

# SYNTHETIC CUSTOMER PROFILING DATA

This report analyzes the customer profiling scenario of retail company for e.g. what age os customers has highest score , is there any relationship between customer age and their purchases ? etc.

## DATASET OVERVIEW

- **ROWS:** 200
- **Columns:** 5
- **Columns Names:** ID, Age, Income in Rs., Score, Purchases
  - **Id** =Unique identifier for each customers
  - **Age** = age of customers
  - **Income in Rs.** =estimated annual income of each customers
  - **Score** =a score from 0 to 100 indicating customer's loyalty towards company
  - **Purchases** =number of purchases by each customers in last 6 months or 1 year

## DATA CLEANING SUMARRY

- Filled 2 empty cells from column- 'Income', 'Score'
- Correct an outlier from 'Score' column,code used-
  - `print(np.where((df['Score']>100) | (df['Score']<0)))`
  - `df.loc[75, 'Score']=100`
  - `df['Score'].fillna(df['Score'].mean().round(2),inplace=True)`
- Converted negative values of 'Age' into positive value
- Fixed data types to int

## EXPLORATORY DATA ANALYSIS

- AVG age of customers are 39,while the minimum age is 18 and maximum age is 59
- AVG annual income of customers are Rs. 52691.97, where around half of the total customers have income around Rs.42266.75
- AVG scores of customers are 50 , with a max score of 100 and minimum of 1
- AVG purchases made in 6 months or 1 year by each customer is around 5 units

```
In [5]: import pandas as pd, numpy as np, matplotlib.pyplot as plt ,seaborn as sns
df=pd.read_csv('cleaned_file.csv')
df.drop('Unnamed: 0',axis=1,inplace=True)
print(df.describe().round(2))
```

	ID	Age	Income in Rs.	Score	Purchases
count	200.00	200.00	200.00	200.00	200.00
mean	100.50	38.72	52691.97	50.44	4.99
std	57.88	12.57	14729.71	31.02	2.20
min	1.00	18.00	20075.00	1.00	1.00
25%	50.75	28.00	42266.75	20.75	3.00
50%	100.50	40.00	52057.50	51.50	5.00
75%	150.25	49.25	63876.25	79.00	7.00
max	200.00	59.00	91016.00	100.00	11.00

- There is no correlation between any of the variables

```
In [6]: print(df.corr().round())
```

	ID	Age	Income in Rs.	Score	Purchases
ID	1.0	0.0	-0.0	-0.0	0.0
Age	0.0	1.0	-0.0	0.0	0.0
Income in Rs.	-0.0	-0.0	1.0	-0.0	0.0
Score	-0.0	0.0	-0.0	1.0	0.0
Purchases	0.0	0.0	0.0	0.0	1.0

- Group wise stats
  - grouping by 'Age' and calculated mean and maximum value of variables 'Income in Rs.','Score'
  - grouping by 'Score' and calculated minimum and maximum value of variables 'Purchases'

```
In [7]: group_1=df.groupby('Age')[['Income in Rs.','Score']].agg(['mean','max']).round(2)
print(group_1)
group_2=df.groupby('Score')['Purchases'].agg(['max','min'])
print("\n\n",group_2)
```



