions from random variations.
- **Confirmation Bias:**  
Avoid looking only for evidence that supports your hypotheses.
- **Overcomplicating:**  
Start simple and build up complexity as needed.
- **Neglecting Documentation:**  
Keep a clear record of your analysis process.

## 10. Future Trends in EDA.

EDA is evolving with advancements in technology. Some future trends include:

- **Automated EDA Tools:**  
Using AI to automate and enhance EDA processes.
- **Interactive Visualizations:**  
More powerful and user-friendly visualization tools.
- **Integration with Machine Learning:**  
Seamless integration with predictive modeling.
- **Enhanced Collaboration:**  
Tools that make it easier to share and collaborate on data analysis.
- **Big Data Handling:**  
Improved capabilities for handling and analyzing large datasets.

## **Exploratory Data Analysis Complete Guide**

Chapter 2: Setting Up The Computer For EDA

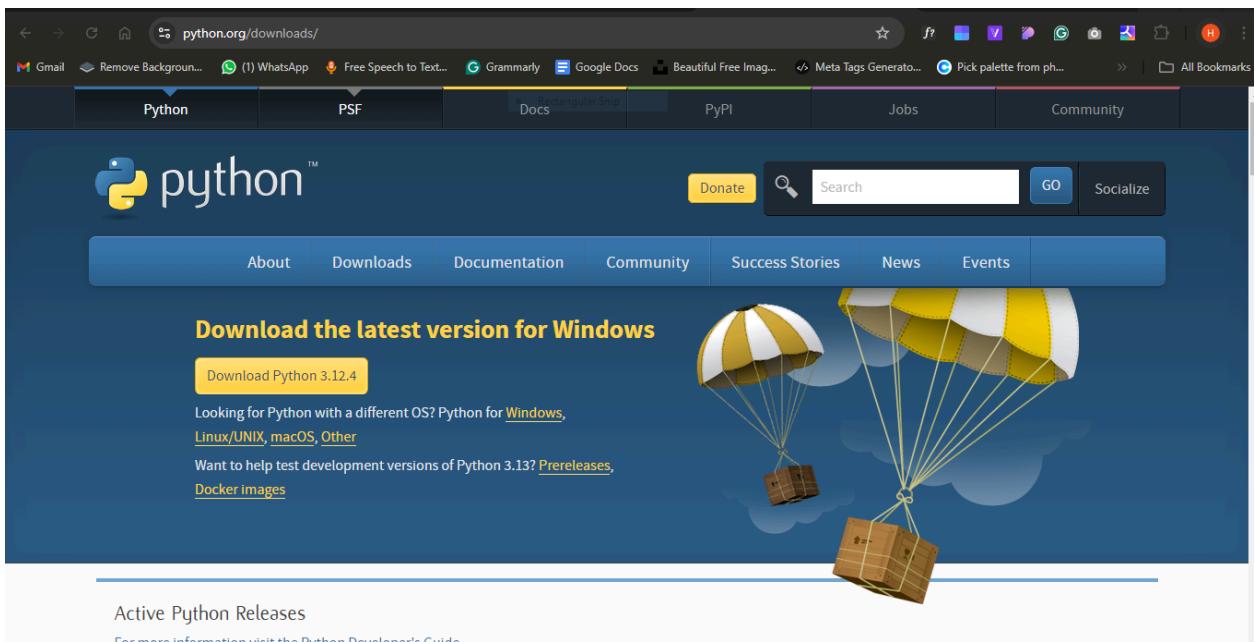
# Preparing Your Computer for EDA.

To start with Exploratory Data Analysis (EDA) in Python, you'll need to set up your computer with the necessary software, extensions, and libraries. Here's a step-by-step guide to help you get everything ready.

## Step 1: Install Python.

1. **Download Python:** Go to the [Python official website](https://www.python.org/downloads/) and download the latest version of Python.

3



2. **Install Python:** Run the installer and make sure to check the box that says "Add Python to PATH". Follow the instructions to complete the installation.

