

```
import pandas as pd
import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
from math import sqrt
from sklearn.linear_model import LinearRegression
```

Milestone 1

Customer retention is a crucial factor in the success of commercial insurance businesses. Insurers constantly seek to understand the factors contributing to customer churn and develop strategies to retain their clients. The objective of this project is to predict customer retention in the commercial insurance sector using a dataset containing information on business policies. By developing a predictive model, I aim to help our insurance companies proactively identify and address potential cancellations, ultimately improving customer retention and increasing profitability.

The business problem I will address is the high rate of customer churn in the commercial insurance industry. The goal is to develop a machine learning model that can predict the likelihood of policy cancellation based on the features available in the dataset, including policy label, effective date, expiration date, cancellation date, product code, subline code, street address, property city, property county name, property zip code, property state, written premium amount, total loss payments, total loss reserve, and total reserve.

The target variable for our model is the "Cancellation Indicator," a binary variable derived from the cancellation date column. If a policy has a cancellation date, the indicator will be 1, signifying a cancellation. Otherwise, it will be 0, indicating that the policy is still active or expired naturally.

I will start by cleaning and preprocessing the dataset, handling missing values, and removing duplicate records. Categorical variables will be converted into numerical variables using appropriate encoding techniques. I will create new features, such as policy duration, time to cancellation, claim frequency, and loss ratio (total loss payments divided by written premium amount), which can provide additional insights into customer retention.

An extensive exploratory data analysis will be performed to understand the relationships between different variables and cancellations. I will identify patterns and trends in the data, such as geographical clusters of cancellations, seasonal variations, and correlations between specific product lines and cancellations.

```
In [37]: # Loading the data and showing the first 10 rows.
df = pd.read_excel('C:/Users/ndyoung/OneDrive - Bankers Financial Corporation/Desktop/Bellevue/DSC550_T102_2231/df.head(10)
```

