

Classification on Imbalanced Data with Python: A strategy for Enhanced Predictive Performance

Introduction: In the realm of machine learning, encountering imbalanced data sets presents a significant challenge, particularly in classification tasks. Imbalanced data occurs when the distribution of classes within a dataset is heavily skewed, with one class dominating over the others. This scenario can compromise the effectiveness of predictive models, as they tend to favor the majority class, resulting in subpar performance, particularly on the minority class, which often holds greater significance. In this article, we delve into the intricacies of addressing imbalanced data in classification tasks using Python, exploring essential techniques and strategies to enhance model efficacy.

Understanding Imbalanced Data: Imbalanced data poses a formidable hurdle in classification endeavors, demanding meticulous attention to various facets including data preprocessing, resampling methodologies, model selection, and evaluation metrics. To embark on this journey effectively, it's imperative to adhere to a structured approach:

1. Analyzing Class Distribution: Begin by scrutinizing the distribution of classes within the dataset to gauge the extent of imbalance. This preliminary step provides crucial insights into the disparity between different classes, laying the groundwork for subsequent actions.

2. Assessing Class Importance: Determine the significance of each class within the context of the problem at hand. Understanding the relative importance of different classes aids in devising targeted strategies to mitigate imbalance and prioritize model performance on critical classes.

3. Mitigating Imbalance: Address the class imbalance by augmenting the number of instances in the minority class. Replication techniques can be employed to rebalance the class distribution, thereby fostering a more equitable representation of all classes within the dataset.

4. Leveraging Robust Algorithms: Certain machine learning algorithms, such as tree-based methods, exhibit resilience to class imbalance. Consider leveraging these algorithms or ensemble techniques like Random Forest or Gradient Boosted Trees, which inherently mitigate the adverse effects of imbalance, ensuring more robust model performance.

5. Adopting Informative Metrics: Beyond conventional accuracy metrics, prioritize evaluation measures that are tailored for imbalanced datasets. Metrics such as Precision, Recall, F1 Score, and the Area Under the Receiver Operating Characteristic (AUROC) curve offer nuanced insights into model performance across different classes, enabling a more comprehensive assessment of predictive efficacy.

Navigating the intricacies of classification on imbalanced data necessitates a multifaceted approach, encompassing diligent data preprocessing, strategic model selection, and meticulous evaluation. By adhering to the prescribed methodologies and leveraging the rich ecosystem of Python libraries and tools, practitioners can unlock the full potential of their classification models, achieving superior predictive performance even in the face of challenging imbalanced datasets.

```
In [2]: import pandas as pd
# Load the dataset
data = pd.read_csv("C:/Users/anike/OneDrive/Desktop/Projects/Imbalance Data/Insurance
print(data.head())
```



	policy_id	subscription_length	vehicle_age	customer_age	region_code	\		
0	POL045360	9.3	1.2	41	C8			
1	POL016745	8.2	1.8	35	C2			
2	POL007194	9.5	0.2	44	C8			
3	POL018146	5.2	0.4	44	C10			
4	POL049011	10.1	1.0	56	C13			
	region_density	segment	model	fuel_type	max_torque	...	is_brake_assist	\
0	8794	C2	M4	Diesel	250Nm@2750rpm	...	Yes	
1	27003	C1	M9	Diesel	200Nm@1750rpm	...	No	
2	8794	C2	M4	Diesel	250Nm@2750rpm	...	Yes	
3	73430	A	M1	CNG	60Nm@3500rpm	...	No	
4	5410	B2	M5	Diesel	200Nm@3000rpm	...	No	
	is_power_door_locks	is_central_locking	is_power_steering	\				
0	Yes	Yes	Yes					
1	Yes	Yes	Yes					
2	Yes	Yes	Yes					
3	No	No	Yes					
4	Yes	Yes	Yes					
	is_driver_seat_height_adjustable	is_day_night_rear_view_mirror	is_ecw	\				
0	Yes	No	Yes					
1	Yes	Yes	Yes					
2	Yes	No	Yes					
3	No	No	No					
4	No	No	Yes					
	is_speed_alert	ncap_rating	claim_status					
0	Yes	3	0					
1	Yes	4	0					
2	Yes	3	0					
3	Yes	0	0					
4	Yes	5	0					

[5 rows x 41 columns]

In [3]: data.info()



