

Exploring Dataset

Using Pandas



Presented by

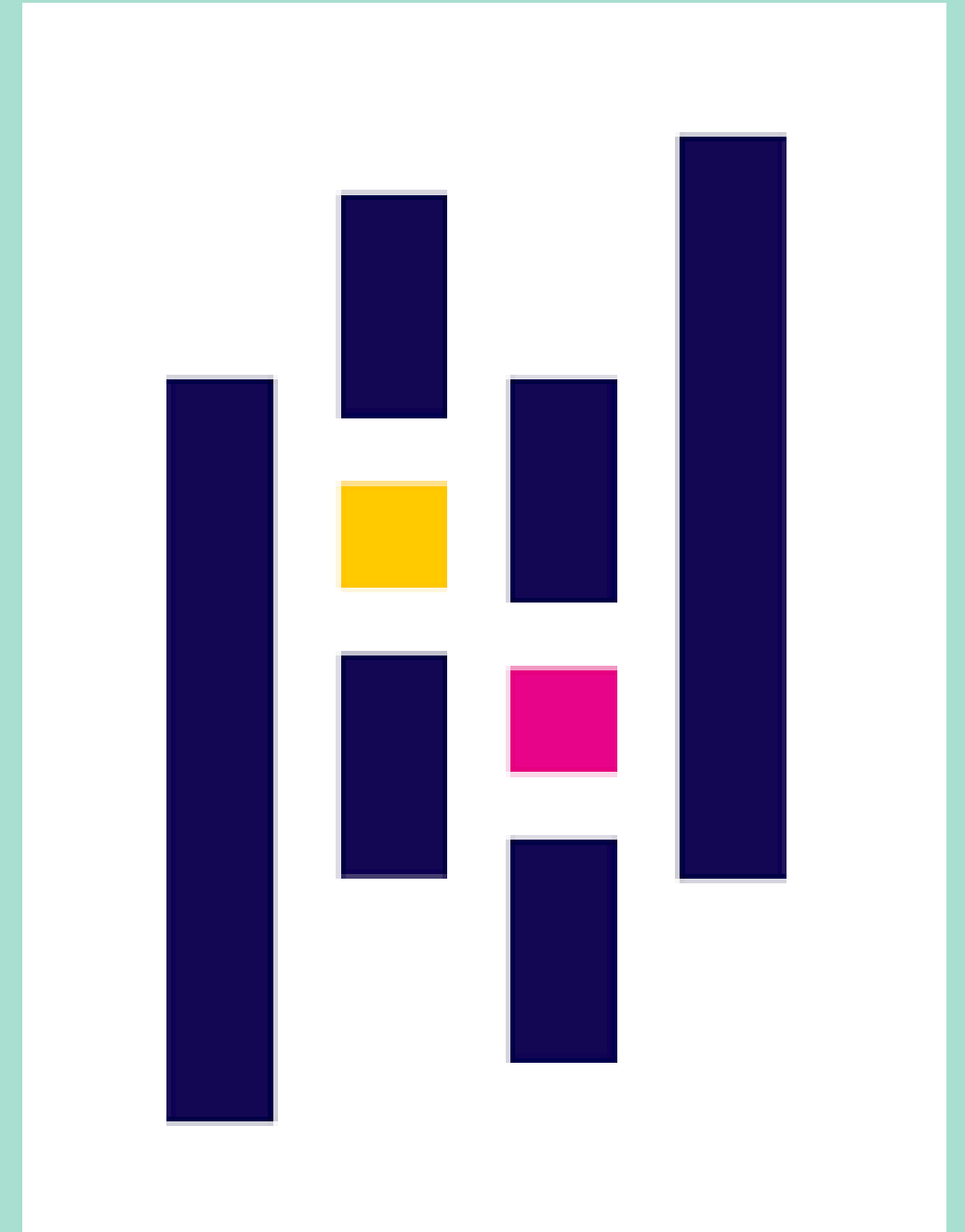
Rajesh Kanna R.
99220041074

What is Pandas?

Pandas is a Python library renowned for its versatility in data manipulation and analysis.