## Step 2: Install Visual Studio Code (VS Code).

1. **Download VS Code:** Visit the [Visual Studio Code website](https://code.visualstudio.com/) and download the installer for your operating system.
2. **Install VS Code:** Run the installer and follow the instructions to complete the installation.

The screenshot shows the official Visual Studio Code download page. At the top, there's a navigation bar with links to Visual Studio Code, Docs, Updates, Blog, API, Extensions, FAQ, Learn, a search bar, and a 'Download' button. Below the navigation, a message says 'Version 1.92 is now available! Read about the new features and fixes from July.' A large section titled 'Download Visual Studio Code' follows, with the subtext 'Free and built on open source. Integrated Git, debugging and extensions.' Below this, there are three main download sections: Windows (with icons for User Installer, System Installer, and .zip), Linux (.deb and .rpm), and Mac (.zip, Intel chip, Apple silicon, Universal). Each section includes specific file names like 'Windows 10, 11', 'Debian, Ubuntu', 'Red Hat, Fedora, SUSE', and 'macOS 10.15+'. At the bottom left, a link to 'https://code.visualstudio.com/blogs' is visible.

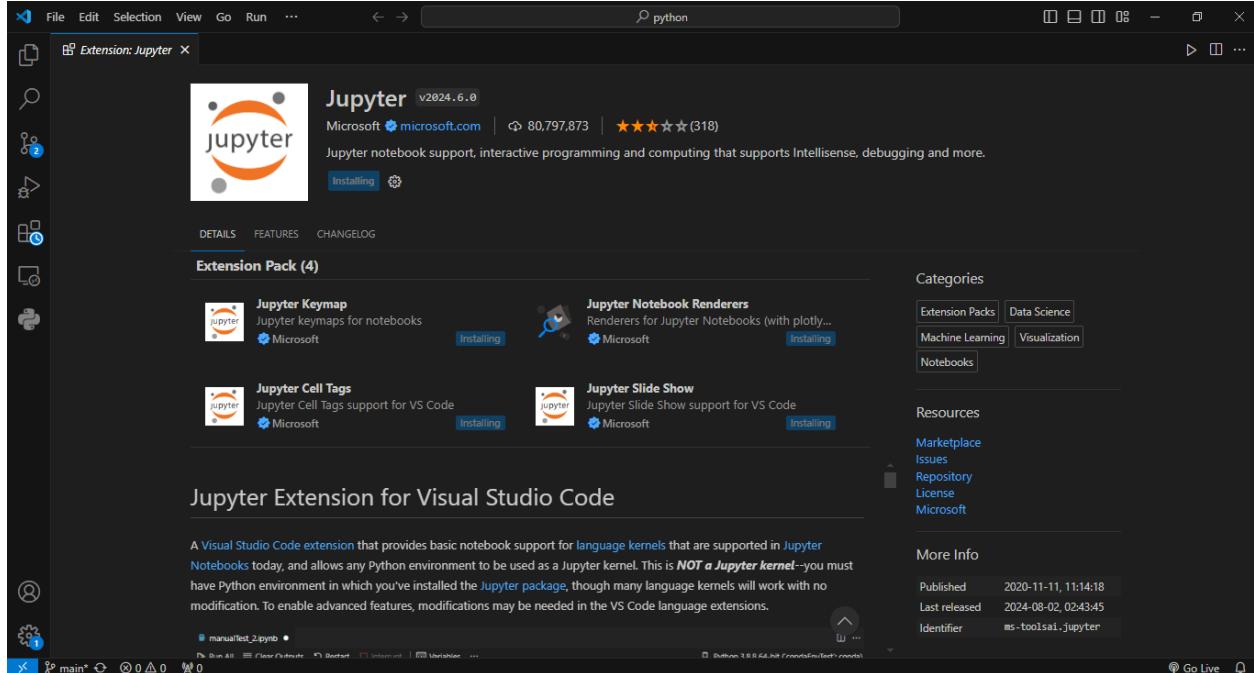
## Step 3: Install Python Extension for VS Code.

