

Start coding or [generate](#) with AI.

Dataset uploading

```
# prompt: DATASET UPLOAD CODE

from google.colab import files
uploaded = files.upload()
# prompt: handle missing values for an csv file in googl

import pandas as pd

# Load the CSV file into a pandas DataFrame
df = pd.read_csv('customer.csv')
```

 Choose Files customer.csv

- **customer.csv**(text/csv) - 3055 bytes, last modified: 5/8/2025
 - 100% done

Saving customer.csv to customer.csv

Handling Missing Values

```
# prompt: check for missing values and fill the missing

# Check for missing values
print(df.isnull().sum())

# Fill missing values with mean for numerical columns
numerical_cols = df.select_dtypes(include=['number']).columns
df[numerical_cols] = df[numerical_cols].fillna(df[numerical_cols].mean())

# Fill missing values with mode for categorical columns
categorical_cols = df.select_dtypes(include=['object']).columns
df[categorical_cols] = df[categorical_cols].fillna(df[categorical_cols].mode().iloc[0])

# Verify if missing values are filled
print(df.isnull().sum())
```

 Year 0

 Date 0

 Question_Number 0

 Question 0

 Number_of_Respondents 0

 Very_Satisfied 0

 Satisfied 0

 Neutral 0

 Dissatisfied 0

ValueError NameError ...

 Please explain the error:

ValueError: Could not interpret value

Hypothesis

The error message "ValueError: Could not interpret value Geography for x. Value is a string, but data was not passed." arises because the sns.boxplot function is called without specifying the data parameter, and you're trying to use column names ('Geography', 'CreditScore') directly. When you don't provide the data argument, seaborn expects the x and y arguments to be arrays or lists of data, not column names.

In this case, you're likely intending to use the original_df DataFrame (as mentioned in your code comments) for the visualization, but you haven't passed it to the sns.boxplot function. As a result, seaborn doesn't know where to find the 'Geography' and 'CreditScore' columns.

Suggested Changes

```
# Example 3: Boxplot for a sp
# (using the original_df to a
plt.figure(figsize=(10, 6))
# Pass the original_df DataFr
sns.boxplot(x='Geography', y=
```

```
Very_Dissatisfied      0
Very_Satisfied_or_Satisfied 0
ObjectId                0
dtype: int64
Year                     0
Date                      0
Question_Number           0
Question                  0
Number_of_Respondents     0
Very_Satisfied            0
Satisfied                 0
Neutral                   0
Dissatisfied              0
Very_Dissatisfied          0
Very_Satisfied_or_Satisfied 0
ObjectId                  0
dtype: int64
```

Duplicate records

```
# prompt: check for the duplicate records and remove them

# Check for duplicate rows
duplicate_rows = df[df.duplicated()]

# Print the duplicate rows
print("Duplicate Rows:")
print(duplicate_rows)

# Remove duplicate rows
df = df.drop_duplicates()

# Print the DataFrame after removing duplicates
print("\nDataFrame after removing duplicates:")
df
```

```
plt.title("CreditScore Distri
plt.xticks(rotation=45, ha='r
plt.show()
```

```
# Example 4: Countplot for a
plt.figure(figsize=(8, 6))
sns.countplot(x='Gender', dat
plt.title("Count of Gender")
plt.show()
```

```
# Example 5: Distribution of
plt.figure(figsize=(10, 6))
sns.countplot(x='Geography',
plt.title('Distribution of Ex
plt.xticks(rotation=45, ha='r
plt.show()
```

```
# Example 6: Histogram of a N
plt.figure(figsize=(10, 6))
sns.histplot(data=original_df
plt.title('Distribution of Ag
plt.show()
```

[Use code with caution](#)

Rate this answer

→ Duplicate Rows:
 Empty DataFrame
 Columns: [Year, Date, Question_Number, Question, Num
 Index: []

DataFrame after removing duplicates:

	Year	Date	Question_Number	Question	Num
0	2017	2017/10/31 07:00:00+00	7-13	Overall quality of customer service	
1	2016	2016/10/31 07:00:00+00	26	Overall quality of customer service	
2	2015	2015/10/31 07:00:00+00	Survey Not Conducted	Survey Not Conducted	
3	2014	2014/10/31 07:00:00+00	10b	How easy was the City to contact	
4	2014	2014/10/31 07:00:00+00	10c	The way you were treated	
5	2014	2014/10/31 07:00:00+00	10d	The accuracy of the information you were given	
6	2014	2014/10/31 07:00:00+00	10e	How quickly staff responded to your request	
7	2014	2014/10/31 07:00:00+00	10f	How well your issue was handled	
8	2013	2013/10/31 07:00:00+00	14b	How easy was the City to contact	
9	2013	2013/10/31 07:00:00+00	14c	The way you were treated	

				The accuracy of the information you were given
10	2013	2013/10/31 07:00:00+00	14d	How quickly staff responded to your request
11	2013	2013/10/31 07:00:00+00	14e	How well your issue was handled
12	2013	2013/10/31 07:00:00+00	14f	How easy was the city to contact
13	2012	2012/10/31 07:00:00+00	14b	The way you were treated
14	2012	2012/10/31 07:00:00+00	14c	The accuracy of the information you were given
15	2012	2012/10/31 07:00:00+00	14d	How quickly responded to request
16	2012	2012/10/31 07:00:00+00	14e	How well issue was handled
17	2012	2012/10/31 07:00:00+00	14f	How easy was the contact
18	2011	2011/10/31 07:00:00+00	14b	The way you were treated
19	2011	2011/10/31 07:00:00+00	14c	Accuracy of information you were given
20	2011	2011/10/31 07:00:00+00	14d	

				How quickly staff responded to your request
21	2011	2011/10/31 07:00:00+00	14e	How well your issue was handled
22	2011	2011/10/31 07:00:00+00	14f	How easy was the city to contact
23	2010	2010/10/31 07:00:00+00	14B	The way you were treated
24	2010	2010/10/31 07:00:00+00	14C	Accuracy of information you were given
25	2010	2010/10/31 07:00:00+00	14D	How quickly staff responded to your request
26	2010	2010/10/31 07:00:00+00	14E	How well your issue was handled
27	2010	2010/10/31 07:00:00+00	14F	How quickly staff responded to your request

Next steps:

[Generate code with df](#)

[View recommended plots](#)

Outliers

```
# prompt: check for the outliers in the above dataset

import pandas as pd
import numpy as np

# Assuming 'df' is your DataFrame with numerical feature

def find_outliers_iqr(data):
    Q1 = np.percentile(data, 25)
    Q3 = np.percentile(data, 75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    outliers = data[(data < lower_bound) | (data > upper_bound)]
    return outliers

numerical_features = df.select_dtypes(include=np.number)
for col in numerical_features:
    outliers = find_outliers_iqr(df[col])
    print(f"Outliers in {col}:")
    print(outliers)
    print("-" * 20)
```

→ Outliers in Year:
Series([], Name: Year, dtype: int64)

Outliers in Number_of_Respondents:
0 882
1 1202
2 99999
Name: Number_of_Respondents, dtype: int64

Outliers in Very_Satisfied:
0 22.16
1 22.39
2 99999.00
Name: Very_Satisfied, dtype: float64

Outliers in Satisfied:
2 99999.0
Name: Satisfied, dtype: float64

Outliers in Neutral:
0 25.03
1 23.37
2 99999.00
Name: Neutral, dtype: float64

```
-----  
Outliers in Dissatisfied:  
2    99999.0  
Name: Dissatisfied, dtype: float64  
-----  
Outliers in Very_Dissatisfied:  
2    99999.0  
Name: Very_Dissatisfied, dtype: float64  
-----  
Outliers in Very_Satisfied_or_Satisfied:  
2    99999.0  
Name: Very_Satisfied_or_Satisfied, dtype: float64  
-----  
Outliers in ObjectId:  
Series([], Name: ObjectId, dtype: int64)  
-----
```

Standardization

```
# prompt: standardize the above dataset  
  
from sklearn.preprocessing import StandardScaler  
  
# Assuming 'df' is your DataFrame with numerical feature  
  
# Create a StandardScaler object  
scaler = StandardScaler()  
  
# Select numerical columns for standardization  
numerical_cols = df.select_dtypes(include=np.number).col  
  
# Fit and transform the numerical columns  
df[numerical_cols] = scaler.fit_transform(df[numerical_c  
  
# Print the standardized DataFrame  
print("\nStandardized DataFrame:")  
df
```



Standardized DataFrame:

	Year	Date	Question_Number	Question
0	2.483682	2017/10/31 07:00:00+00	7-13	Overall quality customer service
1	1.940376	2016/10/31 07:00:00+00	26	Overall quality customer service
2	1.397071	2015/10/31 07:00:00+00	Survey Not Conducted	Survey Not Conducted
3	0.853766	2014/10/31 07:00:00+00	10b	How easily was the City to contact
4	0.853766	2014/10/31 07:00:00+00	10c	The way you were treated
5	0.853766	2014/10/31 07:00:00+00	10d	The accuracy of the information you were given
6	0.853766	2014/10/31 07:00:00+00	10e	How quickly staff responded to your request
7	0.853766	2014/10/31 07:00:00+00	10f	How well your issue was handled
8	0.310460	2013/10/31 07:00:00+00	14b	How easily was the City to contact
9	0.310460	2013/10/31 07:00:00+00	14c	The way you were treated
10	0.310460	2013/10/31 07:00:00+00	14d	The accuracy of the information

				you were given
11	0.310460	2013/10/31 07:00:00+00	14e	How quickly staff responded to your requests
12	0.310460	2013/10/31 07:00:00+00	14f	How we handled your issue
13	-0.232845	2012/10/31 07:00:00+00	14b	How easily was the city to contact
14	-0.232845	2012/10/31 07:00:00+00	14c	The way you were treated
15	-0.232845	2012/10/31 07:00:00+00	14d	The accuracy of the information you were given
16	-0.232845	2012/10/31 07:00:00+00	14e	How quickly responded to requests
17	-0.232845	2012/10/31 07:00:00+00	14f	How we handled your issue
18	-0.776151	2011/10/31 07:00:00+00	14b	How easily was the city to contact
19	-0.776151	2011/10/31 07:00:00+00	14c	The way you were treated
20	-0.776151	2011/10/31 07:00:00+00	14d	Accuracy of the information you were given
21	-0.776151	2011/10/31 07:00:00+00	14e	How quickly staff responded to your requests

				How we your issu wa handle
22	-0.776151	2011/10/31 07:00:00+00	14f	
23	-1.319456	2010/10/31 07:00:00+00	14B	How eas was thi city to contac
24	-1.319456	2010/10/31 07:00:00+00	14C	The wa you were treat
25	-1.319456	2010/10/31 07:00:00+00	14D	Accurac c information you were give
26	-1.319456	2010/10/31 07:00:00+00	14E	How quickl stat responde to you reques
27	-1.319456	2010/10/31 07:00:00+00	14F	How we your issu wa handle

Next
steps:

[Generate code with df](#)[!\[\]\(4cafc60cd39da821525d7c6589540296_img.jpg\) View recommended plots](#)

EDA

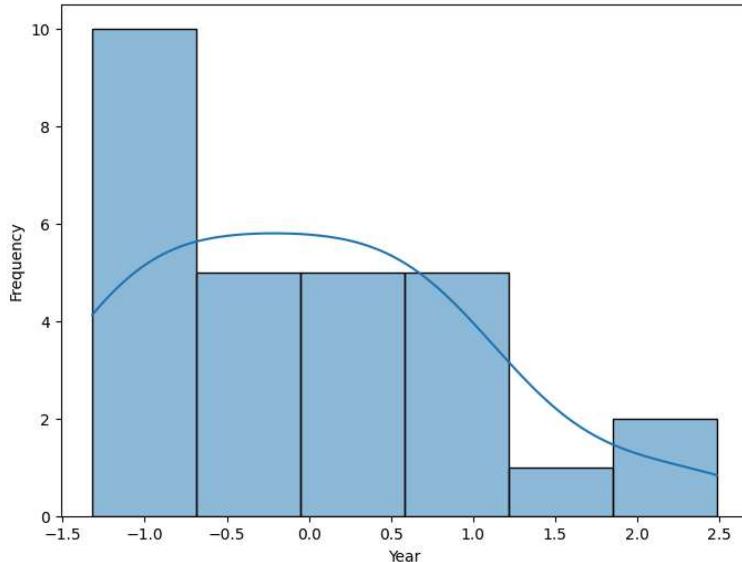
```
# prompt: visualize the dataset using eda by univariate,  
  
import matplotlib.pyplot as plt  
import seaborn as sns  
  
# Univariate Analysis  
  
# Histograms for numerical features  
for col in numerical_features:  
    plt.figure(figsize=(8, 6))  
    sns.histplot(df[col], kde=True)  
    plt.title(f'Distribution of {col}')  
    plt.xlabel(col)  
    plt.ylabel('Frequency')  
    plt.show()  
  
# Box plots for numerical features  
for col in numerical_features:  
    plt.figure(figsize=(8, 6))  
    sns.boxplot(df[col])  
    plt.title(f'Box Plot of {col}')  
    plt.xlabel(col)  
    plt.show()  
  
# Count plots for categorical features  
for col in categorical_cols:  
    plt.figure(figsize=(8, 6))  
    sns.countplot(x=col, data=df)  
    plt.title(f'Count Plot of {col}')  
    plt.xlabel(col)  
    plt.ylabel('Count')  
    plt.xticks(rotation=45, ha='right') # Rotate x-axis  
    plt.show()  
  
# Bivariate Analysis  
  
# Scatter plots for numerical features  
for col1 in numerical_features:  
    for col2 in numerical_features:  
        if col1 != col2:  
            plt.figure(figsize=(8, 6))  
            sns.scatterplot(x=col1, y=col2, data=df)  
            plt.title(f'Scatter Plot of {col1} vs {col2}')  
            plt.xlabel(col1)  
            plt.ylabel(col2)  
            plt.show()
```

```
# Correlation Heatmap
plt.figure(figsize=(12, 10))
sns.heatmap(df[numerical_cols].corr(), annot=True, cmap=
plt.title('Correlation Heatmap of Numerical Features')
plt.show()

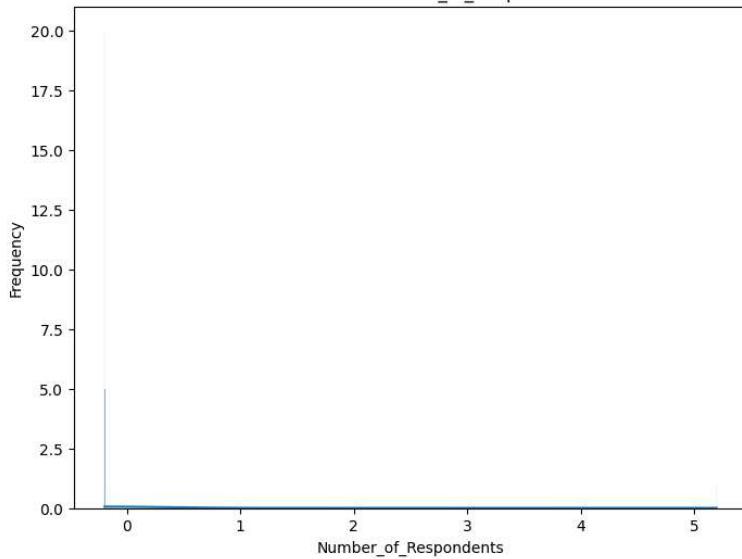
# Box plots for numerical features grouped by a category
for col in numerical_features:
    for cat_col in categorical_cols:
        plt.figure(figsize=(10,6))
        sns.boxplot(x = cat_col, y = col, data=df)
        plt.title(f"Box plot of {col} grouped by {cat_col}")
        plt.show()
```



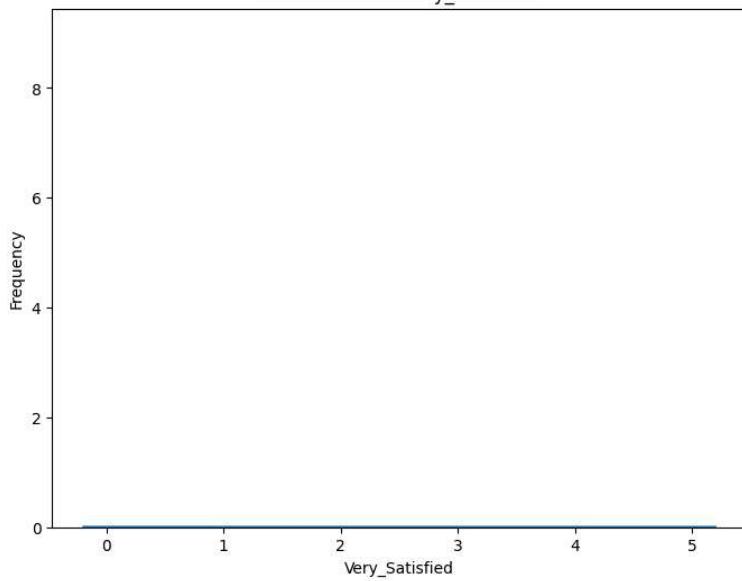
Distribution of Year



Distribution of Number_of_Respondents

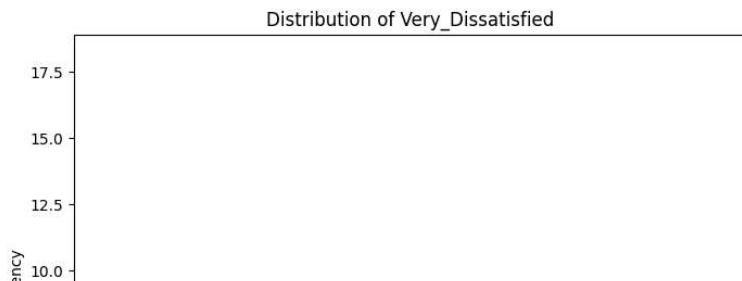
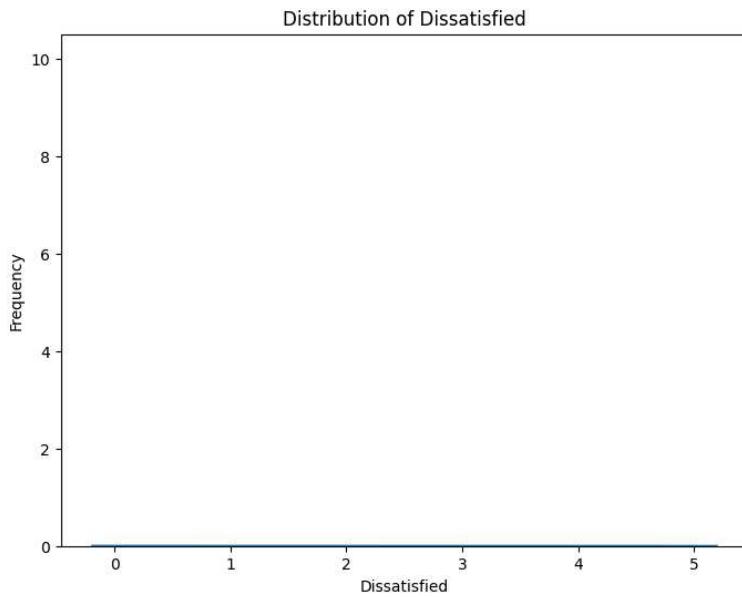
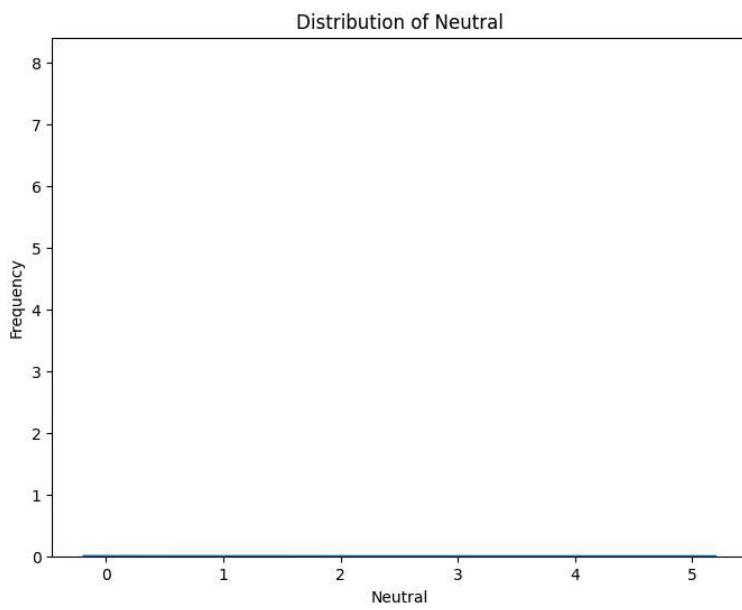
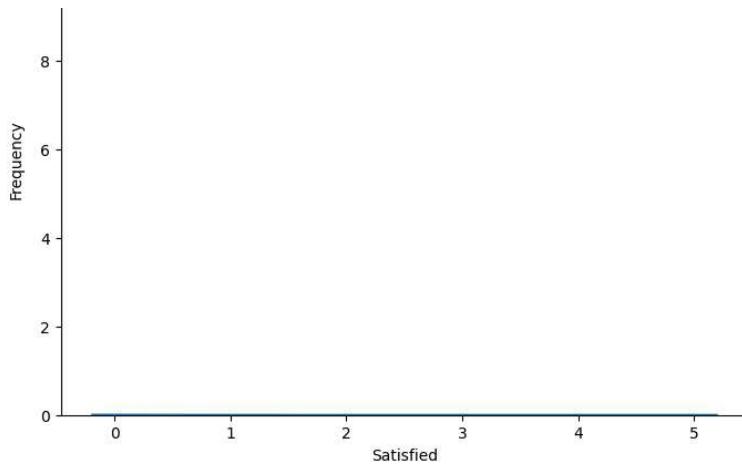


