

THỰC HÀNH: PHÂN TÍCH KHÁM PHÁ DỮ LIỆU

1.1. THỐNG KÊ MÔ TẢ

1.1.1. Ôn tập lý thuyết

1.1.2. Bài làm mẫu

1.1.3. Bài tập thực hành 1

Thực hiện thống kê mô tả trên tập dữ liệu về phân loại chất lượng rượu đỏ. Dữ liệu lấy tại

<https://www.kaggle.com/code/eisgandar/red-wine-quality-eda-classification>

Nhiệm vụ 1: Khám phá dữ liệu chất lượng rượu vang đỏ

```
Mean: 10.422983114446529
Median: 10.2
Mode: 9.5 (count: 139 )
Variance: 1.1356473950004693
Standard Deviation: 1.0656675818473926
Max: 14.9 Min: 8.4
60th Percentile: 10.5
Q3 (0.75 quantile): 11.1
Interquartile Range (IQR): 1.5999999999999996
```

Nhiệm vụ 2: Loại bỏ dữ liệu trùng lặp

```
Out[4]:
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	sulphates	alcohol	quality
0	7.4	0.700	0.00	1.9	0.076	11.0	34.0	0.99780	0.56	9.4	5
1	7.8	0.880	0.00	2.6	0.098	25.0	67.0	0.99680	0.68	9.8	5
2	7.8	0.760	0.04	2.3	0.092	15.0	54.0	0.99700	0.65	9.8	5
3	11.2	0.280	0.56	1.9	0.075	17.0	60.0	0.99800	0.58	9.8	6
4	7.4	0.700	0.00	1.9	0.076	11.0	34.0	0.99780	0.56	9.4	5
...
1594	6.2	0.600	0.08	2.0	0.090	32.0	44.0	0.99490	0.58	10.5	5
1595	5.9	0.550	0.10	2.2	0.062	39.0	51.0	0.99512	0.76	11.2	6
1596	6.3	0.510	0.13	2.3	0.076	29.0	40.0	0.99574	0.75	11.0	6
1597	5.9	0.645	0.12	2.0	0.075	32.0	44.0	0.99547	0.71	10.2	5
1598	6.0	0.310	0.47	3.6	0.067	18.0	42.0	0.99549	0.66	11.0	6

1599 rows × 11 columns

Nhiệm vụ 3: Thay thế dữ liệu và thay đổi định dạng của dữ liệu

```
In [5]: # thay giá trị 3,4,5 của quality thành low, Low, medium
wine_data['quality_replaced'] = wine_data['quality'].replace([3,4,5], ['low','low','medium'])
# điền giá trị 0 cho các ô trống trong cột alcohol
wine_data['alcohol'] = wine_data['alcohol'].fillna(0)
# chuyển kiểu dữ liệu alcohol sang int
wine_data['alcohol_changed'] = wine_data['alcohol'].astype(int)
```

Nhiệm vụ 4: Xử lý dữ liệu thiếu

```
In [6]: # đếm số giá trị thiếu của từng cột
wine_data.isnull().sum()
# xóa các hàng có giá trị thiếu
wine_data_withoutna = wine_data.dropna(how='any')
# xem kích thước dữ liệu sau khi xóa
wine_data_withoutna.shape
```

Out[6]: (1599, 14)

1.1.4. Bài tập thực hành 2

Thực hiện thống kê mô tả trên tập dữ liệu về bệnh tiểu đường. Dữ liệu lấy tại

<https://www.kaggle.com/code/vincentlugat/pima-indians-diabetes-eda-prediction-0-906>

Nhiệm vụ 1: Khám phá dữ liệu về bệnh tiểu đường

```
Glucose - Mean: 120.89453125
Glucose - Median: 117.0
Glucose - Mode: 99 (count: 17 )
Glucose - Variance: 1022.2483142519557
Glucose - Std: 31.97261819513622
Glucose - Min: 0 Max: 199 Range: 199
Glucose - 60th percentile: 125.0
Glucose - Q1: 99.0 Q3: 140.25 IQR: 41.25
```

Nhiệm vụ 2: Loại bỏ dữ liệu trùng lặp

```
Out[8]:
```

	Pregnancies	Glucose	BloodPressure	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
0	6	148	72	0	33.6	0.627	50	1
1	1	85	66	0	26.6	0.351	31	0
2	8	183	64	0	23.3	0.672	32	1
3	1	89	66	94	28.1	0.167	21	0
4	0	137	40	168	43.1	2.288	33	1
...
763	10	101	76	180	32.9	0.171	63	0
764	2	122	70	0	36.8	0.340	27	0
765	5	121	72	112	26.2	0.245	30	0
766	1	126	60	0	30.1	0.349	47	1
767	1	93	70	0	30.4	0.315	23	0

768 rows × 8 columns

Nhiệm vụ 3: Thay thế dữ liệu và thay đổi định dạng của dữ liệu

```
In [9]: # Thay giá trị 0/1 của 'Outcome' thành chuỗi 'no_diabetes'/'diabetes'
diabetes['Outcome_replaced'] = diabetes['Outcome'].replace([0,1], ['no_diabetes','diabetes'])
# Điền giá trị thiếu trong cột 'Insulin' bằng 0
diabetes['Insulin'] = diabetes['Insulin'].fillna(0)
# Chuyển kiểu dữ liệu cột 'Insulin' sang int
diabetes['Insulin_changed'] = diabetes['Insulin'].astype(int)
```

Nhiệm vụ 4: Xử lý dữ liệu thiếu

```
Out[10]: (768, 11)
```

1.2. XỬ LÝ VÀ TRỰC QUAN HÓA DỮ LIỆU

1.2.1. Ôn tập lý thuyết

1.2.2. Bài làm mẫu

1.2.3 Bài tập thực hành 1

+ Thực hiện trực quan hóa dữ liệu trên tập dữ liệu về phân loại chất lượng rượu đỏ.

Dữ liệu lấy tại <https://www.kaggle.com/code/eisgandar/red-wine-quality-eda-classification>

Nhiệm vụ 1: Chuẩn bị dữ liệu cho trực quan hóa dữ liệu

1. Chuẩn bị dữ liệu cho trực quan hóa dữ liệu

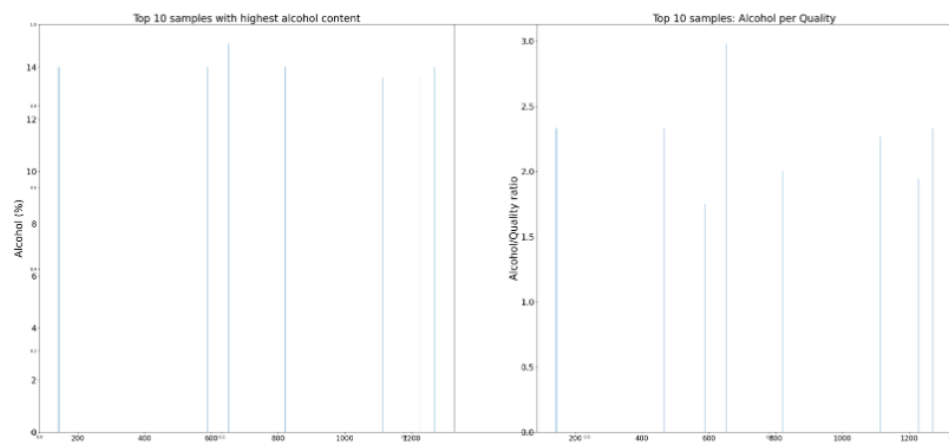
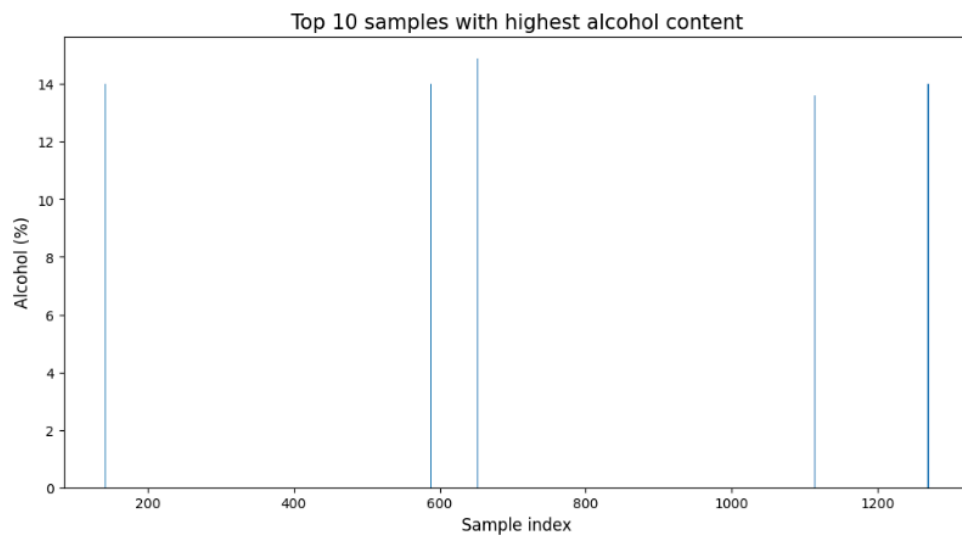
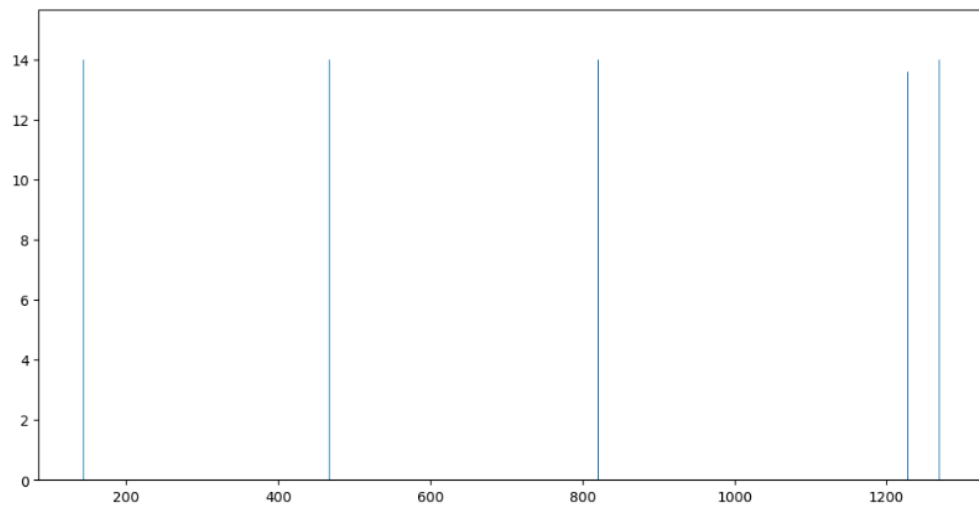
```
In [11]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Đọc file rượu vang đỏ (đặt file csv cùng thư mục làm việc hoặc thay đường dẫn cho đúng)
wine_data = pd.read_csv("winequality-red.csv")