- Age vs INCOME (scatter)
- AGE vs Score (scatter)
- Score vs Purchase (scatter)
- Age Distribution (frequency distribution)
- Age vs purchase (scatter)

```
In [8]: print('Age vs Income in Rs.')
x=df.sort_values(by='Age') # sorting by age in a variable
x.plot(kind='scatter',x='Age',y='Income in Rs.' , xlabel= 'AGE',ylabel='INCOME IN RS.')
plt.title('Age vs Income in Rs.')
plt.show()

print('Age vs Score')
x.plot(kind='scatter',x='Age',y='Score' , xlabel= 'AGE',ylabel='SCORE')
plt.title('Age vs Score')

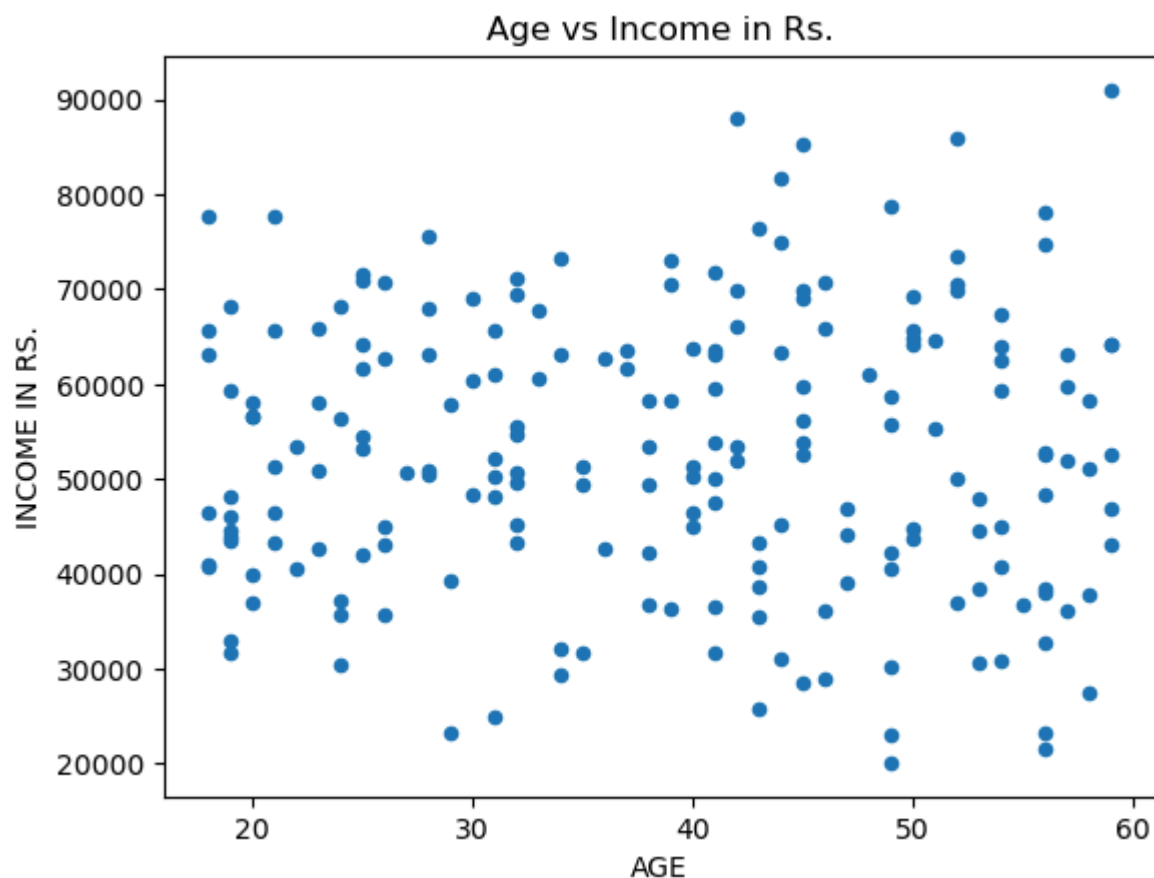