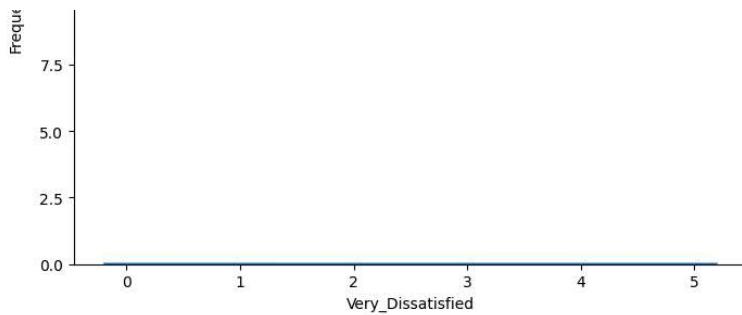
Distribution of Very_Satisfied



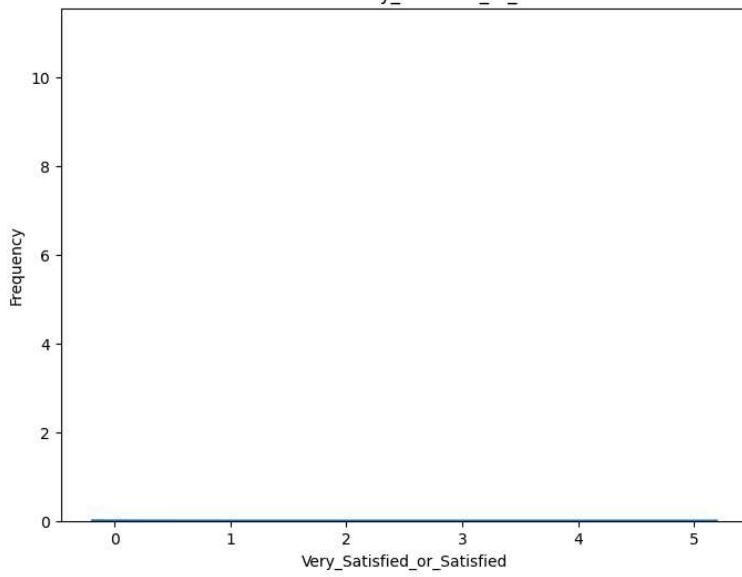
Distribution of Satisfied



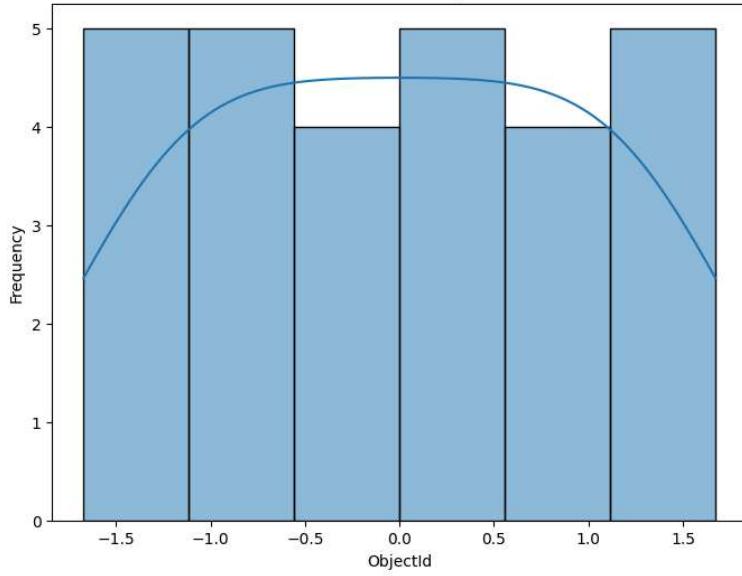




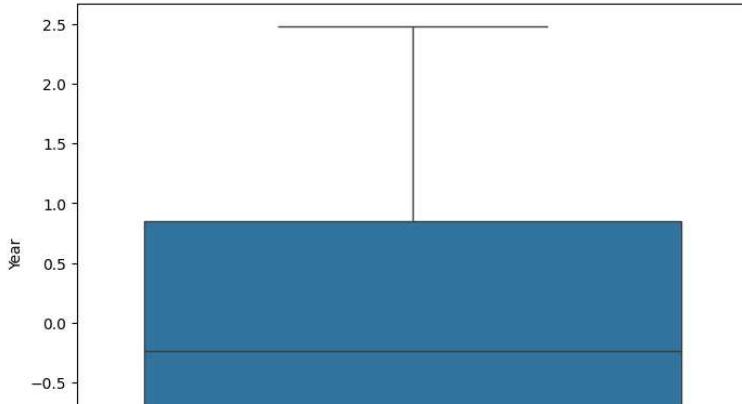
Distribution of Very_Satisfied_or_Satisfied

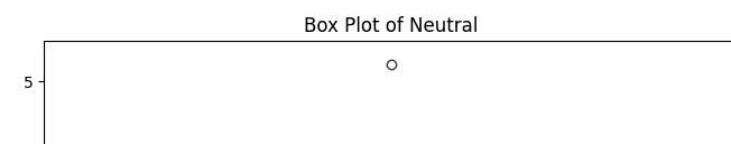
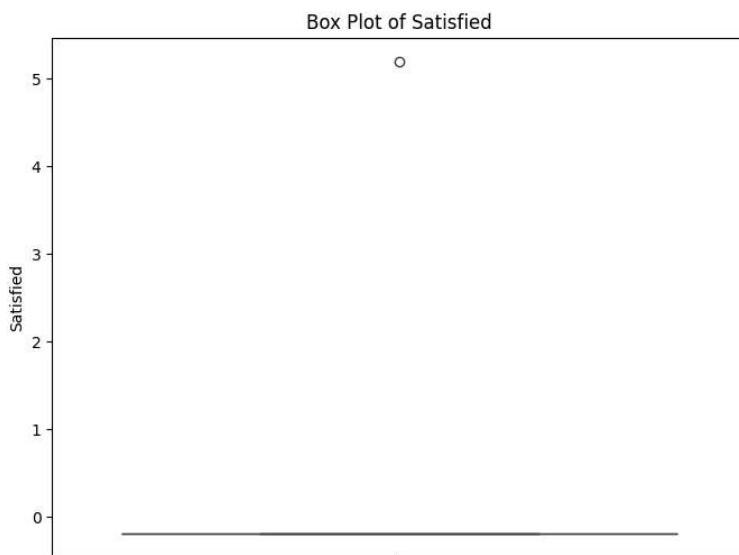
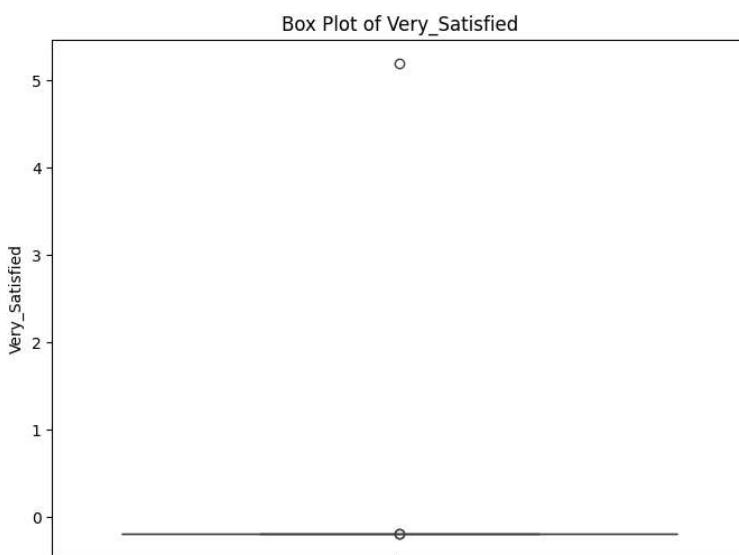
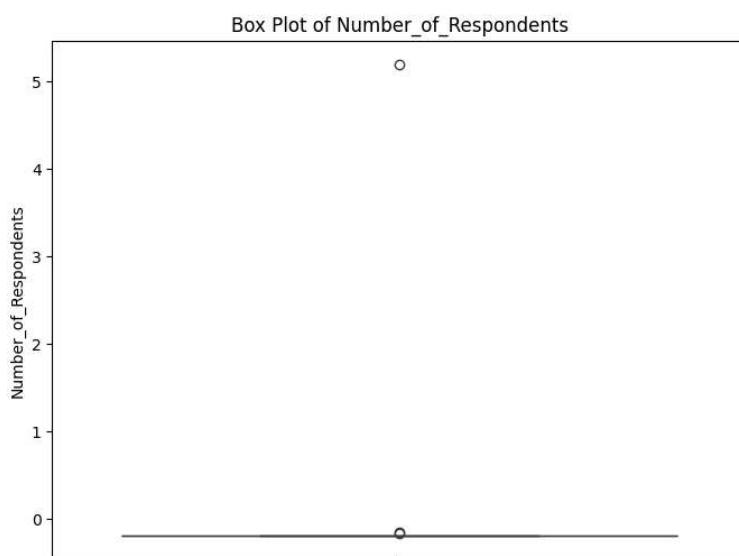
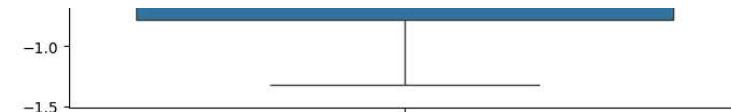


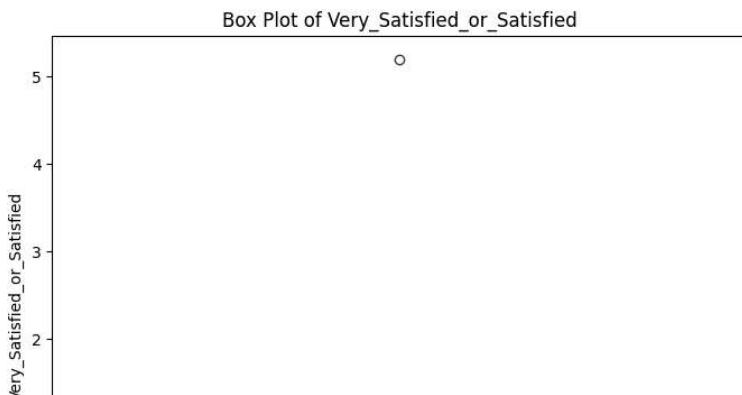
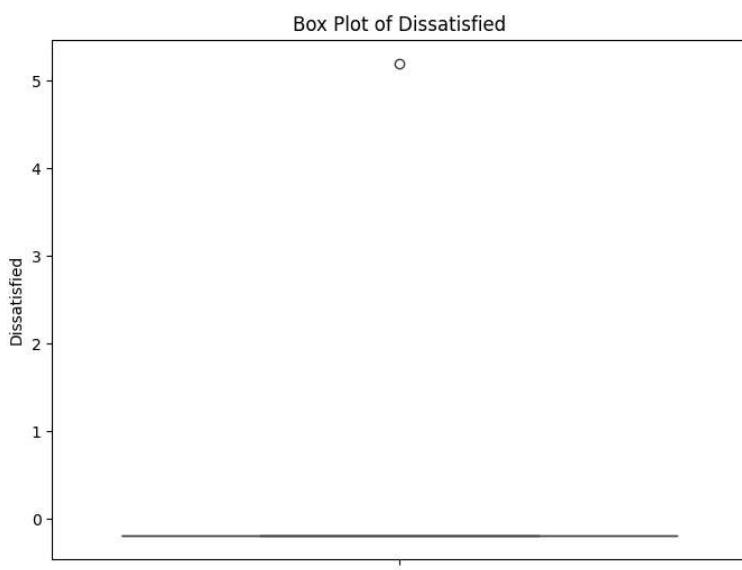
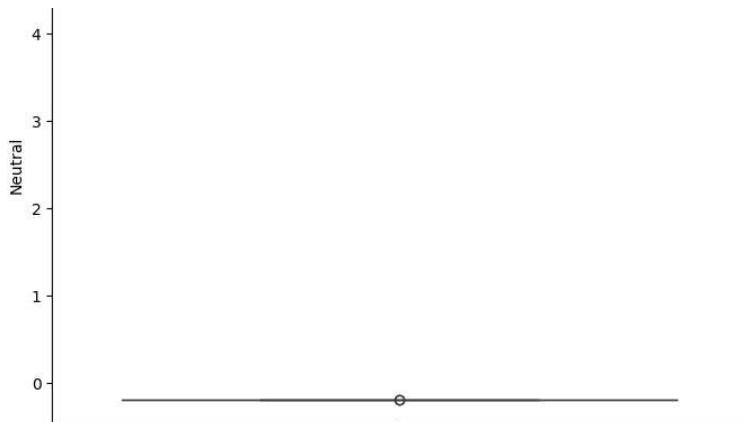
Distribution of ObjectId

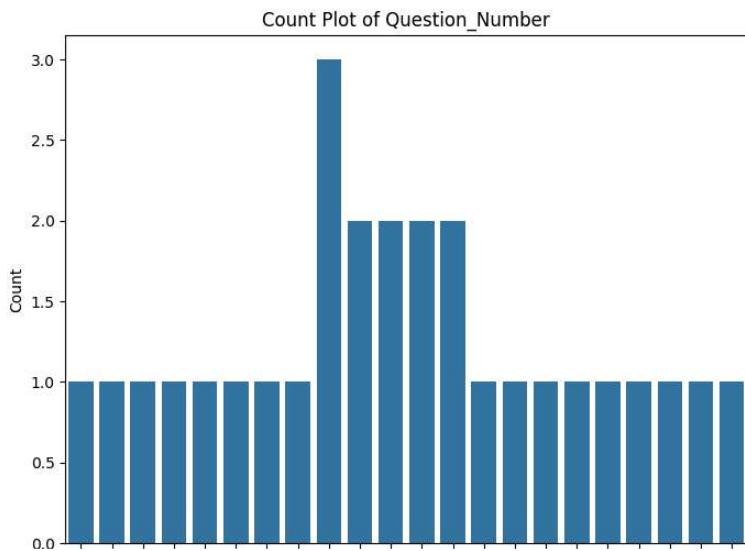
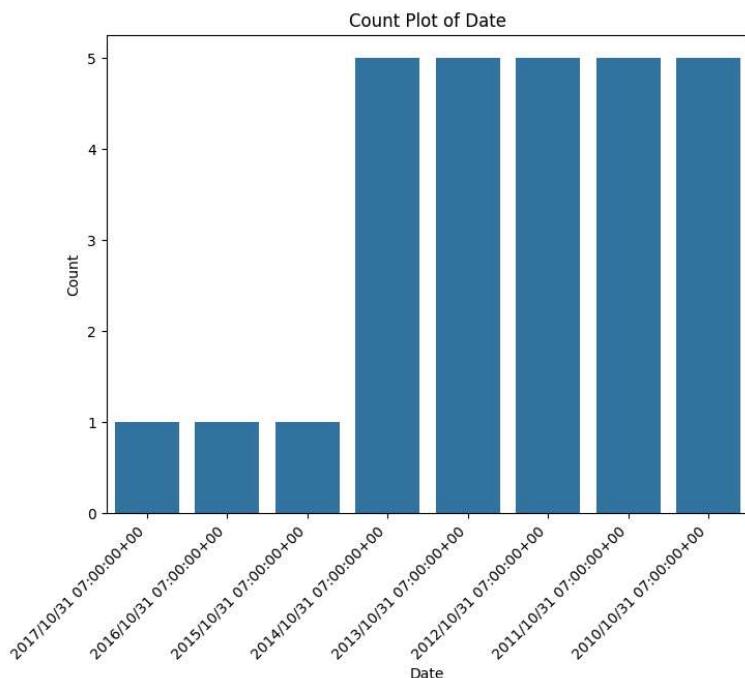
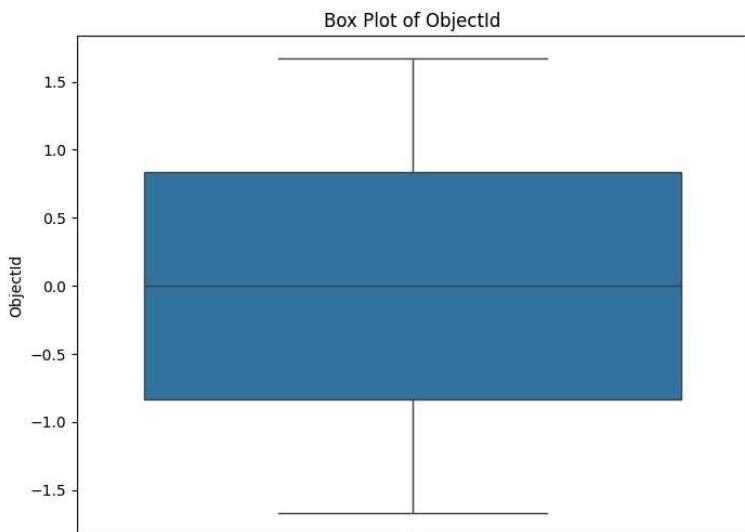


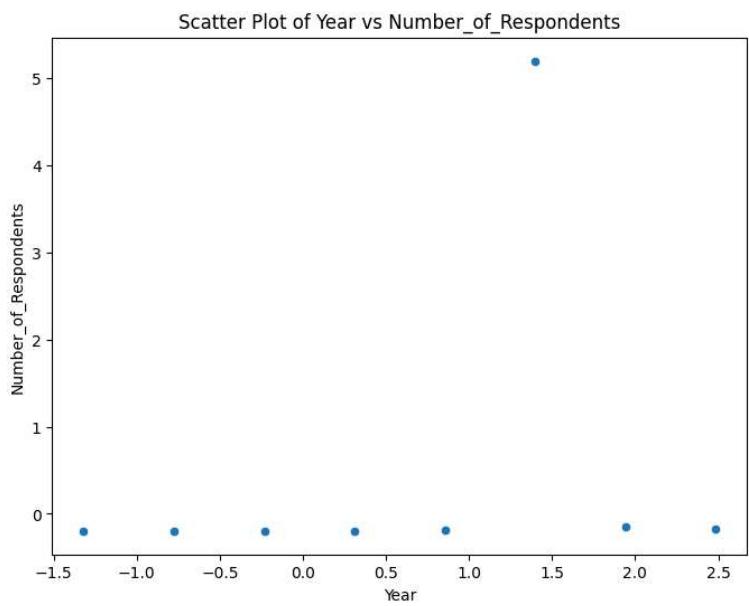
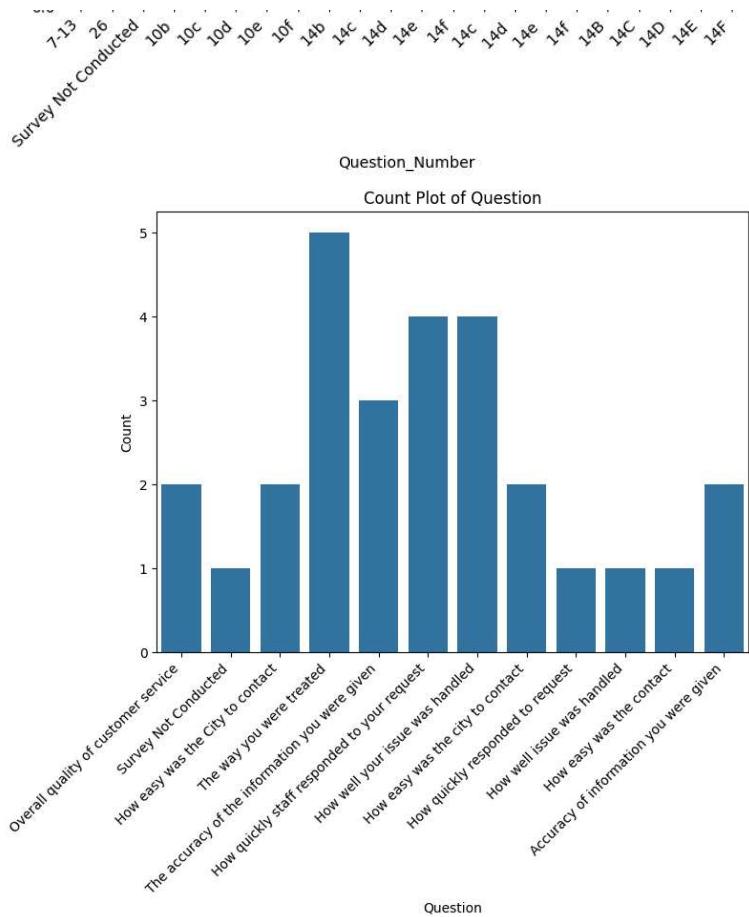
Box Plot of Year