# Chọn các cột cần thiết: alcohol (nồng độ cồn) và quality (điểm chất lượng)
wine_data = wine_data[['alcohol', 'quality']]

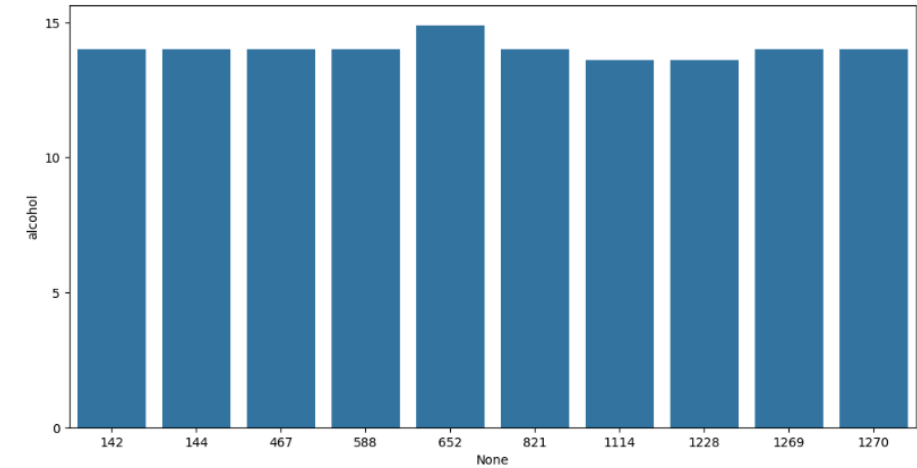
# Tạo biến alcohol_per_quality (nồng độ cồn chia cho điểm chất lượng) để so sánh
wine_data['alcohol_per_quality'] = wine_data['alcohol'] / wine_data['quality']
```

Nhiệm vụ 2: Trực quan hóa dữ liệu với thư viện Matplotlib

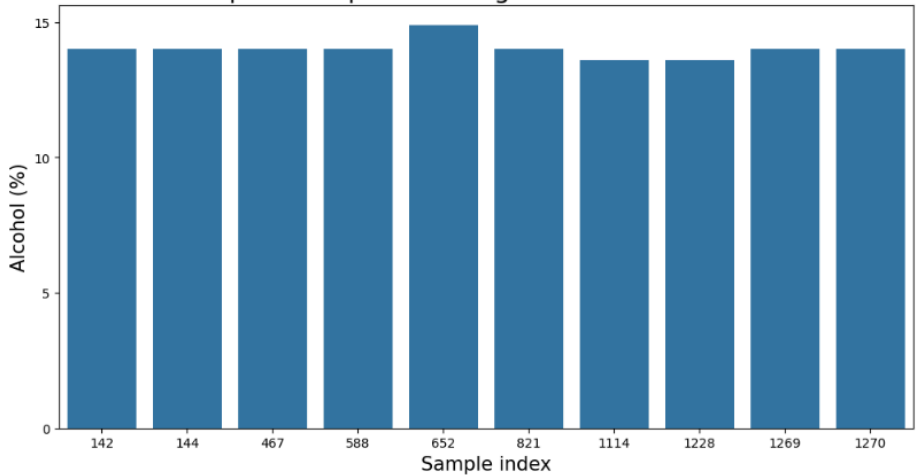


Nhiệm vụ 3: Trắc quan hóa dữ liệu với thư viện Seaborn

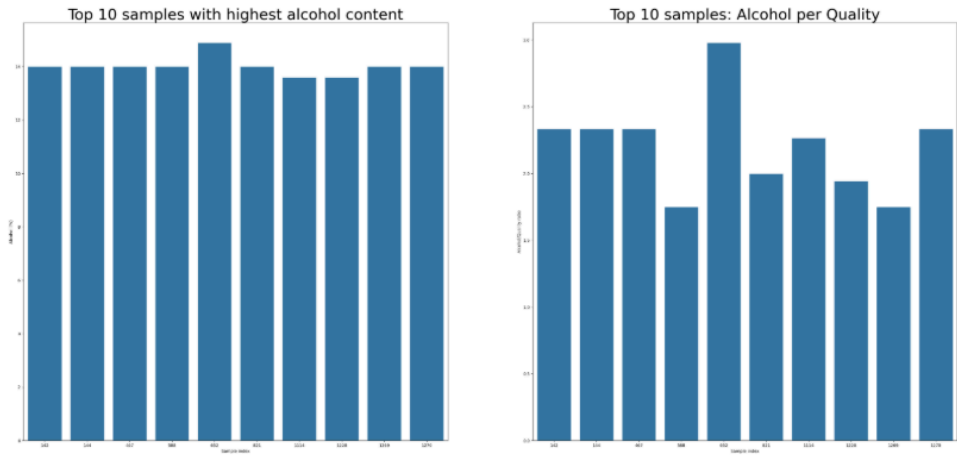
Out[13]: Text(0.5, 1.0, 'Top 10 samples: Alcohol per Quality')



Top 10 samples with highest alcohol content



Sample index



1.2.3. Bài tập thực hành 2

+ Thực hiện trực quan hóa dữ liệu trên tập dữ liệu về bệnh tiểu đường. Dữ liệu lấy tại <https://www.kaggle.com/code/vincentlugat/pima-indians-diabetes-eda-prediction-0-906>

Nhiệm vụ 1: Chuẩn bị dữ liệu cho trực quan hóa dữ liệu

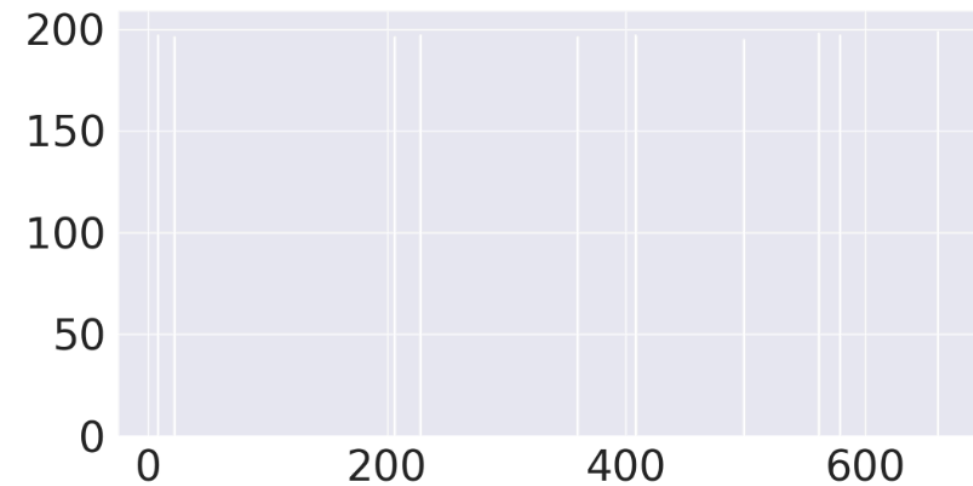
```
In [14]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Đọc file dữ liệu bệnh tiểu đường (đặt file 'diabetes.csv' cùng thư mục làm việc)
diabetes_data = pd.read_csv("diabetes.csv")