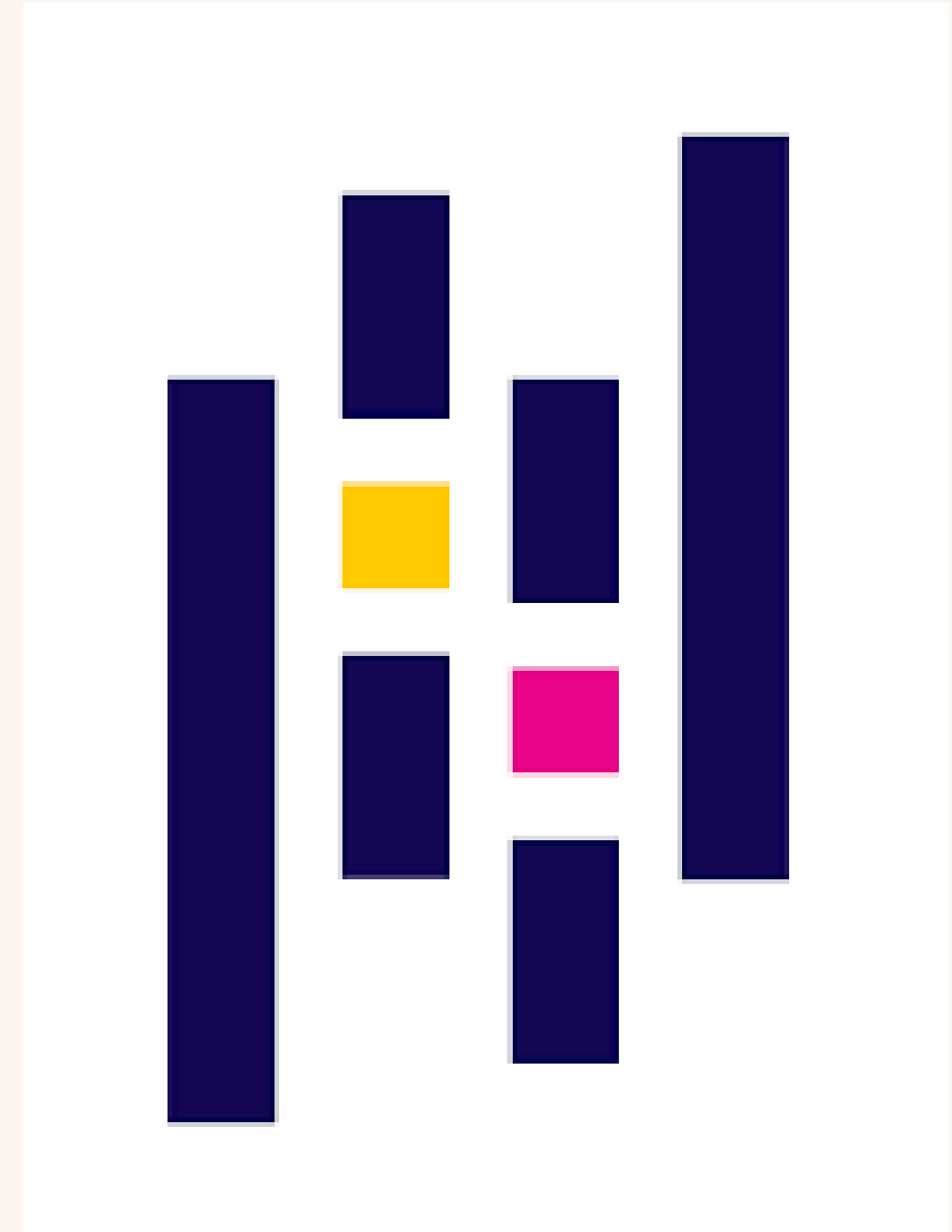
With intuitive data structures like DataFrame and Series, Pandas simplifies tasks such as cleaning, transforming, and analyzing structured data.

Its seamless integration with other libraries enhances its capabilities, making it indispensable for efficient data processing and exploration.



Objective

The objective of this project is to conduct an exploratory data analysis (EDA) of the 'E-Commerce Purchase' dataset using Pandas, with a focus cleaning and preparing the data, visualizing the data, visualizing key findings, and identifying patterns and trends in customer demographics, purchasing behavior, and preferences.