- Open VS Code.**
- Go to Extensions:** Click on the Extensions view icon on the Sidebar or press **Ctrl+Shift+X**.
- Search for Python:** Type "Python" in the search bar.
- Install Python Extension:** Click on the **Install** button next to the Python extension by Microsoft. You can find it [here](#).

The screenshot shows the Visual Studio Code Extensions Marketplace. On the left, the sidebar has icons for File, Edit, Selection, View, Go, Run, and others. The main area shows the Python extension by Microsoft, version v2024.12.1, with a rating of 5 stars (595 reviews). It provides support for IntelliSense (Pylance), Debugging (Python Debugger), linting, formatting, refactoring, and more. Below the extension details, there's a 'Python extension for Visual Studio Code' section with a description, support for vscode.dev, installed extensions (PyLance and Python Debugger), and optional dependencies. To the right, there are sections for Categories (Programming Languages, Debuggers, Other, Data Science, Machine Learning), Resources (Marketplace, Issues, Repository, License, Microsoft), and More Info (Published, Last released, Last updated, Identifier).

## Step 4: Install Jupyter Extension for VS Code.

1. **Open Extensions:** Click on the Extensions view icon on the Sidebar or press **Ctrl+Shift+X**.
2. **Search for Jupyter:** Type "Jupyter" in the search bar.
3. **Install Jupyter Extension:** Click on the **Install** button next to the Jupyter extension by Microsoft. You can find it [here](#).



## Step 5: Install Jupyter via the Command Line.

1. **Open Command Prompt or Terminal.**
2. **Install Jupyter:** Type the following command and press Enter:

```
pip install notebook
```

```
Command Prompt - pip install notebook
Microsoft Windows [Version 10.0.19045.4651]
(c) Microsoft Corporation. All rights reserved.

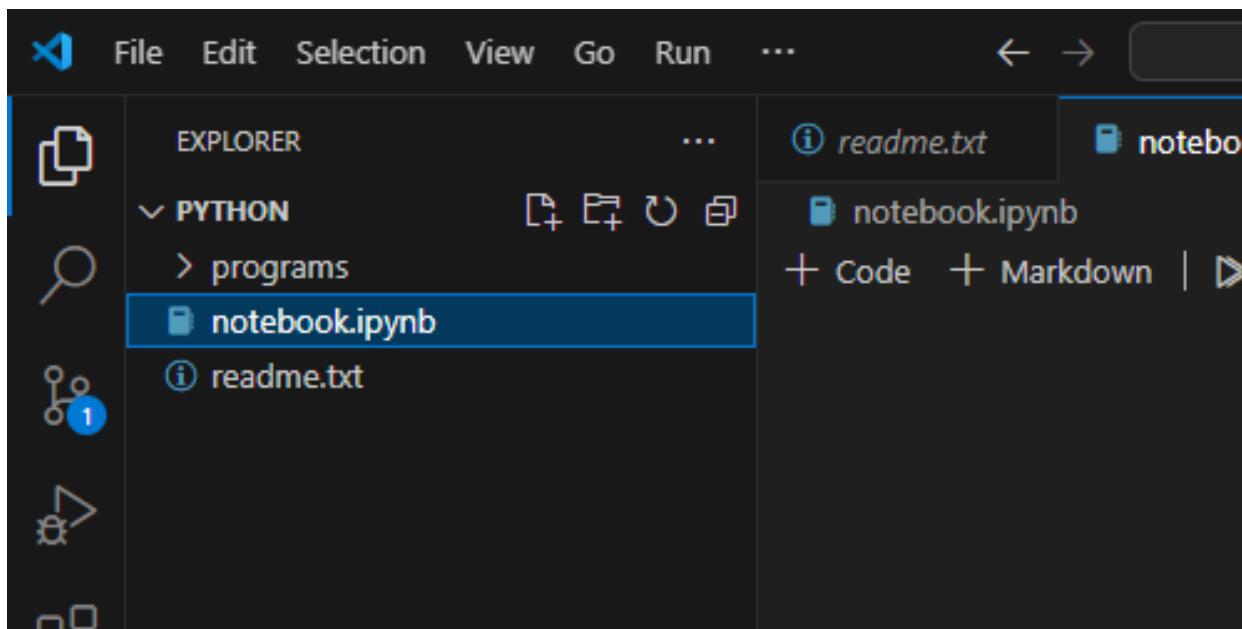
C:\Users\ASAD>pip install notebook
Collecting notebook
  Downloading notebook-7.2.1-py3-none-any.whl.metadata (10 kB)
Collecting jupyter-server<3,>=2.4.0 (from notebook)
  Downloading jupyter_server-2.14.2-py3-none-any.whl.metadata (8.4 kB)
Collecting jupyterlab-server<3,>=2.27.1 (from notebook)
  Downloading jupyterlab_server-2.27.3-py3-none-any.whl.metadata (5.9 kB)
Collecting jupyterlab<4.3,>=4.2.0 (from notebook)
  Downloading jupyterlab-4.2.4-py3-none-any.whl.metadata (16 kB)
Collecting notebook-shim<0.3,>=0.2 (from notebook)
```