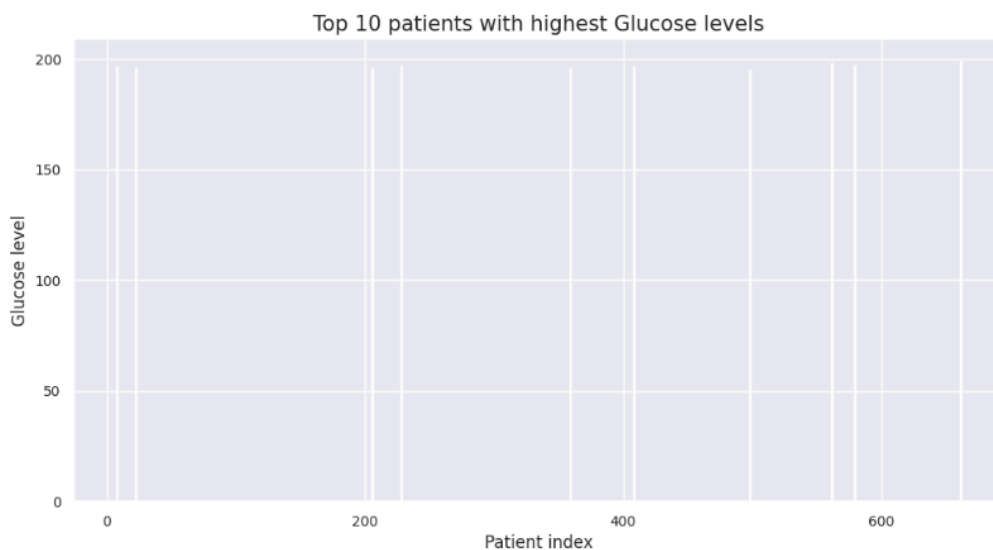
# Chọn các cột quan trọng: Glucose và BMI
diabetes_data = diabetes_data[['Glucose', 'BMI']]

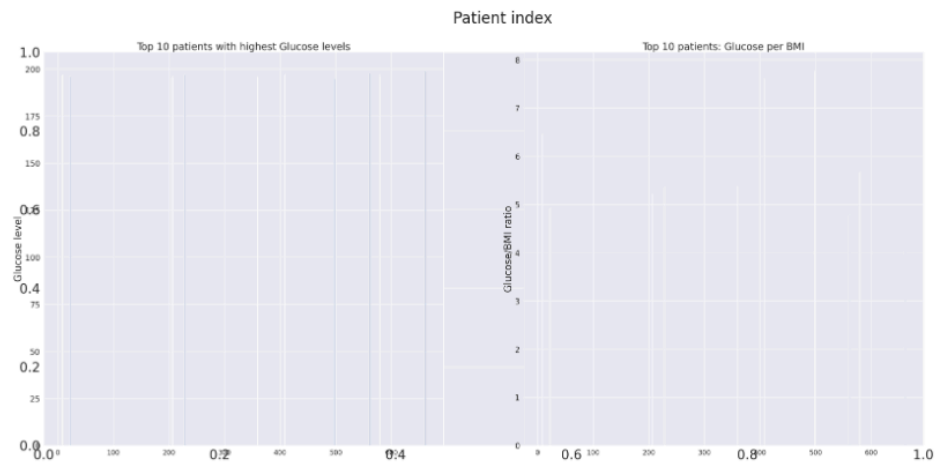
# Tạo biến Glucose_per_BMI (tỉ lệ Glucose/BMI) để so sánh
diabetes_data['Glucose_per_BMI'] = diabetes_data['Glucose'] / diabetes_data['BMI']
```

Nhiệm vụ 2. Trực quan hóa dữ liệu với thư viện Matplotlib



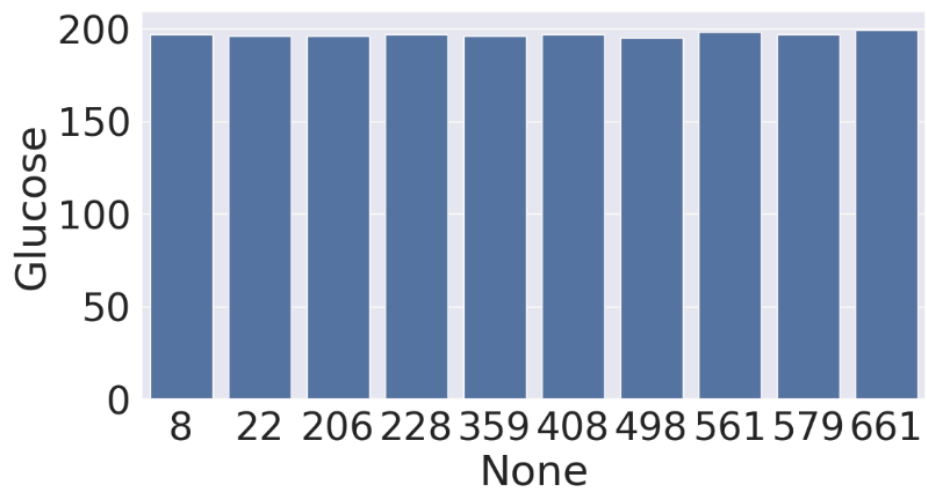
Top 10 patients with highest Glucose levels

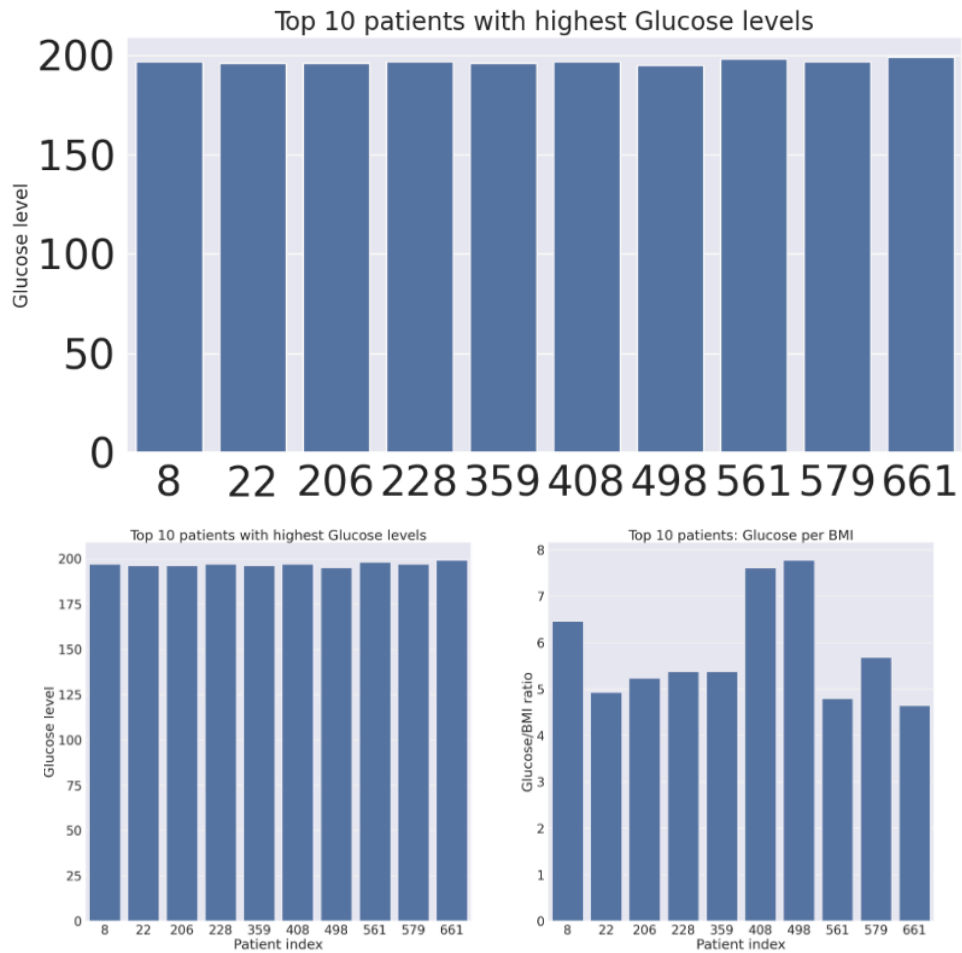




3. Trực quan hóa dữ liệu với thư viện Seaborn

Out[16]: Text(0.5, 1.0, 'Top 10 patients: Glucose per BMI')





+ Thực hiện EDA trên tập dữ liệu mua sắm tại siêu thị. Tập dữ liệu lấy từ <https://www.kaggle.com/code/rajatkumar30/eda-online-retail>


```

# 1. Chuẩn bị dữ liệu cho trực quan hóa
# =====
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Đọc dữ liệu Excel
df = pd.read_excel("Online Retail.xlsx")