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 58592 entries, 0 to 58591
Data columns (total 41 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   policy_id                                58592 non-null  object
1   subscription_length                      58592 non-null  float64
2   vehicle_age                             58592 non-null  float64
3   customer_age                             58592 non-null  int64
4   region_code                             58592 non-null  object
5   region_density                           58592 non-null  int64
6   segment                                  58592 non-null  object
7   model                                    58592 non-null  object
8   fuel_type                               58592 non-null  object
9   max_torque                              58592 non-null  object
10  max_power                                58592 non-null  object
11  engine_type                              58592 non-null  object
12  airbags                                  58592 non-null  int64
13  is_esc                                   58592 non-null  object
14  is_adjustable_steering                   58592 non-null  object
15  is_tpms                                  58592 non-null  object
16  is_parking_sensors                       58592 non-null  object
17  is_parking_camera                       58592 non-null  object
18  rear_brakes_type                         58592 non-null  object
19  displacement                             58592 non-null  int64
20  cylinder                                 58592 non-null  int64
21  transmission_type                       58592 non-null  object
22  steering_type                            58592 non-null  object
23  turning_radius                           58592 non-null  float64
24  length                                   58592 non-null  int64
25  width                                    58592 non-null  int64
26  gross_weight                             58592 non-null  int64
27  is_front_fog_lights                     58592 non-null  object
28  is_rear_window_wiper                    58592 non-null  object
29  is_rear_window_washer                   58592 non-null  object
30  is_rear_window_defogger                 58592 non-null  object
31  is_brake_assist                          58592 non-null  object
32  is_power_door_locks                     58592 non-null  object
33  is_central_locking                      58592 non-null  object
34  is_power_steering                       58592 non-null  object
35  is_driver_seat_height_adjustable         58592 non-null  object
36  is_day_night_rear_view_mirror            58592 non-null  object
37  is_ecw                                   58592 non-null  object
38  is_speed_alert                           58592 non-null  object
39  ncap_rating                             58592 non-null  int64
40  claim_status                            58592 non-null  int64
dtypes: float64(3), int64(10), object(28)
memory usage: 18.3+ MB
```

```
In [4]: data.isnull().sum()
```



```
Out[4]: policy_id      0
subscription_length  0
vehicle_age         0
customer_age        0
region_code         0
region_density      0
segment            0
model              0
fuel_type          0
max_torque          0
max_power           0
engine_type        0
airbags            0
is_esc             0
is_adjustable_steering 0
is_tpms            0
is_parking_sensors 0
is_parking_camera  0
rear_brakes_type   0
displacement       0
cylinder           0
transmission_type  0
steering_type      0
turning_radius     0
length             0
width              0
gross_weight       0
is_front_fog_lights 0
is_rear_window_wiper 0
is_rear_window_washer 0
is_rear_window_defogger 0
is_brake_assist    0
is_power_door_locks 0
is_central_locking 0
is_power_steering  0
is_driver_seat_height_adjustable 0
is_day_night_rear_view_mirror 0
is_ecw            0
is_speed_alert     0
ncap_rating        0
claim_status       0
dtype: int64
```

Exploratory Data Analysis of Insurance Claim Frequency Prediction Dataset

Before delving into the intricacies of classification on imbalanced data, it's imperative to gain a comprehensive understanding of the underlying dataset. In this section, we embark on an exploratory data analysis (EDA) journey, unraveling key insights and patterns within the insurance claim frequency prediction dataset.

Dataset Overview: The dataset comprises 58,592 entries and 41 columns, encompassing a diverse array of features pertinent to insurance claim prediction. Notable attributes include policy-related metrics such as `subscription_length` and `customer_age`, categorical descriptors like `region_code` and `vehicle segment`, as well as specifications regarding vehicle safety and convenience. Crucially, the target variable, `claim_status`, serves as the focal point for predictive modeling, indicating whether a claim was made (1) or not (0).

Key Features: Before delving deeper into the dataset's nuances, let's highlight some key features:

- **policy_id:** Unique identifier for insurance policies
- **subscription_length, vehicle_age, customer_age:** Numeric attributes reflecting policy, vehicle, and customer characteristics
- **region_code, segment, model, fuel_type:** Categorical attributes delineating geographical regions, vehicle segments, models, and fuel types
- **max_torque, max_power, engine_type:** Specifications pertaining to vehicle engines
- **airbags, is_esc, is_adjustable_steering:** Features indicative of vehicle safety and convenience



- **claim_status:** Target variable denoting claim occurrence (1) or absence (0)

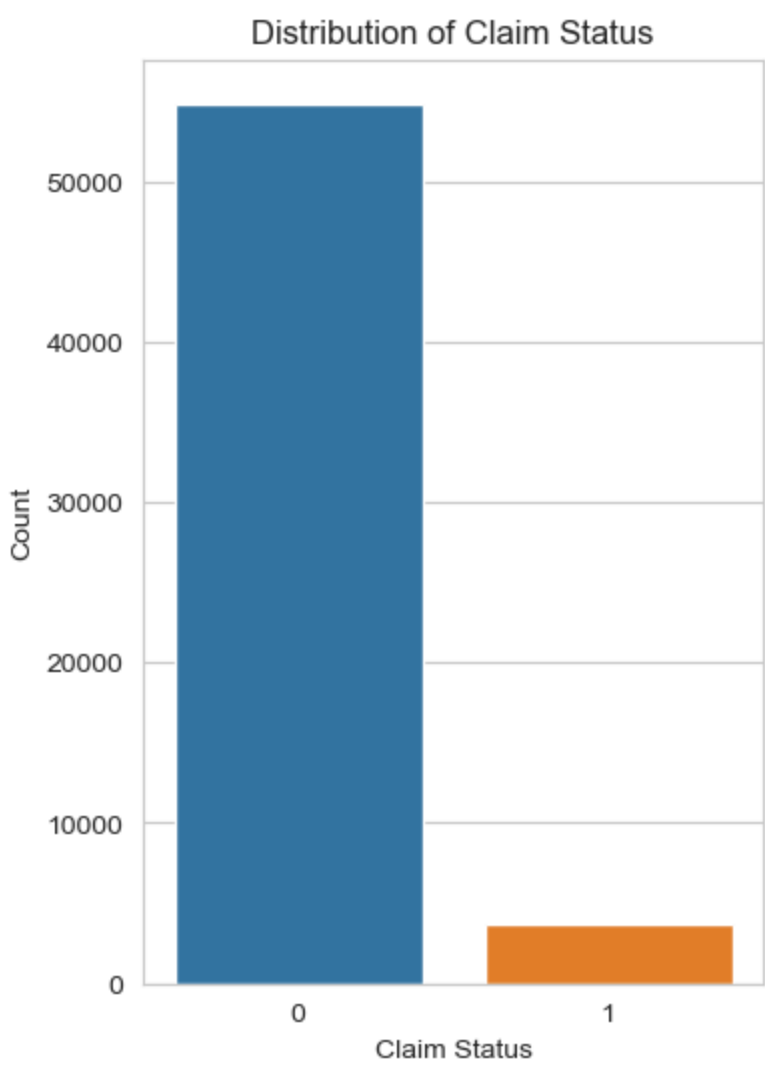
Visualizing Class Balance: A fundamental aspect of grappling with imbalanced data is understanding the distribution of the target variable. Let's initiate our exploration by visualizing the distribution of claim_status:

The preliminary exploration sets the stage for our subsequent endeavors in classification on imbalanced data. By unraveling the nuances of feature distributions and class imbalances, we pave the way for informed preprocessing, modeling, and evaluation strategies, aimed at fostering robust predictive performance in the realm of insurance claim frequency prediction.

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

sns.set_style("whitegrid")

# Plot the distribution of the target variable 'claim_status'
plt.figure(figsize=(4, 6)) # Adjusted figure size
sns.countplot(x='claim_status', data=data)
plt.title('Distribution of Claim Status')
plt.xlabel('Claim Status')
plt.ylabel('Count')
plt.show()
```



