

A person wearing a dark suit is seated at a wooden desk, working on a silver laptop. The laptop screen displays a dashboard with several data visualizations, including bar charts and a line graph. The person's hands are visible, with one hand pointing at the screen and the other resting on the keyboard. A smartphone is lying on the desk next to the laptop. The background is slightly blurred, showing a modern office environment. The image is partially overlaid by a large purple triangle on the right side, which contains the main title.

**Data Analysis with Python**

# **E-COMMERCE PURCHASE CASE STUDY**

# Import the file

```
: import pandas as pd
import zipfile

# Unzipping and loading the dataset from Kaggle
with zipfile.ZipFile('C:\\Users\\HP\\Downloads\\Ecommerce Purchases (1).zip', 'r') as zip_ref:
    zip_ref.extractall('Downloads//')

# Load the CSV file after extracting
df = pd.read_csv('C:\\Users\\HP\\Downloads\\Ecommerce Purchases (1).zip')

# Previewing the data
df.head()
```

|   | Address   | Lot      | AM<br>or<br>PM | Browser<br>Info   | Company             | Credit Card      | CC<br>Exp<br>Date | CC<br>Security<br>Code | CC<br>Provider  | Email             | Job  | IP Address     |
|---|---|----------|----------------|---|---------------------|------------------|-------------------|------------------------|-----------------|-------------------|--|----------------|
| 0 | 16629 Pace Camp<br>Apt.<br>448\\nAlexisborough,<br>NE 77... | 46<br>in | PM             | Opera/9.56.<br>(X11; Linux<br>x86_64; sl-<br>SI)<br>Presto/2... | Martinez-<br>Herman | 6011929061123406 | 02/20             | 900                    | JCB 16<br>digit | pdunlap@yahoo.com | Scientist,<br>product/process<br>development | 149.146.147.20 |
|   | 9374 Jasmine Spurs  |          |                | Opera/8.93.<br>(Windows   | Fletcher,           |                  |                   |                        |                 |                   |  |                |

# Problem

DISPLAY TOP 10 ROWS OF THE DATA SET

```
df.head(10)
```

|   | Address   | Lot      | AM<br>or<br>PM | Browser Info   | Company                                  | Credit Card      | CC<br>Exp<br>Date | CC<br>Security<br>Code | CC<br>Provider  |
|---|---|----------|----------------|--|--|------------------|-------------------|------------------------|-----------------|
| 0 | 16629 Pace Camp Apt.<br>448\nAlexisborough, NE<br>77... | 46<br>in | PM             | Opera/9.56.(X11;<br>Linux x86_64; sl-SI)<br>Presto/2...  | Martinez-<br>Herman                      | 6011929061123406 | 02/20             | 900                    | JCB 16<br>digit |
| 1 | 9374 Jasmine Spurs Suite<br>508\nSouth John, TN 8...    | 28<br>rn | PM             | Opera/8.93.<br>(Windows 98; Win<br>9x 4.90; en-US) Pr... | Fletcher,<br>Richards<br>and<br>Whitaker | 3337758169645356 | 11/18             | 561                    | Mastercard      |

# Problem

check last 10 rows of the data set

```
df.tail(10)
```

|      | Address   | Lot      | AM<br>or<br>PM | Browser Info   | Company                             | Credit Card      |
|------|---|----------|----------------|--|-------------------------------------|------------------|
| 9990 | 75731 Molly<br>Springs\nWest<br>Danielle, VT 96934-<br>5102 | 93<br>ty | PM             | Mozilla/5.0<br>(Macintosh; Intel<br>Mac OS X<br>10_7_4;... | Pace,<br>Vazquez<br>and<br>Richards | 869968197049750  |
| 9991 | PSC 8165, Box<br>8498\nAPO AP 60327-<br>0346                | 50<br>dA | AM             | Mozilla/5.0<br>(compatible; MSIE<br>8.0; Windows NT        | Snyder<br>Inc                       | 4221582137197481 |

# Problem

check data type of each column.

```
df.dtypes
```

|                  |         |
|------------------|---------|
| Address          | object  |
| Lot              | object  |
| AM or PM         | object  |
| Browser Info     | object  |
| Company          | object  |
| Credit Card      | int64   |
| CC Exp Date      | object  |
| CC Security Code | int64   |
| CC Provider      | object  |
| Email            | object  |
| Job              | object  |
| IP Address       | object  |
| Language         | object  |
| Purchase Price   | float64 |
| dtype:           | object  |