| Policy Label | Effective Date | Expiration Date | Cancellation Date | Product Code | Subline Code | Street Address | Property City | Property County Name | Property Zip | Property State | Written Premium Amount |
|--------------|-----------------|-----------------|-------------------|--------------|--------------|--------------------------|------------------|----------------------|--------------|----------------|------------------------|
| 0 | 0004932602-6-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 603 7TH ST S STE 590 | SAINT PETERSBURG | PINELLAS | 33701 | FL | 357 |
| 1 | 0004932602-6-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 101 NORMANDY RD | CASSELBERRY | SEMINOLE | 32707 | FL | 1749 |
| 2 | 0004932602-6-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 10281 SW 72ND ST STE 101 | MIAMI | MIAMI-DADE | 33173 | FL | 699 |
| 3 | 00049327-7-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 1206 MANATEE AVE W | BRADENTON | MANATEE | 34205 | FL | 4098 |
| 4 | 0004932602-6-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 5110 N HABANA AVE STE 2 | TAMPA | HILLSBOROUGH | 33614 | FL | 494 |
| 5 | 0004932602-6-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 1215 JACARANDA BLVD | VENICE | SARASOTA | 34292 | FL | 525 |
| 6 | 0004932602-6-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 1305 S FORT HARRISON AVE | CLEARWATER | PINELLAS | 33756 | FL | 820 |
| 7 | 0004932342-5-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 33 SW 2ND AVE STE 101 | MIAMI | MIAMI-DADE | 33130 | FL | 525 |
| 8 | 0004932602-6-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 525 PINE ISLAND RD STE A | FORT MYERS | LEE | 33903 | FL | 628 |
| 9 | 0004932181-1-02 | 2020-01-01 | 2021-01-01 | NaT | BBOP | 7627 LITTLE RD | NEW PORT RICHEY | PASCO | 34654 | FL | 2153 |

```
In [38]: # Data types and missing values
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 69567 entries, 0 to 69566
Data columns (total 14 columns):
 # Column                    Non-Null Count  Dtype
---  --
 0 Policy Label              69567 non-null object
 1 Effective Date            69567 non-null datetime64[ns]
 2 Expiration Date           69567 non-null datetime64[ns]
 3 Cancellation Date         4456 non-null  datetime64[ns]
 4 Product Code              69567 non-null object
 5 Subline Code              69567 non-null object
 6 Street Address            69567 non-null object
 7 Property City             69567 non-null object
 8 Property County Name      69567 non-null object
 9 Property Zip Code         69567 non-null int64
10 Property State            69567 non-null object
11 Written Premium Amount    69567 non-null float64
12 Total Loss Payments       69567 non-null float64
13 Total Reserve             69567 non-null float64
dtypes: datetime64[ns](3), float64(2), int64(2), object(7)
memory usage: 7.4+ MB

In [39]: # Summary statistics
df.describe()

In [40]: # Calculate the Cancellation Indicator
df['Cancellation_Indicator'] = df['Cancellation Date'].notna().astype(int)

In [41]: # Class distribution
print(df['Cancellation_Indicator'].value_counts(normalize=True))

Class distribution:
0    0.933072
1    0.066928
Name: Cancellation_Indicator, dtype: float64

In [42]: # To plot both graphs (all policies and canceled policies) together, I created a grouped bar chart
# with grouped bars using pandas and matplotlib, displaying the count of all policies and canceled
# policies for each state:

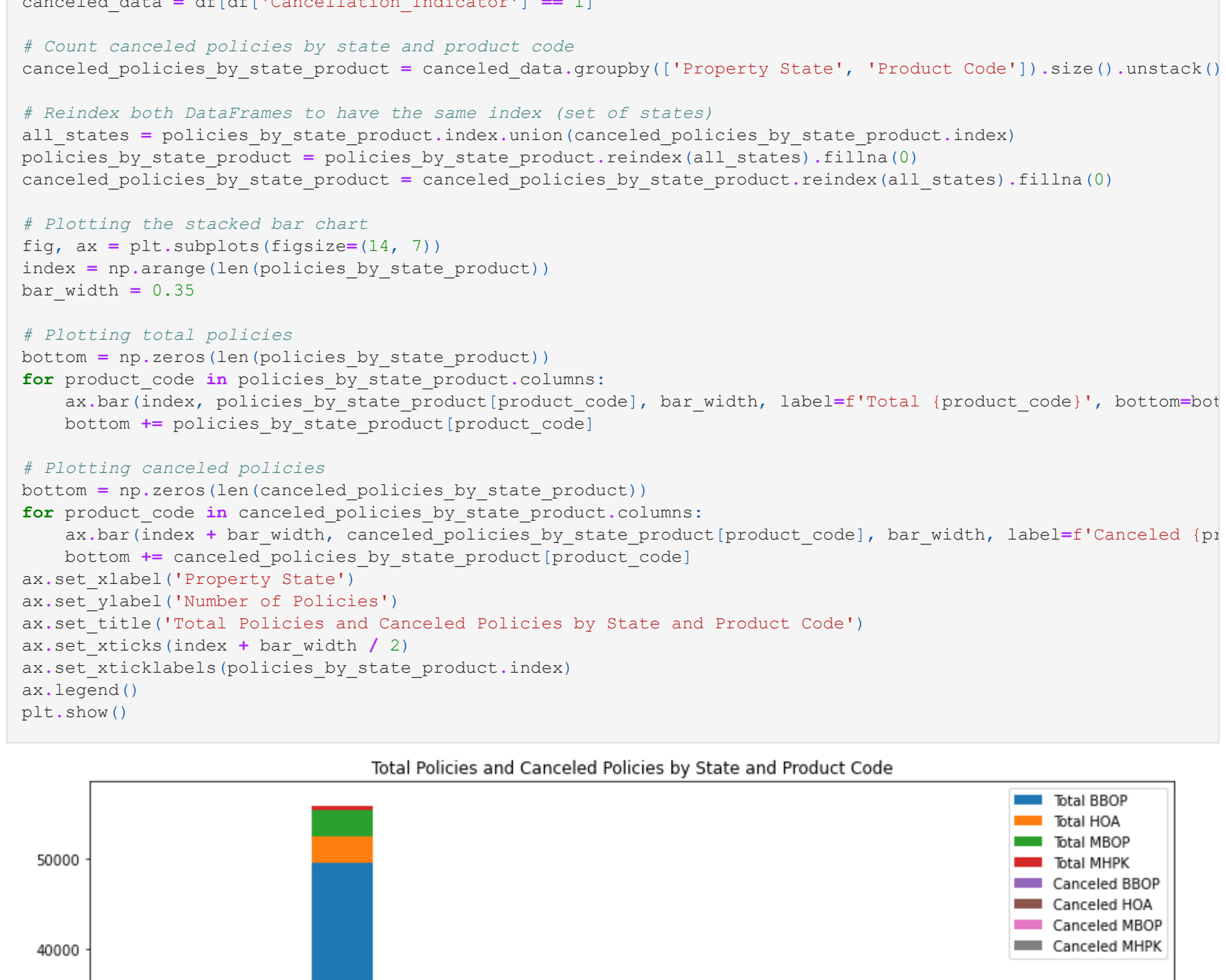
# Counting policies by state
policies_by_state = df['Property State'].value_counts()

# Filter rows with a cancellation
canceled_data = df[df['Cancellation_Indicator'] == 1]

# Count canceled policies by state
canceled_policies_by_state = canceled_data['Property State'].value_counts()

# Combine both Series into a single DataFrame
combined_data = pd.concat([policies_by_state, canceled_policies_by_state], axis=1, keys=['Total Policies', 'Canceled Policies'])
combined_data.fillna(0, inplace=True)

# Plotting the grouped bar chart
fig, ax = plt.subplots(figsize=(14, 7))
index = np.arange(len(combined_data))
rects1 = ax.bar(index, combined_data['Total Policies'], bar_width, label='Total Policies')
rects2 = ax.bar(index + bar_width, combined_data['Canceled Policies'], bar_width, label='Canceled Policies')
ax.set_xlabel('Property State')
ax.set_ylabel('Number of Policies')
ax.set_title('Total Policies vs. Canceled Policies by State')
ax.set_ticks(index + bar_width / 2)
ax.set_xticklabels(combined_data.index)
ax.legend()
plt.show()
```



```
In [43]: # Count policies by state and product code
policies_by_state_product = df.groupby(['Property State', 'Product Code']).size().unstack().fillna(0)

# Filter rows with a cancellation
canceled_data = df[df['Cancellation_Indicator'] == 1]

# Count canceled policies by state and product code
canceled_policies_by_state_product = canceled_data.groupby(['Property State', 'Product Code']).size().unstack().fillna(0)

# Reindex both DataFrames to have the same index (set of states)
all_states = policies_by_state_product.index.union(canceled_policies_by_state_product.index)
policies_by_state_product = policies_by_state_product.reindex(all_states).fillna(0)
canceled_policies_by_state_product = canceled_policies_by_state_product.reindex(all_states).fillna(0)

# Plotting the stacked bar chart
fig, ax = plt.subplots(figsize=(14, 7))
index = np.arange(len(policies_by_state_product))
bar_width = 0.35

# Plotting total policies
bottom = np.zeros(len(policies_by_state_product))
for product_code in policies_by_state_product.columns:
    ax.bar(index, policies_by_state_product[product_code], bar_width, label=f'Total {product_code}', bottom=bottom)
    bottom += policies_by_state_product[product_code]

# Plotting canceled policies
bottom = np.zeros(len(canceled_policies_by_state_product))
for product_code in canceled_policies_by_state_product.columns:
    ax.bar(index + bar_width, canceled_policies_by_state_product[product_code], bar_width, label=f'Canceled {product_code}', bottom=bottom)
    bottom += canceled_policies_by_state_product[product_code]

ax.set_xlabel('Property State')
ax.set_ylabel('Number of Policies')
ax.set_title('Total Policies and Canceled Policies by State and Product Code')
ax.set_ticks(index + bar_width / 2)
ax.set_xticklabels(policies_by_state_product.index)
ax.legend()
plt.show()
```