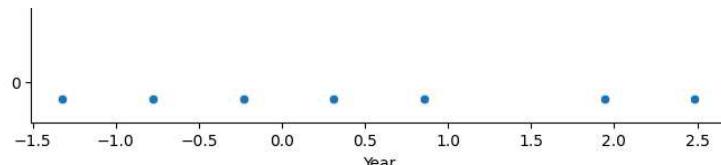




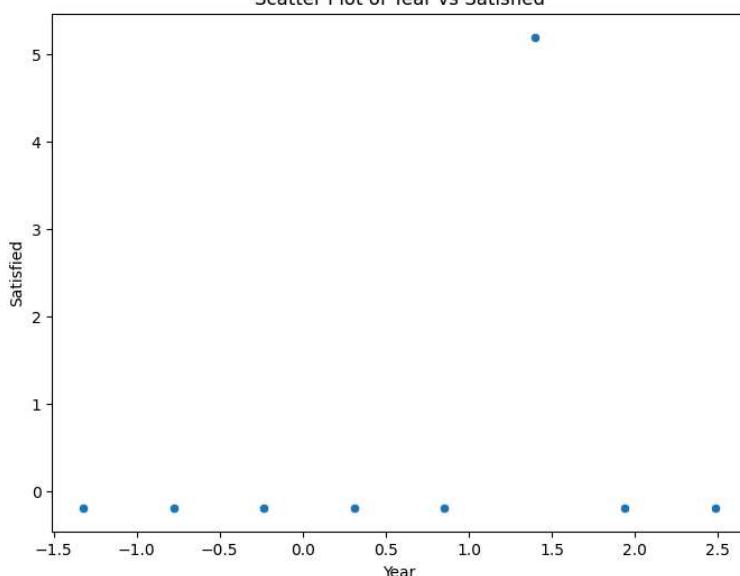




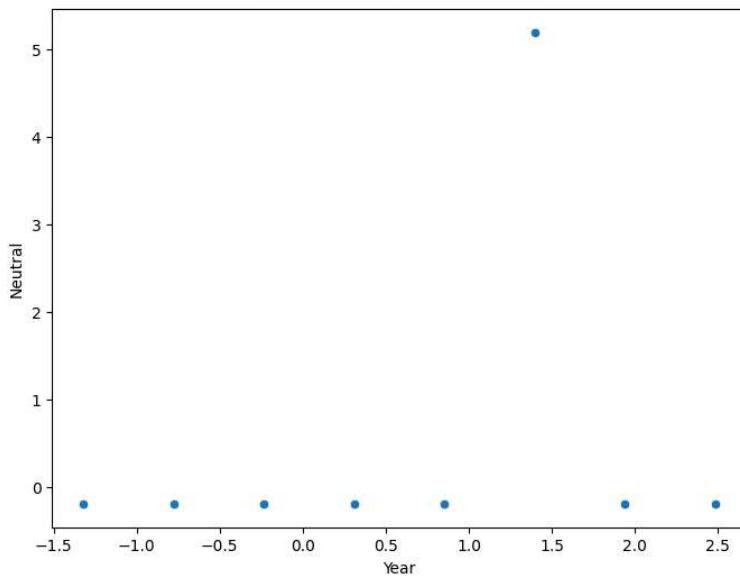




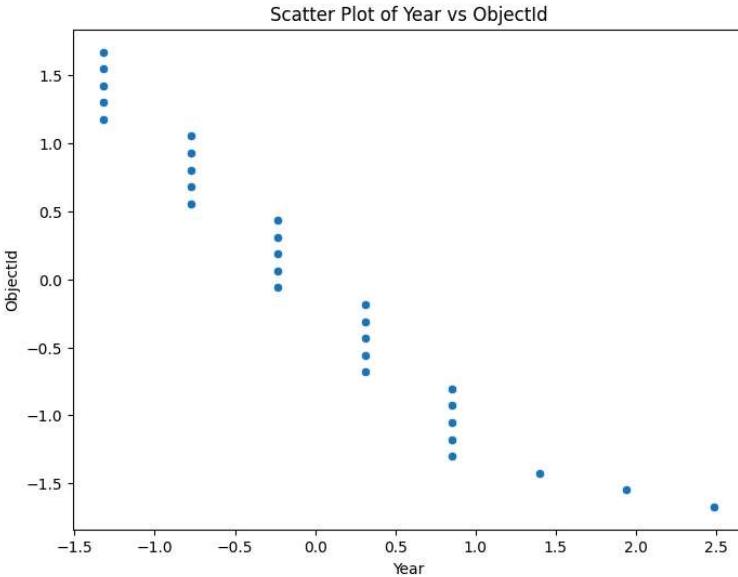
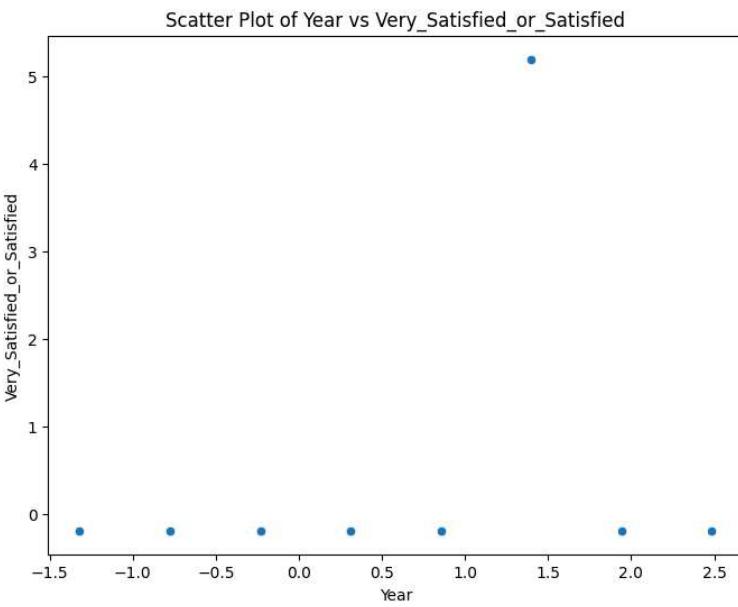
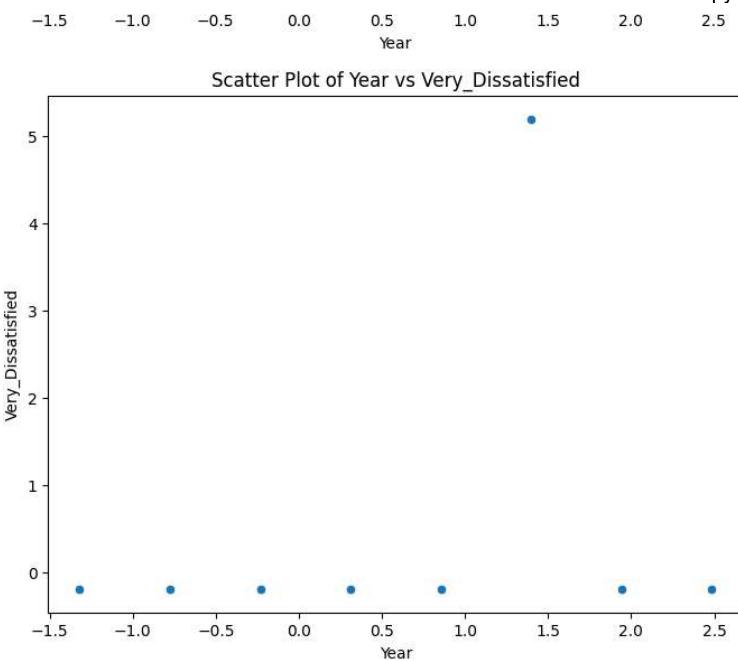
Scatter Plot of Year vs Satisfied



Scatter Plot of Year vs Neutral

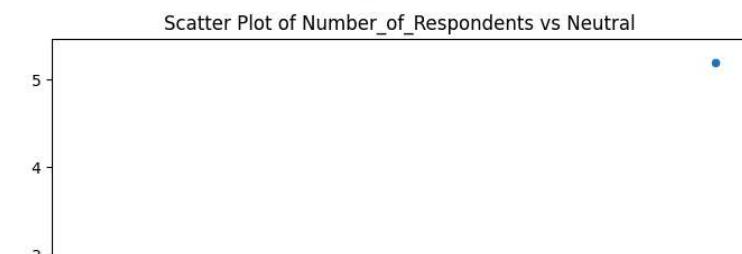
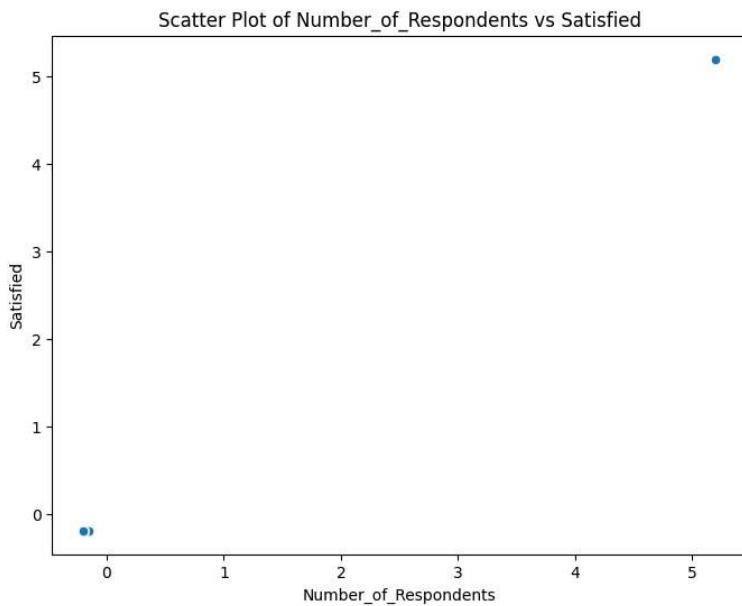
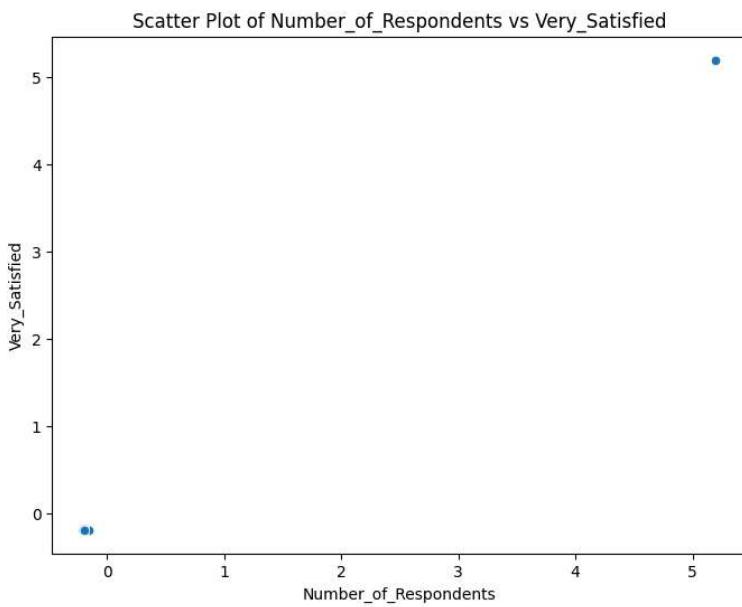
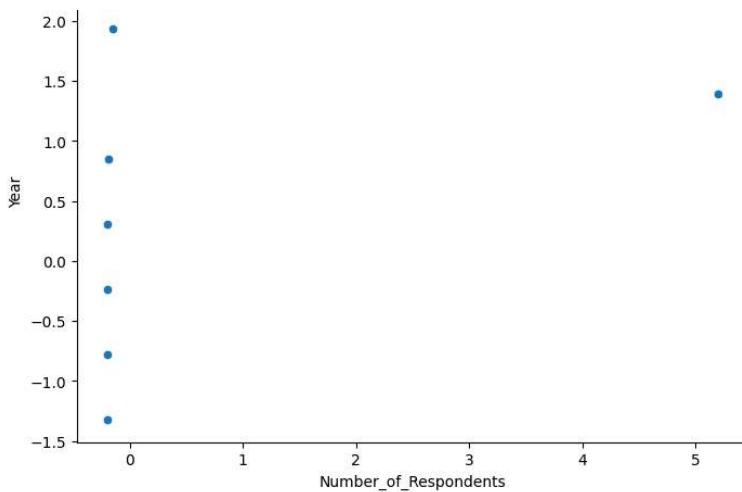


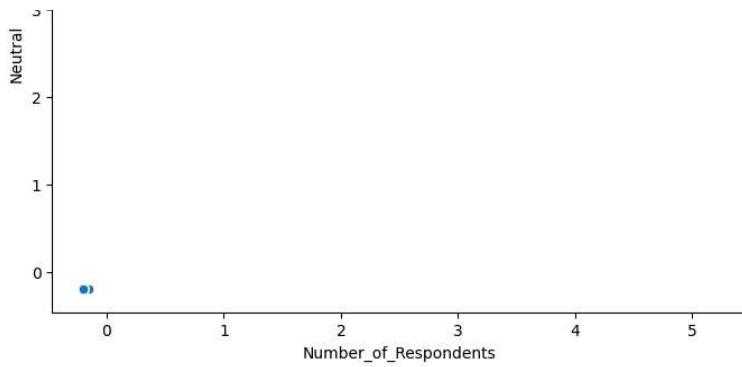
Scatter Plot of Year vs Dissatisfied



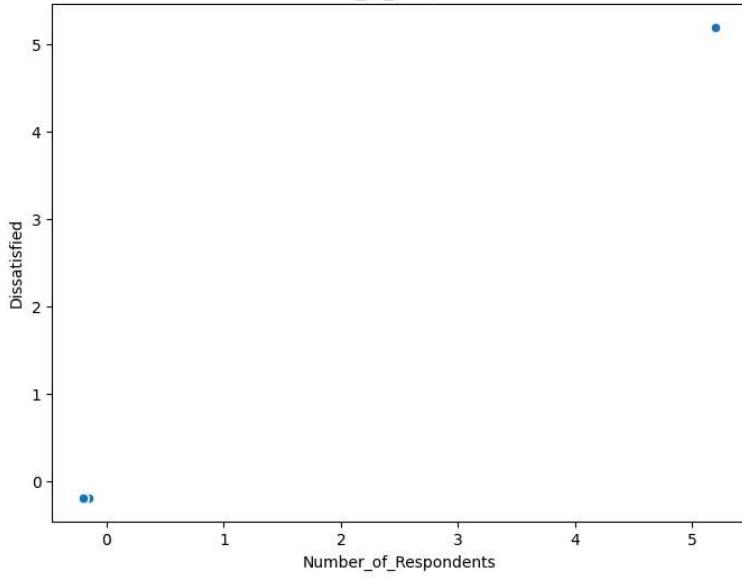
Scatter Plot of Number_of_Respondents vs Year



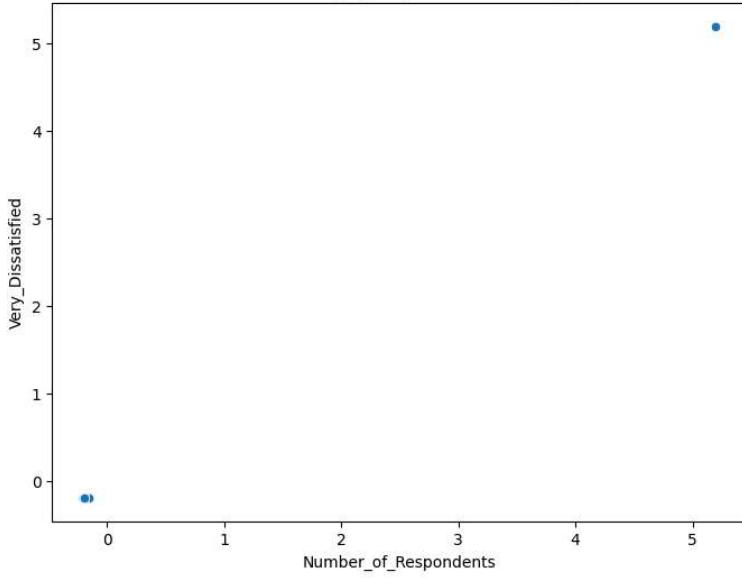




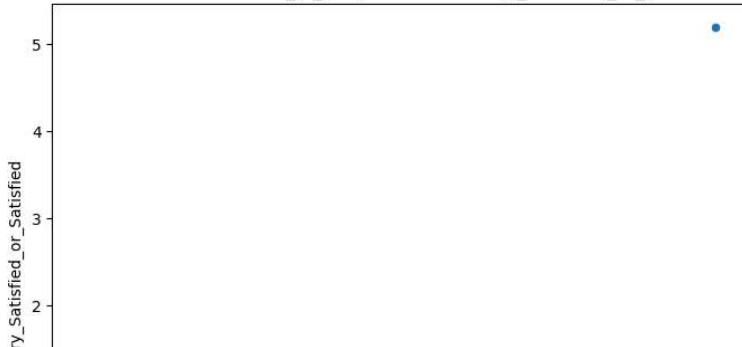
Scatter Plot of Number_of_Respondents vs Dissatisfied

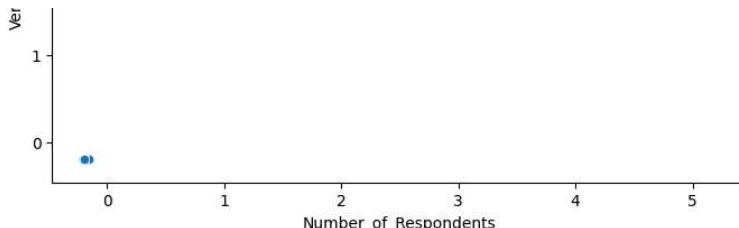


Scatter Plot of Number_of_Respondents vs Very_Dissatisfied

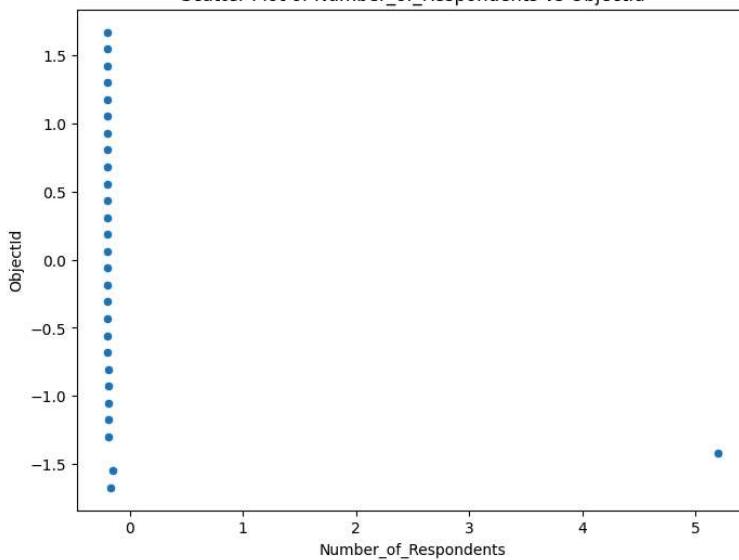


Scatter Plot of Number_of_Respondents vs Very_Satisfied_or_Satisfied

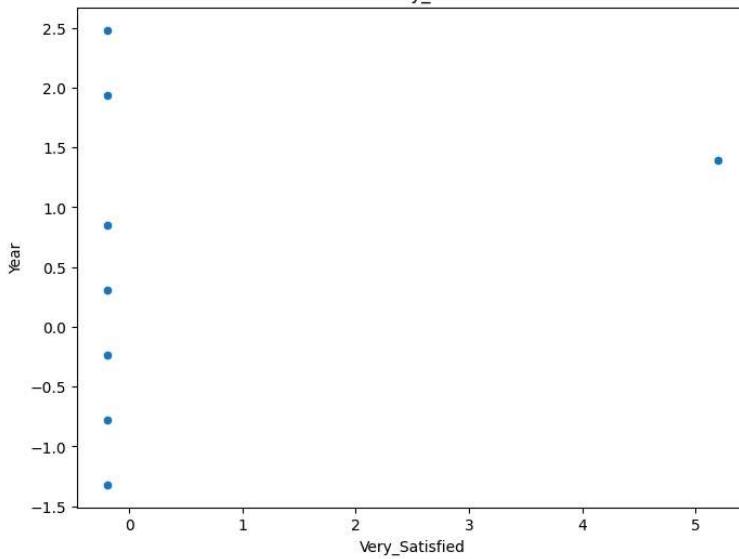




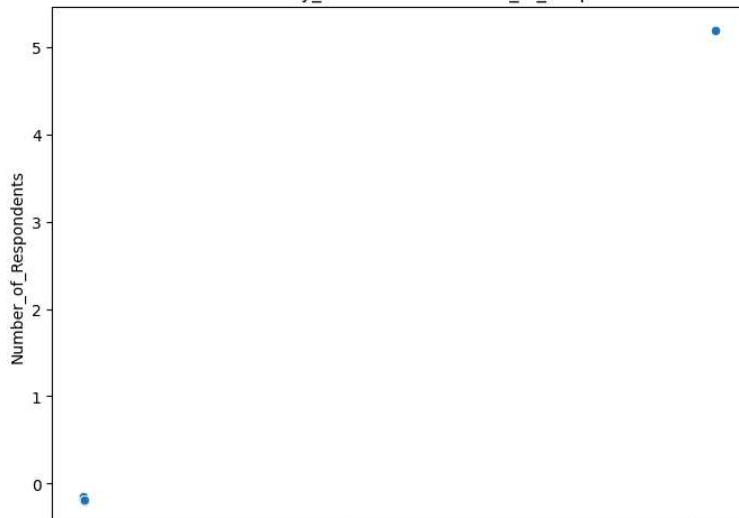
Scatter Plot of Number_of_Respondents vs ObjectId