Unveiling Feature Insights: Analysis of Numerical Features in Relation to Claim Status

The exploration of numerical features within the insurance claim frequency prediction dataset serves as a pivotal step in elucidating their distributions and discerning potential relationships with the target variable, claim_status. By scrutinizing key numerical attributes such as subscription_length, vehicle_age, and customer_age, we endeavor to unravel pertinent insights that inform subsequent modeling endeavors.

Distribution Analysis: Let's delve into the distributions of pivotal numerical features to glean insights into their characteristics and potential implications on claim occurrence:



1. Subscription Length:

- **Meaning:** Duration of insurance policy subscription.
- **Distribution:** Visualize the spread of subscription lengths across the dataset.

2. Vehicle Age:

- **Meaning:** Age of the insured vehicle.
- **Distribution:** Explore the distribution of vehicle ages and discern any discernible patterns.

3. Customer Age:

- **Meaning:** Age of the policyholder/customer.
- **Distribution:** Investigate the age distribution of customers and its potential impact on claim likelihood.

Visualization: Utilizing appropriate visualizations such as histograms or box plots, we aim to illuminate the nuances of these numerical features while concurrently examining their relationship with claim_status.

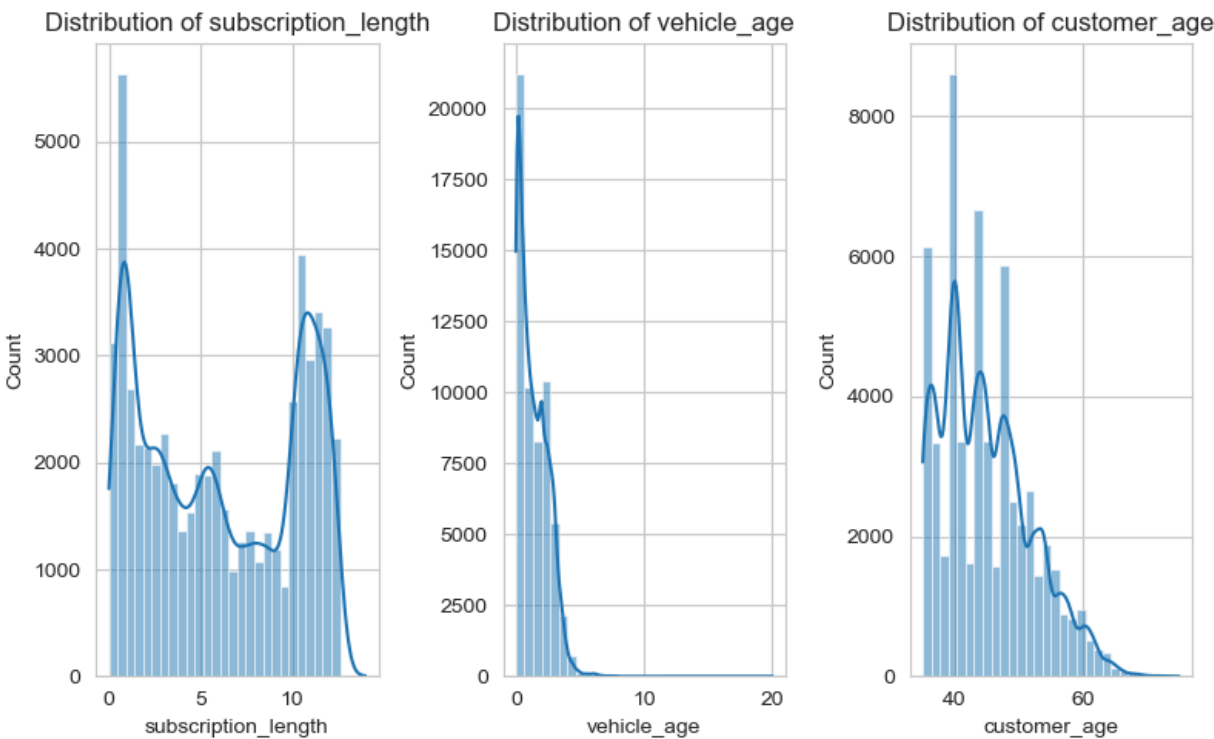
The analysis of numerical features lays the groundwork for a nuanced understanding of their distributions and potential implications on claim occurrence. By discerning patterns and relationships with the target variable, claim_status, we equip ourselves with invaluable insights to guide subsequent modeling endeavors, fostering the development of robust predictive models tailored to the intricacies of insurance claim frequency prediction.