Packages installed

```
•[2]: import pandas as pd  
import matplotlib.pyplot as plt
```



Implementation of the project



Basic things

- Reading the first 10 column in the Dataset

```
[13]: # basic
# 1.Display Top 10 rows in the Dataset
data.head(10)
```

	Address	Lot	AM or PM	Browser Info	Company	Credit Card	CC Exp Date	CC Security Code	CC Provider	Email	Job
0	16629 Pace Camp Apt 448/nAlexandriaborough, HI 77...	46 in	PM	Opera/9.56(X11; Linux x86_64; sr 50) Presto/2...	Martinez- Herman	6011829061123400	02/20	900	JCB 16 digit	gdonlap@yahoo.com	Scientist, product/process development
1	9974 Jasmine Spurs Suite 508/south John, TN E...	28 on	PM	Opera/8.83. (Windows 98; Win 9x 4.90; en-US) Pr...	Fletcher, Richards and Whitaker	3337758169645356	11/18	561	Mastercard	anthony41@meed.com	Drilling engineer
2	Unit 0065 Box 5062/nDPO AP 27450	94 ve	PM	Mozilla/5.0 (compatible; MSIE 8.0; Windows NT...	Simpson, Williams and Plam	675957666125	08/18	689	JCB 16 digit	amymiller@monales- hamilton.com	Customer service manager
3	7780 Julia Ford/nNew Stacy, WA 45788	36 on	PM	Mozilla/5.0 (Macintosh; Intel Mac OS X 10.8.0 ...	Williams, Marshall and Buchanan	6011578584430710	02/24	384	Discover	brent16@olson-robinson.info	Drilling engineer
4	23012 Munoz Drive Suite 337/nNew Cynthis, TX S...	20 VE	AM	Opera/9.58(X11; Linux x86_64; sr 17) Presto/2...	Brown, Watson and Andreas	6011456623207998	10/25	678	Diners Club / Carte Blanche	christophenwright@gmail.com	Fire artist
5	7502 Powell Mission Apt. 768/nTravisland, VA E...	21 XT	PM	Mozilla/5.0 (Macintosh; U; PPC Mac OS X 10.8.5...	Siva- Anderson	30246185196287	07/25	7168	Discover	ynsguyen@gmail.com	Fish farm manager
6	93971 Conway Causeway/nAndersonburgh, AZ 75107	96 M	AM	Mozilla/5.0 (compatible; MSIE 7.0; Windows NT...	Gibson and Sons	6011388782555569	07/24	714	VISA 16 digit	olivia04@yahoo.com	Dancer
7	260 Rachel Plains Suite 360/nCastroberg, WV 24...	96 pG	PM	Mozilla/5.0 (X11; Linux i686) AppleWebKit/5350...	Marshall- Collins	561252141908	06/25	256	VISA 13 digit	phillip48@parks.info	Event organizer
8	2129 Dylan Burg/nNew Michelle, ME 28650	45 M	PM	Mozilla/5.0 (Macintosh; U; Intel Mac OS X 10.7...	Gaffney and Sons	180041795790001	04/24	895	JCB 16 digit	kdavis@raumussen.com	Financial manager
9	2795 Dawson Extensions/nLake Tinsfort, ID 88738	15 Ug	AM	Mozilla/5.0 (X11; Linux i686) rv:1.9.2.20; Ge...	Rivers, Buthanan and Ramirez	4396283918171	01/17	931	American Express	spoleman@hunt-huerta.com	Forensic scientist

```
[14]: # basic
# 1.Display last 10 rows in the Dataset
```

Data Cleaning

- Handling missing values in the dataset

```
•[15]: # basic
      # 2.Check Null values in the dataset

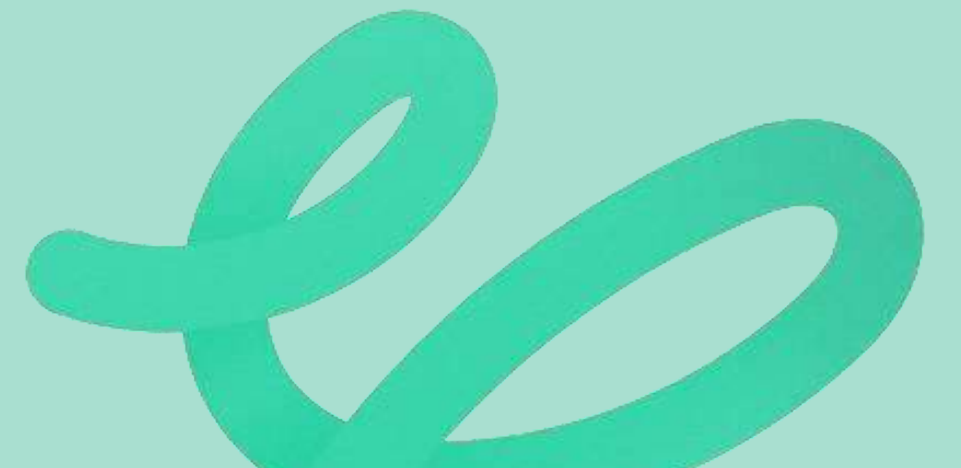
      data.isnull().sum()
```

```
[15]: Address      0
      Lot          0
      AM or PM     0
      Browser Info  0
      Company      0
      Credit Card   0
      CC Exp Date   0
      CC Security Code 0
      CC Provider   0
      Email         0
      Job           0
      IP Address    0
      Language      0
      Purchase Price 0
      dtype: int64
```