# Lọc các cột quan trọng
retail_data = df[['Country', 'Quantity', 'UnitPrice']].copy()

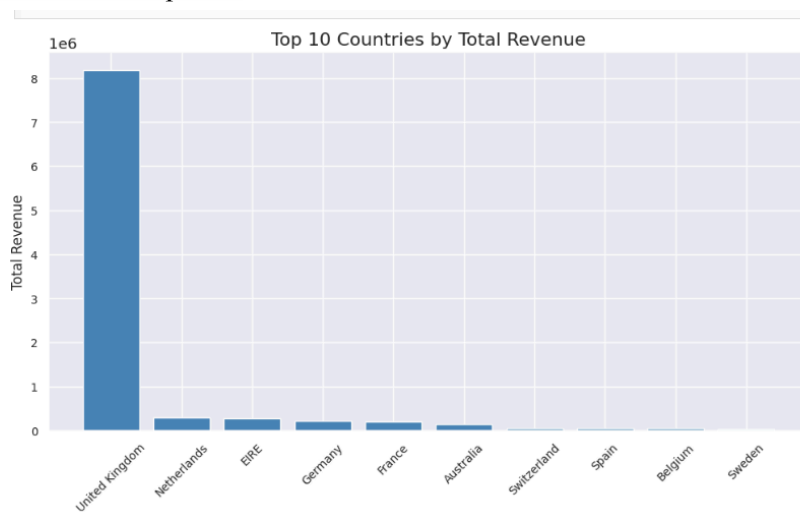
# Tạo cột doanh thu = Quantity * UnitPrice
retail_data['TotalPrice'] = retail_data['Quantity'] * retail_data['UnitPrice']

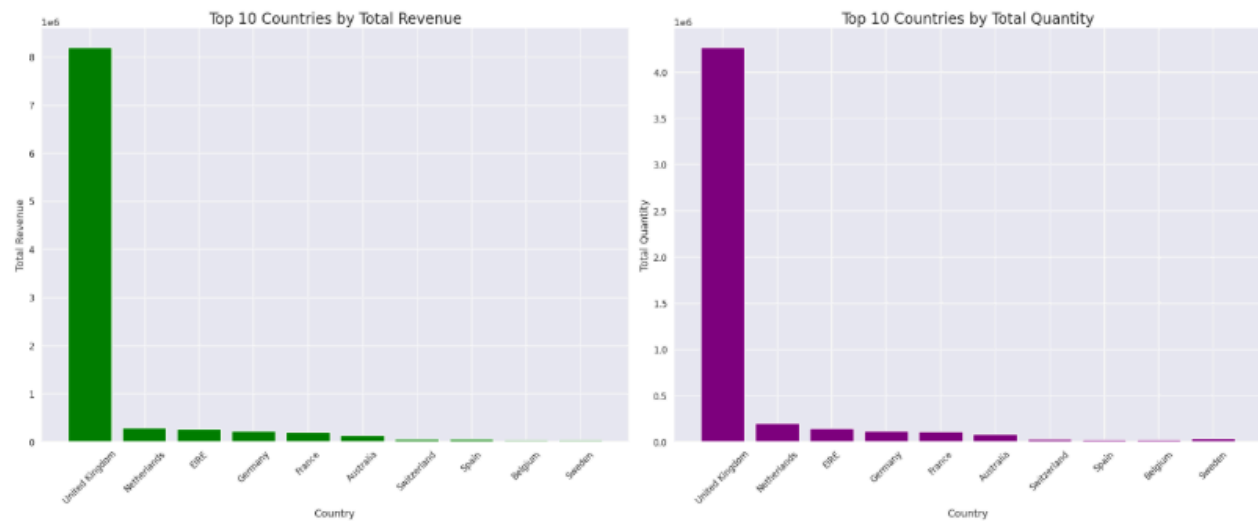
# Gom nhóm theo quốc gia để tính tổng doanh thu và tổng số lượng
country_revenue = retail_data.groupby('Country', as_index=False)['TotalPrice'].sum()
country_revenue = country_revenue.sort_values('TotalPrice', ascending=False)

country_quantity = retail_data.groupby('Country', as_index=False)['Quantity'].sum()
country_quantity = country_quantity.sort_values('Quantity', ascending=False)

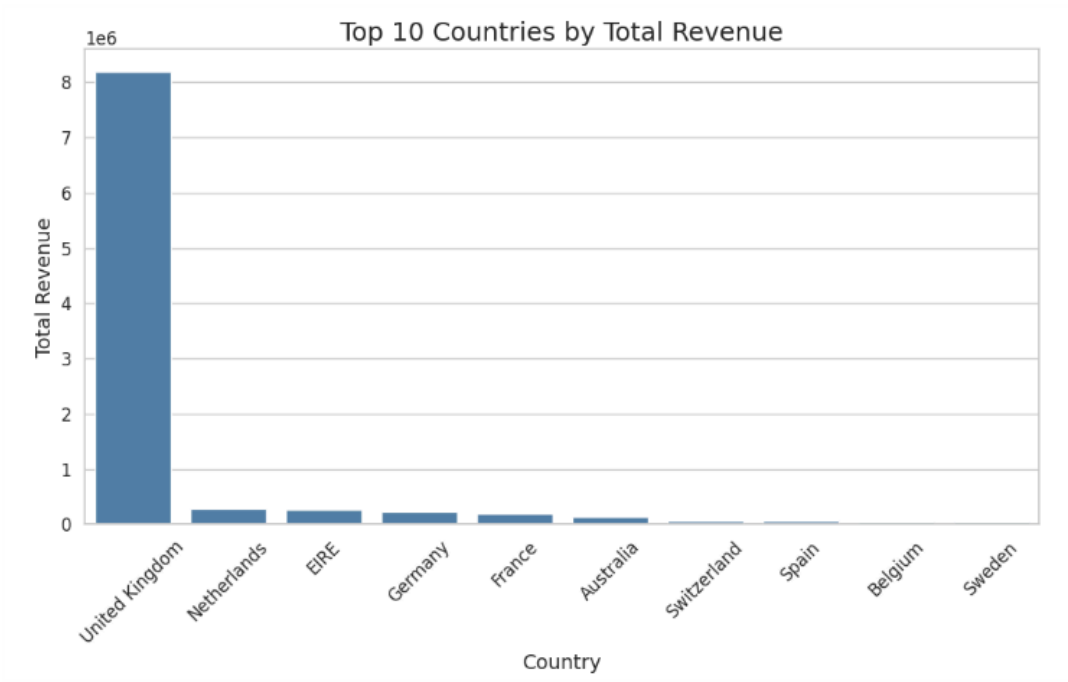
```

2. Trực quan hóa dữ liệu với Matplotlib





3. Trực quan hóa dữ liệu với Seaborn





1.3. Phân tích đơn biến và hai biến

1.3.1. Ôn lý thuyết

1.3.2. Bài làm mẫu

1.3.3. Bài tập thực hành 1

Tìm hiểu các tính năng và cách sử dụng sản phẩm SweetViz (<https://pypi.org/project/sweetviz>) áp dụng trên tập dữ liệu Marketing Campaign

Đọc file (dataset Kaggle dùng tab '\t' làm separator)

```
marketing_data = pd.read_csv("marketing_campaign.csv", sep="\t")
```

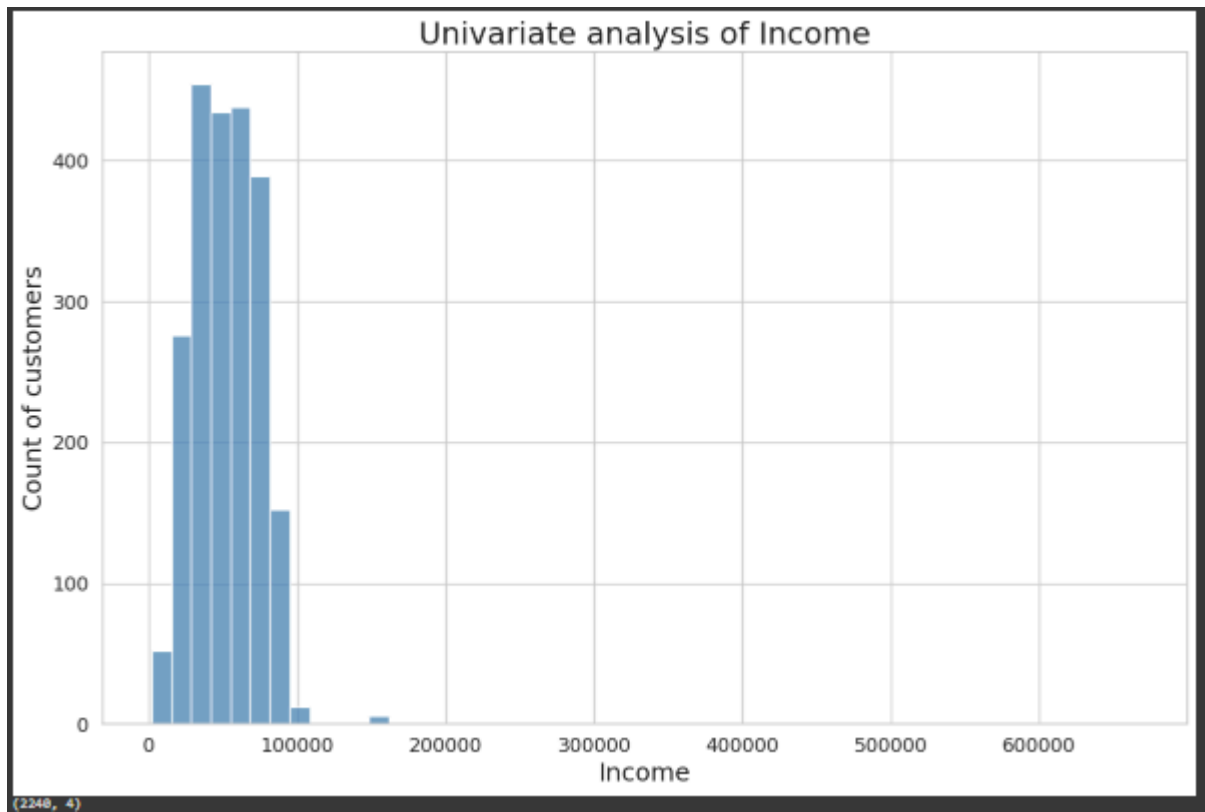
Chọn một số cột tiêu biểu để phân tích