plt.show()

print('Score vs Purchases')
y=df.sort_values(by='Score') # sorting by Score in a variable
y.plot(kind='scatter',x='Purchases',y='Score' , xlabel= 'PURCHASES',ylabel='SCORE')
plt.title('Score vs Purchases')
plt.show()

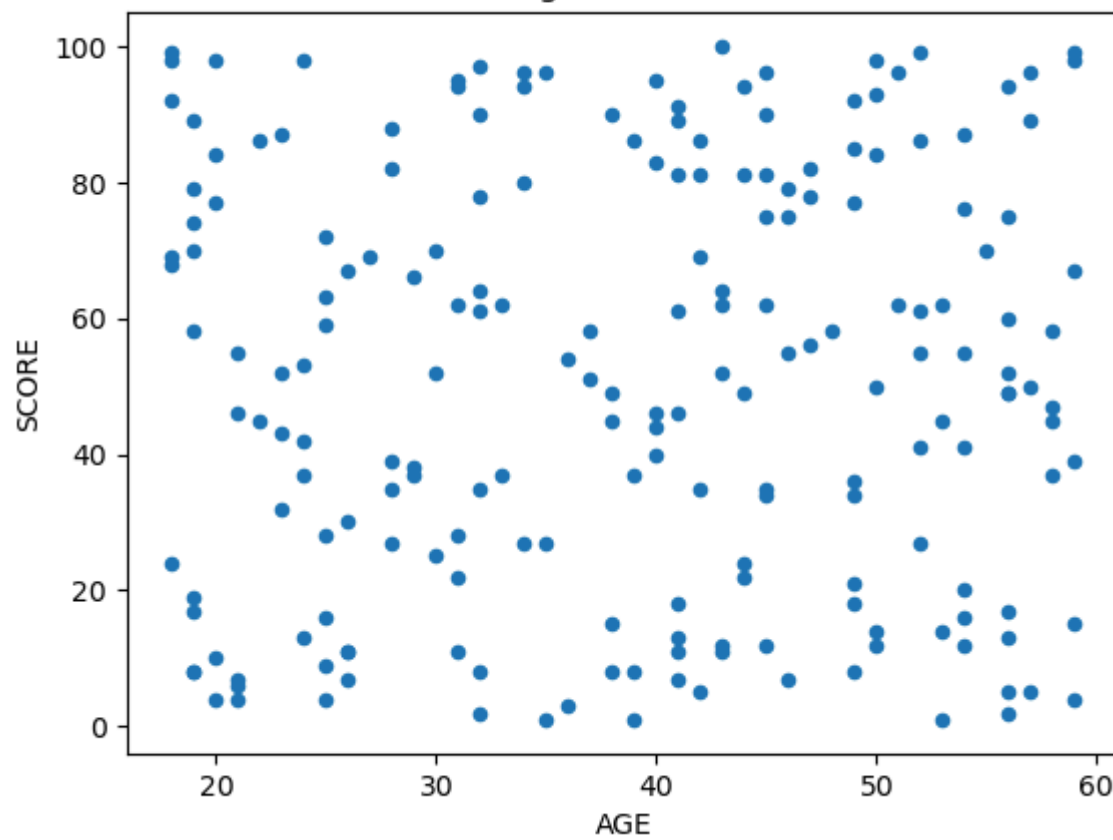
print('Age vs Purchases')
x.plot(kind='scatter',x='Age',y='Purchases' , xlabel= 'AGE',ylabel='Purchases')
plt.title('Age vs Purchases')
plt.show()

print('Age Distribution')
df['Age'].plot(kind='kde',xlabel= 'AGE',ylabel='FREQUENCY')
plt.title('Age Distribution')
plt.show()
```