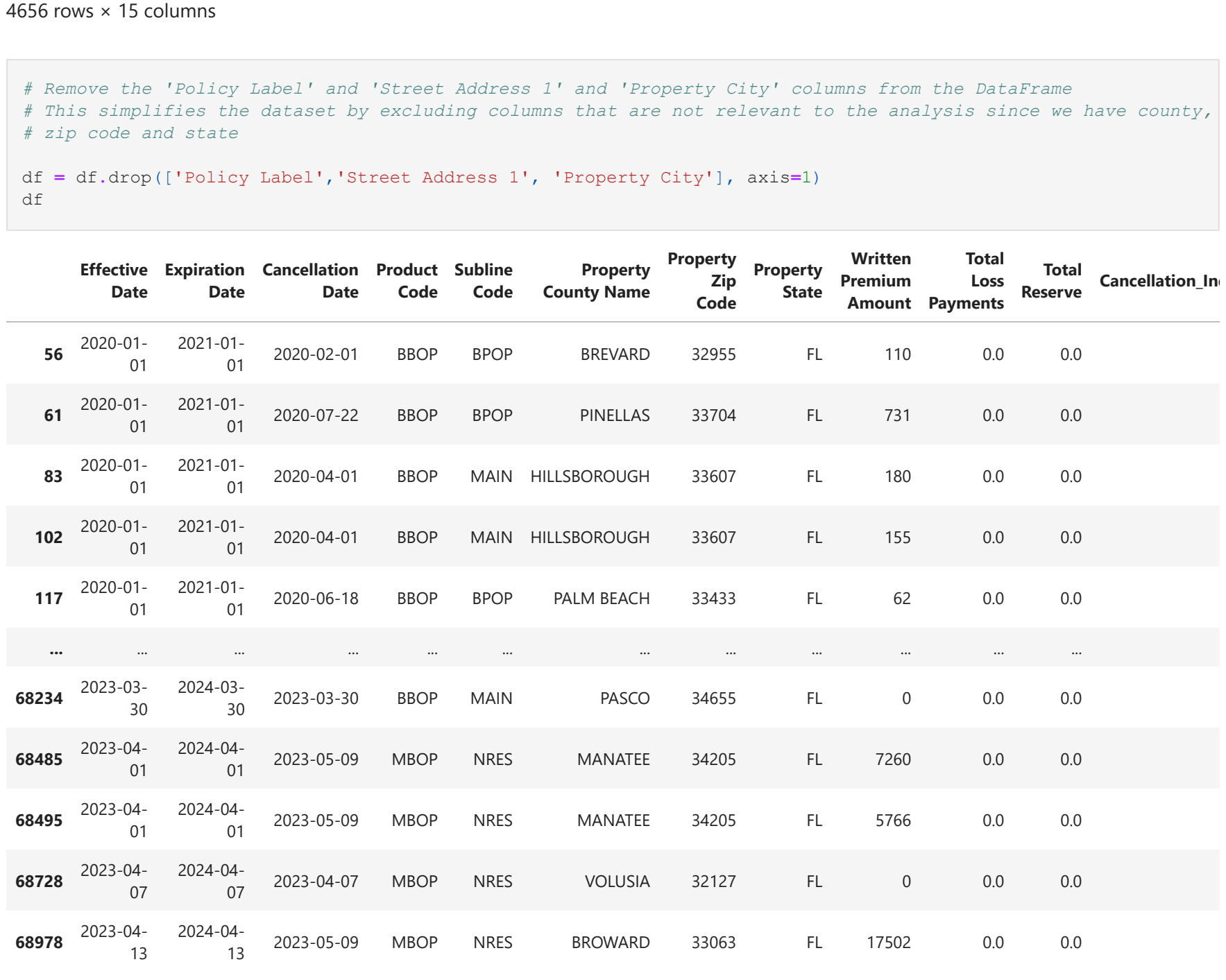


The code filtered the rows with cancellations, counted the policies and canceled policies by state, and created a grouped bar chart displaying the results. We can see that the majority of the policies written and canceled were in the states of Florida and Louisiana.

```
In [44]: # Creating a scatter plot of Written Premium Amount vs. Total Loss Payments
plt.scatter(df['Written Premium Amount'], df['Total Loss Payments'], c=df['Cancellation_Indicator'], cmap='viridis')
plt.xlabel('Written Premium Amount')
plt.ylabel('Total Loss Payments')
plt.title('Written Premium Amount vs. Total Loss Payments')
plt.colorbar(label='Cancellation Indicator')
plt.show()

# Calculating the written premium total and written premium canceled
written_premium_total = df['Written Premium Amount'].sum()
written_premium_canceled = df[df['Cancellation_Indicator'] == 1]['Written Premium Amount'].sum()

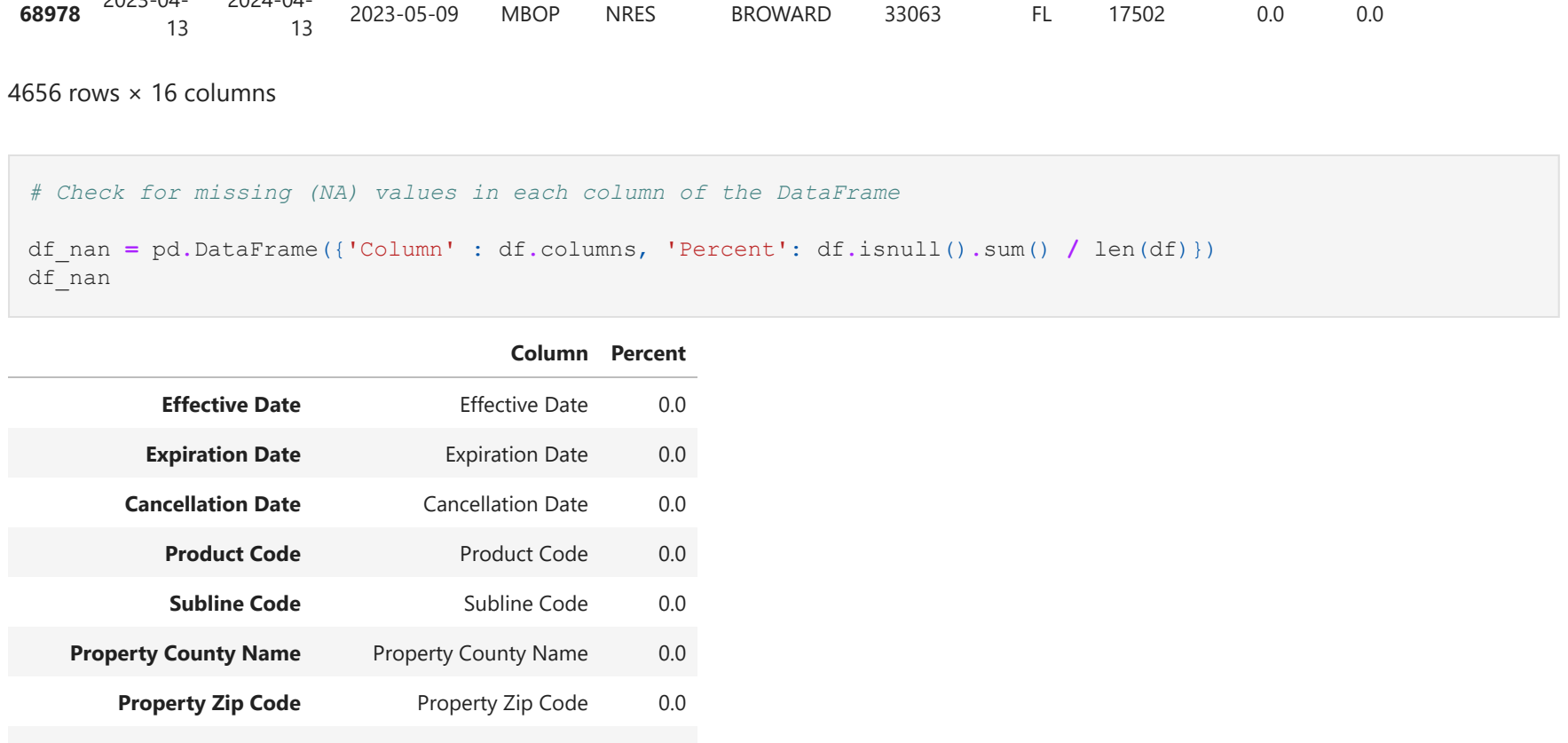
print('Written Premium Total:', written_premium_total)
print('Written Premium Canceled:', written_premium_canceled)
```