```
marketing_data = marketing_data[['Education', 'Income', 'Kidhome', 'Teenhome']]
```

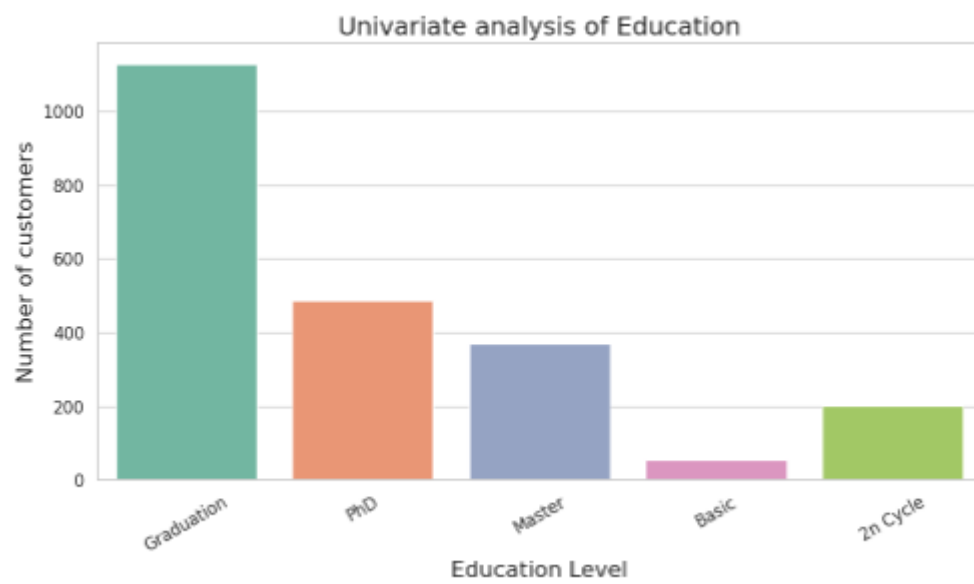
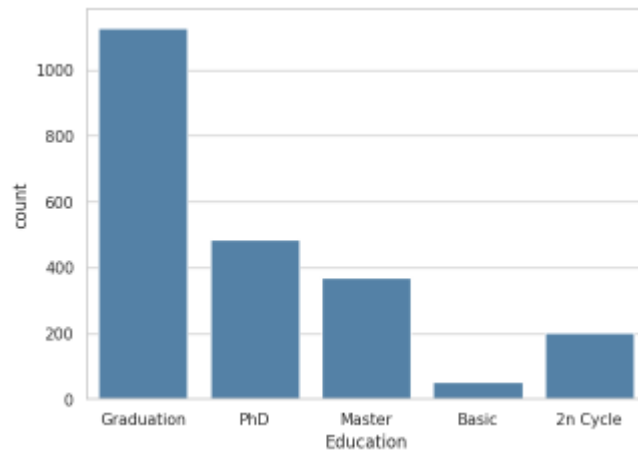
```
marketing_data.head()
```

	Education	Income	Kidhome	Teenhome
0	Graduation	58138.0	0	0
1	Graduation	46344.0	1	1
2	Graduation	71613.0	0	0
3	Graduation	26646.0	1	0
4	PhD	58293.0	1	0

Phân tích đơn biến bằng Histogram (liên tục – Biến thu nhập)

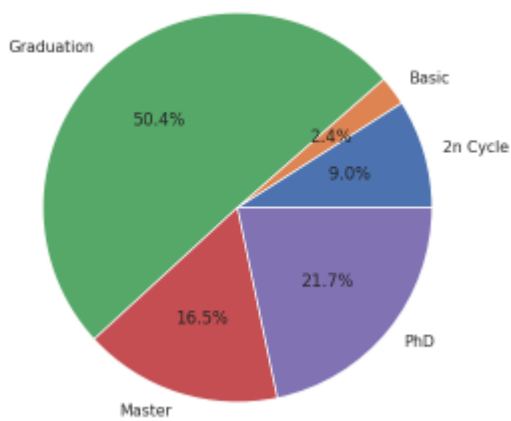


Phân tích đơn biến bằng Bar Chart (rời rạc – Biến education)

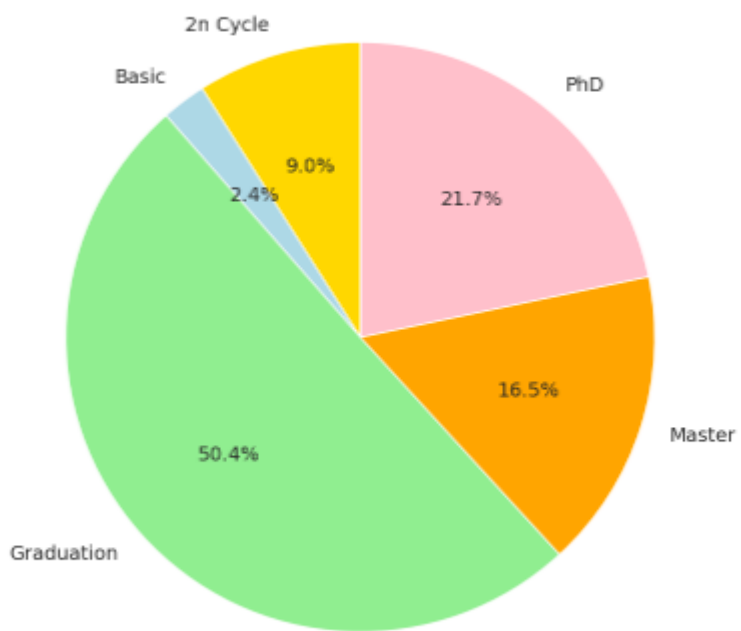


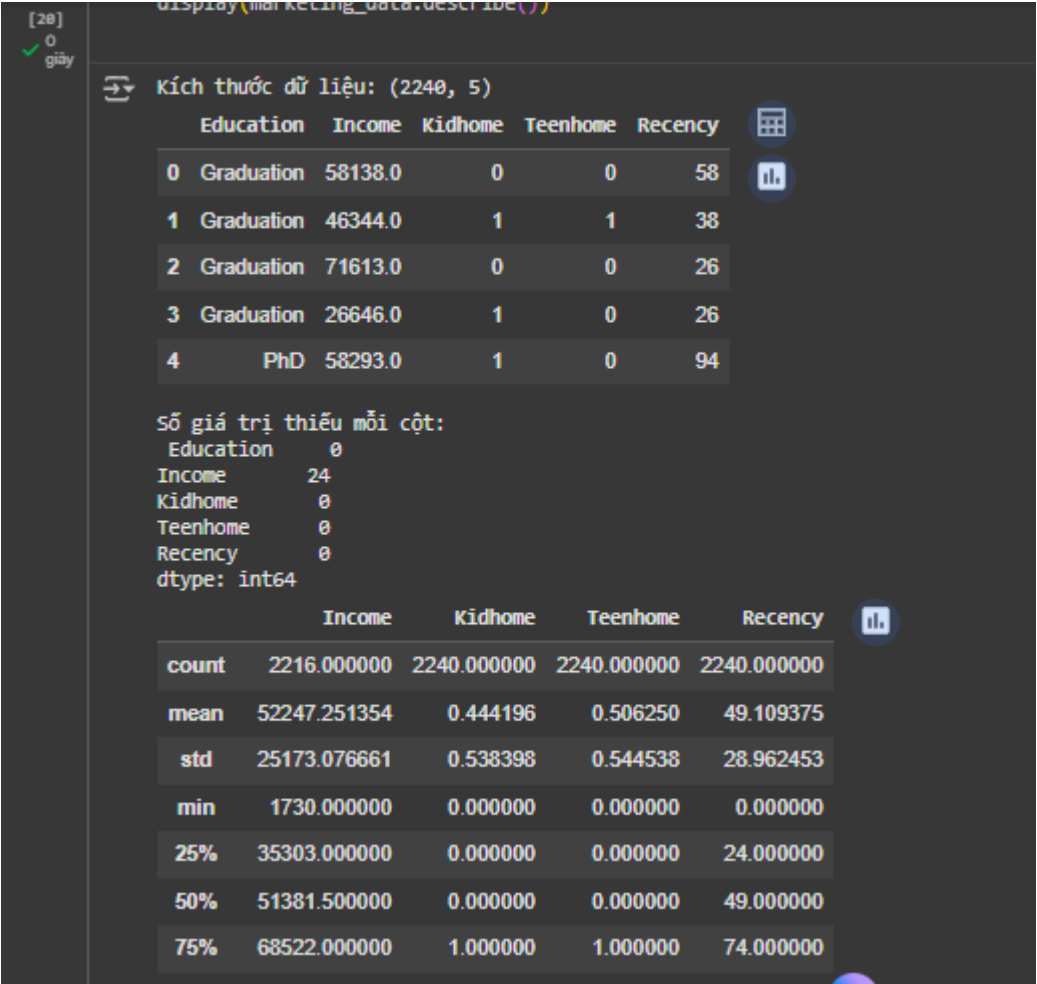
(2248, 4)

Phân tích đơn biến bằng Biểu đồ tròn (Pie Chart – Education)



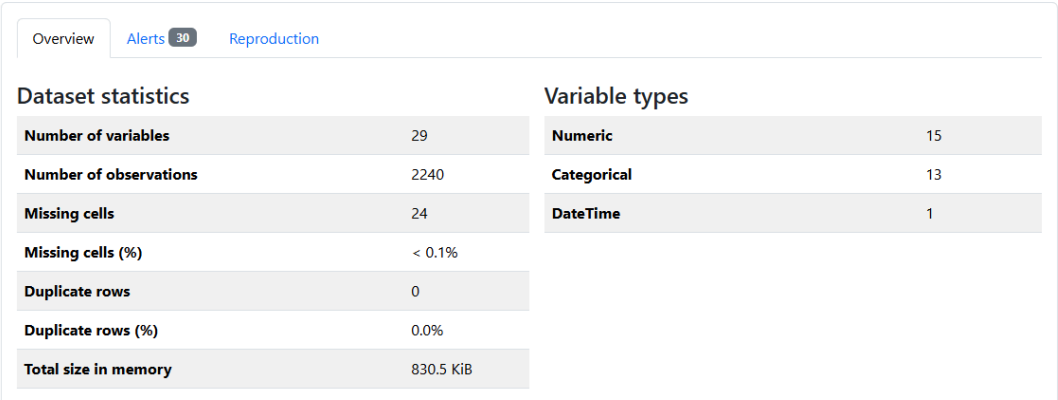
Distribution of Education Levels

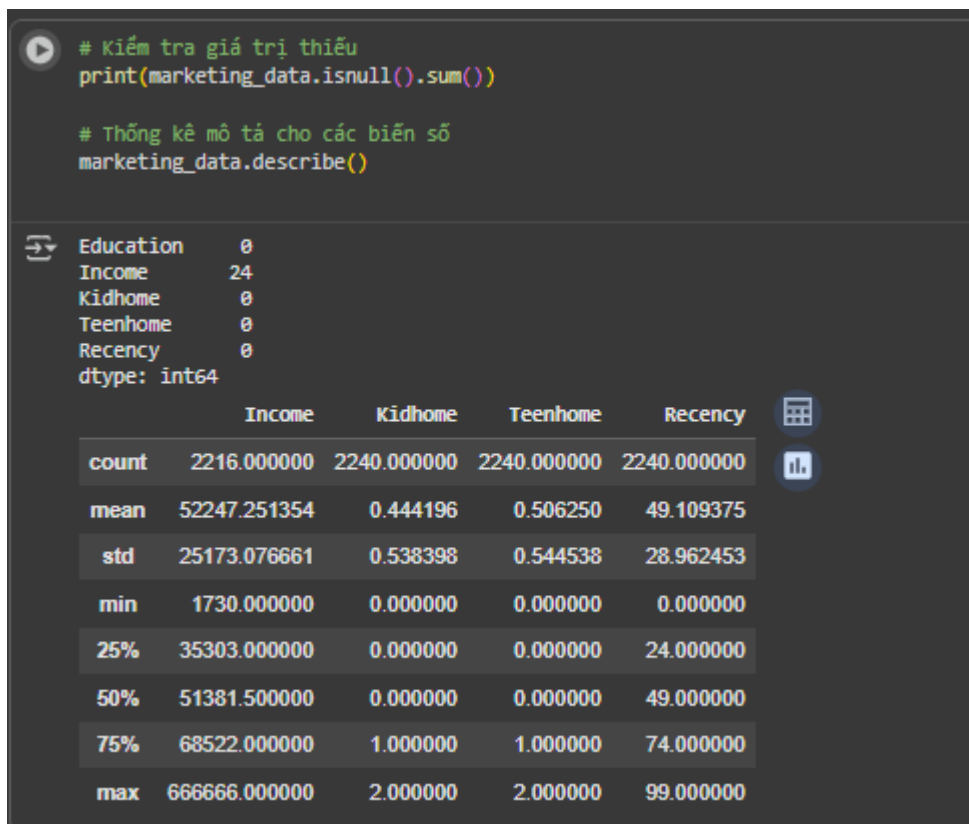




Overview

Brought to you by YData