## Step 6: Install Essential Python Libraries for EDA.

1. **Open Command Prompt or Terminal.**
2. **Install Libraries:** Use `pip` to install the essential libraries for EDA. Run the following commands:  
`pip install pandas`  
`pip install numpy`  
`pip install matplotlib`  
`pip install seaborn`  
`pip install plotly`  
`pip install scipy`

## Step 8: Configure Jupyter Notebook in VS Code.

1. **Open VS Code.**
2. **Create or Open a Jupyter Notebook:** You can create a new Jupyter Notebook file by clicking on the `File` menu, selecting `New File`, and then changing the file extension to `.ipynb`. Alternatively, you can open an existing `.ipynb` file.



3. **Select Interpreter:** When you open a Jupyter Notebook file in VS Code for the first time, it will prompt you to select a Python interpreter. Choose the interpreter you installed in Step 1.
4. **Start Coding:** You should now see the Jupyter Notebook interface within VS Code, allowing you to create and run code cells.

By following these steps, your computer will be fully set up for performing Exploratory Data Analysis (EDA) using Python and Visual Studio Code.

# **Exploratory Data Analysis Complete Guide**

## Chapter 3: Exploring The Libraries

# Exploring The Libraries And Tools.

In the world of data science, having the right tools and libraries at your disposal can significantly affect how efficiently and effectively you can analyze data. Exploratory Data Analysis (EDA) is no exception. In this section, we'll introduce you to some of the most essential libraries and tools used in EDA. These libraries will help you clean, manipulate, visualize, and analyze your data with ease, enabling you to uncover insights and make data-driven decisions.

## 1. Pandas

- **Description:** [Pandas](#) is a powerful data manipulation and analysis library for Python. It provides data structures like DataFrames, which are essential for handling and analyzing structured data.
- **Uses in EDA:**
  - Loading data from various file formats (CSV, Excel, SQL, etc.).
  - Cleaning and preprocessing data (handling missing values, filtering, etc.).
  - Aggregating and summarizing data.
  - Merging and joining datasets.

## 2. NumPy

- **Description:** [NumPy](#) is a fundamental library for numerical computing in Python. It provides support for arrays, matrices, and a large collection of mathematical functions to operate on these data structures.
- **Uses in EDA:**
  - Performing mathematical and statistical operations on arrays and matrices.
  - Efficient numerical computations.
  - Generating random samples and performing random sampling.

## 3. Matplotlib

- **Description:** [Matplotlib](#) is a comprehensive library for creating static, animated, and interactive visualizations in Python. It is highly customizable and provides a wide range of plotting functionalities.
- **Uses in EDA:**
  - Creating basic plots like line charts, bar charts, and histograms.
  - Customizing plot aesthetics (titles, labels, legends, etc.).
  - Visualizing distributions and trends in data.

## 4. Seaborn

- **Description:** [Seaborn](#) is built on top of Matplotlib and provides a high-level interface for drawing attractive and informative statistical graphics. It simplifies complex visualizations and is particularly useful for statistical plots.