In [22]:

```
# selecting numerical columns for analysis
numerical_columns = ['subscription_length', 'vehicle_age', 'customer_age']

# plotting distributions of numerical features
plt.figure(figsize=(8, 5))
for i, column in enumerate(numerical_columns, 1):
    plt.subplot(1, 3, i)
    sns.histplot(data[column], bins=30, kde=True)
    plt.title(f'Distribution of {column}')

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



Categorical Features in Relation to Claim Status



Transitioning from the exploration of numerical attributes, our analysis now pivots towards unraveling insights embedded within categorical features within the insurance claim frequency prediction dataset. By scrutinizing pivotal categorical attributes such as region_code, segment,

and fuel_type, we aim to elucidate their variation and discern potential relationships with the target variable, claim_status.

Characteristics Analysis: Let's delve into the characteristics of key categorical features to glean insights into their variation and potential implications on claim occurrence:

1. Region Code:

- **Meaning:** Geographic region associated with insurance policies.
- **Variation Analysis:** Examine the distribution of region codes and discern any geographic patterns.

2. Segment:

- **Meaning:** Categorization of vehicles into specific segments.
- **Variation Analysis:** Explore the diversity within vehicle segments and ascertain their impact on claim likelihood.

3. Fuel Type:

- **Meaning:** Type of fuel utilized by insured vehicles.
- **Variation Analysis:** Investigate the distribution of fuel types and discern potential relationships with claim_status.

Visualization: Utilizing visualizations such as bar plots or pie charts, we aim to elucidate the variation within these categorical features while concurrently assessing their relationship with claim_status.

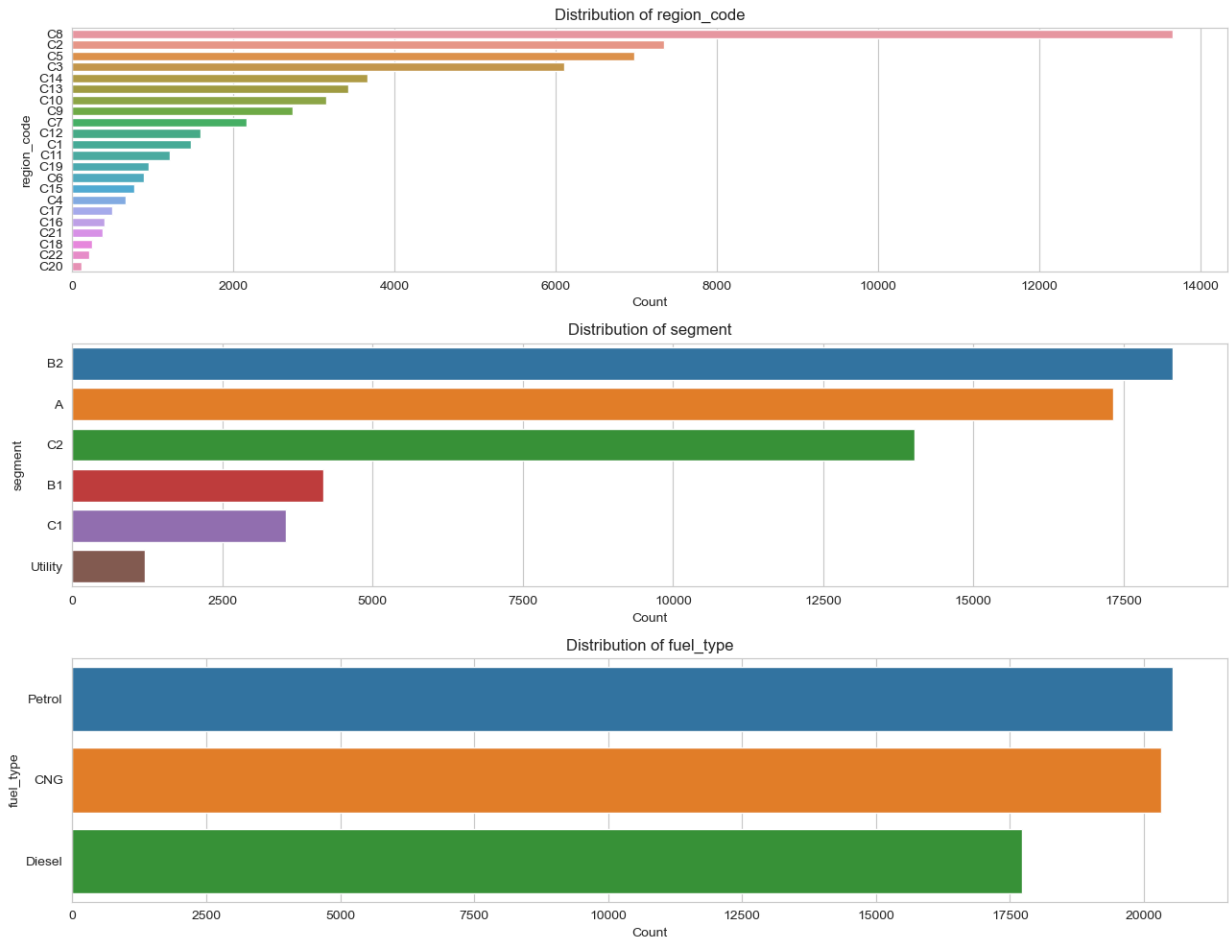
The analysis of categorical features unveils valuable insights into their variation and potential impact on claim occurrence. By discerning patterns and relationships with the target variable, claim_status, we equip ourselves with crucial knowledge to inform subsequent modeling endeavors, facilitating the development of predictive models adept at navigating the complexities of insurance claim frequency prediction.