```
# Kiểm tra giá trị thiếu
print(marketing_data.isnull().sum())
```

```
# Thống kê mô tả cho các biến số
marketing_data.describe()
```


Tìm hiểu các tính năng và cách sử dụng sản phẩm SweetViz

1. Import thư viện và nạp dữ liệu

+ Mã

+ Văn bản

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Đọc file (dataset Kaggle dùng tab '\t' làm separator)
marketing_data = pd.read_csv("marketing_campaign.csv", sep="\t")

# Chọn một số cột tiêu biểu để phân tích
marketing_data = marketing_data[['Education', 'Income', 'Kidhome', 'Teenhome']]
marketing_data.head()
```



	Education	Income	Kidhome	Teenhome
0	Graduation	58138.0	0	0
1	Graduation	46344.0	1	1
2	Graduation	71613.0	0	0
3	Graduation	26646.0	1	0
4	PhD	58293.0	1	0



Các bước tiếp theo: [Tạo mã bằng marketing_data](#) [New interactive sheet](#)

Nhiệm vụ 1: Tạo báo cáo tự động với pandas_profiling / ydata_profiling

Bước 1 – Cài đặt thư viện

[2]
51
giây

```
!pip install ydata-profiling
import pandas as pd
from ydata_profiling import ProfileReport

# Đọc dữ liệu Customer Personality Analysis
marketing_data = pd.read_csv("marketing_campaign.csv", sep="\t")

# Tạo báo cáo EDA tự động
profile = ProfileReport(marketing_data,
                        title="Customer Personality Analysis - Profiling Report",
                        explorative=True)

# Xuất báo cáo ra file HTML
profile.to_file("profile_output.html")
```

[Hiện kết quả đã ẩn](#)

Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.12/dist-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)
[Upgrade to ydata-sdk](#)

Improve your data and profiling with ydata-sdk, featuring data quality scoring, redundancy detection, outlier identification, text validation, and more.

Summarize dataset: 100% 263/263 [00:26<00:00, 6.14it/s, Completed]

0%| | 0/29 [00:00<?, ?it/s]
17%| | 5/29 [00:00<00:00, 34.22it/s]
100%| | 29/29 [00:00<00:00, 76.56it/s]

Generate report structure: 100% 1/1 [00:00<00:00, 6.75s/it]

Render HTML: 100% 1/1 [00:01<00:00, 1.52s/it]

Export report to file: 100% 1/1 [00:00<00:00, 27.45it/s]

1.3.4. Bài tập thực hành

2. Tìm hiểu các tính năng và cách sử dụng sản phẩm AutoViz (<https://pypi.org/project/autoviz/>) áp dụng trên tập dữ liệu Marketing Campaign

#Cài đặt & import

#Nạp dữ liệu

#	Column	Non-Null Count	Dtype
0	ID	2240 non-null	int64
1	Year_Birth	2240 non-null	int64
2	Education	2240 non-null	object
3	Marital_Status	2240 non-null	object
4	Income	2216 non-null	float64
5	Kidhome	2240 non-null	int64
6	Teenhome	2240 non-null	int64
7	Dt_Customer	2240 non-null	object
8	Recency	2240 non-null	int64
9	MntWines	2240 non-null	int64
10	MntFruits	2240 non-null	int64
11	MntMeatProducts	2240 non-null	int64
12	MntFishProducts	2240 non-null	int64
13	MntSweetProducts	2240 non-null	int64
14	MntGoldProds	2240 non-null	int64
15	NumDealsPurchases	2240 non-null	int64
16	NumWebPurchases	2240 non-null	int64
17	NumCatalogPurchases	2240 non-null	int64
18	NumStorePurchases	2240 non-null	int64
19	NumWebVisitsMonth	2240 non-null	int64
20	AcceptedCmp3	2240 non-null	int64
21	AcceptedCmp4	2240 non-null	int64
22	AcceptedCmp5	2240 non-null	int64
23	AcceptedCmp1	2240 non-null	int64
24	AcceptedCmp2	2240 non-null	int64
25	Z_CostContact	2240 non-null	int64
26	Z_Revenue	2240 non-null	int64
27	Response	2240 non-null	int64

dtypes: float64(1), int64(25), object(3)
memory usage: 507.6+ KB
Shape of your Data Set loaded: (2240, 29)

CLASSIFYING VARIABLES #####

Classifying variables in data set...
Number of Numeric Columns = 1
Number of Integer-Categorical Columns = 15
Number of String-Categorical Columns = 2
Number of Factor-Categorical Columns = 0
Number of String-Boolean Columns = 0
Number of Numeric-Boolean Columns = 7
Number of Discrete String Columns = 1
Number of NLP String Columns = 0
Number of Date Time Columns = 0
Number of ID Columns = 1
Number of Columns to Delete = 2
29 Predictors classified...
3 variable(s) removed since they were ID or low-information variables
List of variables removed: ['ID', 'Z_CostContact', 'Z_Revenue']
To fix these data quality issues in the dataset, import FixDQ from autoviz...
All variables classified into correct types.