- **Uses in EDA:**
  - Creating advanced visualizations like violin plots, box plots, and heatmaps.
  - Visualizing the distribution of data and relationships between variables.
  - Enhancing the aesthetics of Matplotlib plots.

## 5. Plotly

- **Description:** [Plotly](#) is a library for creating interactive, web-based visualizations. It supports a wide variety of plots and is particularly useful for creating dashboards and sharing interactive plots online.
- **Uses in EDA:**
  - Creating interactive plots that allow for zooming, panning, and hovering.
  - Building dashboards for exploratory analysis.
  - Visualizing complex datasets in an interactive manner.

## 6. SciPy

- **Description:** [SciPy](#) is a library used for scientific and technical computing. It builds on NumPy and provides a large number of higher-level functions for optimization, integration, interpolation, eigenvalue problems, and other advanced mathematical operations.
- **Uses in EDA:**
  - Performing statistical analysis and hypothesis testing.
  - Computing advanced mathematical functions.
  - Conducting signal processing and image processing.

## 7. Jupyter Notebooks

- **Description:** [Jupyter](#) Notebooks is an open-source web application that allows you to create and share documents containing live code, equations, visualizations, and narrative text. It is widely used in data science for exploratory analysis, data cleaning, and sharing results.
- **Uses in EDA:**
  - Writing and running Python code in an interactive environment.
  - Documenting the analysis process with markdown cells.
  - Visualizing data inline with the code.
  - Sharing notebooks with others for collaboration.

By understanding and utilizing these libraries and tools, you'll be well-equipped to perform comprehensive Exploratory Data Analysis (EDA). Pandas and NumPy will handle your data manipulation and numerical operations, while Matplotlib, Seaborn, and Plotly will help you create insightful visualizations. SciPy will enable advanced statistical analysis, and Jupyter Notebooks will provide an interactive platform to document and share your findings. Together, these tools form a robust toolkit for uncovering the hidden patterns and insights within your data.

# **Exploratory Data Analysis Complete Guide**

Chapter 4: Pandas Tutorial

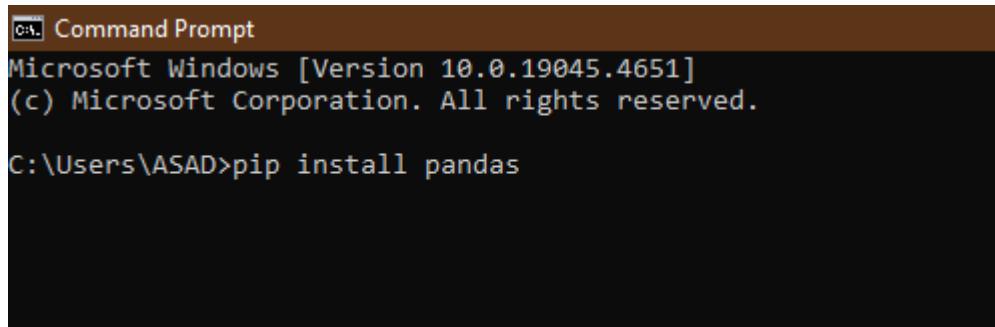
# Pandas Cheat Sheet.

Let's dive into the world of Pandas and see what we can do with it!

## 1. Installing Pandas

If you have Python installed, you can use the following command to install Pandas:

```
pip install pandas
```

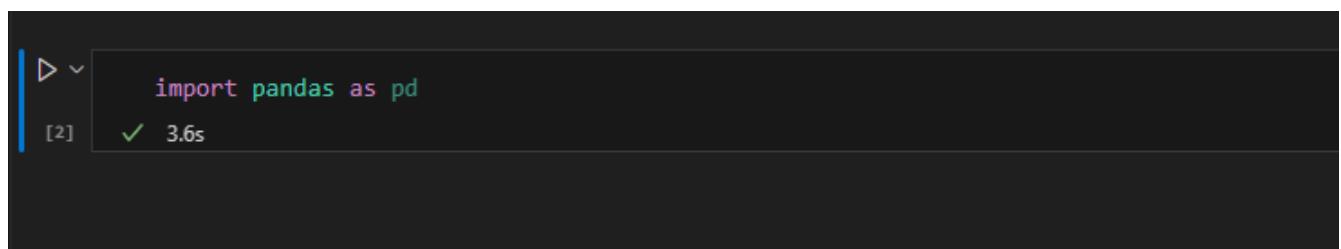


A screenshot of a Windows Command Prompt window. The title bar says "Command Prompt". The window shows the following text:  
Microsoft Windows [Version 10.0.19045.4651]  
(c) Microsoft Corporation. All rights reserved.  
C:\Users\ASAD>pip install pandas