Exploratory Data Analysis

- Exploratory Data Analysis (EDA) is an essential step in the data analysis process, as it helps us understand the structure of the dataset, identify patterns, and generate hypotheses for further analysis.
- In this project, we conducted a thorough EDA of the 'E-Commerce Purchase' dataset, using Pandas to perform basic and intermediate analyses.



Basic Analysis

- We counted the number of people with French as their language and job titles containing "Engineer."

```
•[26]: # basic  
# 9.How many people have French 'fr' as their language ?
```

```
len(data[data['Language']=='fr'])
```

```
[26]: 1097
```

```
•[35]: # basic  
# 10.How many people's job title contains Engineer ?
```

```
len(data[data['Job'].str.contains('engineer',case=False)])
```

```
[35]: 984
```

Intermediate Analysis

- We identified people with Mastercard as their credit card provider who made purchases above a certain threshold.
- We counted the number of people with a credit card that expires in 2020.
- We identified the top 5 most popular email providers.

```
•[22]: # Intermediate
# 1.How many people have Mastercard as their credit card provider and made a purchase above 50 ?

len(data[(data['CC Provider']=="Mastercard") & (data['Purchase Price']>50)])

[22]: 405

•[34]: # Intermediate
# 2.How many people have a credit card that expires in 2020 ?

len(data[data['CC Exp Date'].apply(lambda x:x[3:]=='20')])

[34]: 988

•[35]: # Intermediate
# 3.Top 5 most populare Email providers (e.g. gamil.com , yahoo.com etc...)

list1=[]
for email in data['Email']:
    list1.append(email.split('@')[1])

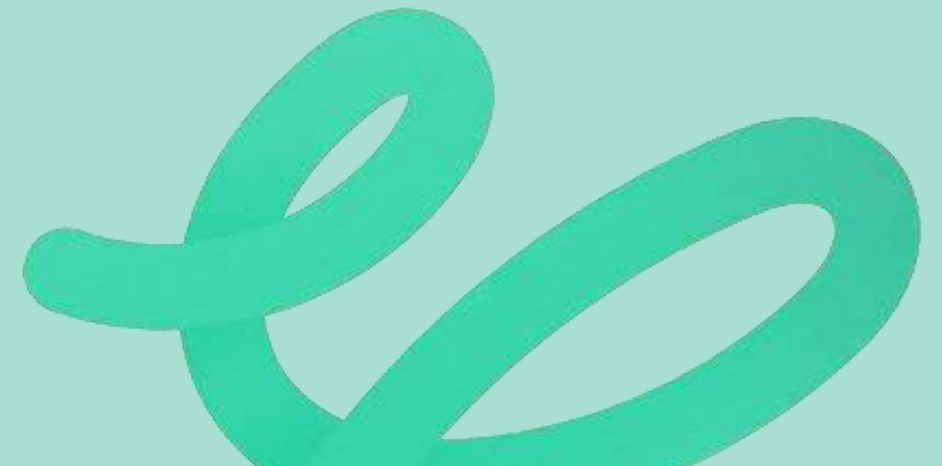
[36]: data['temp']=list1 # creating a new column in dataset

[38]: data['temp'].value_counts().head()

[38]: temp
hotmail.com    1638
yahoo.com      1616
gmail.com      1605
smith.com       42
williams.com    37
Name: count, dtype: int64
```