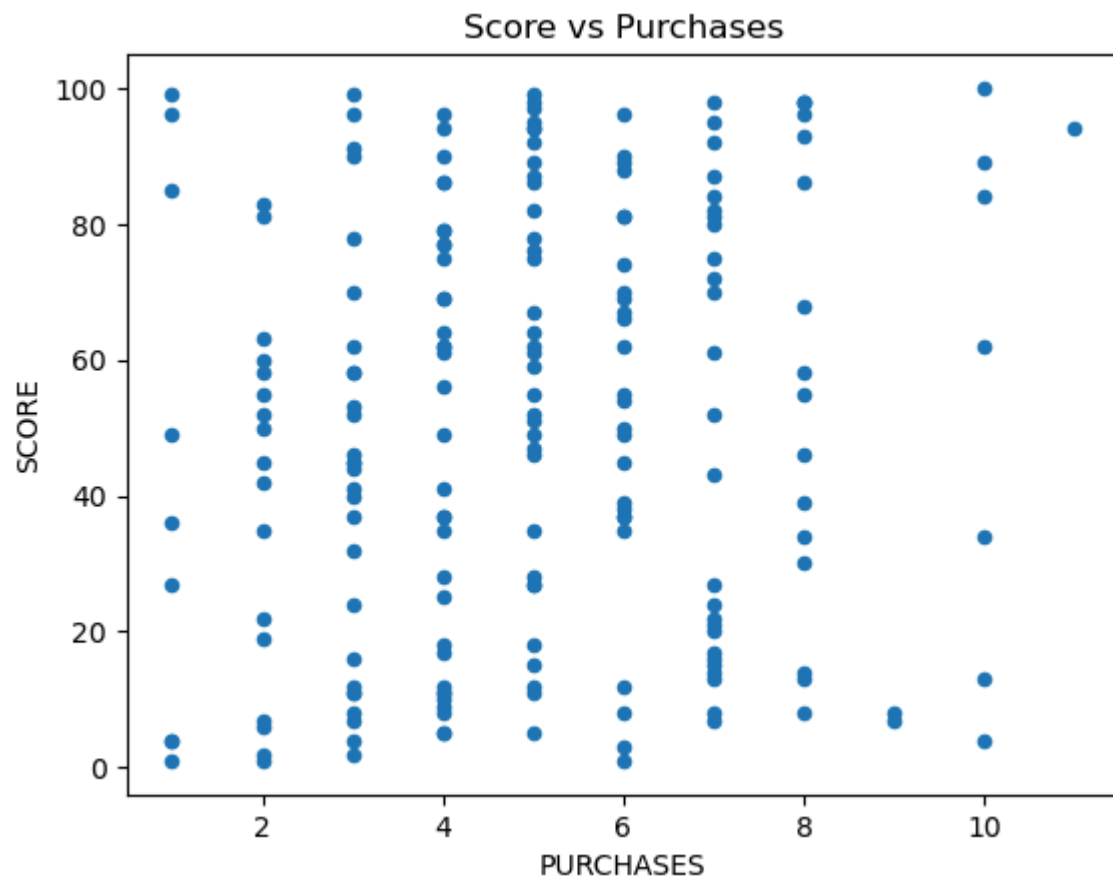
Age vs Income in Rs.



