

FINLATICS DATA SCIENCE

CASE PROJECT

- Name: Ullas G

TOPIC: Banking

Problem Statement:

The goal of this data science project is to **analyze and predict client behavior** in a banking dataset related to direct marketing campaigns. Specifically, **the objective is to explore correlations and patterns among various client attributes** and determine the likelihood of clients subscribing to a term deposit.

Libraries Used:

- 1) Matplotlib – Version: 3.9.0
- 2) Seaborn – Version: 0.13.2
- 3) Pandas - Version: 2.2.2
- 4) Numpy – Version: 1.26.4

Questions:

- 1) What is the distribution of age among the clients?

Solution:

```

q1.py > ...
1  #Q1-- What is the distribution of age among the clients?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df = pd.read_csv('banking_data.csv')
8  fig,ax = plt.subplots(ncols=2,figsize=(12,6))
9  #Histogram of the given age distribution
10 sns.histplot(df['age'],bins = 25,kde = True,ax = ax[0])
11 ax[0].set_title('Age Distribution of Clients')
12 ax[0].set_xlabel('Age')
13 ax[0].set_ylabel('Frequency')
14
15 #Box plot for the age of clients to learn about the measures of central tendency
16 min_age = df['age'].min()
17 q1_age = df['age'].quantile(0.25)
18 median_age = df['age'].median()
19 q3_age = df['age'].quantile(0.75)
20 max_age = df['age'].max()
21
22 sns.boxplot(x=df['age'],ax=ax[1])
23 ax[1].annotate(f'Min: {min_age}', xy=(min_age, 0.5), xytext=(min_age - 10, 0.5),
24               arrowprops=dict(facecolor='black', arrowstyle="->")),
25 ax[1].annotate(f'Q1: {q1_age}', xy=(q1_age, 0.5), xytext=(q1_age - 10, 0.6),
26               arrowprops=dict(facecolor='black', arrowstyle="->")),
27 ax[1].annotate(f'Median: {median_age}', xy=(median_age, 0.5), xytext=(median_age - 10, 0.4),
28               arrowprops=dict(facecolor='black', arrowstyle="->")),
29 ax[1].annotate(f'Q3: {q3_age}', xy=(q3_age, 0.5), xytext=(q3_age - 10, 0.6),
30               arrowprops=dict(facecolor='black', arrowstyle="->")),
31 ax[1].annotate(f'Max: {max_age}', xy=(max_age, 0.5), xytext=(max_age - 10, 0.5),
32               arrowprops=dict(facecolor='black', arrowstyle="->")),
33 ax[1].set_title('Box plot for age of client')
34 ax[1].set_xlabel('Age')
35 plt.tight_layout()
36 plt.show()

```

2) How does the job type vary among the clients?

Solution:

```

q2.py > ...
1  #Q2-How does the job type vary among the clients?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df = pd.read_csv('banking_data.csv')
8  fig,ax = plt.subplots(ncols=2,figsize = (12,6))
9
10 #print(df.shape[0])
11 #print(df[df['job'].notnull()].shape[0]) , no missing values
12
13 #Count plot for distribution of clients among different jobs
14 sns.countplot(y='job', data=df, order=df['job'].value_counts().index, ax=ax[0])
15 ax[0].set_title('Distribution of Job Types Among Clients')
16 ax[0].set_xlabel('Count')
17 ax[0].set_ylabel('Job Type')
18
19 # Pie chart of job type distribution
20 df['job'].value_counts().plot.pie(autopct='%1.1f%%', startangle=140, ax=ax[1])
21 ax[1].set_title('Proportion of Job Types Among Clients')
22 ax[1].set_ylabel('')
23 plt.show()

```

3) What is the marital status distribution of the clients?

Solution:

```

q3.py > ...
1  #Q3- What is the marital status distribution of the clients?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df = pd.read_csv('banking_data.csv')
8
9  #print(df.shape[0])
10 #print(df[df['marital'].notnull()].shape[0])
11 #only 3 missing values out of 45216 hence no need to worry about null values
12
13 fig,ax = plt.subplots(ncols = 2,figsize=(12,6))
14 #Count plot for marital status distribution
15 sns.countplot(y='marital',data=df,order=df['marital'].value_counts().index,ax = ax[0])
16 ax[0].set_title('Distribution of Marital Status Among Clients')
17 ax[0].set_xlabel('Count')
18 ax[0].set_ylabel('Marital Status')
19
20 #Pie chart of marital status distribution
21 df['marital'].value_counts().plot.pie(autopct='%1.1f%%', startangle=140, ax=ax[1])
22 ax[1].set_title('Proportion of Marital Status Among Clients')
23 ax[1].set_ylabel('')
24 plt.show()
25

```

4) What is the level of education among the clients?