The above scatter plot shows the relationship between 'Written Premium Amount' and 'Total Loss Payments'. This scatter plot can help you visualize the correlation between the premium amounts and the loss payments for each policy, potentially revealing patterns or trends related to cancellations.

```
In [45]: # Calculating the cancellation rates by product code
cancellation_rate_by_product = df.groupby('Product Code')['Cancellation_Indicator'].mean()

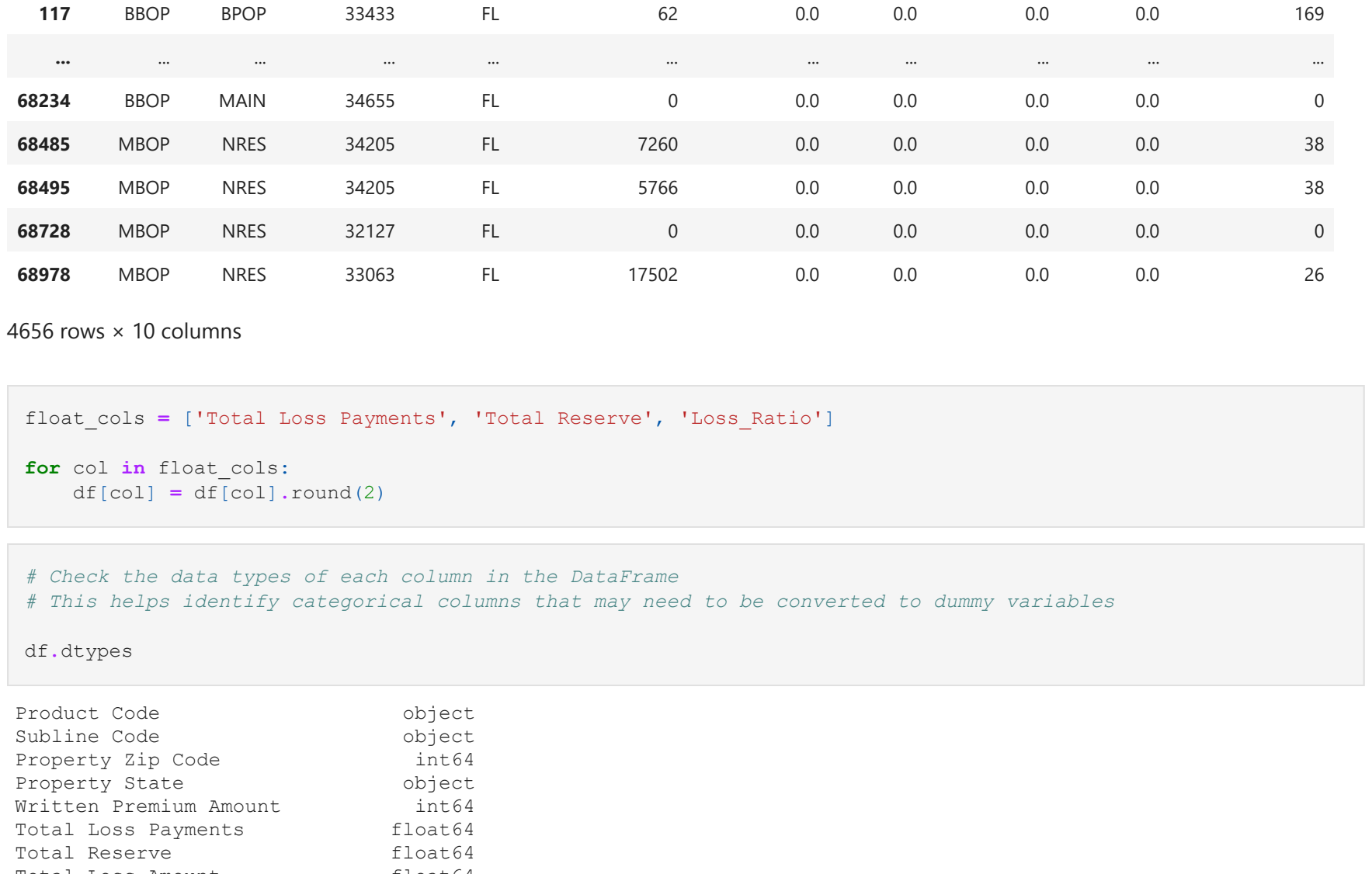
# Plotting the bar chart
plt.figure(figsize=(12, 6))
cancellation_rate_by_product.plot(kind='bar')
plt.xlabel('Product Code')
plt.ylabel('Cancellation Rate')
plt.title('Cancellation Rate by Product Code')
plt.show()
```



This bar chart shows the cancellation rate by product code. This visualization can reveal patterns or trends related to specific product lines, which could be helpful in understanding the factors that influence cancellations and improve customer retention strategies. This visualization can provide insights into the specific product lines that have higher or lower cancellation rates, allowing you to focus on improving the products or customer segments that need more attention. As we can see here, MBOP and BBOP products have the highest potential for investigation.

```
In [46]: # Calculating the cancellation rates by product code and property state
cancellation_rate_by_product_state = df.groupby(['Product Code', 'Property State'])['Cancellation_Indicator'].mean()

# Plotting the data in a heatmap
plt.figure(figsize=(16, 8))
cancellation_rate_by_product_state.heatmap(cmap='coolwarm', annot=True, fmt='.2f', linewidths=.5, cbar_kws={'label': 'Cancellation Rate'})
plt.xlabel('Property State')
plt.ylabel('Product Code')
plt.title('Cancellation Rate by Product Code and Property State')
plt.show()
```



This heatmap is displaying the cancellation rate across different product codes and property states. This visualization can help you identify geographic trends in cancellations as well as any interactions between product code and location. This visualization can help you understand the cancellation patterns across different product lines and locations, allowing you to focus on specific combinations of product codes and property states that exhibit higher cancellation rates. This insight can guide targeted improvements in products or customer segments, ultimately contributing to better customer retention.