```
In [24]: # selecting some relevant categorical columns for analysis
categorical_columns = ['region_code', 'segment', 'fuel_type']

# plotting distributions of categorical features
plt.figure(figsize=(13, 10))
for i, column in enumerate(categorical_columns, 1):
    plt.subplot(3, 1, i)
    sns.countplot(y=column, data=data, order = data[column].value_counts().index)
    plt.title(f'Distribution of {column}')
    plt.xlabel('Count')
    plt.ylabel(column)

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





Class Imbalance: Balancing the Dataset Using Oversampling

As we delve deeper into the analysis of the insurance claim frequency prediction dataset, one glaring challenge that surfaces is the pronounced class imbalance within the target variable, claim_status. To mitigate the adverse effects of this imbalance on predictive modeling, we resort to oversampling techniques to ensure equitable representation of both claim and non-claim instances, thereby fostering robust model training and performance.

Class Imbalance Analysis: Before delving into the oversampling process, it's imperative to acknowledge the extent of class imbalance within the dataset:

- **Claim Status Distribution:** The distribution of claim_status reveals a significant disproportion between claim (1) and non-claim (0) instances, with the latter vastly outnumbering the former.

Oversampling Strategy: In light of the observed class imbalance, we opt for oversampling to rectify the disparity between claim and non-claim instances. By replicating minority class instances, we strive to create a more balanced dataset that facilitates equitable representation of both classes during model training.

Implementation: Utilizing Python libraries such as `imbalanced-learn`, we employ oversampling techniques such as Random Oversampling to augment the minority class instances, thereby achieving a balanced dataset conducive to robust predictive modeling.

Addressing class imbalance through oversampling is a crucial step in fortifying the predictive efficacy of classification models, particularly in scenarios characterized by pronounced class disproportion. By ensuring equitable representation of both classes, we pave the way for more robust and accurate predictions, thereby enhancing the utility of our models in the realm of insurance claim frequency prediction.



```
In [25]: from sklearn.utils import resample

# separate majority and minority classes
majority = data[data.claim_status == 0]
minority = data[data.claim_status == 1]
```



```
# oversample the minority class
minority_oversampled = resample(minority,
                                replace=True,
                                n_samples=len(majority),
                                random_state=42)

# combine majority class with oversampled minority class
oversampled_data = pd.concat([majority, minority_oversampled])

# check the distribution of undersampled and oversampled datasets
oversampled_distribution = oversampled_data.claim_status.value_counts()

oversampled_distribution
```

Out[25]: 0 54844
1 54844
Name: claim_status, dtype: int64

After conducting oversampling on the minority class, both categories now contain an equal number of entries, totaling 54,844 each. Let's examine some essential variables to gain insights into the balanced dataset's composition. An alternative approach to visualize the distribution of 'customer_age', 'vehicle_age', and 'subscription_length' with respect to 'claim_status' could be using stacked histograms. Stacked histograms provide a clear comparison of the distributions for each class within the same plot. Here's how you can do it:

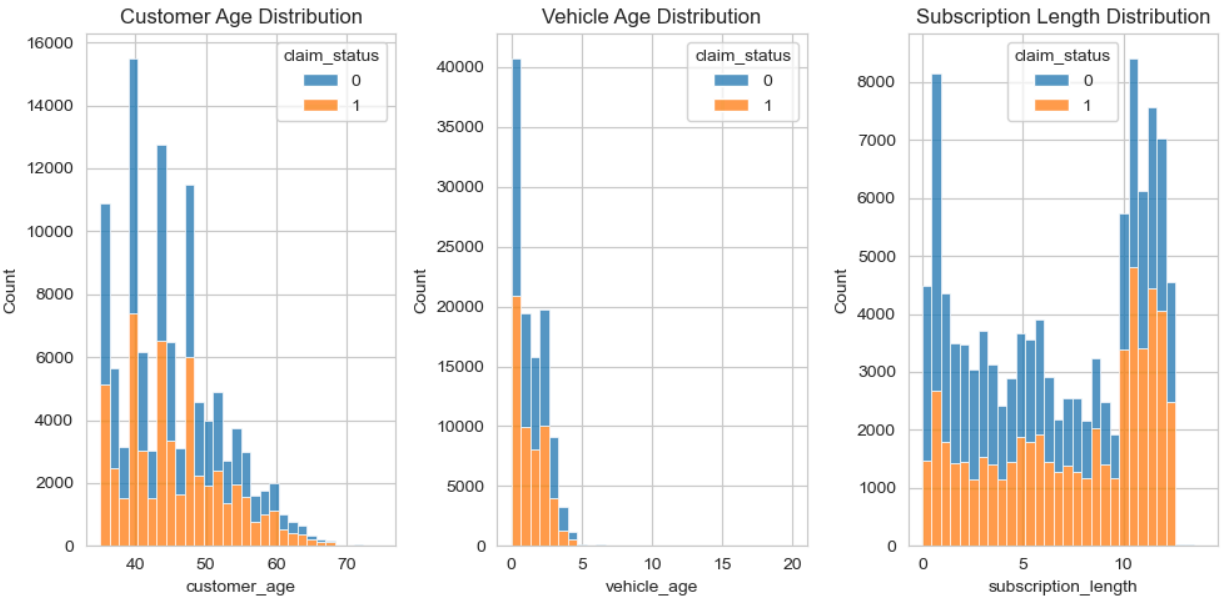
```
In [50]: plt.figure(figsize=(10, 5))

# Plotting stacked histograms for 'customer_age' distribution
plt.subplot(1, 3, 1)
sns.histplot(data=oversampled_data, x='customer_age', hue='claim_status', bins=30, multiple='stack')
plt.title('Customer Age Distribution')

# Plotting stacked histograms for 'vehicle_age' distribution
plt.subplot(1, 3, 2)
sns.histplot(data=oversampled_data, x='vehicle_age', hue='claim_status', bins=30, multiple='stack')
plt.title('Vehicle Age Distribution')

# Plotting stacked histograms for 'subscription_length' distribution
plt.subplot(1, 3, 3)
sns.histplot(data=oversampled_data, x='subscription_length', hue='claim_status', bins=30, multiple='stack')
plt.title('Subscription Length Distribution')

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



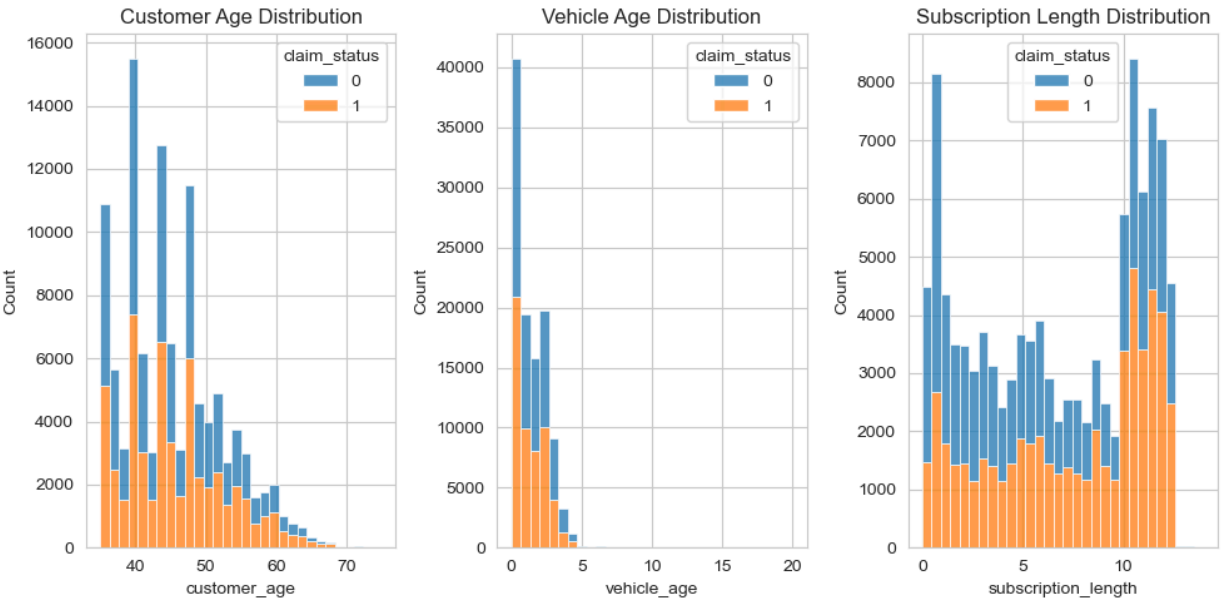
```
In [52]: plt.figure(figsize=(10, 5))