Solution:

```

q4.py > ...
1  #Q4 - - What is the level of education among the clients?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8  #print(df.shape[0])
9  #print(df[df['education'].notnull()].shape[0])
10 #Only 3 missing values, hence no need to worry about the null values
11 fig, ax = plt.subplots(ncols=2, figsize=(12,6))
12
13 # Count plot of education level distribution
14 #The plot also depicts the number of clients who have credit in each educational category
15 sns.countplot(x='education', data=df, order=df['education'].value_counts().index, ax=ax[0],hue=df['default'])
16 ax[0].set_title('Distribution of Education Levels Among Clients')
17 ax[0].set_xlabel('Education Level')
18 ax[0].set_ylabel('Count')
19 ax[0].tick_params(axis='x', rotation=45)
20 handles, labels = ax[0].get_legend_handles_labels()
21 ax[0].legend(handles, labels, title='Credit')#Changing the title of the legend
22
23 # Pie chart of education level distribution
24 df['education'].value_counts().plot.pie(autopct='%1.1f%%', startangle=140, ax=ax[1])
25 ax[1].set_title('Proportion of Education Levels Among Clients')
26 ax[1].set_ylabel('')
27 plt.show()

```

5) What proportion of clients have credit in default?

Solution:

```

q5.py > ...
1  #Q5-What proportion of clients have credit in default?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8
9  #calculating the proportion of clients who have credit
10 proportion = df['default'].value_counts(normalize=True) * 100
11 print(proportion)
12
13 #Pie chart depicting proportion of clients who have credit
14 fig,ax = plt.subplots(ncols=2,figsize=(12,6))
15 proportion.plot(kind='pie', autopct='%1.1f%%', startangle=140, ax=ax[0])
16 ax[0].set_title('Proportion of Clients with Credit in Default')
17 ax[0].set_ylabel('')
18
19
20 #View of Age group vs count of clients who have credit
21 bins = [10, 20, 30, 40, 50, 60, 70, 80, 90, 100]
22
23 # Grouping by age and default to get counts
24 age_default_counts = df.groupby(pd.cut(df['age'], bins=bins))['default'].value_counts().unstack(fill_value=0)
25
26 age_default_counts.plot(kind='bar', ax=ax[1])
27 ax[1].set_title('Distribution of Credit Status by Age Groups')
28 ax[1].set_xlabel('Age Groups')
29 ax[1].set_ylabel('Count')
30 ax[1].legend(title='Credit', labels=['No', 'Yes'])
31 ax[1].tick_params(axis='x', rotation=45)
32 plt.show()

```

6) What is the distribution of average yearly balance among the clients?

Solution:

```

q6.py > ...
1  #Q6-What is the distribution of average yearly balance among the clients?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8  #print(df.shape[0])
9  #print(df[df['balance'].notnull()].shape[0]) , -no missing values
10
11  # Plotting the distribution of average yearly balance
12  plt.figure(1,figsize=(10, 6))
13  sns.histplot(df['balance'], bins=40,kde=True)
14  plt.title('Distribution of Average Yearly Balance Among Clients')
15  plt.xlabel('Average Yearly Balance (euros)')
16  plt.ylabel('Frequency')
17  plt.tight_layout()
18  plt.show()
19

```

7) How many clients have housing loans?

Solution:

```

q7.py > ...
1  #Q7-How many clients have housing loans?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8  #print(df.shape[0])
9  #print(df[df['housing'].notnull()].shape[0]) -no missing values
10
11  #printing the number of clients who have housing loans
12  housing_loans = df[df['housing'] == 'yes'].shape[0]
13  housing_counts = df['housing'].value_counts()
14  print(f'Number of clients with housing loans:{housing_loans}')
15
16  fig,ax = plt.subplots(ncols = 2,figsize=(12,6))
17  #pie chart representing clients with an without housing loans
18  ax[0].pie(housing_counts, labels=housing_counts.index, autopct='%1.1f%%',colors=['skyblue', 'lightcoral'])
19  ax[0].set_title('Proportion of Clients with and without Housing Loans')
20  ax[0].set_ylabel('')
21
22
23  #bar chart representing clients with and without housing loans
24  housing_counts.plot(kind='bar')
25  ax[1].set_title('Number of Clients with and without Housing Loans')
26  ax[1].set_xlabel('Housing Loan')
27  ax[1].set_ylabel('Count')
28  plt.show()

```

8) How many clients have personal loans?