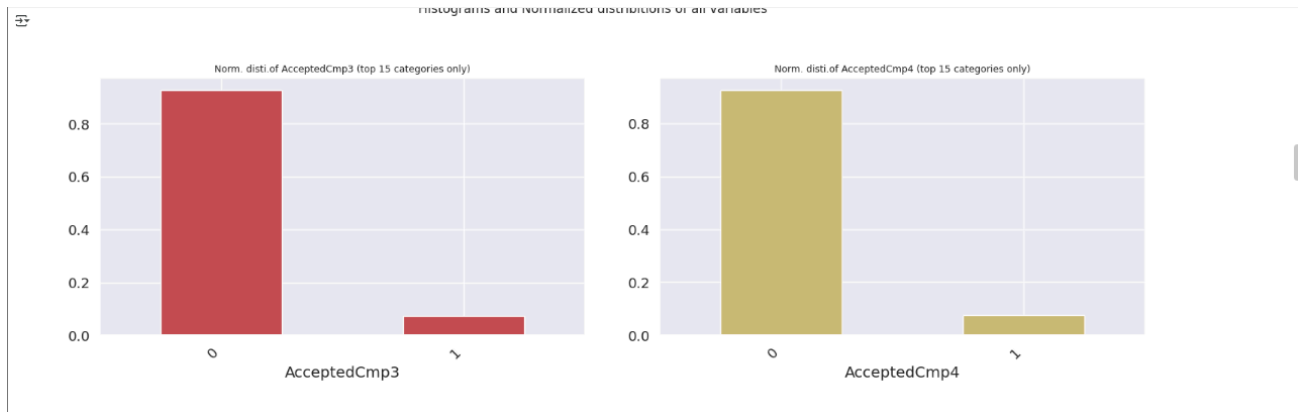
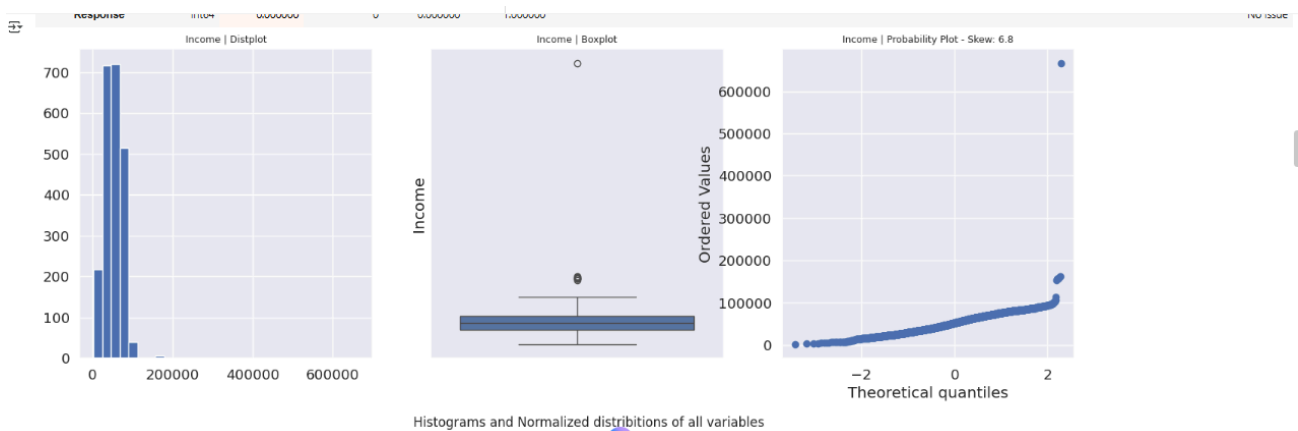
Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value
-----------	-----------------	----------------	---------------	---------------

3 variable(s) removed since they were ID or low-information variables
List of variables removed: ['ID', 'Z_CostContact', 'Z_Revenue']
To fix these data quality issues in the dataset, import FixDQ from autoviz...
All variables classified into correct types.

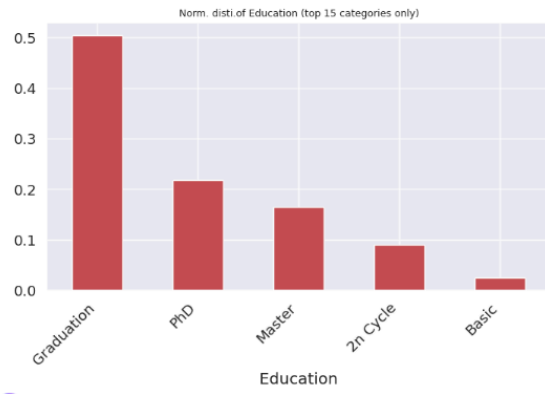
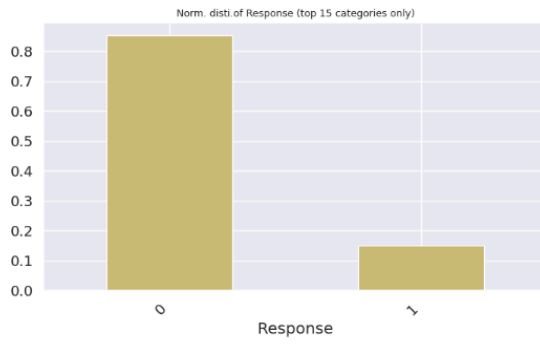
	Data Type	Missing Values%	Unique Values%	Minimum Value	Maximum Value	DQ Issue
ID	int64	0.000000	100	0.000000	11191.000000	Possible ID column: drop before modeling step.
Year_Birth	int64	0.000000	2	1893.000000	1996.000000	Column has 3 outliers greater than upper bound (2004.00) or lower than lower bound(1932.00). Cap them or remove them.
Education	object	0.000000	0			No issue
Marital_Status	object	0.000000	0			3 rare categories: ['Alone', 'Absurd', 'YOLO']. Group them into a single category or drop the categories.
Income	float64	1.071429	NA	1730.000000	666666.000000	24 missing values. Impute them with mean, median, mode, or a constant value such as 123. Column has 6 outliers greater than upper bound (118350.50) or lower than lower bound(-14525.50). Cap them or remove them.
Kidhome	int64	0.000000	0	0.000000	2.000000	No issue
Teenhome	int64	0.000000	0	0.000000	2.000000	No issue
Dt_Customer	object	0.000000	29			Possible high cardinality column with 663 unique values: Use hash encoding or text embedding to reduce dimension.
Recency	int64	0.000000	4	0.000000	99.000000	No issue
MntWines	int64	0.000000	34	0.000000	1493.000000	Column has 35 outliers greater than upper bound (1225.00) or lower than lower bound(-697.00). Cap them or remove them.
MntFruits	int64	0.000000	7	0.000000	199.000000	Column has 227 outliers greater than upper bound (81.00) or lower than lower bound(-47.00). Cap them or remove them.

Cửa sổ dòng lệnh

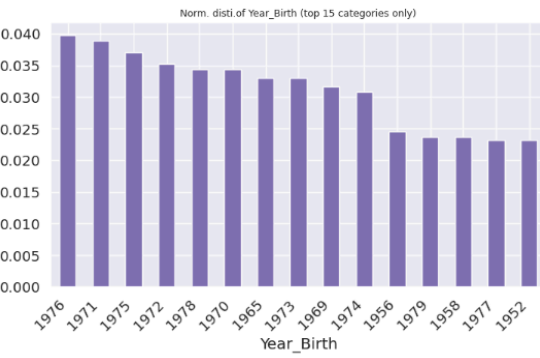
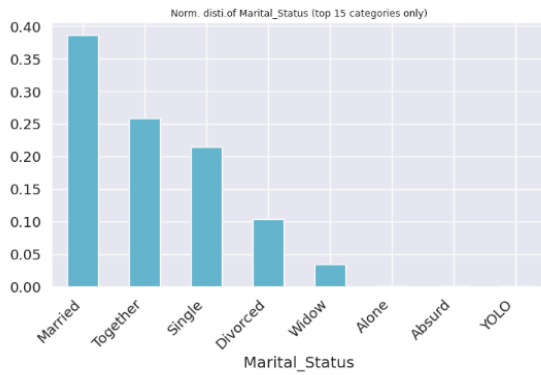
#Sinh báo cáo tự động



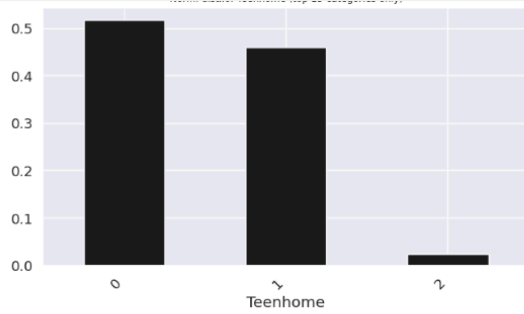
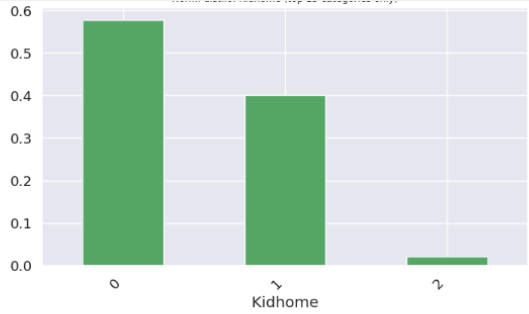
43