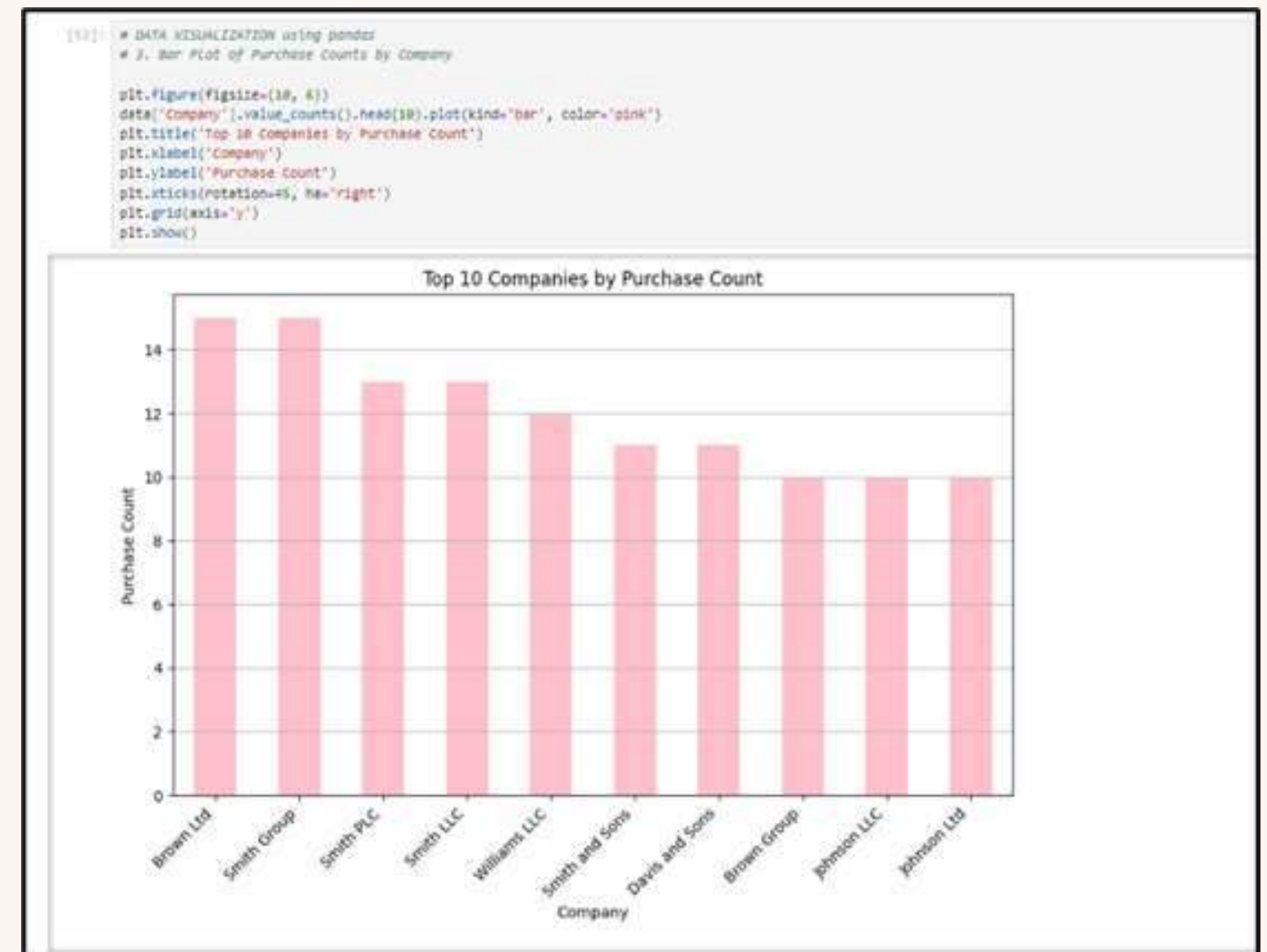


Data Visualization

- Data visualization plays a crucial role in understanding complex datasets, as it allows us to visually explore patterns, trends, and relationships that may not be apparent from the raw data alone.
 - In this project, we used various data visualization techniques to enhance our understanding of the 'E-Commerce Purchase' dataset and communicate our findings effectively.
- 