Solution:

```
q8.py > ...
1  #Q8-How many clients have personal loans?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8
9  #printing the number of clients who have personal loans
10 personal_loans = df[df['loan'] == 'yes'].shape[0]
11 personal_counts = df['loan'].value_counts()
12 print(f'Number of clients with personal loans:{personal_loans}')
13
14 fig,ax = plt.subplots(ncols = 2,figsize=(12,6))
15 #pie chart representing clients with an without personal loans
16 ax[0].pie(personal_counts, labels=personal_counts.index, autopct='%1.1f%%',colors=['skyblue', 'lightcoral'])
17 ax[0].set_title('Proportion of Clients with and without Personal Loans')
18 ax[0].set_ylabel('')
19
20
21 #Histogram representing clients with and without personal loans
22 personal_counts.plot(kind='bar')
23 ax[1].set_title('Number of Clients with and without Personal Loans')
24 ax[1].set_xlabel('Personal Loan')
25 ax[1].set_ylabel('Number of Clients')
26
27
```

9) What are the communication types used for contacting clients during the campaign?

Solution:


```

q9.py > ...
1  #Q9-What are the communication types used for contacting clients during the campaign
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8  #print(df.shape[0])
9  #print(df[df['contact'].notnull()].shape[0])    , - no missing values
10
11  types_of_contact = df['contact'].value_counts()
12  print('The different Types of Communication used to contact the Client are:')
13  for i in types_of_contact.index:
14      if(i=='unknown'):
15          continue
16      print(i)
17  print('The communication type for some of the clients are unknown')
18  #Printing the types of communication used
19
20
21  #Bar plot representing the types of communication used by the Bank to contact the Clients
22  fig = plt.figure(figsize=(12,6))
23  types_of_contact.plot(kind='bar',color='skyblue',title='Types Of Communication Used By the Bank')
24  for i, count in enumerate(types_of_contact):
25      plt.text(i, count + 100, str(count), ha='center', va='bottom', fontsize=10)
26  plt.xlabel('Communication Type')
27  plt.ylabel('Count')
28  plt.show()
29
30

```

10) What is the distribution of the last contact day of the month?

Solution:

```

q10.py > ...
1  #Q10-What is the distribution of the last contact day of the month?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8
9  # Plotting the distribution of last contact day of the month
10 plt.figure(figsize=(10, 6))
11 sns.histplot(df['day'], bins=31, kde=False)
12 plt.title('Distribution of Last Contact Day of the Month')
13 plt.xlabel('Day of the Month')
14 plt.ylabel('Frequency')
15
16
17 df['month'] = pd.Categorical(df['month'], categories=['jan', 'feb', 'mar', 'apr', 'may', 'jun', 'jul', 'aug', 'sep', 'oct', 'nov', 'dec'])
18
19 #HEAT MAP
20 #To represent the distribution of last contact day of the month for each month
21 pivot = df.pivot_table(index='day', columns='month', values='day_month', aggfunc='count')
22 plt.figure(figsize=(10,6))
23 sns.heatmap(pivot, cmap='YlGnBu', annot=True, fmt='g', linewidths=.5)
24 plt.title('Distribution of Last Contact Day of the Month')
25 plt.xlabel('Month')
26 plt.ylabel('Last Contact Day of the Month')
27 plt.tight_layout()
28 plt.show()
29
30

```

11) How does the last contact month vary among the clients?

Solution:

```
q11.py > ...
1  #Q11-How does the last contact month vary among the clients?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8  last_contact_month = df['month'].value_counts()
9  #PIE CHART
10 plt.pie(last_contact_month.values,labels=last_contact_month.index,autopct='%1.1f%%')
11 plt.title('Proportion of last contact month of Clients')
12 plt.ylabel('')
13
14 #COUNT PLOT WITH PREVIOUS OUTCOME
15 g = sns.catplot(x='month', col='poutcome', data=df, kind='count', palette='viridis')
16 g.set_xlabel('Month')
17 g.set_ylabel('Number of Clients')
18 g.set_titles('previous outcome: {col_name}')
19 g.fig.suptitle('Distribution Of Last Contact Month Of Clients with Previous Outcome', y=1)
20
21 for ax in g.axes.flat:
22     for p in ax.patches:
23         height = p.get_height()
24         ax.text(p.get_x() + p.get_width() / 2., height + 50, f'{int(height)}', ha='center', va='bottom', fontsize=10)
25 #COUNT PLOT WITH JOB TYPE
26 l = sns.catplot(x='month', col='job', data=df, kind='count', palette='viridis', col_wrap=3,sharex=True)
27 l.set_xlabel('Month')
28 l.set_ylabel('Number of Clients')
29 l.set_titles('Job: {col_name}',y=0.90)
30 l.fig.suptitle('Distribution Of Last Contact Month Of Clients with Job Type', y=1)
31
32 for ax in l.axes.flat:
33     for p in ax.patches:
34         height = p.get_height()
35         ax.text(p.get_x() + p.get_width() / 2., height + 50, f'{int(height)}', ha='center', va='bottom', fontsize=10)
36 plt.tight_layout()
37 plt.show()
```