## 2. Importing Pandas

Once Pandas is installed, you can import it into your Python script or Jupyter Notebook using the following import statement:

```
import pandas as pd
```



A screenshot of a Jupyter Notebook cell. The cell contains the code: `import pandas as pd`. To the left of the code, there is a blue dropdown arrow icon. Below the code, the text "[2]" is followed by a green checkmark and the time "3.6s".

## 3. Reading/writing Files

- **Read the .csv file:**

```
pd.read_csv('filename.csv')
```

- **Saving the .csv file:**

```
df.to_csv('xyz.csv')
```

- **Read the Sheet1 of the Excel file 'xyz.xls' :**

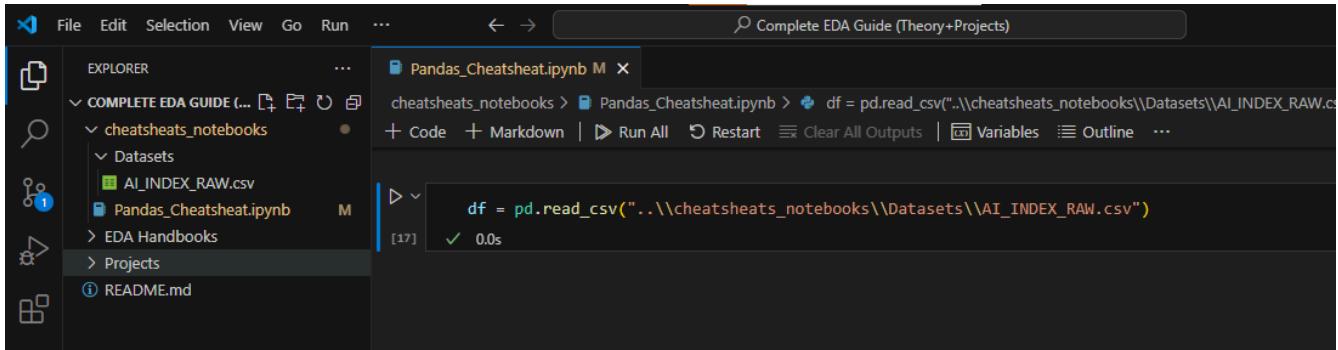
```
pd.read_excel(filename, 'Sheet1')
```

- **Saving the Sheet1 of the Excel file 'xyz.xls': `df.to_excel('filename.xlsx', sheet_name='Sheet1')`**