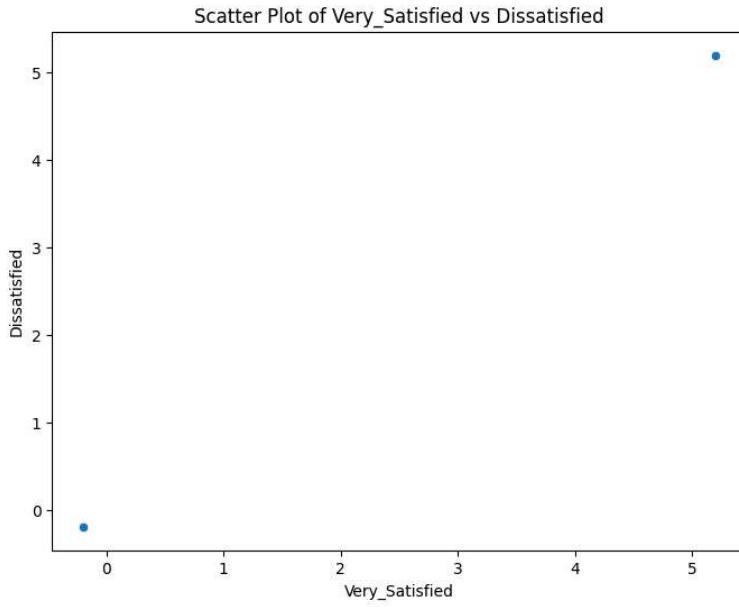
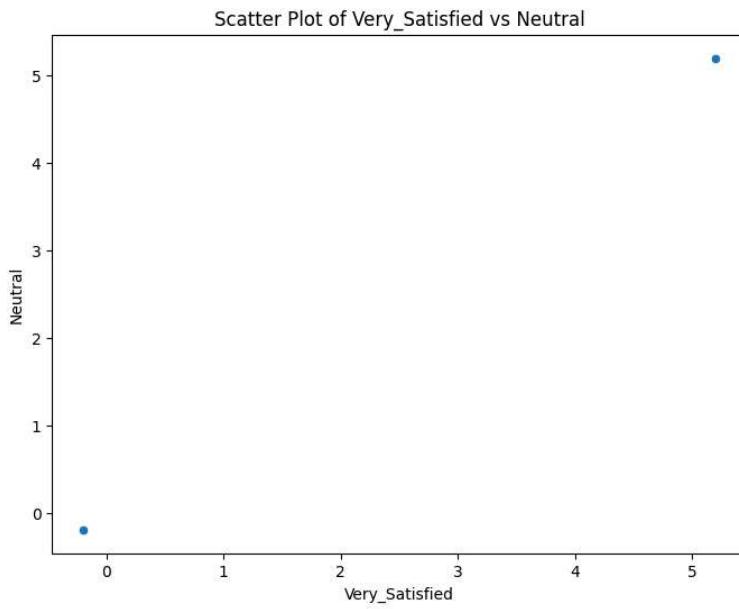
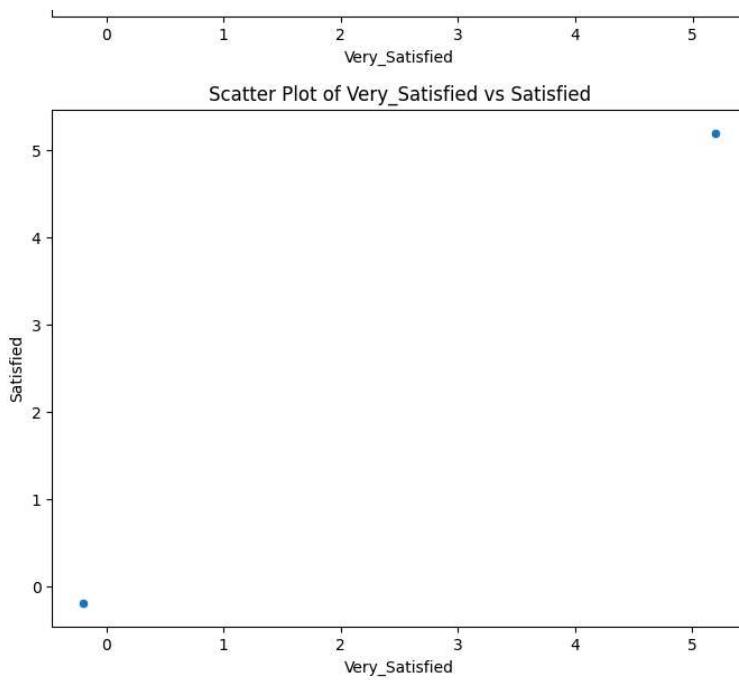


Scatter Plot of Very_Satisfied vs Year

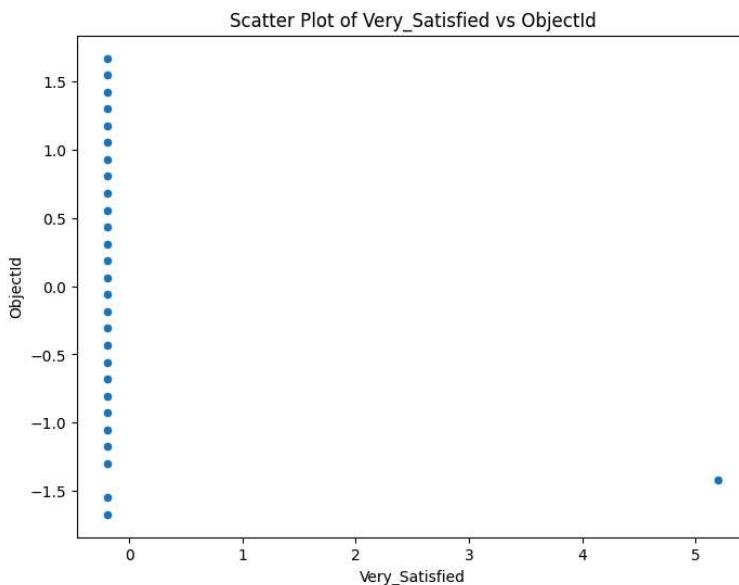
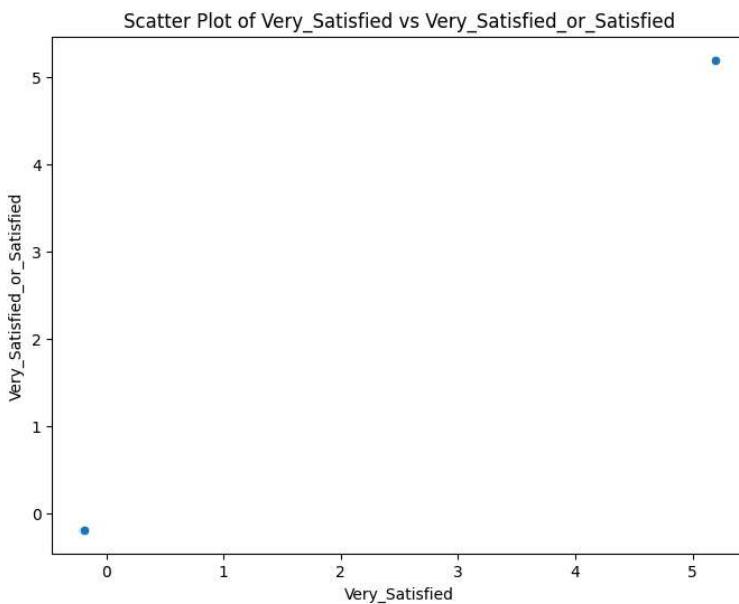
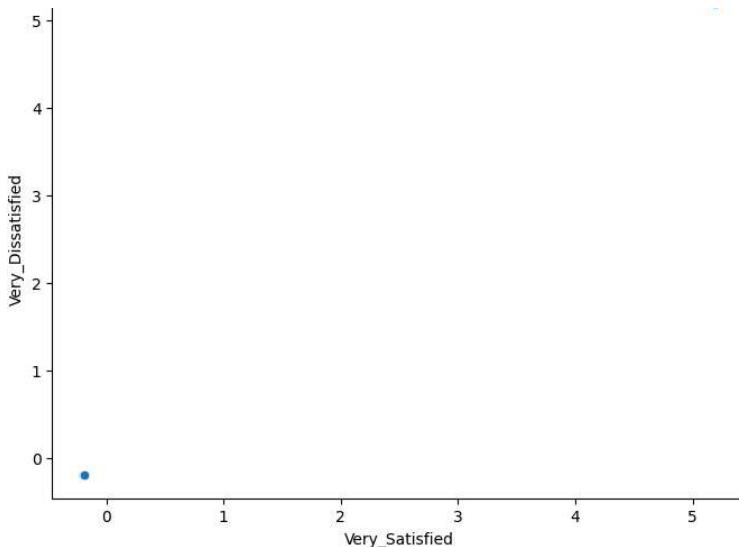


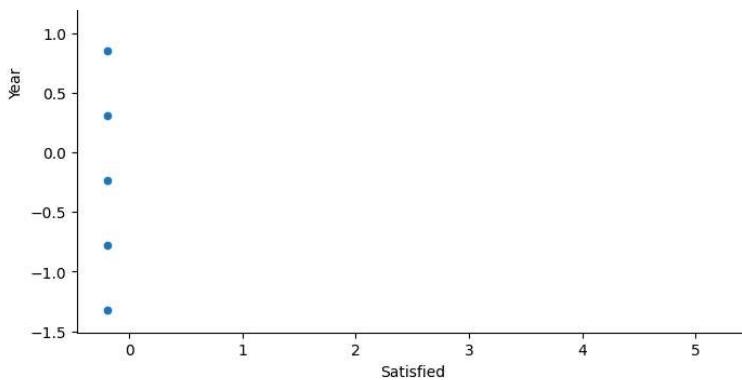
Scatter Plot of Very_Satisfied vs Number_of_Respondents



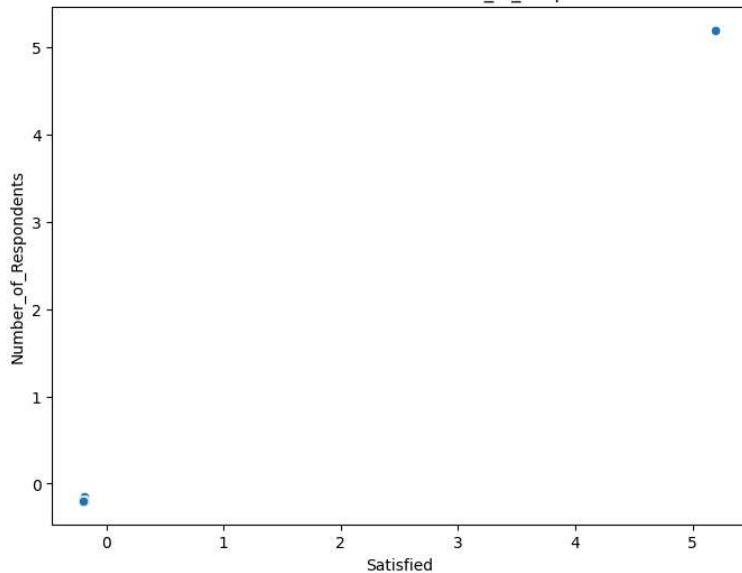


Scatter Plot of Very_Satisfied vs Very_Dissatisfied

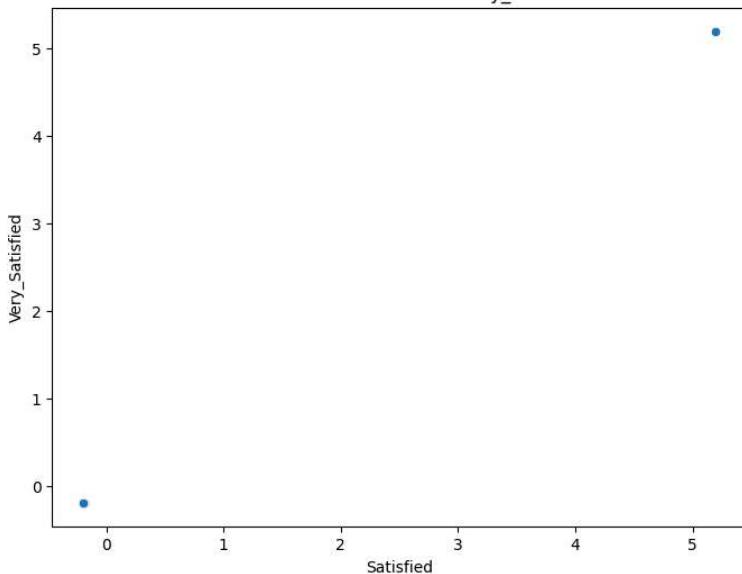




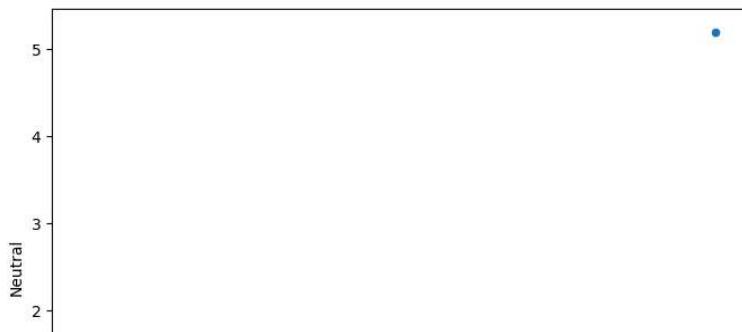
Scatter Plot of Satisfied vs Number_of_Respondents

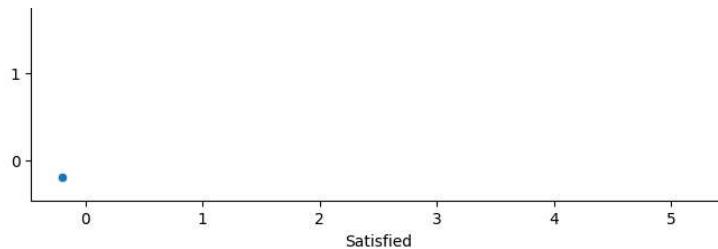


Scatter Plot of Satisfied vs Very_Satisfied

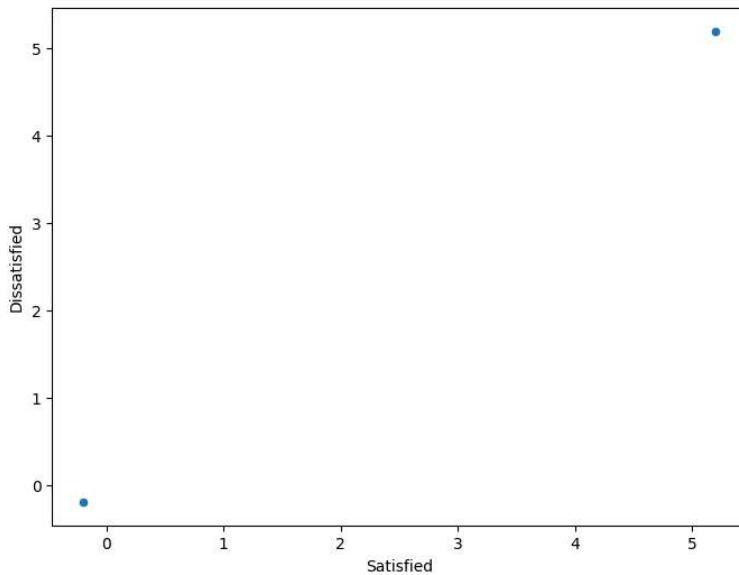


Scatter Plot of Satisfied vs Neutral

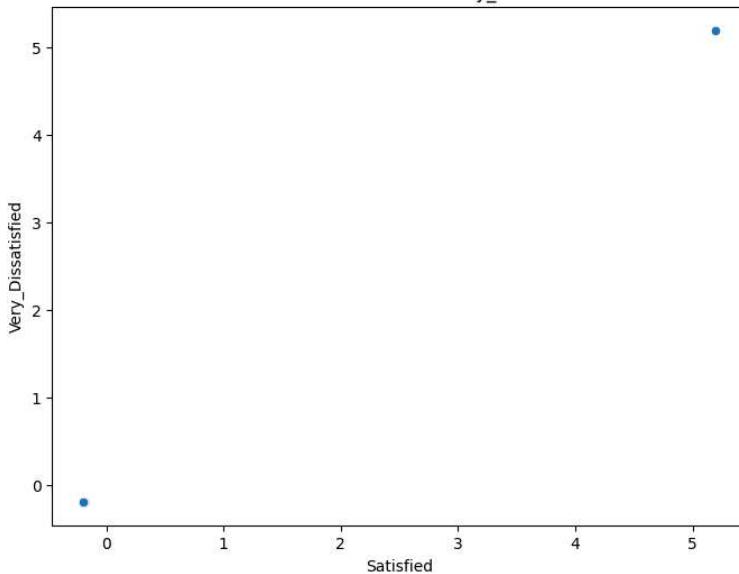




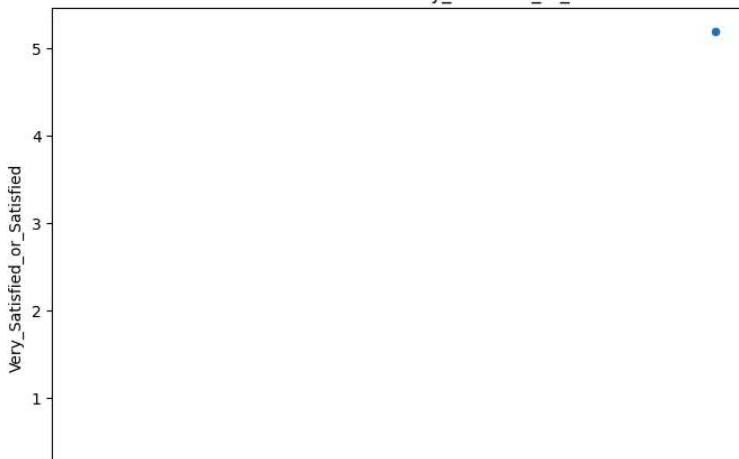
Scatter Plot of Satisfied vs Dissatisfied