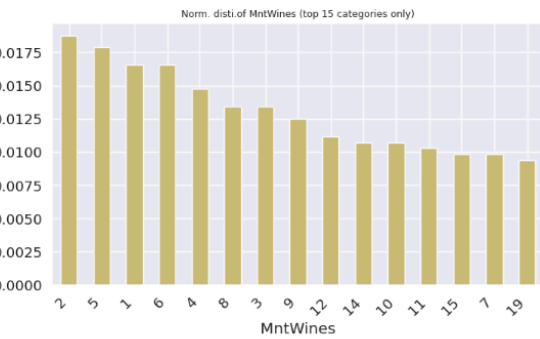
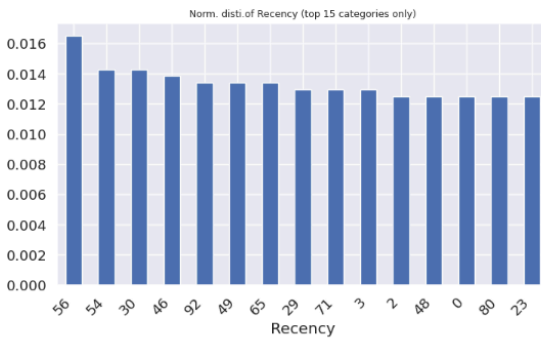
44



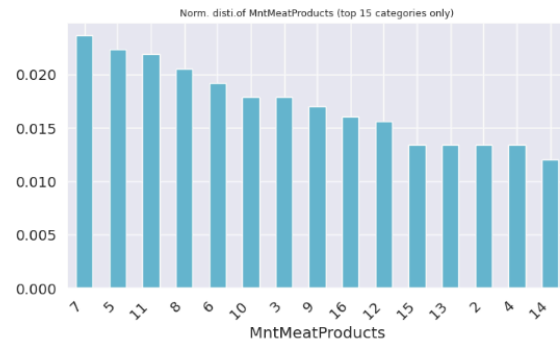
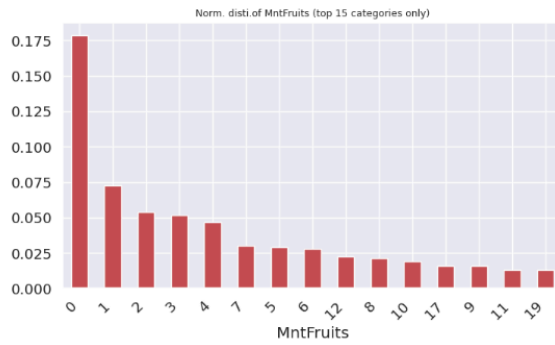
45



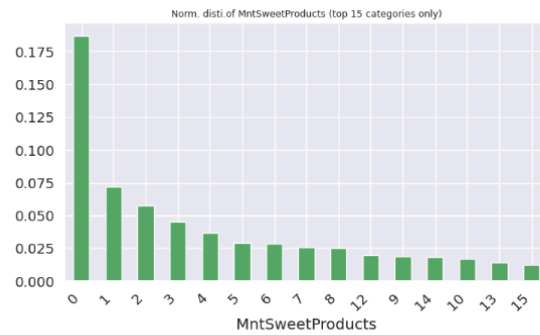
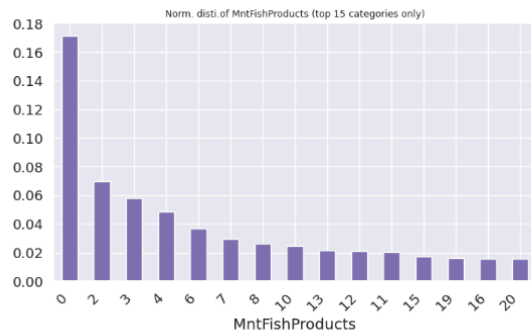
46



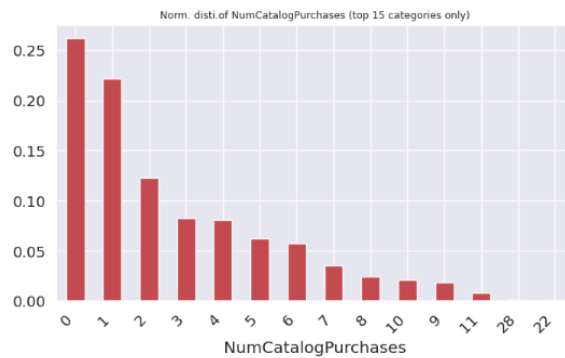
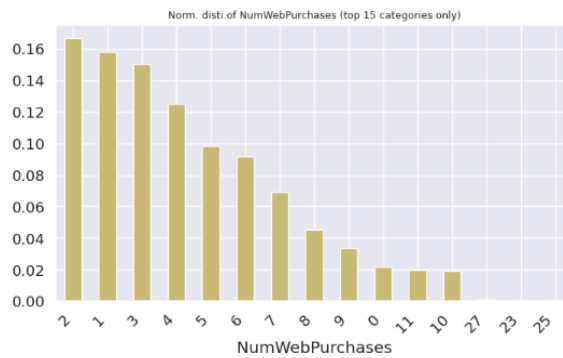
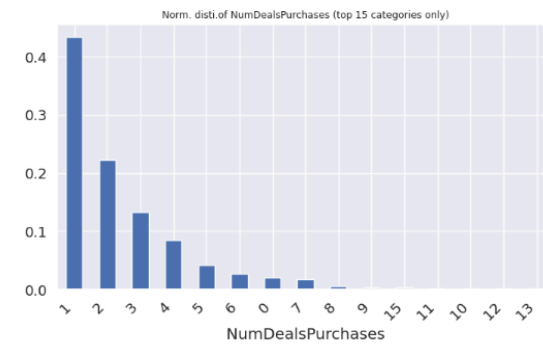
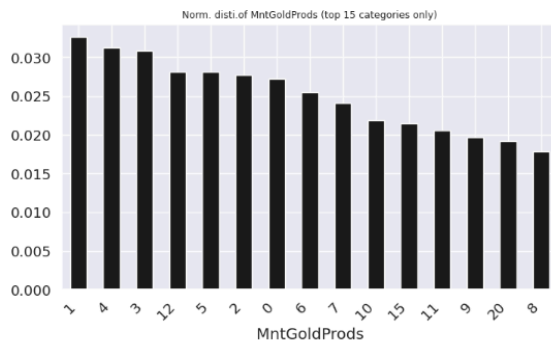
13



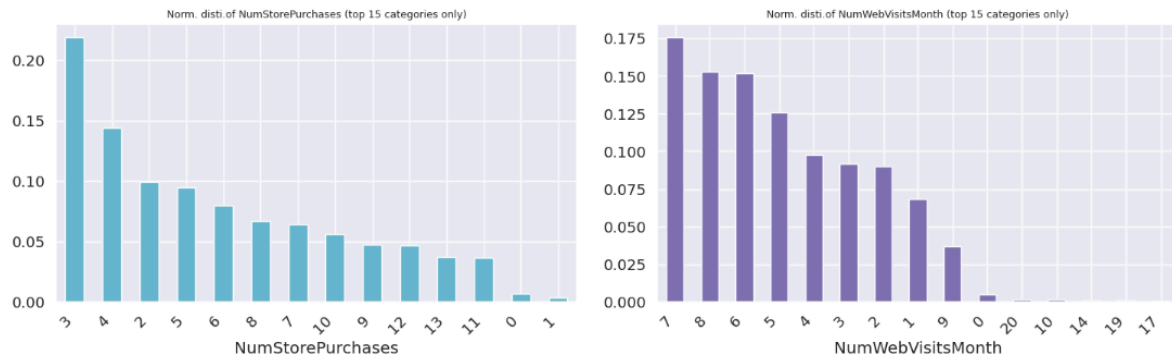
14



15

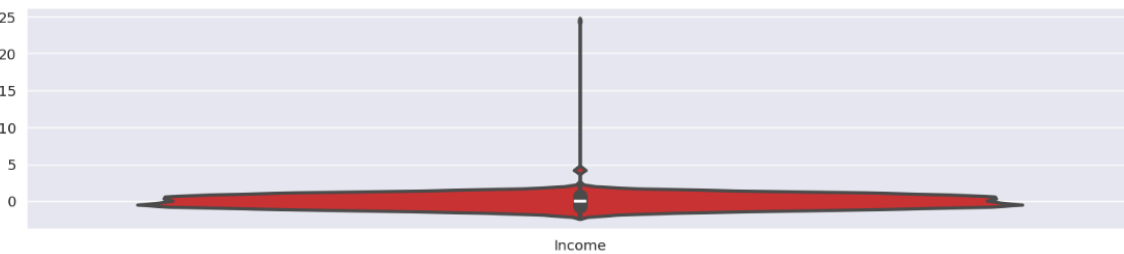


13



#Lưu thêm biểu đồ quan trọng

14



Heatmap of all Numeric Variables including target:

#Đồ thị heatmap

15



4)