12) What is the distribution of the duration of the last contact?

Solution:

```

q12.py > ...
1  #Q12-What is the distribution of the duration of the last contact?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7  df=pd.read_csv('banking_data.csv')
8  sns.histplot(df['duration'], bins=30, kde=True, color='skyblue')
9  plt.title('Distribution of Duration of Last Contact')
10 plt.xlabel('Duration (seconds)')
11 plt.ylabel('Frequency')
12 plt.xlim(right=2000)
13 plt.grid()
14 plt.show()

```

13) How many contacts were performed during the campaign for each client?

Solution:

```

q13.py > ...
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7
8  df=pd.read_csv('banking_data.csv')
9  print('Summary Statistics for Camapign:')
10 print(df['campaign'].describe())
11
12 #Handling the missing values by filling them up with median
13 df['campaign'].fillna(df['campaign'].median(), inplace=True)
14
15 #Hist plot to display the the distribution of number of contacts made in the campaign
16 plt.figure(figsize=(10,6))
17 sns.histplot(df['campaign'], bins=30, kde=False, color='skyblue')
18 plt.title('Distribution of Number of Contacts Performed During Campaign')
19 plt.xlabel('Number of Contacts')
20 plt.ylabel('Frequency')
21 plt.grid(True)
22 plt.xlim(right=20)
23
24
25 #Box plot to Observe the statistics of the number of contacts made
26 plt.figure(figsize=(10,6))
27 sns.boxplot(y='campaign', data=df, palette='viridis')
28 plt.title('Number of Contacts During Campaign')
29 plt.ylabel('Number of Contacts')
30 plt.ylim()
31
32 median = df['campaign'].median()
33 min_val=df['campaign'].min()
34 max_val=df['campaign'].max()
35 plt.text(0, median, f'Median: {median}', ha='center', va='bottom', color='black', fontsize=10)
36 plt.text(0, min_val, f'Min: {min_val}', ha='center', va='bottom', color='black', fontsize=10)
37 plt.text(0, max_val, f'Max: {max_val}', ha='center', va='bottom', color='black', fontsize=10)
38 plt.show()
39

```

14) What is the distribution of the number of days passed since the client was last contacted from a previous campaign?

Solution:

```
q14.py > ...
1  #Q14-What is the distribution of the number of days passed since the client was last contacted from a previous campaign?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7
8  df=pd.read_csv('banking_data.csv')
9  df['pdays'] = df['pdays'].replace(-1, pd.NA)
10 plt.figure(figsize=(10,6))
11 sns.histplot(df['pdays'],bins=10,color='skyblue')
12 plt.title('Distribution of Number of Days Since Last Contact from Previous Campaign')
13 plt.xlabel('Number of Days')
14 plt.ylabel('Frequency')
15 plt.grid()
16
17
18 #Distribution Of number of days since client was previously contacted vs previous outcome with client
19 plt.figure(figsize=(10, 6))
20 sns.barplot(x='poutcome', y='pdays', data=df, estimator='mean', palette='viridis')
21 plt.title('Mean Number of Days Since Last Contact for Each Previous Outcome')
22 plt.xlabel('Previous Outcome')
23 plt.ylabel('Mean Number of Days')
24 plt.show()
25
```

15) How many contacts were performed before the current campaign for each client?

Solution:

```
q15.py > ...
1  #Q15-How many contacts were performed before the current campaign for each client?
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import pandas as pd
6
7
8  df=pd.read_csv('banking_data.csv')
9  df.fillna({'previous':0})
10 num_of_contacts = df['previous'].sum()
11 print(f'The total number of contacts made to each client before the current campaign is:{num_of_contacts}')
12
13
14
15
```

16) What were the outcomes of the previous marketing campaigns?