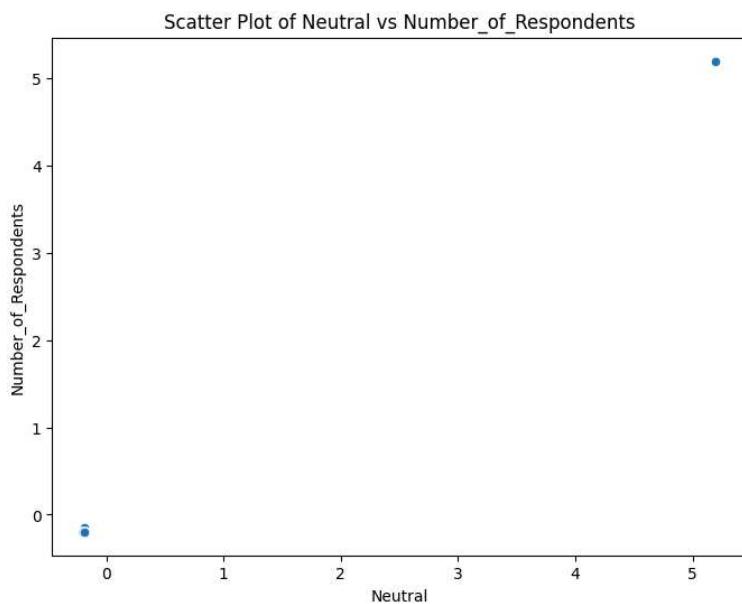
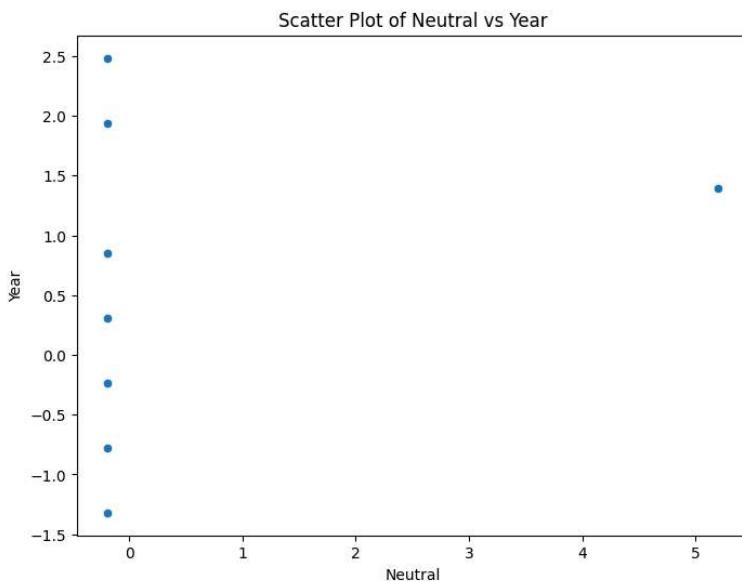
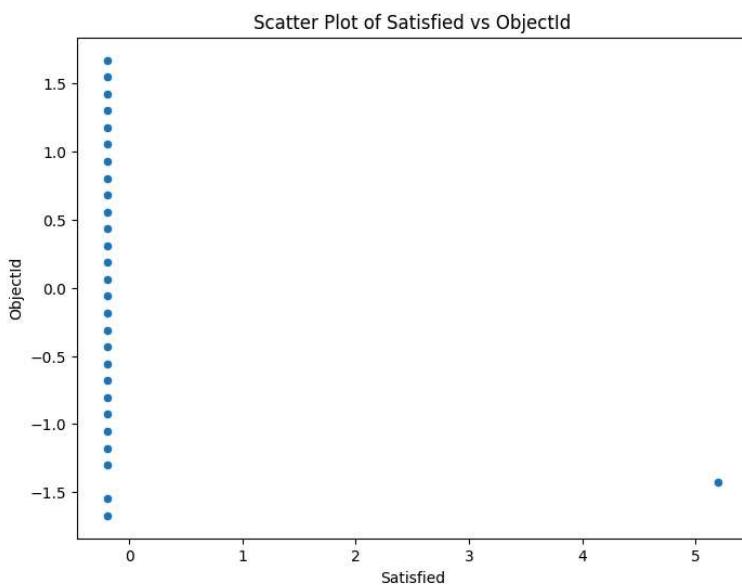


Scatter Plot of Satisfied vs Very_Dissatisfied

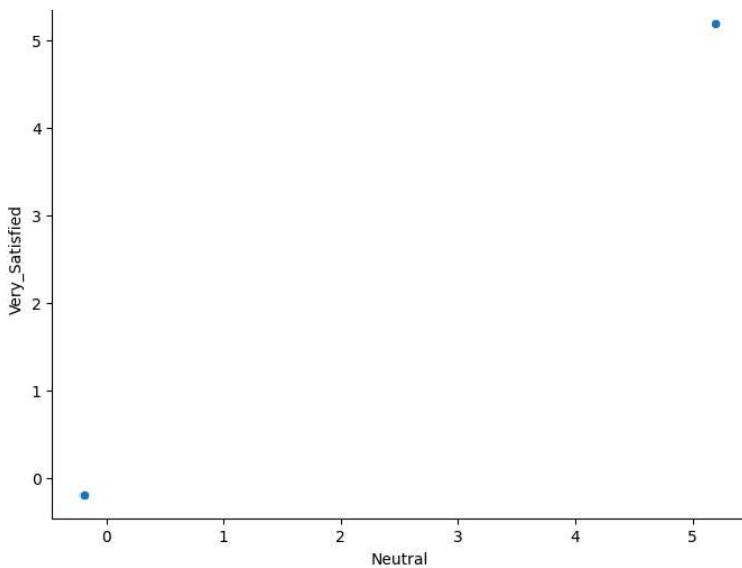


Scatter Plot of Satisfied vs Very_Satisfied_or_Satisfied

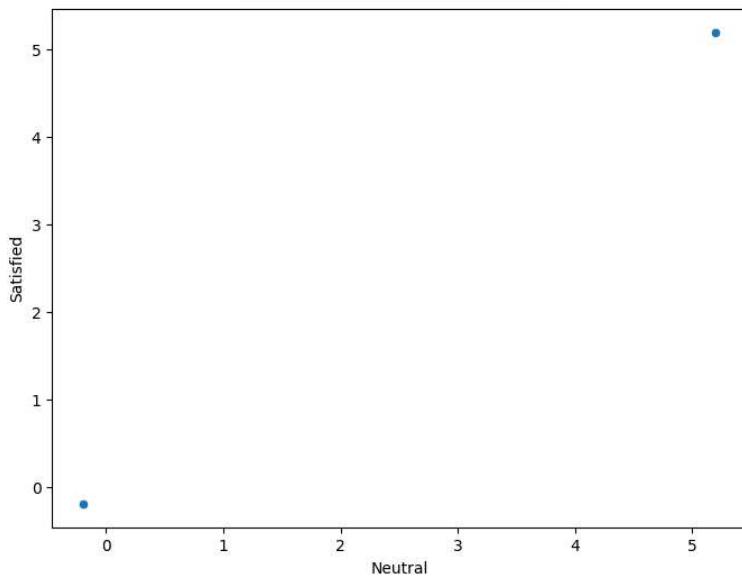




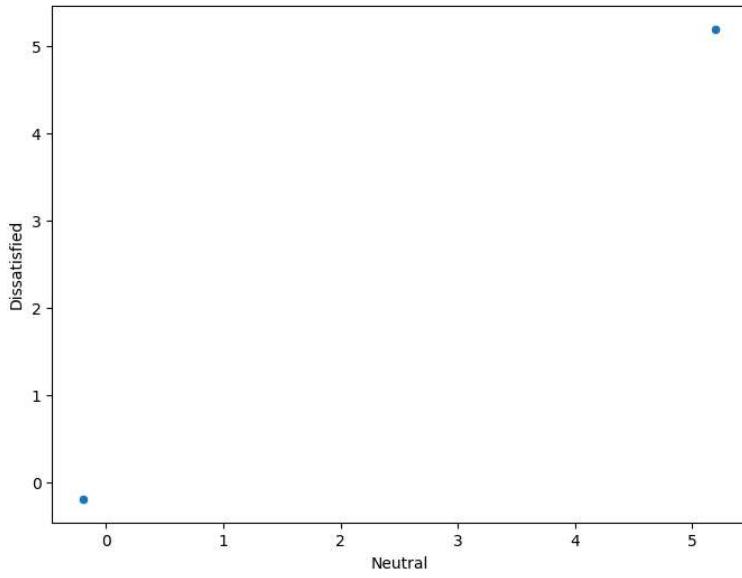
Scatter Plot of Neutral vs Very_Satisfied



Scatter Plot of Neutral vs Satisfied

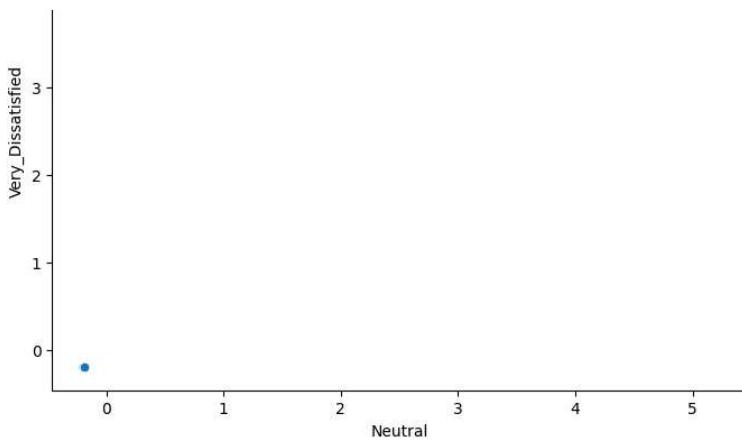


Scatter Plot of Neutral vs Dissatisfied

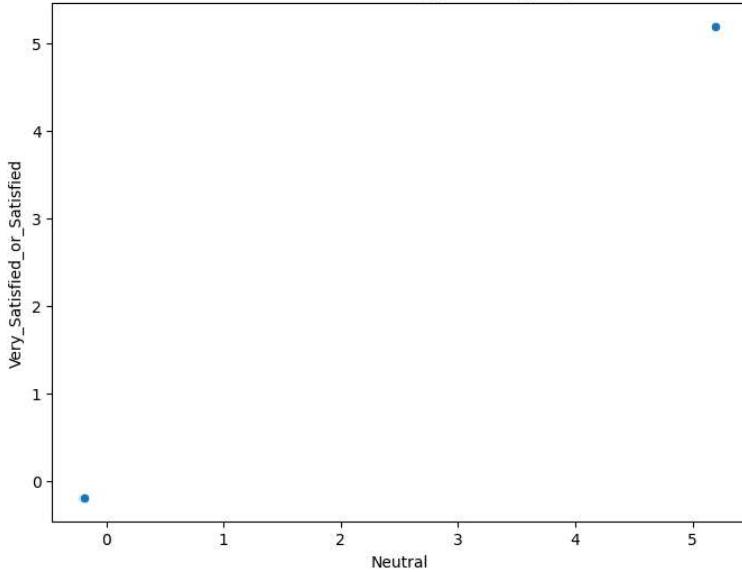


Scatter Plot of Neutral vs Very_Dissatisfied

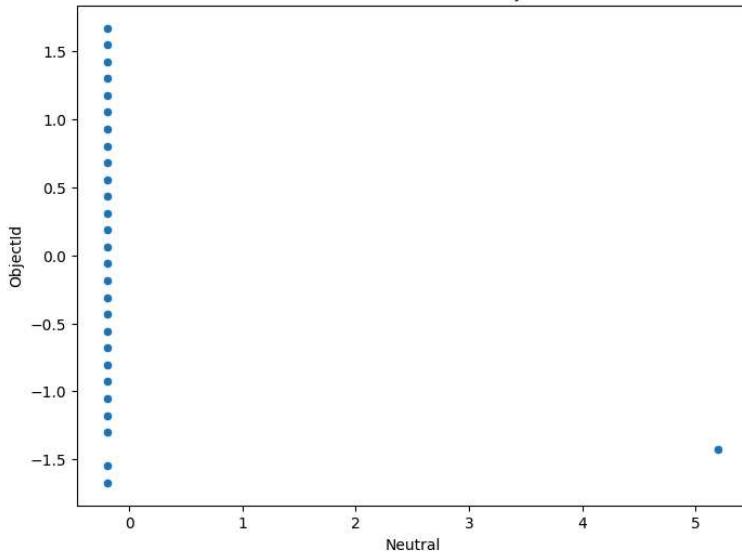




Scatter Plot of Neutral vs Very_Dissatisfied

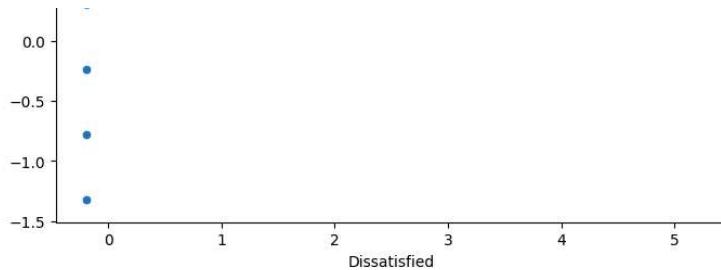


Scatter Plot of Neutral vs Objectid

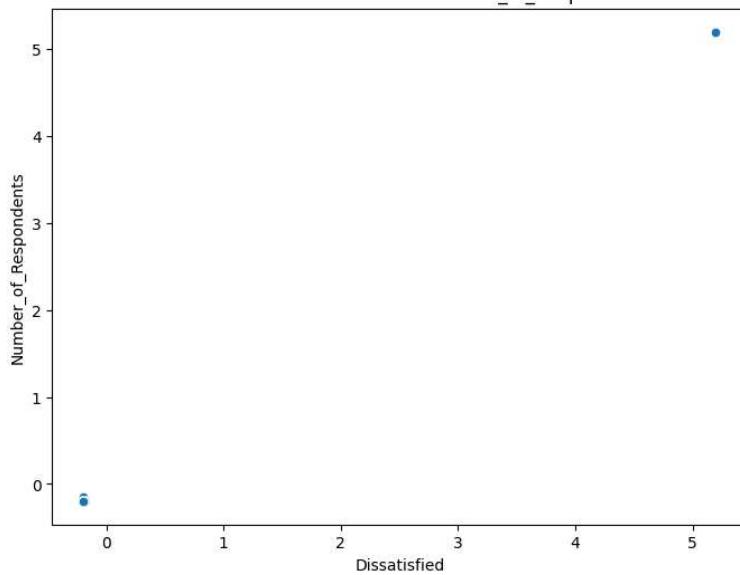


Scatter Plot of Dissatisfied vs Year

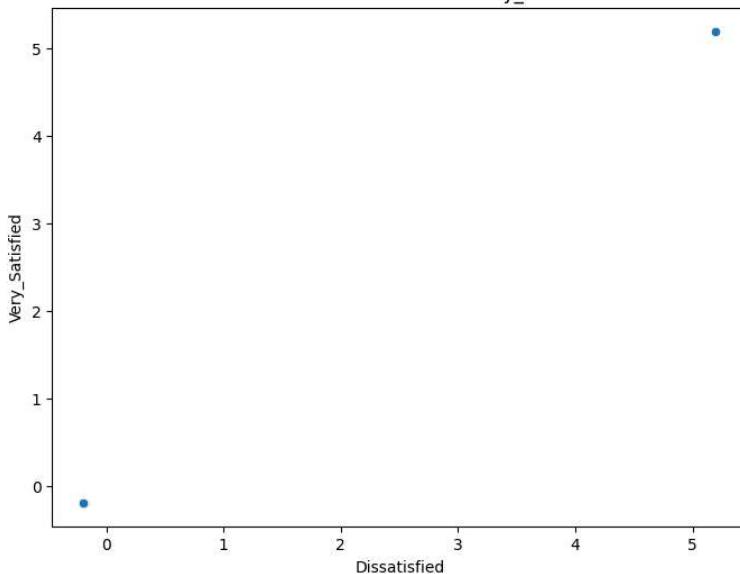




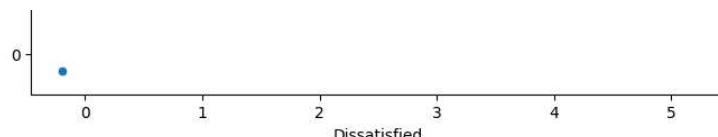
Scatter Plot of Dissatisfied vs Number_of_Respondents



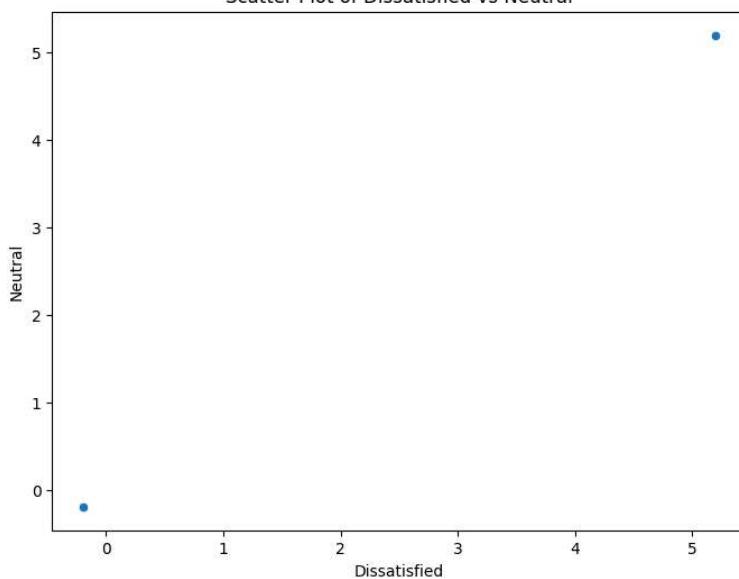
Scatter Plot of Dissatisfied vs Very_Satisfied



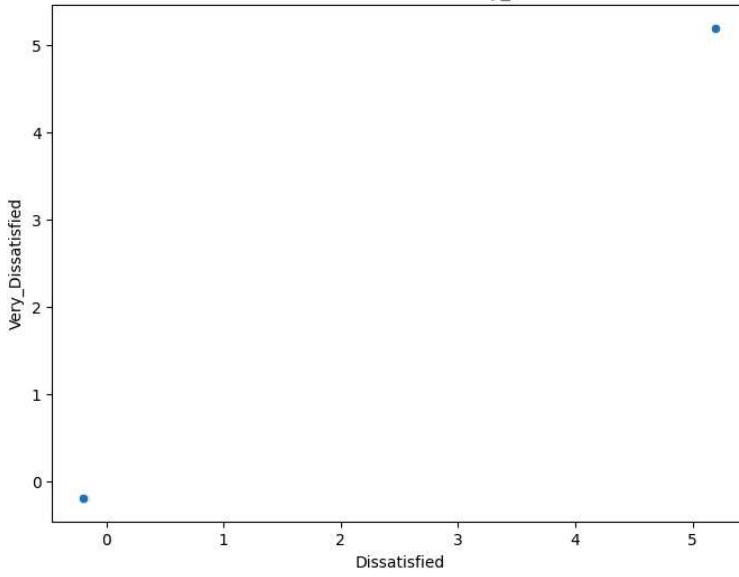
Scatter Plot of Dissatisfied vs Satisfied