From the graphical analysis above we can conclude that the focus of the datasets should be in the states of Florida and Louisiana. The cancellation rates are highest amongst the MBOP and BBOP product lines. Which is understandable since a large majority of the policies sold are in these lines.

Milestone 2

Upon completing the original EDA in Milestone 1, I concluded that I will focus on the policies that were canceled.

```
In [47]: # Filtering the DataFrame df to only include rows where the 'Cancellation_Indicator' is '1' (yes) = canceled
# Prior to expiration
df = df[df['Cancellation_Indicator'] == 1]
df
```

| Policy Label | Effective Date | Expiration Date | Cancellation Date | Product Code | Subline Code | Street Address | Property City | Property County Name | Property Zip | Property State | Written Premium Amount | Total Loss Payments | Total Reserve | Cancellation Indicator |
|--------------|-----------------|-----------------|-------------------|--------------|--------------|-----------------------------|---------------|----------------------|--------------|----------------|------------------------|---------------------|---------------|------------------------|
| 56 | 0004991546-8-07 | 2020-01-01 | 2020-02-01 | BBOP | BPOP | 1362 US HIGHWAY 1 STE 1 | ROCKLEDGE | BREVARD | 32955 | FL | 110 | 0.0 | 0.0 | 1 |
| 61 | 0004991687-6-07 | 2020-01-01 | 2020-07-22 | BBOP | BPOP | 3179 4TH ST N | ST PETERSBURG | PINELLAS | 33704 | FL | 731 | 0.0 | 0.0 | 1 |
| 83 | 0005801748-9-05 | 2020-01-01 | 2020-04-01 | BBOP | MAIN | 2203 N LOIS AVE STE 961 | TAMPA | HILLSBOROUGH | 33607 | FL | 180 | 0.0 | 0.0 | 1 |
| 102 | 0005801748-9-05 | 2020-01-01 | 2020-04-01 | BBOP | MAIN | 2203 N LOIS AVE STE 932 | TAMPA | HILLSBOROUGH | 33607 | FL | 155 | 0.0 | 0.0 | 1 |
| 117 | 0005806008-3-04 | 2020-01-01 | 2020-06-18 | BBOP | BPOP | 7000 W PALMETTO PARK RD STE | BOCA RATON | PALM BEACH | 33433 | FL | 62 | 0.0 | 0.0 | 1 |
| 68234 | 0005816483-7-02 | 2023-03-30 | 2023-03-30 | BBOP | MAIN | 10015 TRINITY BLVD STE 101 | TRINITY | PASCO | 34655 | FL | 0 | 0.0 | 0.0 | 1 |
| 68485 | 0037792044-6-00 | 2023-04-01 | 2023-05-09 | MBOP | NRES | 3021 MANATEE AVE W | BRADENTON | MANATEE | 34205 | FL | 7260 | 0.0 | 0.0 | 1 |
| 68495 | 0037792044-6-00 | 2023-04-01 | 2023-05-09 | MBOP | NRES | 3027 MANATEE AVE W | BRADENTON | MANATEE | 34205 | FL | 5766 | 0.0 | 0.0 | 1 |
| 68728 | 0037792044-1-00 | 2023-04-07 | 2023-04-07 | MBOP | NRES | 1720 DUNLAWTON ORANGE | PORT ORANGE | VOLUSIA | 32127 | FL | 0 | 0.0 | 0.0 | 1 |
| 68978 | 0037792152-2-00 | 2023-04-13 | 2023-05-09 | MBOP | NRES | 603 MELELUKUA DR | MARGATE | BROWARD | 33063 | FL | 17502 | 0.0 | 0.0 | 1 |

4656 rows x 15 columns

```
In [48]: # Remove the 'Policy Label' and 'Street Address' and 'Property City' columns from the DataFrame
# This simplifies the dataset by excluding columns that are not relevant to the analysis since we have county,
# zip code and state
df = df.drop(['Policy Label', 'Street Address', 'Property City'], axis=1)
df
```

26 rows x 26 columns

Cancel Date correlation

correlation['Cancel Date Difference Days']

| | |
|-----------------------------|-----------|
| Property Zip Code | 0.015964 |
| Written Premium Amount | 0.334575 |
| Total Loss Payments | 0.054889 |
| Total Reserve | 0.051515 |
| Total Loss Amount | 0.047554 |
| Loss Ratio | 0.015202 |
| Cancel Date Difference Days | 1.000000 |
| Product Code_BBOP | 0.092356 |
| Product Code_MOA | -0.027407 |
| Product Code_MBOA | -0.093494 |
| Product Code_MHPK | 0.006516 |
| Subline Code_MHPT | -0.002662 |
| Subline Code_BLDG | -0.003782 |
| Subline Code_BPDP | -0.020461 |
| Subline Code_COND | -0.022842 |
| Subline Code_ISO | 0.010624 |
| Subline Code_MAIN | 0.068202 |
| Subline Code_NRES | -0.093494 |
| Subline Code_NRMH | 0.006403 |
| Subline Code_NBSB | -0.030595 |
| Subline Code_RGMH | 0.001431 |
| Property State_FL | -0.022425 |
| Property State_LA | 0.030791 |
| Property State_NC | 0.004579 |
| Property State_SC | 0.009699 |
| Property State_TX | -0.038308 |

Name: Cancel Date Difference Days, dtype: float64

4656 rows x 12 columns

```
In [49]: # Calculate the loss ratio for each policy by dividing 'Total Loss Amount' by
# 'Written Premium Amount'
df['Total Loss Amount'] = df['Total Loss Payments'] + df['Total Reserve']
df['Loss_Ratio'] = df['Total Loss Amount'] / df['Written Premium Amount']
df['Loss_Ratio'].fillna(0, inplace=True)
```