Solution:

```

q16.py > ...
1 #Q16-What were the outcomes of the previous marketing campaigns?
2
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import pandas as pd
6
7
8 df=pd.read_csv('banking_data.csv')
9 df.fillna({'poutcome':'unknown'})
10 prev_out = df['poutcome'].value_counts()
11
12 print('Summary of Outcomes of previous marketing campaigns')
13 print(f'Total number of Reported Success:{df[df['poutcome']=='success'].shape[0]}')
14 print(f'Total number of Reported Failure:{df[df['poutcome']=='failure'].shape[0]}')
15 print(f'Other Results:{df[df['poutcome']=='other'].shape[0]}')
16 print(f'Unknown:{df[df['poutcome']=='unknown'].shape[0]}')
17
18
19 #Proportion of Each Previous Outcome
20 plt.figure(figsize=(10,6))
21 plt.pie(prev_out.values,labels=prev_out.index,autopct='%1.1f%%',colors=['lightgreen', 'lightcoral', 'lightskyblue', 'lightgrey'])
22 plt.title('Distribution of Outcomes from Previous Marketing Campaigns')
23 plt.ylabel('')
24 plt.show()

```

17) What is the distribution of clients who subscribed to a term deposit vs. those who did not?

Solution:

```

q17.py > ...
1 #Q17-What is the distribution of clients who subscribed to a term deposit vs. those who did not?
2
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 df = pd.read_csv('banking_data.csv')
8 subscription_counts = df['y'].value_counts(normalize=True) * 100
9 print(f"Percentage subscribed to term deposit: {subscription_counts['yes']:.1f}%")
10 print(f"Percentage not subscribed to term deposit: {subscription_counts['no']:.1f}%")
11
12 #Distribution of clients subscribing to term deposit
13 plt.figure(figsize=(10, 6))
14 sns.countplot(x='y', data=df, palette=['lightcoral', 'lightskyblue'])
15 plt.title('Distribution of Clients Subscribing to Term Deposit')
16 plt.xlabel('Term Deposit Subscription')
17 plt.ylabel('Number of Clients')
18
19
20 for p in plt.gca().patches:
21     plt.gca().annotate(f'{p.get_height()}', (p.get_x() + p.get_width() / 2., p.get_height()),
22                       ha='center', va='center', xytext=(0, 9), textcoords='offset points', fontsize=12, color='black')
23 plt.legend(labels=['Not Subscribed', 'Subscribed'], loc='upper right')
24 plt.grid()
25 plt.show()

```

18) Are there any correlations between different attributes and the likelihood of subscribing to a term deposit?

Solution:

```

q18.py > ...
2
3 import pandas as pd
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 df = pd.read_csv('banking_data.csv')
7 numerical_columns = ['age', 'balance', 'duration', 'campaign', 'pdays', 'previous']
8 categorical_columns = ['job', 'education', 'marital', 'outcome']
9
10 #Converting some columns to numerical type for correlation
11 df['y_numerical'] = df['y'].apply(lambda x: 1 if x == 'yes' else 0)
12 df['Housing_loan'] = df['housing'].apply(lambda x: 1 if x == 'yes' else 0)
13 df['Personal_loan'] = df['loan'].apply(lambda x: 1 if x == 'yes' else 0)
14 numerical_data = df[numerical_columns + ['Housing_loan'] + ['Personal_loan'] + ['y_numerical']]
15
16 # Correlation Matrix for Numerical Variables
17 corr_matrix = numerical_data.corr()
18 plt.figure(figsize=(12, 6))
19 sns.heatmap(corr_matrix, annot=True, cmap='viridis', fmt=".2f", annot_kws={"size": 10})
20 plt.title('Correlation Matrix of Numerical Variables')
21
22 #Contingency Tables for Categorical Variables
23 for col in categorical_columns:
24     contingency_table = pd.crosstab(df[col], df['y'], margins=True, margins_name='Total')
25     print(f"Contingency Table for {col} vs y:")
26     print(contingency_table)
27     print("\n")
28
29 print('The above Heat map shows how the numerical variables are related with themselves and other variables')
30 print("Summary Of the logically drawn results of correlation Statistics of Numerical variables with the target variable 'y'")
31 print('Variable    Correlation Coefficient')
32 variables_of_interest = ['Housing_loan', 'Personal_loan', 'age', 'balance']
33 for variable in variables_of_interest:
34     correlation = corr_matrix.loc[variable, 'y_numerical']
35     print(f'{variable}    {correlation:.2f}')
36 print("\n\nThe contingency table above shows how the value of 'y' depends on some of the categorical variables")
37 plt.tight_layout()
38 plt.show()
39

```