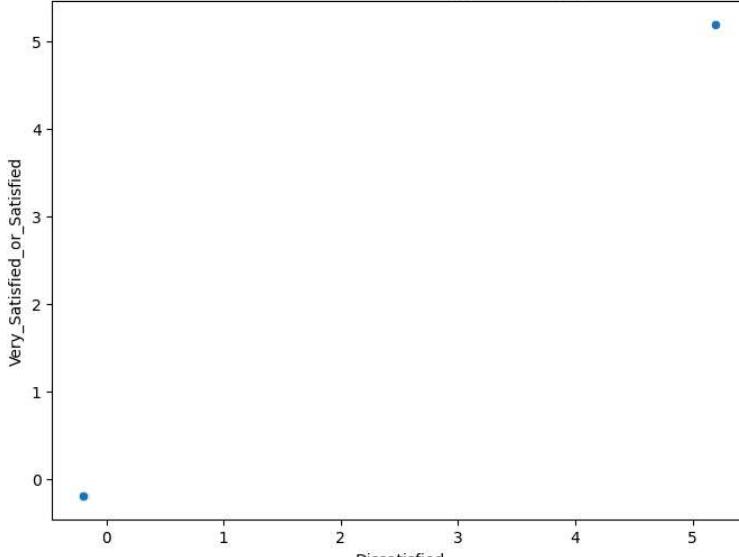
Scatter Plot of Dissatisfied vs Neutral

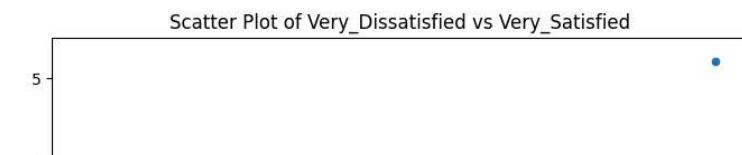
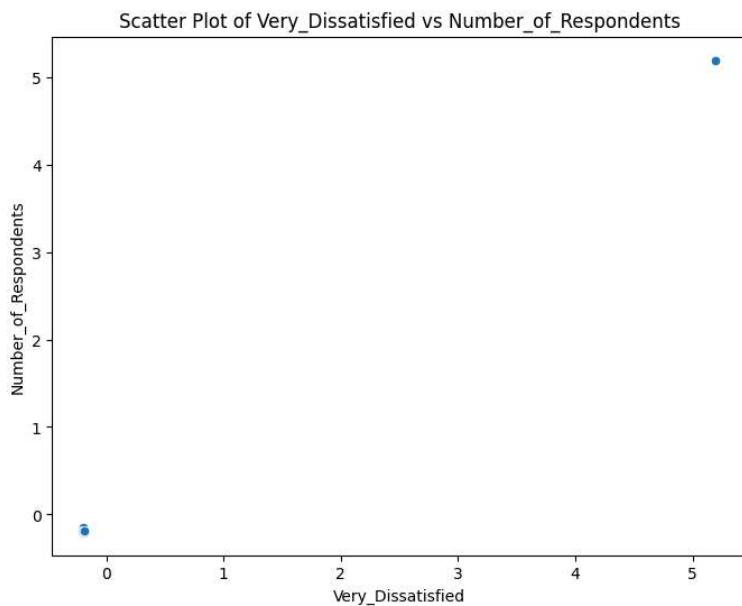
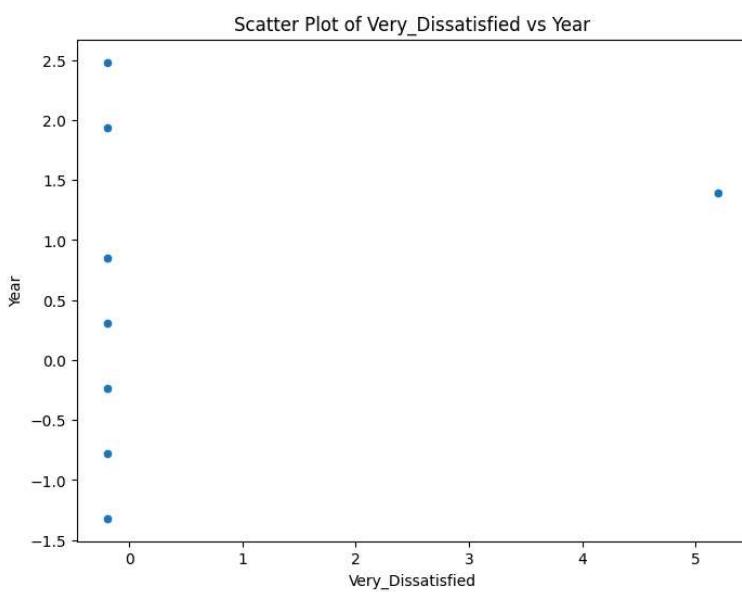
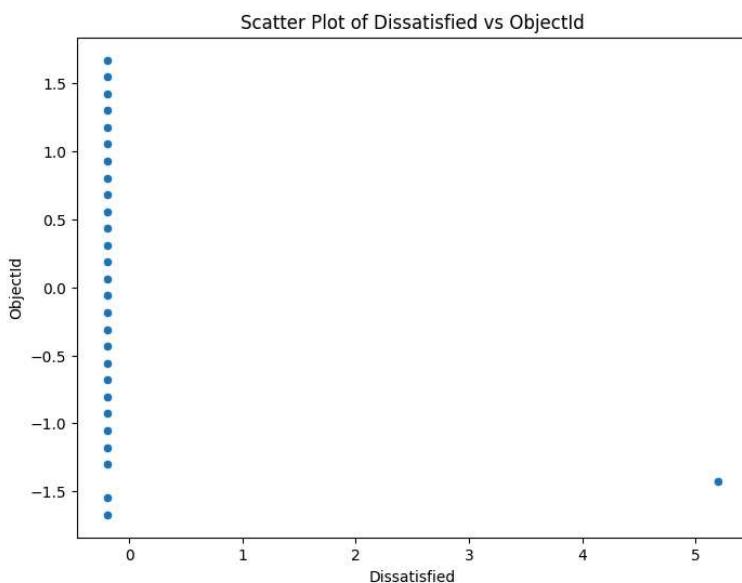


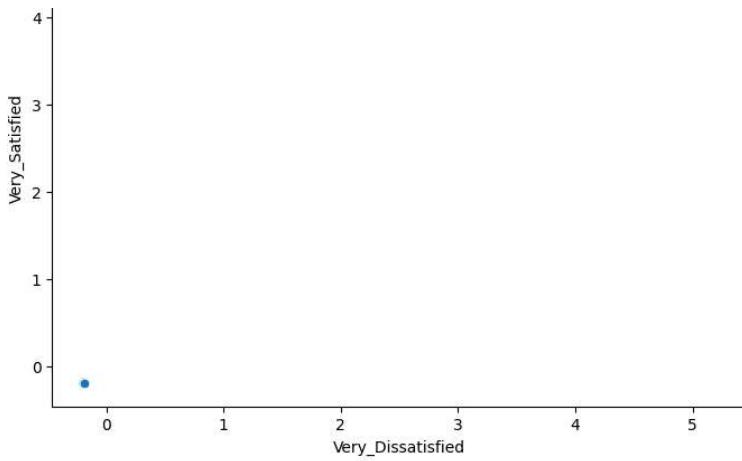
Scatter Plot of Dissatisfied vs Very_Dissatisfied



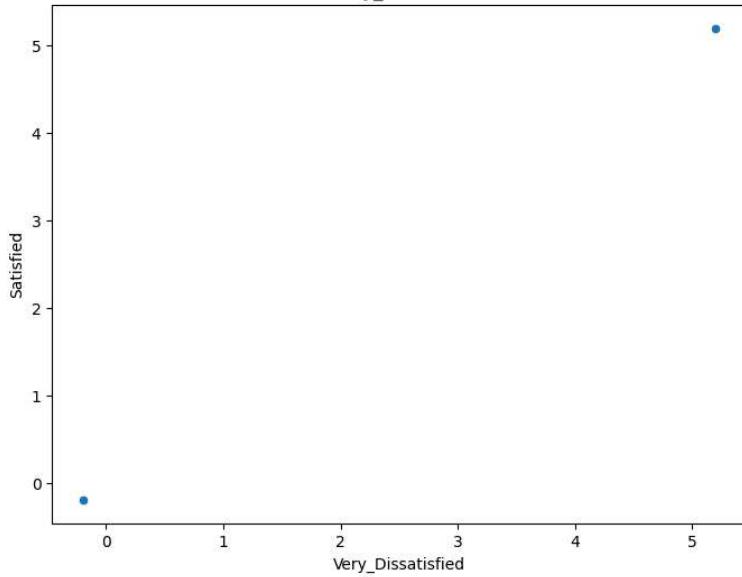
Scatter Plot of Dissatisfied vs Very_Satisfied_or_Satisfied



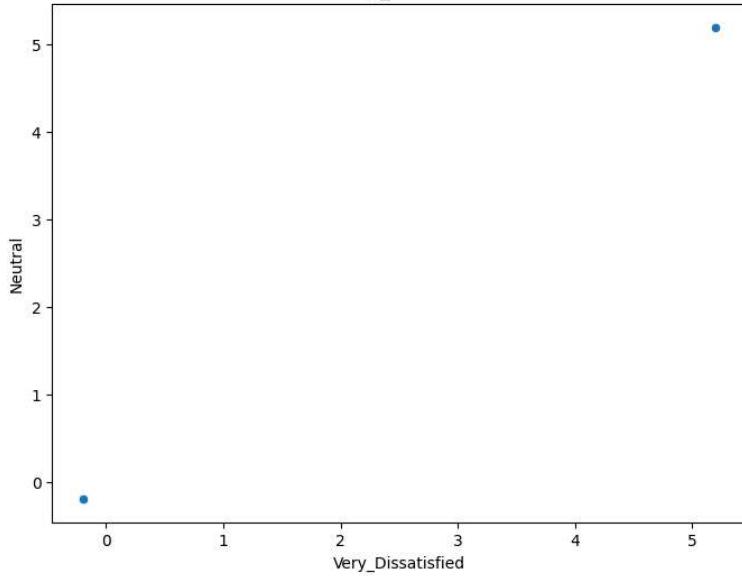




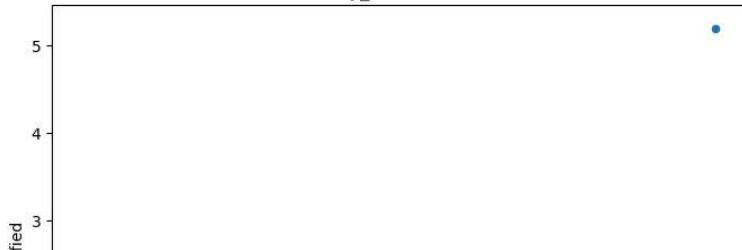
Scatter Plot of Very_Dissatisfied vs Satisfied

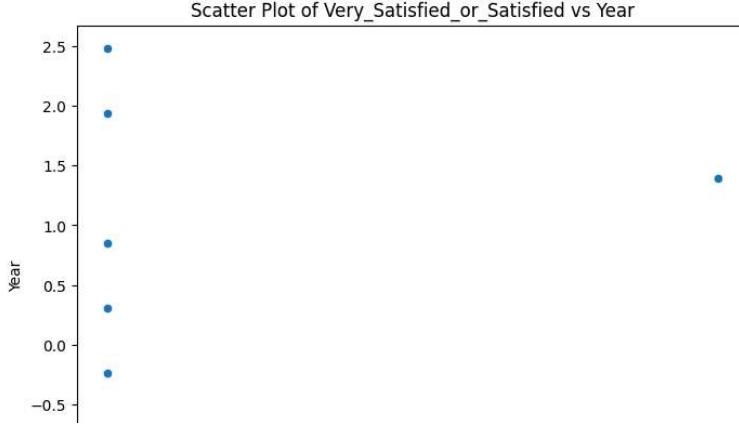
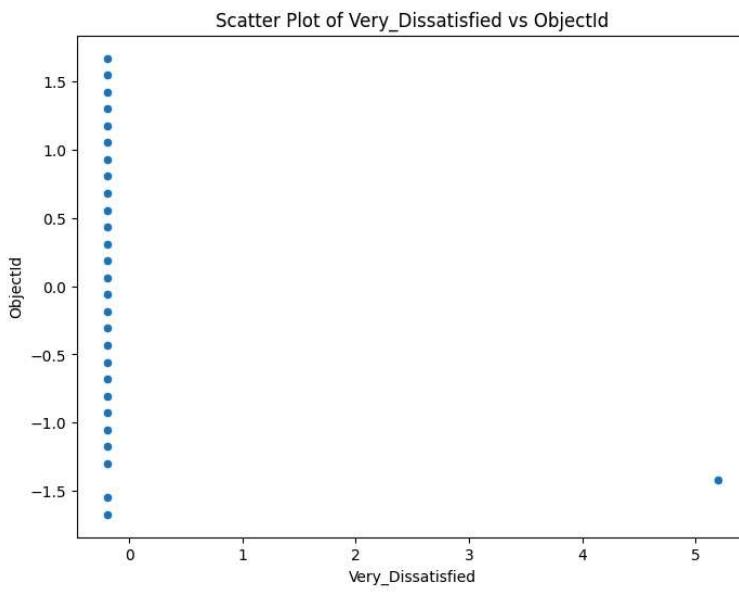
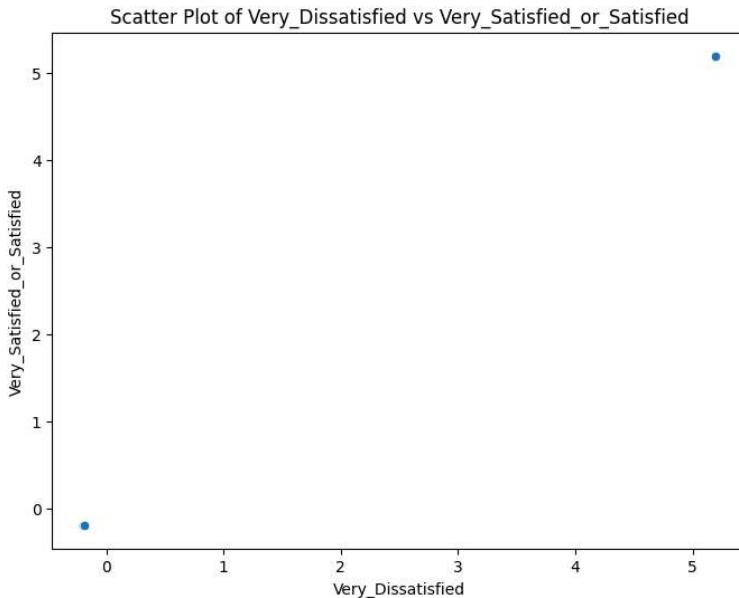
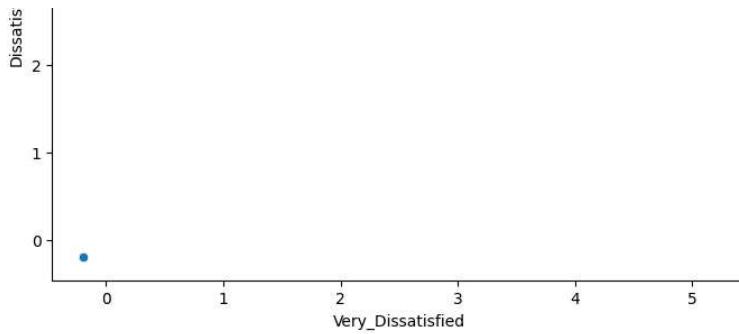


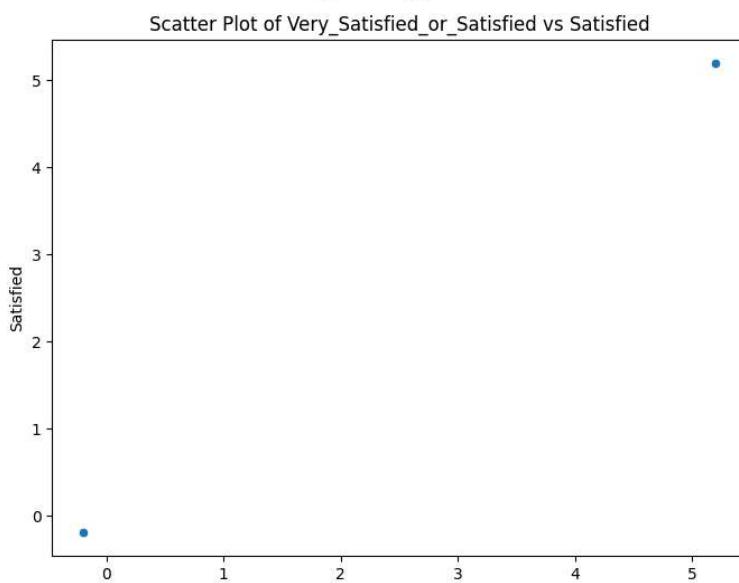
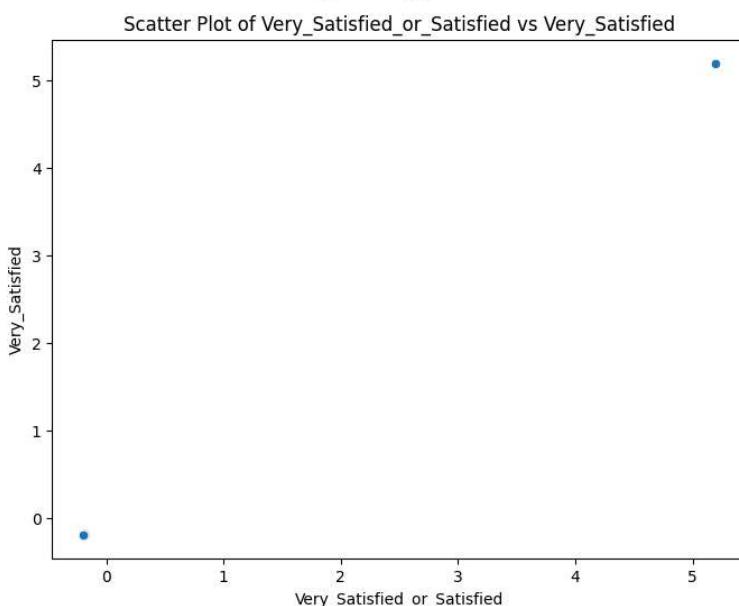
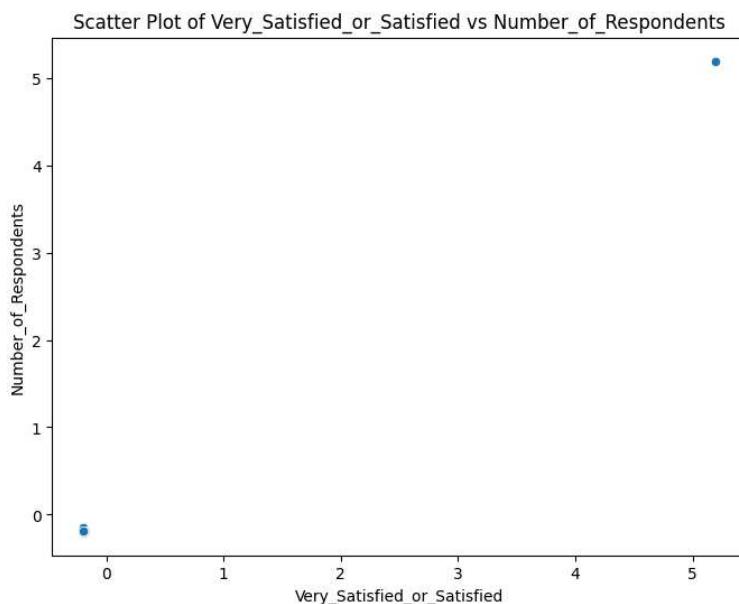
Scatter Plot of Very_Dissatisfied vs Neutral

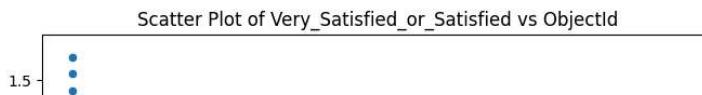
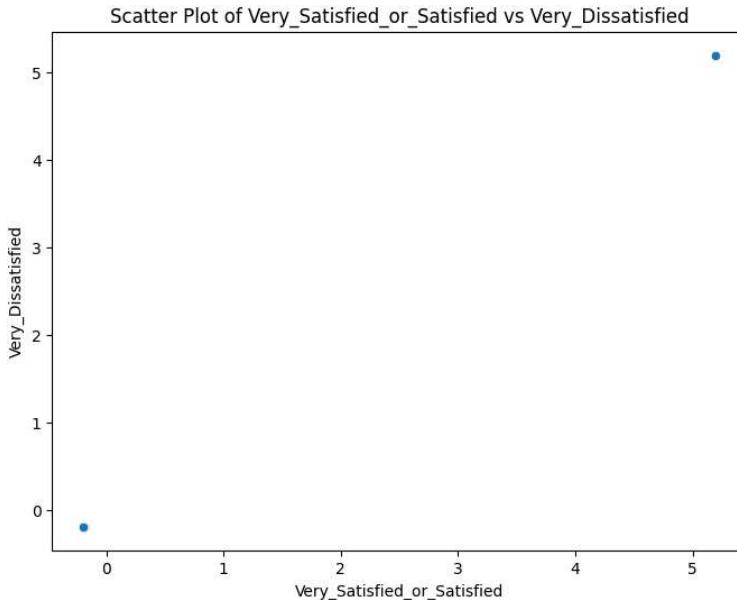
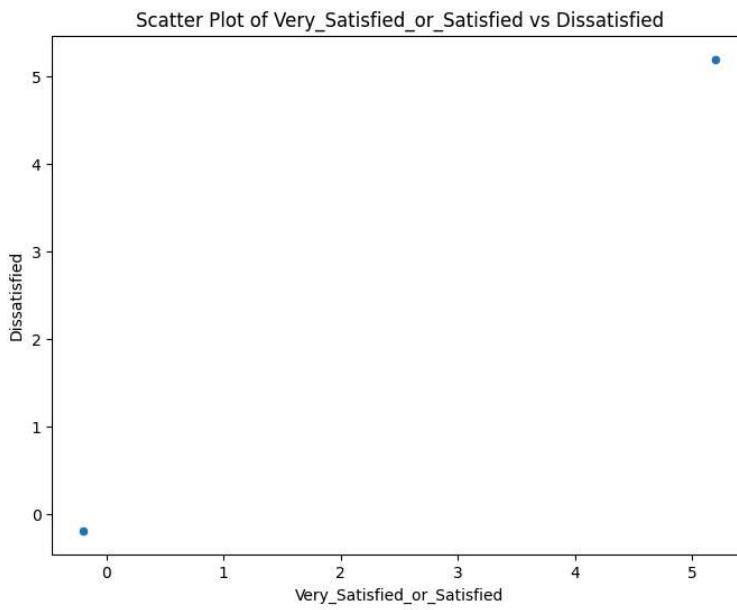
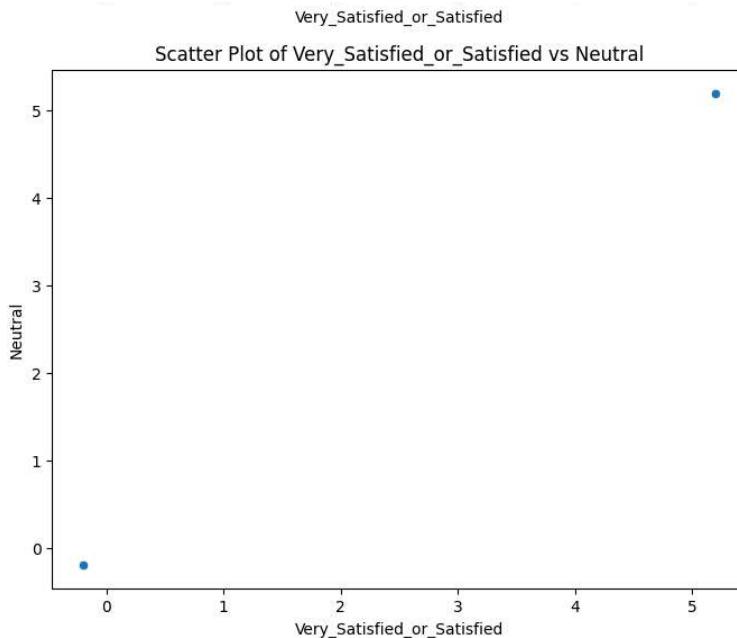