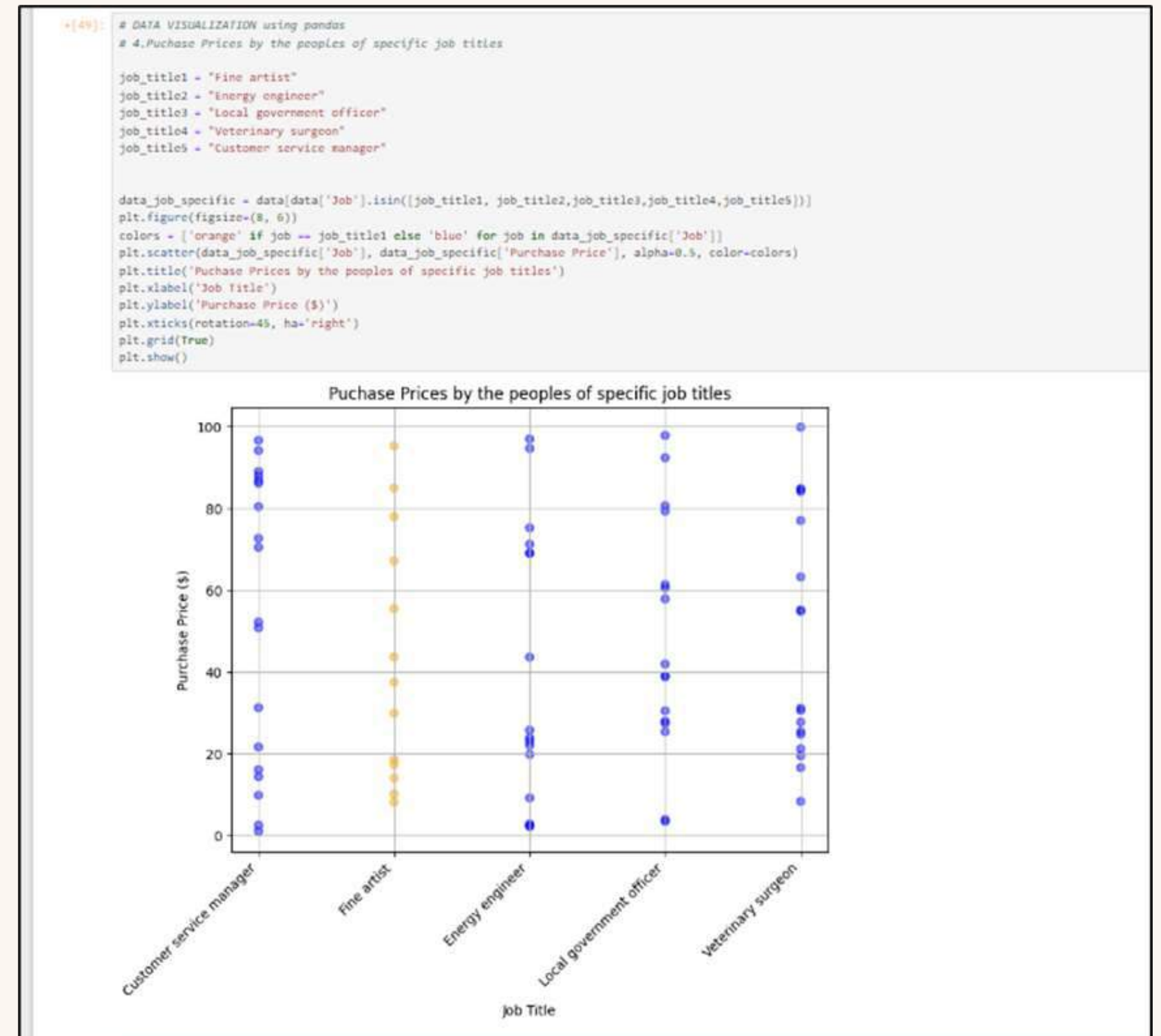
Bar Plot of Purchase Counts by Company

- This visualization allows us to compare the purchase counts across different companies.
- It contributes to the analysis by highlighting the popularity of certain companies among customers.




Scatter Plot of Purchase Prices by Specific Job Titles

- The scatter plot helps us visualize the relationship between purchase prices and specific job titles.
- It contributes to the analysis by identifying any patterns or trends in purchasing behavior based on occupation.





Conclusion

- The analysis of the 'E-Commerce Purchase' dataset using Pandas has provided valuable insights into customer behavior and trends within the e-commerce platform. Through data cleaning, exploratory data analysis (EDA), and data visualization, we were able to uncover key patterns and trends that can inform business decisions and strategies.
 - In conclusion, the analysis of the 'E-Commerce Purchase' dataset has demonstrated the power of data analysis and visualization in extracting meaningful insights from complex datasets. By leveraging these insights, businesses can make informed decisions that drive growth and enhance the customer experience.
- 

Thank
you very
much!