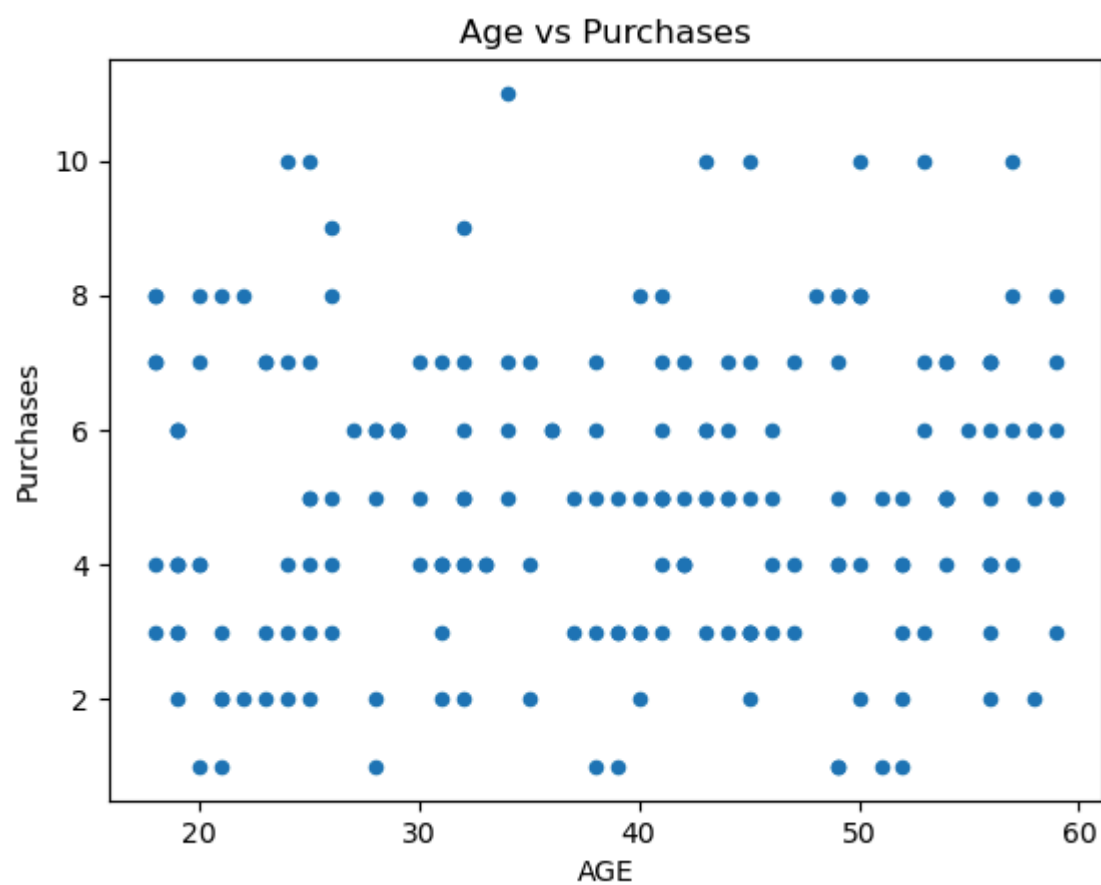
### Age vs Score



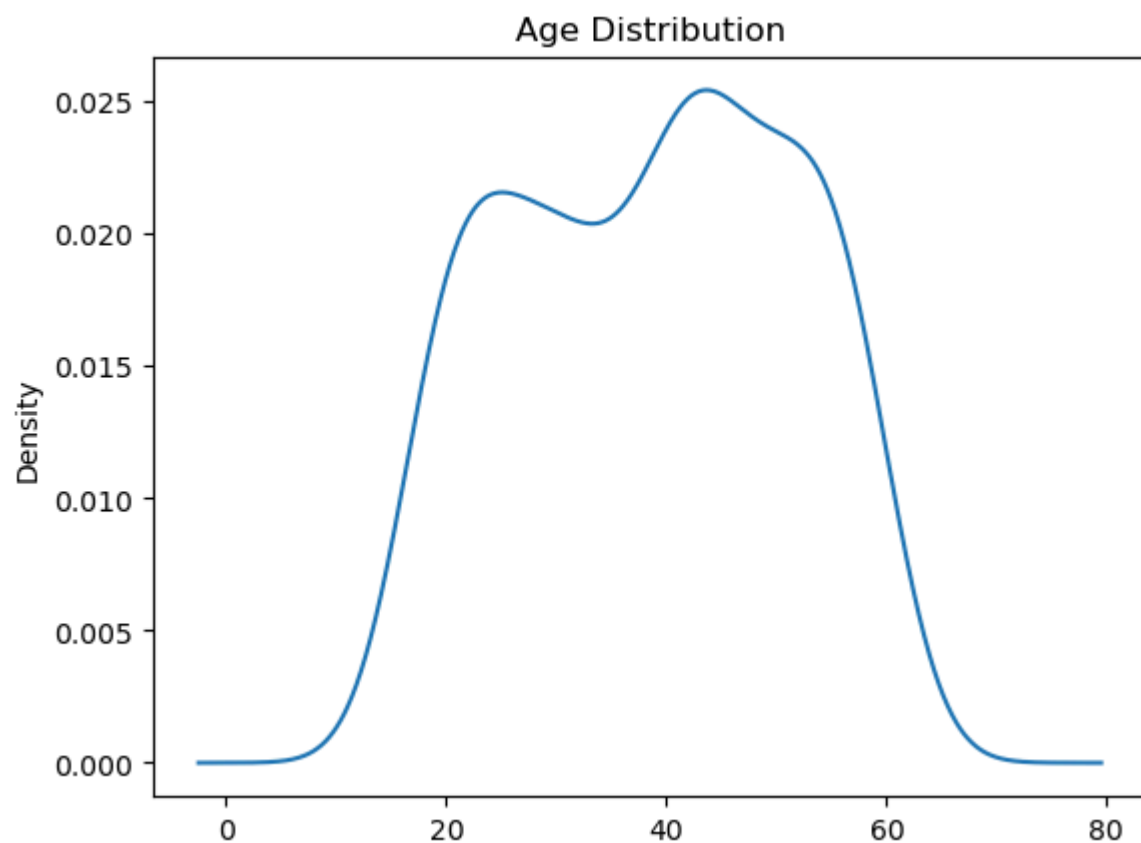
### Score vs Purchases



### Age vs Purchases



Age Distribution



## KEY INSIGHTS

- The company has over half of his customers of age around 40
- The wealthiest customer has annual income of Rs. 91016 , with an avg income of their customers Rs. 52691.97
- On an average customers scored 50 in company's customer loyalty list
- In a year or in 6 months , the maximum sale has gone by a single customer is 11 , while avg sales by

- From the Scatter Graphs and Correlation formula we found that there is no association between any variables