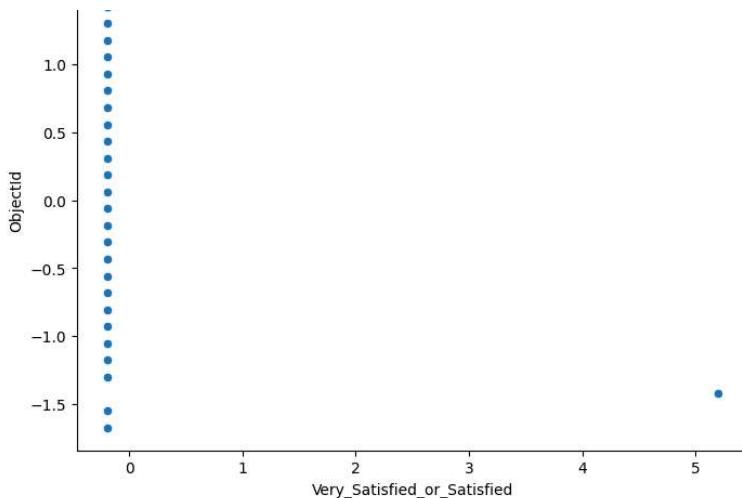
Scatter Plot of Very_Dissatisfied vs Dissatisfied



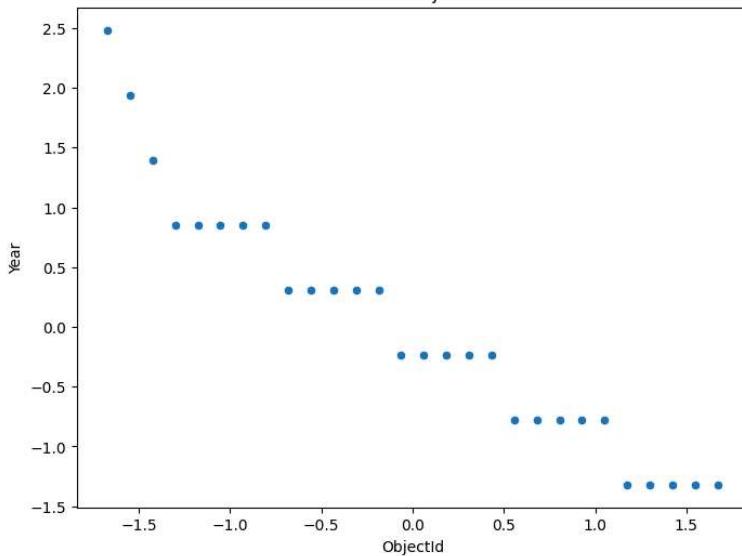




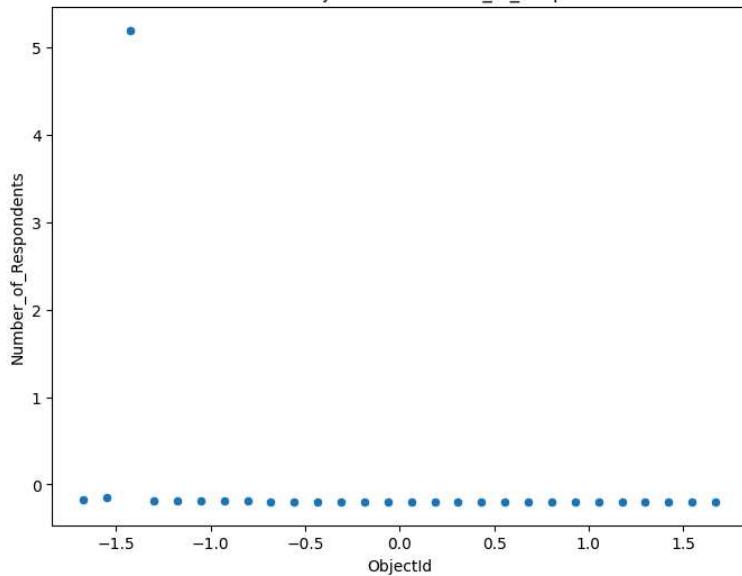




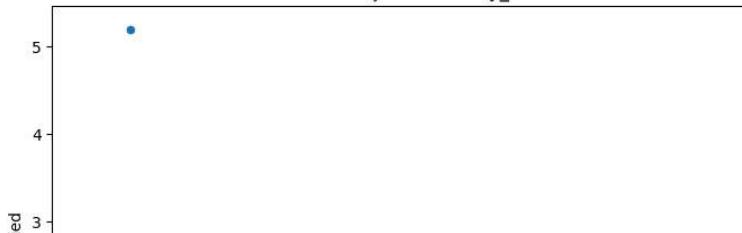
Scatter Plot of Objectid vs Year

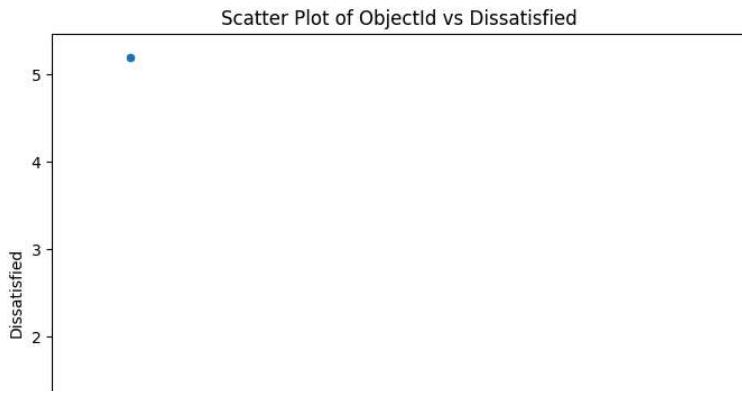
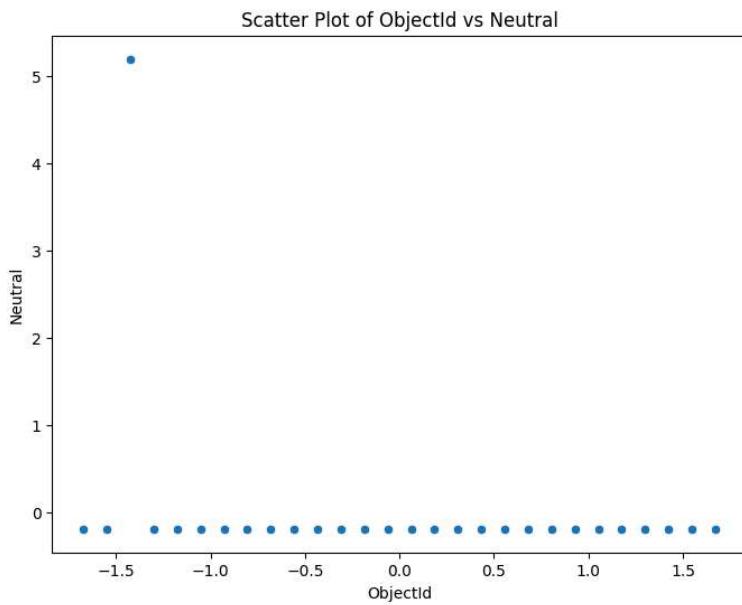
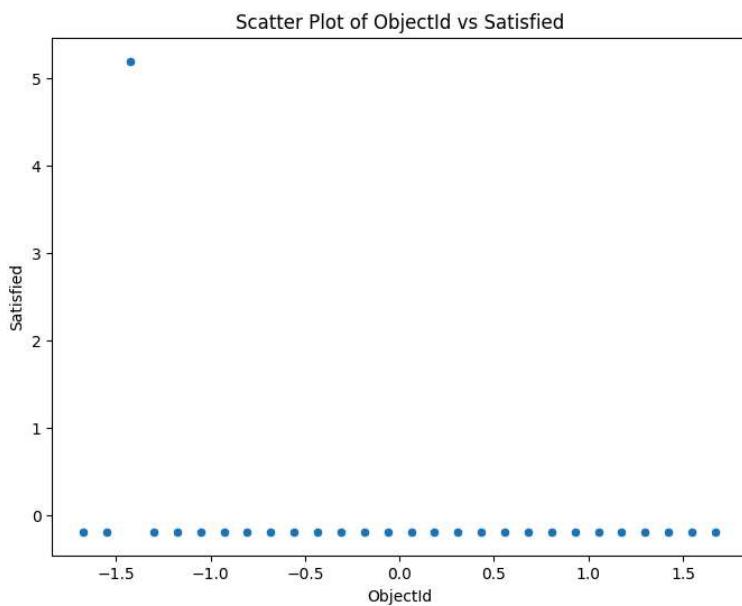
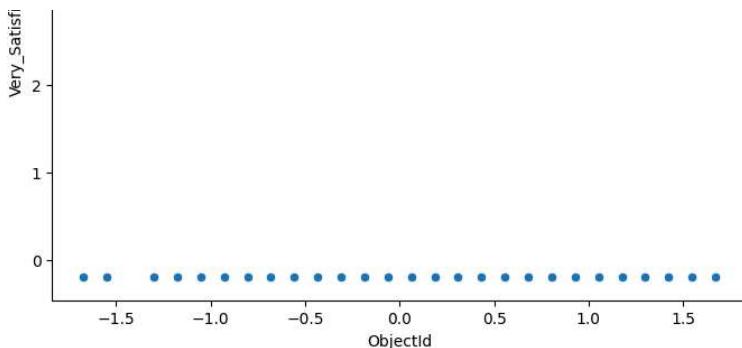


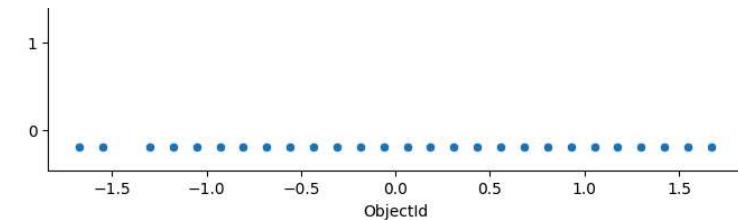
Scatter Plot of Objectid vs Number_of_Respondents



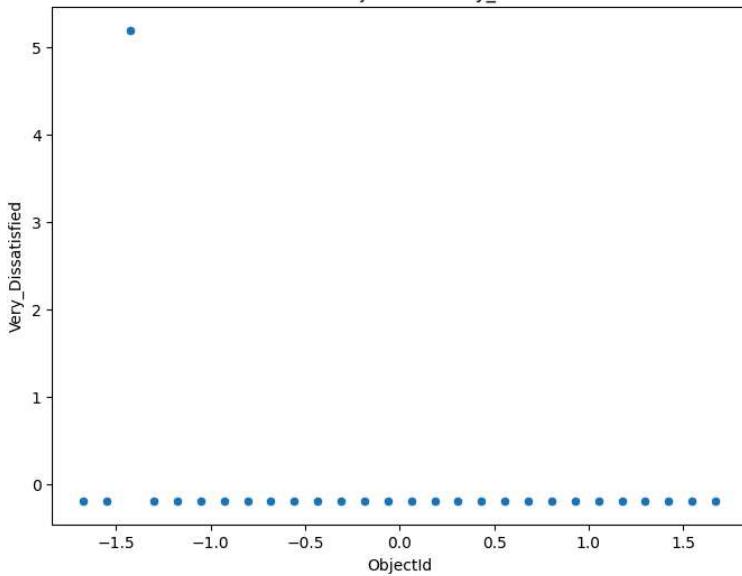
Scatter Plot of Objectid vs Very_Satisfied



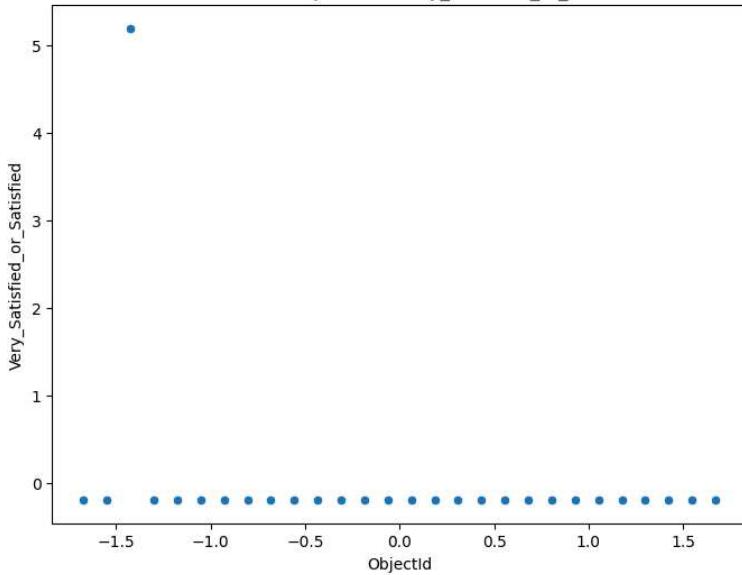




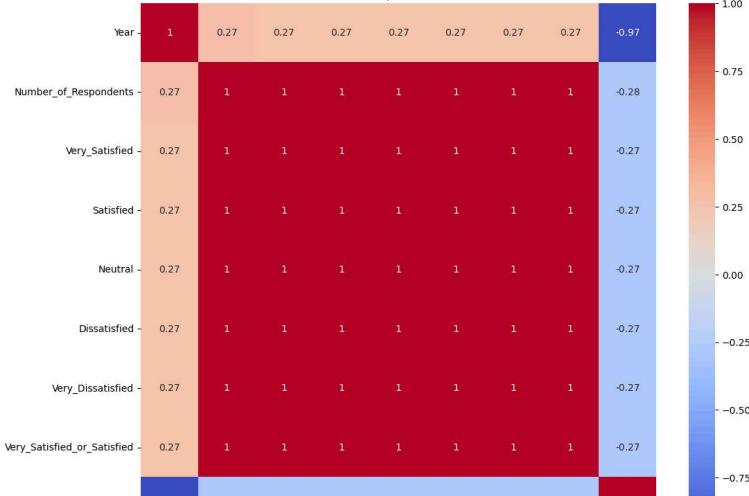
Scatter Plot of Objectid vs Very_Dissatisfied

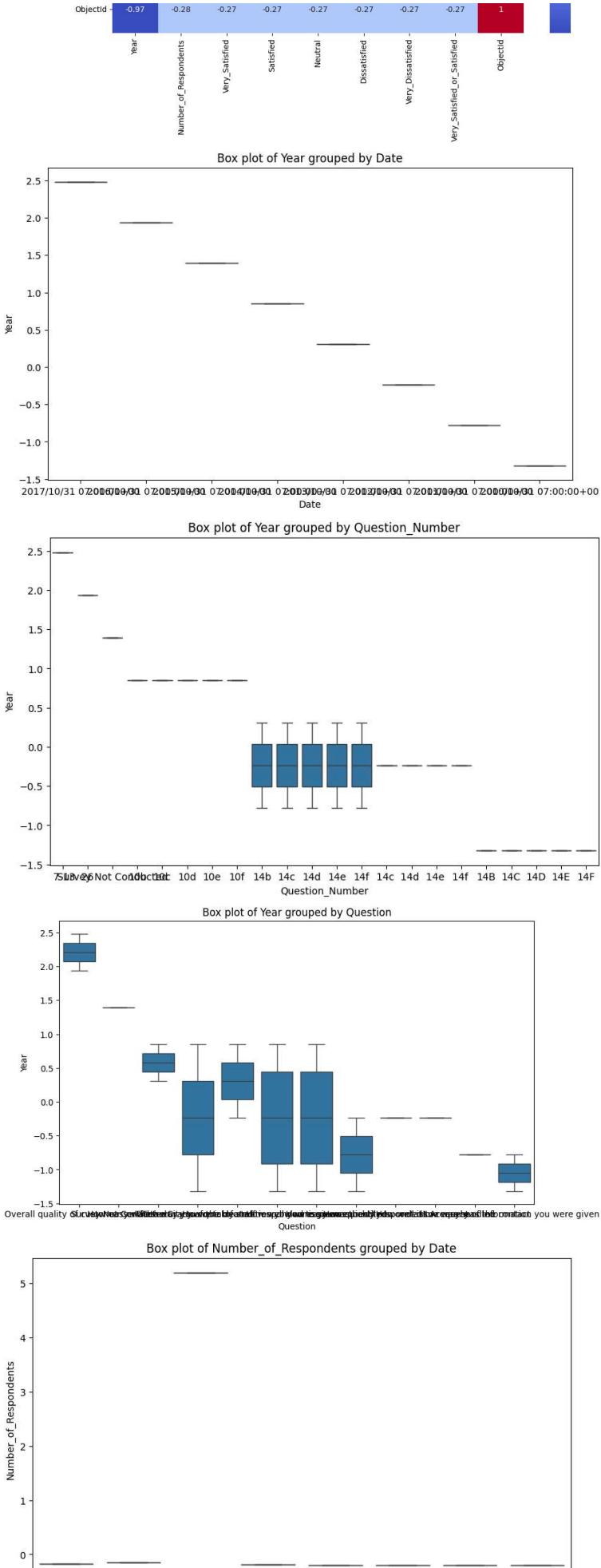


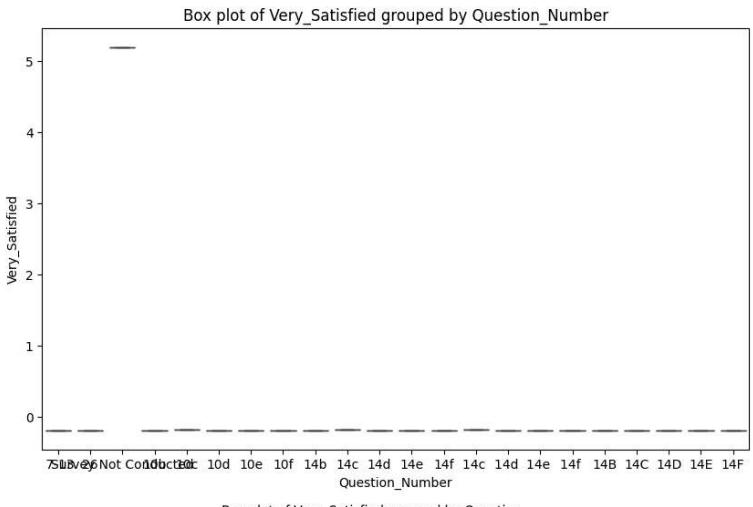
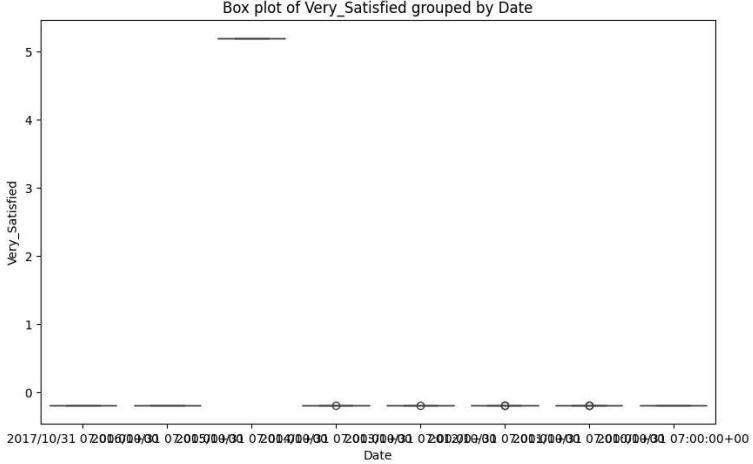
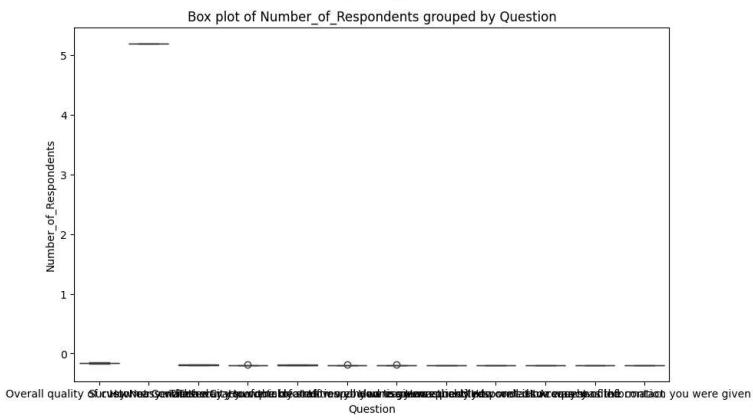
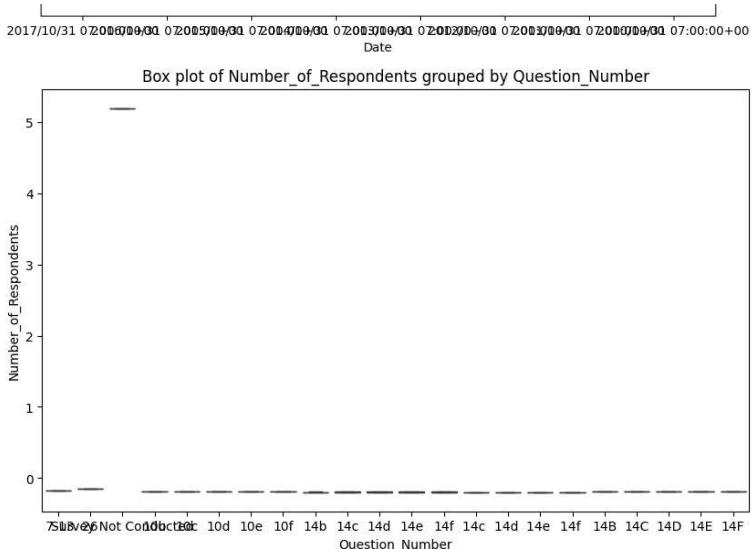
Scatter Plot of Objectid vs Very_Satisfied_or_Satisfied