- **Read the xyz.json file:**  
`pd.read_json('filename.json')`

- **Read the xyz.sql file:**  
`pd.read_sql('filename.sql')`

- **Read the xyz.html file: `pd.read_html(filename.html')`**



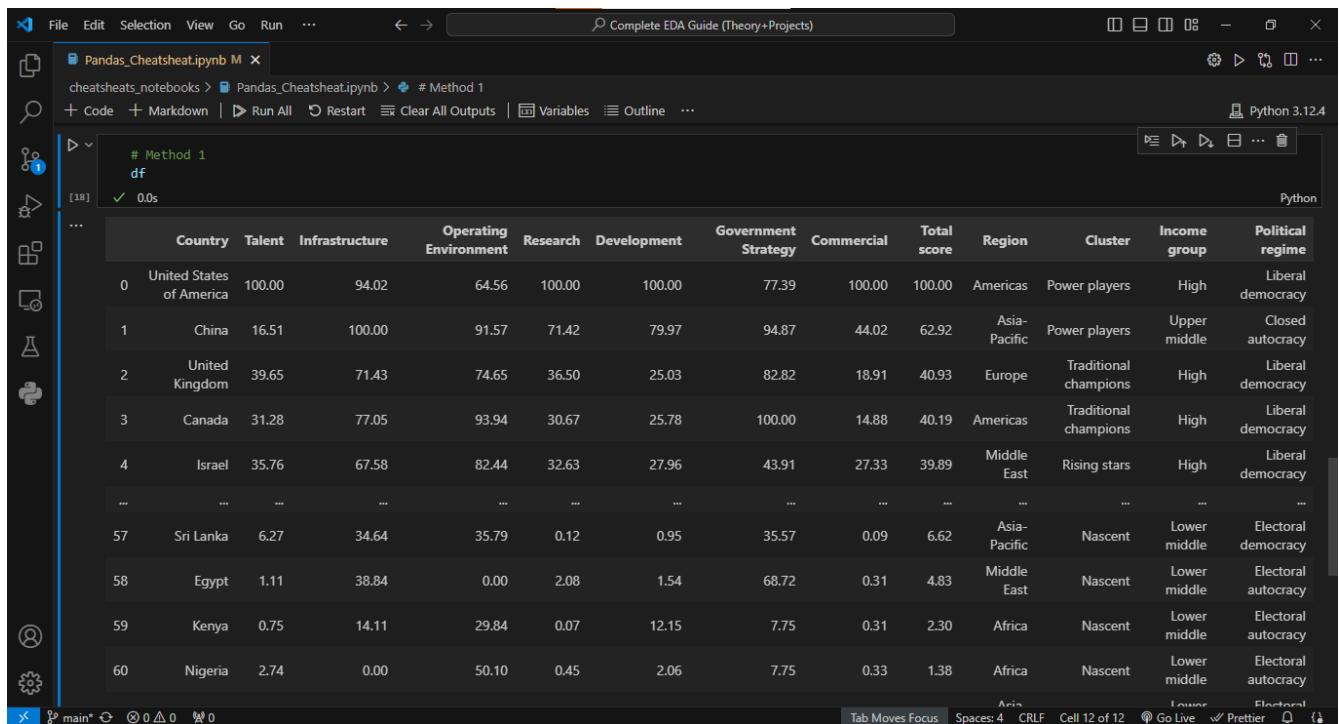
A screenshot of a Jupyter Notebook interface. The left sidebar shows a tree view of files and notebooks, including 'Pandas\_Cheatsheat.ipynb' and 'AI\_INDEX\_RAW.csv'. The main area displays a code cell with the following content:

```
df = pd.read_csv("../cheatsheets_notebooks\Datasets\AI_INDEX_RAW.csv")
[17]    ✓ 0.0s
```

The status bar at the bottom indicates the Python version is 3.12.4.

## 4. Viewing The Data

- **Method 1 - Printing whole Data File: df or dataframe-name**



A screenshot of a Jupyter Notebook interface showing a large data frame output. The code cell contains:

```
# Method 1
df
```

The resulting output is a table with 61 rows and 14 columns. The columns are labeled: Country, Talent, Infrastructure, Operating Environment, Research, Development, Government Strategy, Commercial, Total score, Region, Cluster, Income group, and Political regime. The first few rows of data are:

	Country	Talent	Infrastructure	Operating Environment	Research	Development	Government Strategy	Commercial	Total score	Region	Cluster	Income group	Political regime
0	United States of America	100.00	94.02	64.56	100.00	100.00	77.39	100.00	100.00	Americas	Power players	High	Liberal democracy
1	China	16.51	100.00	91.57	71.42	79.97	94.87	44.02	62.92	Asia-Pacific	Power players	Upper middle	Closed autocracy
2	United Kingdom	39.65	71.43	74.65	36.50	25.03	82.82	18.91	40.93	Europe	Traditional champions	High	Liberal democracy
3	Canada	31.28	77.05	93.94	30.67	25.78	100.00	14.88	40.19	Americas	Traditional champions	High	Liberal democracy
4	Israel	35.76	67.58	82.44	32.63	27.96	43.91	27.33	39.89	Middle East	Rising stars	High	Liberal democracy
...	...	...	...	...	...	...	...	...	...	...	...	...	...
57	Sri Lanka	6.27	34.64	35.79	0.12	0.95	35.57	0.09	6.62	Asia-Pacific	Nascent	Lower middle	Electoral democracy
58	Egypt	1.11	38.84	0.00	2.08	1.54	68.72	0.31	4.83	Middle East	Nascent	Lower middle	Electoral autocracy
59	Kenya	0.75	14.11	29.84	0.07	12.15	7.75	0.31	2.30	Africa	Nascent	Lower middle	Electoral autocracy
60	Nigeria	2.74	0.00	50.10	0.45	2.06	7.75	0.33	1.38	Africa	Nascent	Lower middle	Electoral autocracy