# Problem

## 4.check null values in the data set

```
df.isnull().sum()

]: Address      0
   Lot          0
   AM or PM     0
   Browser Info 0
   Company      0
   Credit Card  0
   CC Exp Date  0
   CC Security Code 0
```

# Problem

how many rows and columns are there in our data set

```
len(df.columns)

: 14

: len(df)

: 10000
```

# Problem

## Highest and Lowest purchase price

```
df.columns
```

```
Index(['Address', 'Lot', 'AM or PM', 'Browser Info', 'Company', 'Credit Card',  
      'CC Exp Date', 'CC Security Code', 'CC Provider', 'Email', 'Job',  
      'IP Address', 'Language', 'Purchase Price'],  
      dtype='object')
```

```
df['Purchase Price'].max()
```

```
np.float64(99.99)
```

```
df['Purchase Price'].min()
```

```
np.float64(0.0)
```

## Average purchase price

```
df['Purchase Price'].mean()
```

```
np.float64(50.347302)
```

# Problem

# Problem

How many people have french 'fr' as their language

```
In: df.columns
Out: Index(['Address', 'Lot', 'AM or PM', 'Browser Info', 'Company', 'Credit Card',
          'CC Exp Date', 'CC Security Code', 'CC Provider', 'Email', 'Job',
          'IP Address', 'Language', 'Purchase Price'],
          dtype='object')

In: len(df[df['Language']=='fr'])
Out: 1097
```

Job title contains engineer

# Problem

df[df['Job'].str.contains('engineer',case=False)]

|   | CC Exp Date | CC Security Code | CC Provider | Email              | Job               | IP Address   | La |
|---|-------------|------------------|-------------|--------------------|-------------------|--------------|----|
| 6 | 11/18       | 561              | Mastercard  | anthony41@reed.com | Drilling engineer | 15.160.41.51 |    |

# Problem

**find email of the person with the following ip address:  
132.207.160.22**

```
df.columns

Index(['Address', 'Lot', 'AM or PM', 'Browser Info', 'Company', 'Credit Card',
      'CC Exp Date', 'CC Security Code', 'CC Provider', 'Email', 'Job',
      'IP Address', 'Language', 'Purchase Price'],
      dtype='object')

df[df['IP Address']=='132.207.160.22']['Email']

2    amymiller@morales-harrison.com
Name: Email, dtype: object
```

**How many people have mastercard as their credit card provider and  
a purchase above 50**

```
df.columns

Index(['Address', 'Lot', 'AM or PM', 'Browser Info', 'Company', 'Credit Card',
      'CC Exp Date', 'CC Security Code', 'CC Provider', 'Email', 'Job',
      'IP Address', 'Language', 'Purchase Price'],
      dtype='object')

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

405
```

# Problem



# Problem

Find email of the person with the following credit number-  
4664825258997302

```
df.columns
```

```
Index(['Address', 'Lot', 'AM or PM', 'Browser Info', 'Company', 'Credit Card',  
      'CC Exp Date', 'CC Security Code', 'CC Provider', 'Email', 'Job',  
      'IP Address', 'Language', 'Purchase Price'],  
      dtype='object')
```

```
df[df['Credit Card']==4664825258997302]['Email']
```

```
9992    bberry@wright.net  
Name: Email, dtype: object
```

How many people purchase during 'am' and how many people purchase during 'pm'

# Problem

```
df.columns
```

```
Index(['Address', 'Lot', 'AM or PM', 'Browser Info', 'Company', 'Credit Card',  
      'CC Exp Date', 'CC Security Code', 'CC Provider', 'Email', 'Job',  
      'IP Address', 'Language', 'Purchase Price'],  
      dtype='object')
```

```
df['AM or PM'].value_counts()
```

```
AM or PM  
PM      5068  
AM      4932  
Name: count, dtype: int64
```

# Problem

How many people have credit card that expires in 2020

```
df.columns
```

```
Index(['Address', 'Lot', 'AM or PM', 'Browser Info', 'Company', 'Credit Card',  
      'CC Exp Date', 'CC Security Code', 'CC Provider', 'Email', 'Job',  
      'IP Address', 'Language', 'Purchase Price'],  
      dtype='object')
```

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

```
988
```

Top 5 most popular in email providers (e.g, gmail.com,yahoo.com etc)

# Problem

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

```
df['temp']=list1
```

```
df['temp'].value_counts().head()
```

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

# created by

**Arpita Chakroborty**

---

**arpitachakroborty.official@gmail.com**

---