The loss ratio is a useful metric in insurance as it helps measure the profitability of an insurance company's policies.

```
In [50]: # Calculate the difference between 'Effective Date' and 'Cancellation Date'
df['Cancel Date Difference Days'] = (df['Cancellation Date'] - df['Effective Date']).dt.days

# Convert the difference to months and store it in a new column named 'Months Between'
df['Cancel Months Between'] = df['Cancel Date Difference Days'] / 30.44

# Display the updated DataFrame with the new 'Months Between' column
df
```

| Effective Date | Expiration Date | Cancellation Date | Product Code | Subline Code | Property County Name | Property Zip | Property State | Written Premium Amount | Total Loss Payments | Total Reserve | Cancellation Indicator | |
|----------------|-----------------|-------------------|--------------|--------------|----------------------|--------------|----------------|------------------------|---------------------|---------------|------------------------|---|
| 56 | 2020-01-01 | 2021-01-01 | 2020-02-01 | BBOP | BPOP | BREVARD | 32955 | FL | 110 | 0.0 | 0.0 | 1 |
| 61 | 2020-01-01 | 2021-01-01 | 2020-07-22 | BBOP | BPOP | PINELLAS | 33704 | FL | 731 | 0.0 | 0.0 | 1 |
| 83 | 2020-01-01 | 2021-01-01 | 2020-04-01 | BBOP | MAIN | HILLSBOROUGH | 33607 | FL | 180 | 0.0 | 0.0 | 1 |
| 102 | 2020-01-01 | 2021-01-01 | 2020-04-01 | BBOP | MAIN | HILLSBOROUGH | 33607 | FL | 155 | 0.0 | 0.0 | 1 |
| 117 | 2020-01-01 | 2021-01-01 | 2020-06-18 | BBOP | BPOP | PALM BEACH | 33433 | FL | 62 | 0.0 | 0.0 | 1 |
| 68234 | 2023-03-30 | 2024-03-30 | 2023-03-30 | BBOP | MAIN | PASCO | 34655 | FL | 0 | 0.0 | 0.0 | 1 |
| 68485 | 2023-04-01 | 2024-04-01 | 2023-05-09 | MBOP | NRES | MANATEE | 34205 | FL | 7260 | 0.0 | 0.0 | 1 |
| 68495 | 2023-04-01 | 2024-04-01 | 2023-05-09 | MBOP | NRES | MANATEE | 34205 | FL | 5766 | 0.0 | 0.0 | 1 |
| 68728 | 2023-04-07 | 2024-04-07 | 2023-04-07 | MBOP | NRES | VOLUSIA | 32127 | FL | 0 | 0.0 | 0.0 | 1 |
| 68978 | 2023-04-13 | 2024-04-13 | 2023-05-09 | MBOP | NRES | BROWARD | 33063 | FL | 17502 | 0.0 | 0.0 | 1 |

4656 rows x 16 columns

```
In [51]: # Check for missing (NA) values in each column of the DataFrame
df_nan = pd.DataFrame({'Column': df.columns, 'Percent': df.isnull().sum() / len(df)})
df_nan
```

| Column | Percent |
|-----------------------------|---------|
| Effective Date | 0.0 |
| Expiration Date | 0.0 |
| Cancellation Date | 0.0 |
| Product Code | 0.0 |
| Subline Code | 0.0 |
| Property County Name | 0.0 |
| Property Zip Code | 0.0 |
| Property State | 0.0 |
| Written Premium Amount | 0.0 |
| Total Loss Payments | 0.0 |
| Total Reserve | 0.0 |
| Cancellation Indicator | 0.0 |
| Total Loss Amount | 0.0 |
| Loss_Ratio | 0.0 |
| Cancel Date Difference Days | 0.0 |
| Cancel Months Between | 0.0 |

4656 rows x 26 columns

Milestone 3

Creating a correlation matrix to identify relationships.

```
In [58]: correlation = df.corr()
correlation
```

| Property Zip Code | Written Premium Amount | Total Loss Payments | Total Reserve | Total Loss Amount | Loss_Ratio | Cancel Date Difference Days | Product Code_BBOP | Product Code_HBA | Product Code_MBOP | Product Code_MPKP | Subline Code_MAIN | Subline Code_MBOP | Subline Code_MPKP |
|-----------------------------|------------------------|---------------------|---------------|-------------------|------------|-----------------------------|-------------------|------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| 1.000000 | 0.074139 | 0.098505 | 0.001372 | 0.099877 | 0.070750 | 0.015964 | 0.096044 | -0.039446 | -0.084333 | ... | 0.10822 | ... | 0.10822 |
| Written Premium Amount | 1.000000 | 0.067000 | 0.003678 | 0.068219 | 0.006498 | 0.034575 | -0.045161 | 0.001188 | 0.029760 | ... | 0.04645 | ... | 0.04645 |
| Total Loss Payments | 0.098505 | 1.000000 | 0.002947 | 0.034720 | 0.034720 | 0.004589 | 0.007168 | -0.006500 | -0.003799 | ... | 0.02660 | ... | 0.02660 |
| Total Reserve | 0.001372 | 0.003678 | 1.000000 | 0.003515 | 0.004586 | 0.015451 | -0.005116 | -0.002117 | 0.007252 | ... | -0.00217 | ... | -0.00217 |
| Total Loss Amount | 0.099877 | 0.068219 | 0.034720 | 1.000000 | 0.023623 | 0.047554 | 0.005084 | -0.006768 | -0.001245 | ... | 0.01798 | ... | 0.01798 |
| Loss_Ratio | 0.070750 | 0.006498 | 0.034720 | 0.004586 | 1.000000 | 0.015202 | -0.001755 | -0.003852 | 0.005923 | ... | 0.02286 | ... | 0.02286 |
| Cancel Date Difference Days | 0.015964 | 0.034575 | 0.004586 | 0.015451 | 0.047554 | 1.000000 | 0.002356 | -0.027607 | -0.093494 | ... | 0.08520 | ... | 0.08520 |
| Product Code_BBOP | 0.096044 | -0.045161 | 0.007168 | -0.005116 | 0.005084 | 0.001755 | 1.000000 | 0.041384 | -0.063137 | ... | 0.36191 | ... | 0.36191 |
| Product Code_HBA | -0.03 | | | | | | | | | | | | |


```
[79]: correlation = pd.Series({
    'Property Zip Code': 0.015964,
    'Written Premium Amount': 0.334575,
    'Total Loss Payments': 0.045489,
    'Total Reserve': 0.015415,
    'Total Loss Amount': 0.047554,
    'Loss_Ratio': 0.013122,
    'Product Code_BBOP': 0.092356,
    'Product Code_HOA': -0.027607,
    'Product Code_MBOP': -0.093494,
    'Product Code_RRPT': 0.000516,
    'Subline Code_BBOP': -0.020461,
    'Subline Code_LSO': 0.010824,
    'Subline Code_MAIN': 0.063202,
    'Subline Code_NRMH': -0.022842,
    'Subline Code_NRMH': 0.006403,
    'Subline Code_NSGO': -0.030595,
    'Property State_FL': -0.022425,
    'Property State_IL': 0.030791,
    'Property State_NC': 0.004579,
    'Property State_SC': 0.009699,
    'Property State_TX': -0.038309
})

# Creating the bar graph
plt.figure(figsize=(14, 10))
bars = plt.bar(correlation.index, correlation.values)
plt.xlabel('Variables')
plt.ylabel('Correlation')
plt.title('Cancel Date Correlation')

# Adjusting the x-axis labels for proper alignment
plt.xticks(rotation=45, ha='right')

# Function to add labels to the bars
def rect_text(x, y, text):
    height = rect.get_height()
    plt.annotate(text, (x, y), xytext=(x, y + rect.get_height() / 2, height),
                xycoords='offset points', ha='center', va='bottom')

# Adding labels to the bars
autolabel(bars)

plt.tight_layout()
plt.show()
```