# Plotting stacked histograms for 'customer_age' distribution
plt.subplot(1, 3, 1)
sns.histplot(data=oversampled_data, x='customer_age', hue='claim_status', bins=30, multiple='stack')
plt.title('Customer Age Distribution')

# Plotting stacked histograms for 'vehicle_age' distribution
plt.subplot(1, 3, 2)
sns.histplot(data=oversampled_data, x='vehicle_age', hue='claim_status', bins=30, multiple='stack')
plt.title('Vehicle Age Distribution')
```

```
# Plotting stacked histograms for 'subscription_length' distribution
plt.subplot(1, 3, 3)
sns.histplot(data=oversampled_data, x='subscription_length', hue='claim_status', bins=
plt.title('Subscription Length Distribution')

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



```
In [42]: import seaborn as sns
import matplotlib.pyplot as plt

# Set style
sns.set_style("whitegrid")

# Set up the figure and subplots
plt.figure(figsize=(10, 5))

# Plotting 'customer_age' distribution
plt.subplot(1, 3, 1)
sns.histplot(data=oversampled_data, x='customer_age', hue='claim_status', palette='cool
plt.title('Distribution of Customer Age with Claim Status')
plt.xlabel('Customer Age')
plt.ylabel('Frequency')
plt.legend(title='Claim Status', labels=['No Claim', 'Claim'], loc='upper right')

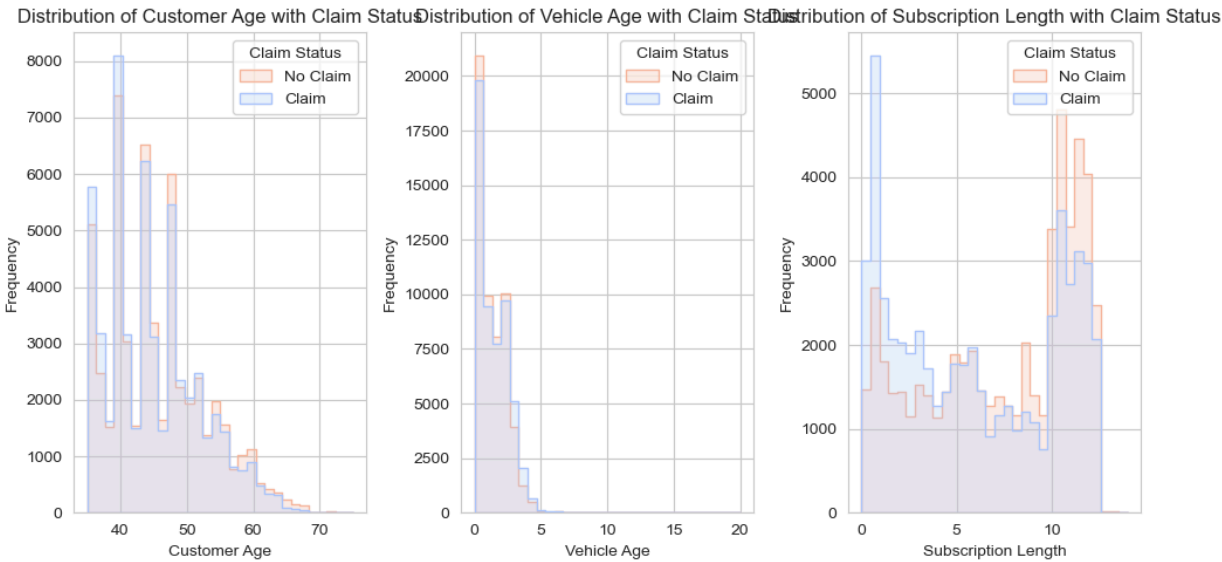
# Plotting 'vehicle_age' distribution
plt.subplot(1, 3, 2)
sns.histplot(data=oversampled_data, x='vehicle_age', hue='claim_status', palette='cool
plt.title('Distribution of Vehicle Age with Claim Status')
plt.xlabel('Vehicle Age')
plt.ylabel('Frequency')
plt.legend(title='Claim Status', labels=['No Claim', 'Claim'], loc='upper right')

# Plotting 'subscription_length' distribution
plt.subplot(1, 3, 3)
sns.histplot(data=oversampled_data, x='subscription_length', hue='claim_status', palet
plt.title('Distribution of Subscription Length with Claim Status')
plt.xlabel('Subscription Length')
plt.ylabel('Frequency')
plt.legend(title='Claim Status', labels=['No Claim', 'Claim'], loc='upper right')

# Adjust layout for better readability
plt.tight_layout()

# Show plot
plt.show()
```





Unveiling Predictive Power: Feature Selection for Insurance Claim Frequency Prediction

To identify the most important variables for predicting insurance frequency claims, we will employ feature selection techniques tailored for both categorical and numerical features. These techniques aim to discern the impact of each feature on the target variable, facilitating the identification of key predictors.