- **Method 2 - using .head() method:**

The `df.head()` method shows rows from top to bottom. By default, it shows the first 5 rows.

We can also change the number of rows by specifying them in the parentheses.

`df.head(8)`: it will show the first 8 rows of the dataframe.

	Country	Talent	Infrastructure	Operating Environment	Research	Development	Government Strategy	Commercial	Total score	Region	Cluster	Income group	Political regime
0	United States of America	100.00	94.02	64.56	100.00	100.00	77.39	100.00	100.00	Americas	Power players	High	Liberal democracy
1	China	16.51	100.00	91.57	71.42	79.97	94.87	44.02	62.92	Asia-Pacific	Power players	Upper middle	Closed autocracy
2	United Kingdom	39.65	71.43	74.65	36.50	25.03	82.82	18.91	40.93	Europe	Traditional champions	High	Liberal democracy
3	Canada	31.28	77.05	93.94	30.67	25.78	100.00	14.88	40.19	Americas	Traditional champions	High	Liberal democracy
4	Israel	35.76	67.58	82.44	32.63	27.96	43.91	27.33	39.89	Middle East	Rising stars	High	Liberal democracy
5	Singapore	39.38	84.30	43.15	37.67	22.55	79.82	15.07	38.67	Asia-Pacific	Rising stars	High	Electoral democracy
6	South Korea	14.54	85.23	68.86	26.66	77.25	87.50	5.41	38.60	Asia-Pacific	Rising stars	High	Liberal democracy
7	The Netherlands	33.83	81.99	88.05	25.54	30.17	62.35	4.97	36.35	Europe	Rising stars	High	Liberal democracy

- **method 3- using tail method:**

The `df.tail()` method is similar to the `.head` method but shows the last 5 rows of the dataframe.

`df.tail(8)` : It will show the last 8 rows of the dataframe.

	Country	Talent	Infrastructure	Operating Environment	Research	Development	Government Strategy	Commercial	Total score	Region	Cluster	Income group	Political regime
54	South Africa	4.61	45.73	58.43	0.83	7.52	0.00	2.03	9.71	Africa	Waking up	Upper middle	Electoral democracy
55	Morocco	3.36	44.88	60.17	1.46	0.05	15.90	0.10	8.87	Africa	Waking up	Lower middle	Closed autocracy
56	Armenia	6.69	37.84	58.40	0.28	0.33	14.40	1.37	8.49	Europe	Waking up	Upper middle	Electoral democracy
57	Sri Lanka	6.27	34.64	35.79	0.12	0.95	35.57	0.09	6.62	Asia-Pacific	Nascent	Lower middle	Electoral democracy
58	Egypt	1.11	38.84	0.00	2.08	1.54	68.72	0.31	4.83	Middle East	Nascent	Lower middle	Electoral autocracy
59	Kenya	0.75	14.11	29.84	0.07	12.15	7.75	0.31	2.30	Africa	Nascent	Lower middle	Electoral autocracy
60	Nigeria	2.74	0.00	50.10	0.45	2.06	7.75	0.33	1.38	Africa	Nascent	Lower middle	Electoral autocracy
61	Pakistan	8.00	2.43	12.48	2.17	1.09	13.92	0.27	0.00	Asia-Pacific	Nascent	Lower middle	Electoral autocracy