'Cancel Date Difference Days' has the highest positive correlation (0.335) with the 'Written Premium Amount'. This indicates that as the 'Written Premium Amount' increases, the length of time until a policy is cancelled also tends to increase. 'Total Loss Amount' and 'Product Code_BBOP' also show weak positive correlations of 0.048 and 0.092, respectively.

Conversely, there are factors like 'Product Code_HOA' and 'Product Code_MBOP' that exhibit weak negative correlations of -0.028 and -0.093, respectively. The correlations for various 'Subline Codes' and 'Property States' range from very weak negative to weak positive, indicating a lack of strong linear relationships between these factors and 'Cancel Date Difference Days'.

```
In [61]: # Drop the specified columns as they are not needed
df = df.drop(columns=['Loss_Ratio', 'Property State_NC', 'Property State_SC', 'Subline Code_RRPT',
                    'Subline Code_NRMH'])
```

```
In [62]: # Cancel Date Difference Days is target variable
X = df.drop('Cancel Date Difference Days', axis=1) # Features
y = df['Cancel Date Difference Days'] # Target variable
```

```
In [63]: # Split the data into a training set and a test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

print('Training data shape:')
print(X_train.shape)
print(y_train.shape)
print('Testing data shape:')
print(X_test.shape)
print(y_test.shape)
```

Training data shape:
(3724, 19)
(3724, 19)
Testing data shape:
(932, 19)
(932, 19)

I will first run a linear regression model to assess the initial performance and gain insights into the relationships between variables before exploring more complex modeling techniques.

```
In [64]: # Initialize the model
model = LinearRegression()

# Train the model
model.fit(X_train, y_train)
```

```
Out[64]: * LinearRegression
LinearRegression()
```

```
In [65]: # Predicting the training dataset
y_predict_train = model.predict(X_train)
```

```
In [66]: # Calculating R2, RMSE and MAE for the training set
r2_train = r2_score(y_train, y_predict_train)
rmse_train = np.sqrt(mean_squared_error(y_train, y_predict_train))
mae_train = mean_absolute_error(y_train, y_predict_train)

print('Training data R2, RMSE, MAE:')
print(f'R2: {r2_train}')
print(f'RMSE: {rmse_train}')
print(f'MAE: {mae_train}')
```

Training data R2, RMSE, MAE:
R2: 0.1281354607918873
RMSE: 100.80535402420925
MAE: 84.3050881731694

```
In [67]: # Predicting the test dataset
y_predict_test = model.predict(X_test)
```

```
# Calculating R2, RMSE and MAE for the test set
r2_test = r2_score(y_test, y_predict_test)
rmse_test = np.sqrt(mean_squared_error(y_test, y_predict_test))
mae_test = mean_absolute_error(y_test, y_predict_test)

print('Test data R2, RMSE, MAE:')
print(f'R2: {r2_test}')
print(f'RMSE: {rmse_test}')
print(f'MAE: {mae_test}')
```

Test data R2, RMSE, MAE:
R2: 0.0486079146431145
RMSE: 104.2255587997634
MAE: 85.12426892072823

Based on these metrics, it appears that the linear regression model may not be capturing the underlying patterns in the data very well. The R2 values are relatively low, indicating that the model explains only a small portion of the variance in the target variable. Additionally, the RMSE and MAE values are relatively high, suggesting that the model's predictions have a significant average error.

Experimenting with other models to get best fit.

```
In [68]: from sklearn.tree import DecisionTreeRegressor

# Initialize the model
model = DecisionTreeRegressor(random_state=42)

# Train the model
model.fit(X_train, y_train)

# Predicting the training dataset
y_predict_train = model.predict(X_train)

# Calculate R2, RMSE, and MAE for the training set
r2_train = r2_score(y_train, y_predict_train)
rmse_train = np.sqrt(mean_squared_error(y_train, y_predict_train))
mae_train = mean_absolute_error(y_train, y_predict_train)

# Print training set performance metrics
print('Training data R2, RMSE, MAE:')
print(f'R2: {r2_train}')
print(f'RMSE: {rmse_train}')
print(f'MAE: {mae_train}')
```

```
# Predicting the test dataset
y_predict_test = model.predict(X_test)

# Calculate R2, RMSE, and MAE for the test set
r2_test = r2_score(y_test, y_predict_test)
rmse_test = np.sqrt(mean_squared_error(y_test, y_predict_test))
mae_test = mean_absolute_error(y_test, y_predict_test)

# Print test set performance metrics
print('Test data R2, RMSE, MAE:')
print(f'R2: {r2_test}')
print(f'RMSE: {rmse_test}')
print(f'MAE: {mae_test}')
```

Training data R2, RMSE, MAE:
R2: 0.997617965117985
RMSE: 1.66628012450182
MAE: 0.660303739396942
Test data R2, RMSE, MAE:
R2: -0.0502739366791396
RMSE: 109.60649297153116
MAE: 78.2043991416309

The high R2 value (close to 1) and low RMSE and MAE values for the training data indicate that the decision tree model fits the training data very well. However, the negative R2 value and relatively high RMSE and MAE values for the test data suggest that the model does not generalize well to unseen data. This indicates potential overfitting, where the model has learned the training data too well and is not performing well on new data.