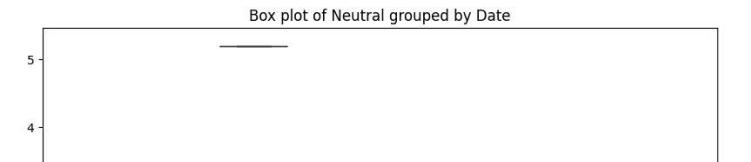
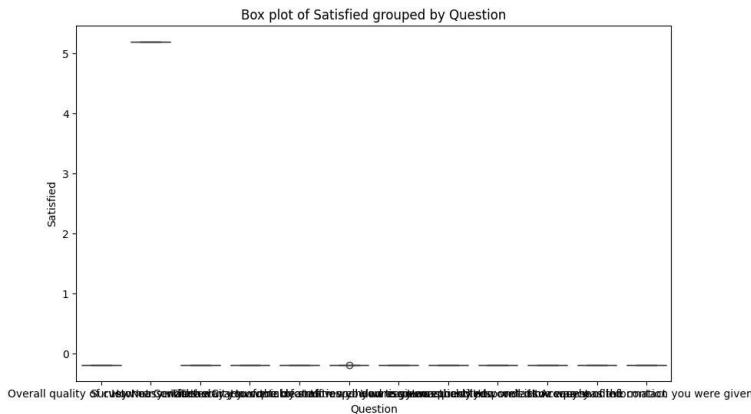
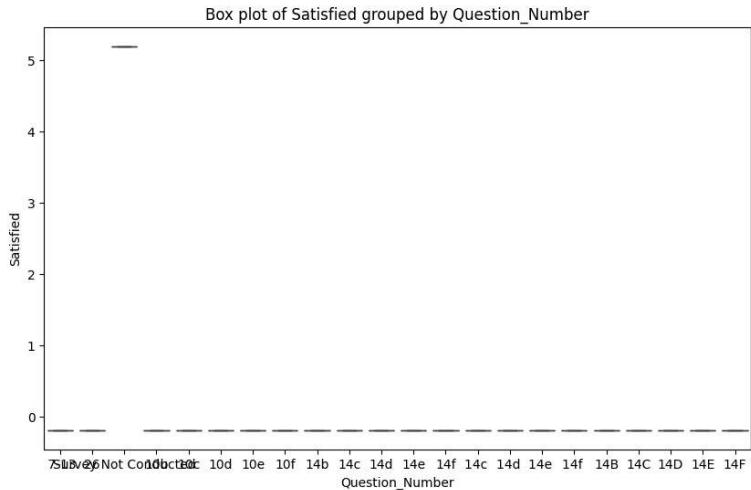
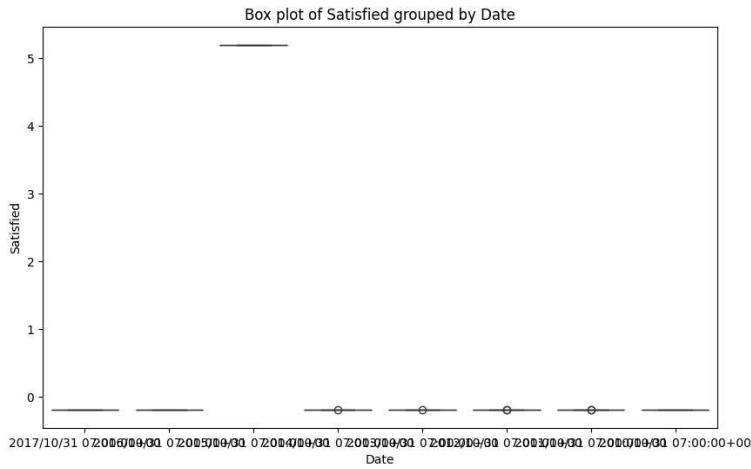
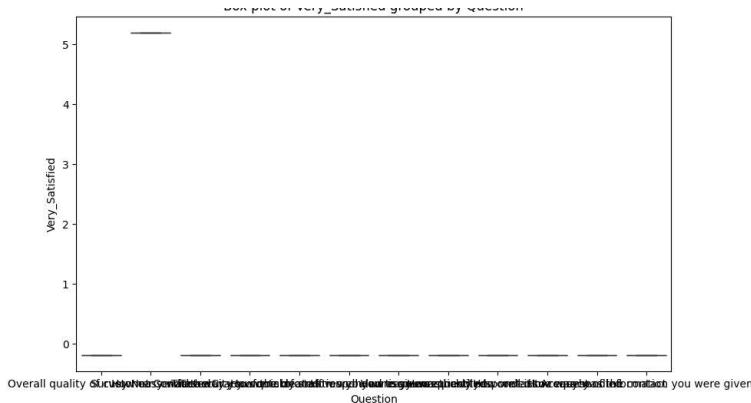
Correlation Heatmap of Numerical Features

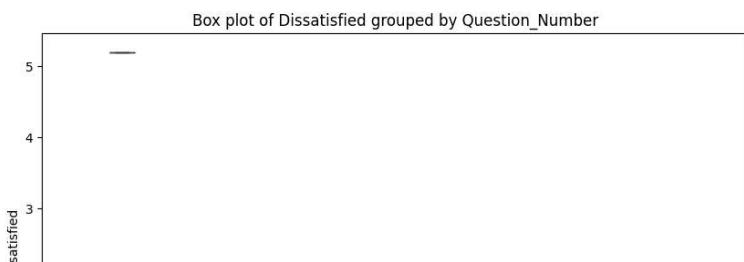
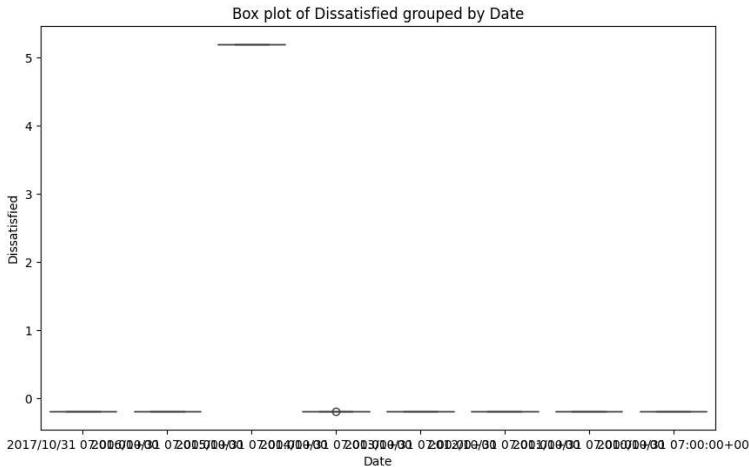
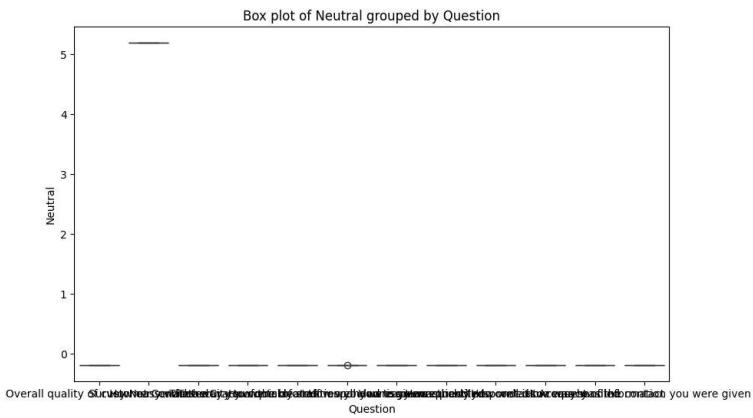
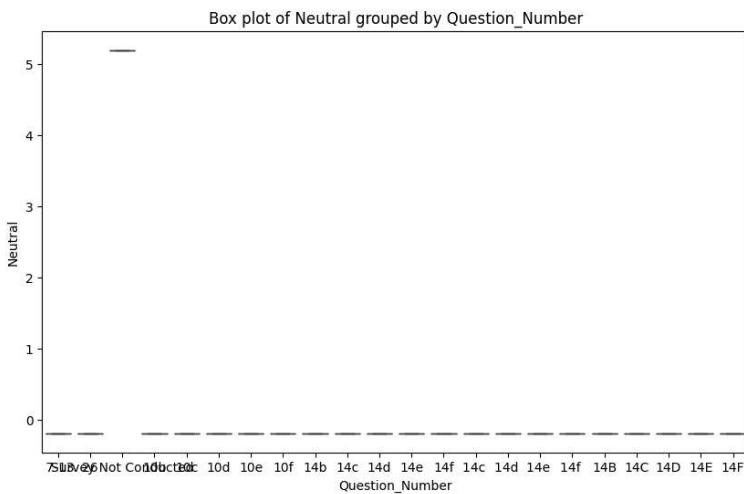
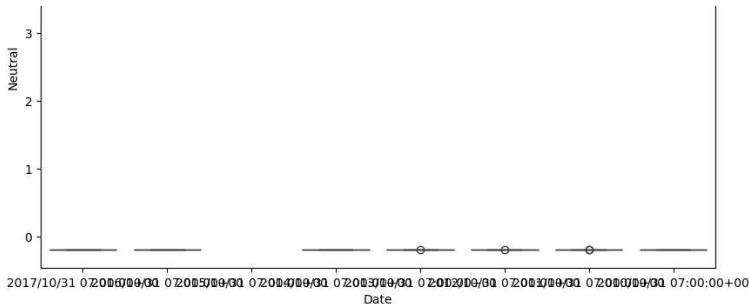


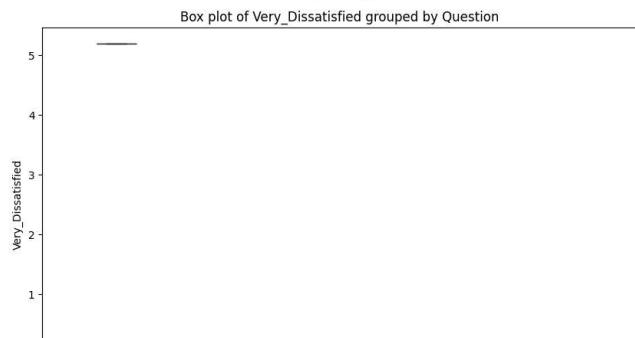
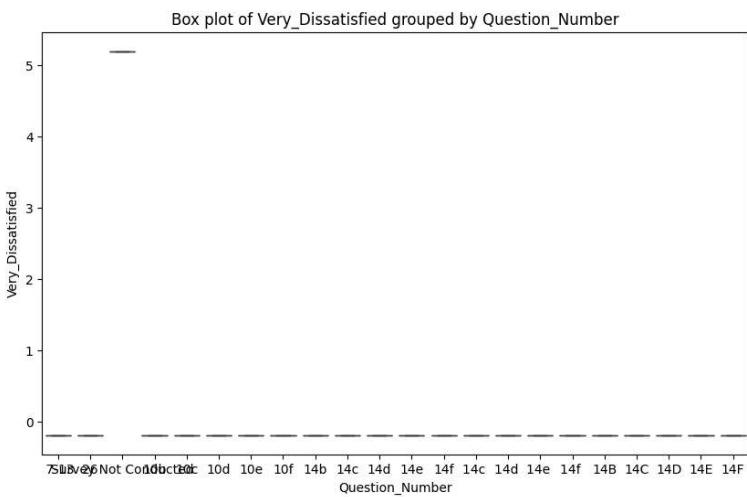
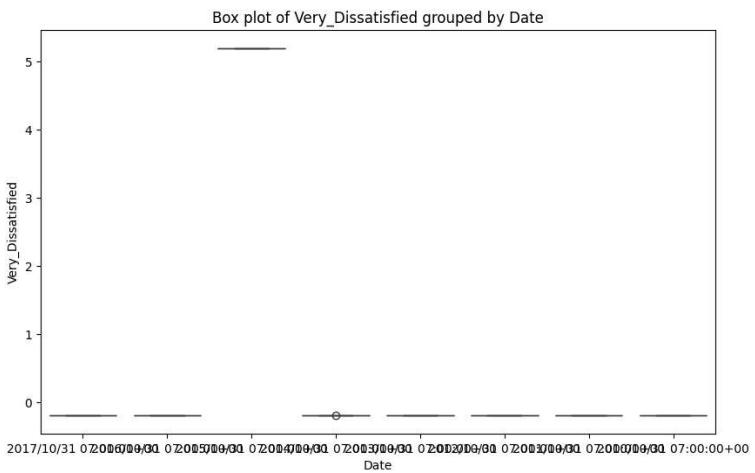
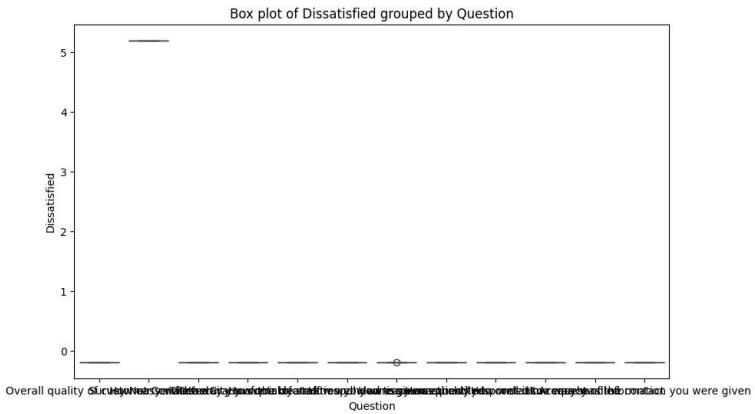
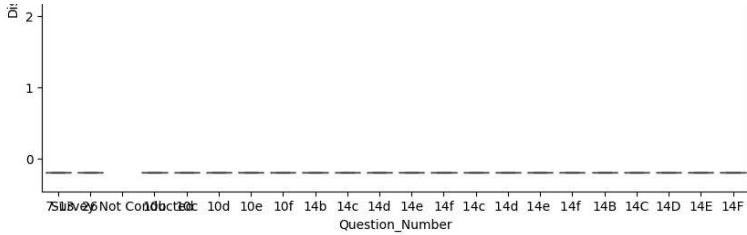


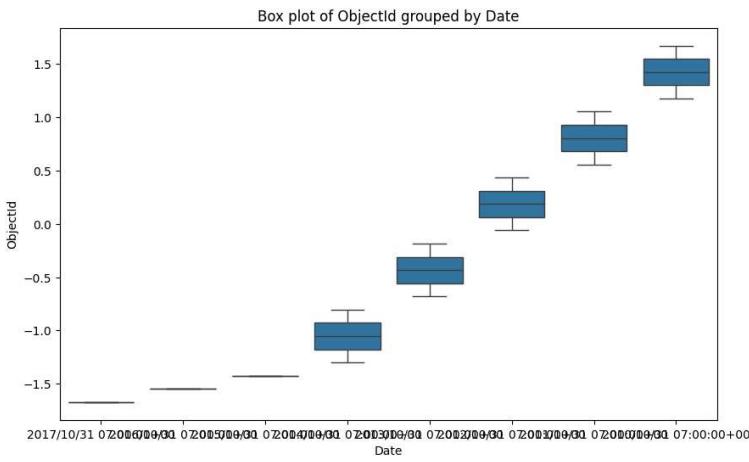
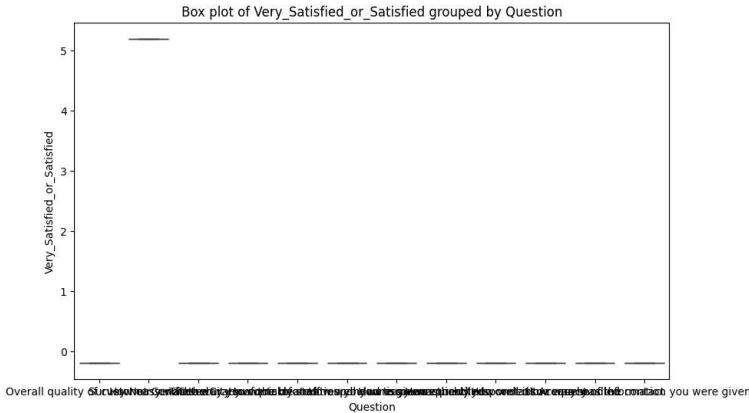
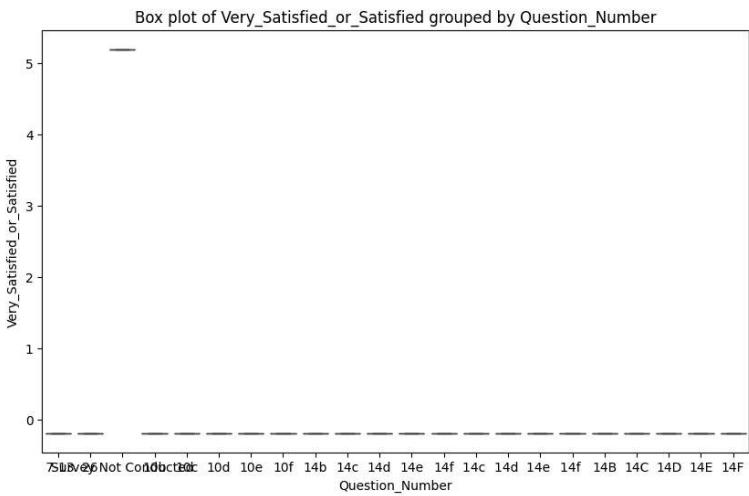
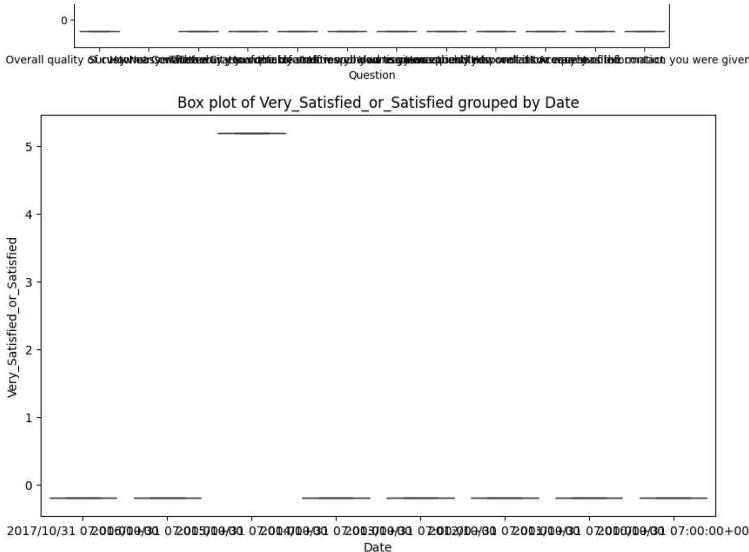


Box plot of Very_Satisfied grouped by Question









Box plot of ObjectId grouped by Question_Number