Adjusting hyperparameters in the decision tree model.

```
In [69]: from sklearn.model_selection import GridSearchCV

# Define the hyperparameter grid
param_grid = {
    'max_depth': [None, 5, 10, 15],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4, 8]
}

# Initialize the model
model = DecisionTreeRegressor(random_state=42)

# Perform grid search
grid_search = GridSearchCV(estimator=model, param_grid=param_grid, scoring='neg_mean_squared_error', cv=5)
grid_search.fit(X_train, y_train)

# Get the best hyperparameters
best_params = grid_search.best_params_
print('Best Hyperparameters:', best_params)

# Initialize a new model with the best hyperparameters
best_model = DecisionTreeRegressor(random_state=42, **best_params)

# Train the best model
best_model.fit(X_train, y_train)

# Predict on the test data using the best model
y_pred_test = best_model.predict(X_test)

# Evaluate the best model's performance
r2_test = r2_score(y_test, y_pred_test)
rmse_test = np.sqrt(mean_squared_error(y_test, y_pred_test))
mae_test = mean_absolute_error(y_test, y_pred_test)

print('Best Model Performance:')
print(f'R2: {r2_test}')
print(f'RMSE: {rmse_test}')
print(f'MAE: {mae_test}')
```

Best Hyperparameters: {'max_depth': 5, 'min_samples_leaf': 4, 'min_samples_split': 2}
Best Model Performance (Training data):
R2: 0.46985001203827836
RMSE: 81.33963394052921
MAE: 59.8852081677907216

These metrics indicate that the model with tuned hyperparameters performs better than the previous model. The R2 score of 0.47 suggests that the model explains 47% of the variance in the target variable. The lower values of RMSE and MAE indicate smaller prediction errors.

```
In [70]: from sklearn.ensemble import RandomForestRegressor

# Initialize the model
model = RandomForestRegressor(random_state=42)

# Train the model
model.fit(X_train, y_train)

# Predicting the training dataset
y_predict_train = model.predict(X_train)

# Calculate R2, RMSE, and MAE for the training set
r2_train = r2_score(y_train, y_predict_train)
rmse_train = np.sqrt(mean_squared_error(y_train, y_predict_train))
mae_train = mean_absolute_error(y_train, y_predict_train)

print('Training data R2, RMSE, MAE:')
print(f'R2: {r2_train}')
print(f'RMSE: {rmse_train}')
print(f'MAE: {mae_train}')
```

```
# Predicting the test dataset
y_predict_test = model.predict(X_test)

# Calculate R2, RMSE, and MAE for the test set
r2_test = r2_score(y_test, y_predict_test)
rmse_test = np.sqrt(mean_squared_error(y_test, y_predict_test))
mae_test = mean_absolute_error(y_test, y_predict_test)

print('Test data R2, RMSE, MAE:')
print(f'R2: {r2_test}')
print(f'RMSE: {rmse_test}')
print(f'MAE: {mae_test}')
```

Training data R2, RMSE, MAE:
R2: 0.9192363650397426
RMSE: 30.680884965184404
MAE: 22.669396053910287
Test data R2, RMSE, MAE:
R2: 0.41774800894869635
RMSE: 81.33963394052921
MAE: 59.8852081677907216

Overall, the Random Forest model shows good performance on the training data, but there is some drop in performance on the test data. This may indicate a degree of overfitting, where the model is capturing the training data patterns too closely and not generalizing well to new data. Further optimization or tuning of the model parameters could potentially improve its performance on the test data.

Adjusting hyperparameters in the random forest model.

```
In [71]: # Define the parameter grid for hyperparameter tuning
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [3, 5, 7],
    'min_samples_split': [2, 4, 6],
    'min_samples_leaf': [1, 2, 4]
}

# Initialize the random forest regressor
model = RandomForestRegressor(random_state=42)

# Create the GridSearchCV object
grid_search = GridSearchCV(model, param_grid, cv=5)

# Fit the GridSearchCV object to the training data
grid_search.fit(X_train, y_train)

# Get the best hyperparameters and best model
best_params = grid_search.best_params_
best_model = grid_search.best_estimator_

# Print the best hyperparameters
print("Best Hyperparameters:", best_params)

# Evaluate the best model on the training data
y_train_pred = best_model.predict(X_train)
r2_train = r2_score(y_train, y_train_pred)
rmse_train = np.sqrt(mean_squared_error(y_train, y_train_pred))
mae_train = mean_absolute_error(y_train, y_train_pred)
print("Best Model Performance (Training data):")
print(f'R2: {r2_train}')
print(f'RMSE: {rmse_train}')
print(f'MAE: {mae_train}')
```

```
# Evaluate the best model on the test data
y_test_pred = best_model.predict(X_test)
r2_test = r2_score(y_test, y_test_pred)
rmse_test = np.sqrt(mean_squared_error(y_test, y_test_pred))
mae_test = mean_absolute_error(y_test, y_test_pred)
print("Best Model Performance (Test data):")
print(f'R2: {r2_test}')
print(f'RMSE: {rmse_test}')
print(f'MAE: {mae_test}')
```

Best Hyperparameters: {'max_depth': 5, 'min_samples_leaf': 4, 'min_samples_split': 2, 'n_estimators': 300}
Best Model Performance (Training data):
R2: 0.918243711817887
RMSE: 77.19207524240134
MAE: 57.37015112400198

The results indicate improved performance compared to the previous iterations.

The best Random Forest model with tuned hyperparameters achieved an R2 score of 0.48 on the test data, indicating that it explains 48% of the variance in the target variable. The RMSE value of 77.19 and MAE value of 57.37 suggest relatively small prediction errors. This model outperformed the previous iterations and shows promise in predicting the "Cancel Date Difference Days" based on the given features.

Incorporating new features can help improve the prediction of cancellation patterns. I might consider adding more customer related demographics, customer behaviors (payment history), and cancellation reasons to correlate and perform textual analysis.