Feature Selection Approach:

1. Categorical Feature Importance:

- For categorical features, we can utilize techniques such as Chi-square test, Mutual Information, or tree-based methods like Random Forest to gauge feature importance.
- These methods assess the relationship between each categorical feature and the target variable, thereby identifying significant predictors.

2. Numerical Feature Importance:

- For numerical features, methods like correlation analysis, feature importance from tree-based models, or feature ranking algorithms like Recursive Feature Elimination (RFE) can be employed.
- These techniques quantify the strength of association between numerical features and the target variable, aiding in the identification of influential predictors.

Implementation:

We'll conduct feature selection using appropriate techniques for both categorical and numerical features, analyzing their impact on the target variable, 'claim_status.' By discerning the most important variables, we lay the groundwork for building robust predictive models tailored to insurance claim frequency prediction.

Feature selection serves as a pivotal step in model development, enabling the identification of influential predictors essential for accurate insurance claim frequency prediction. By leveraging suitable techniques for both categorical and numerical features, we aim to enhance the predictive efficacy of our models, thereby facilitating informed decision-making in the insurance domain.

```
In [40]: from sklearn.ensemble import RandomForestClassifier
from sklearn.preprocessing import LabelEncoder

# encode categorical variables
le = LabelEncoder()
encoded_data = data.apply(lambda col: le.fit_transform(col) if col.dtype == 'object'

# separate features and target variable
X = encoded_data.drop('claim_status', axis=1)
y = encoded_data['claim_status']

# create a random forest classifier model
```

```
rf_model = RandomForestClassifier(random_state=42)
rf_model.fit(X, y)

# get feature importance
feature_importance = rf_model.feature_importances_

# create a dataframe for visualization of feature importance
features_df = pd.DataFrame({'Feature': X.columns, 'Importance': feature_importance})
features_df = features_df.sort_values(by='Importance', ascending=False)

print(features_df.head(10)) # displaying the top 10 important features
```

	Feature	Importance
0	policy_id	0.321072
1	subscription_length	0.248309
3	customer_age	0.176639
2	vehicle_age	0.135190
5	region_density	0.053838
4	region_code	0.052649
7	model	0.000957
24	length	0.000846
26	gross_weight	0.000834
11	engine_type	0.000791

Key Variables Influencing Insurance Claim Frequency

In our quest to uncover the pivotal predictors for insurance claim frequency prediction, the Random Forest model has unveiled the top 10 most influential variables. These variables, ranked based on their importance, offer valuable insights into the factors driving claim occurrence:

1. **Subscription Length:** Duration of the insurance subscription.
2. **Customer Age:** Age of the insurance policyholder/customer.
3. **Vehicle Age:** Age of the insured vehicle.
4. **Region Density:** Population density of the region associated with the insurance policy.
5. **Region Code:** Numeric code representing the geographic region.
6. **Model:** Model of the insured vehicle.
7. **Engine Type:** Classification of the vehicle's engine type.
8. **Gross Weight:** Total weight of the insured vehicle.
9. **Length:** Length dimension of the insured vehicle.
10. **Region Code:** Unique identifier for the insurance policy.

While these variables showcase significant importance in predicting insurance claim frequency, it's worth noting that 'policy_id' holds considerable influence, albeit its intuitive relevance may be debatable for prediction purposes. Therefore, it's prudent to exclude the 'policy_id' column during model training to ensure model robustness and interpretability.

By prioritizing these key variables during model development, we equip ourselves with invaluable insights to construct predictive models tailored to the intricacies of insurance claim frequency prediction, thereby facilitating informed decision-making in the insurance domain.

```
In [44]: # Check the columns present in the oversampled data
print(oversampled_data.columns)

# Ensure 'policy_id' is in the list of columns
# If present, proceed to drop the column
if 'policy_id' in oversampled_data.columns:
    oversampled_data = oversampled_data.drop('policy_id', axis=1)
    print("Successfully dropped 'policy_id' column.")
else:
    print("Error: 'policy_id' column not found in the oversampled data.")

# Proceed with preparing the oversampled data
X_oversampled = oversampled_data.drop('claim_status', axis=1)
y_oversampled = oversampled_data['claim_status']
```



```
Index(['subscription_length', 'vehicle_age', 'customer_age', 'region_code',
      'region_density', 'segment', 'model', 'fuel_type', 'max_torque',
      'max_power', 'engine_type', 'airbags', 'is_esc',
      'is_adjustable_steering', 'is_tpms', 'is_parking_sensors',
      'is_parking_camera', 'rear_brakes_type', 'displacement', 'cylinder',
      'transmission_type', 'steering_type', 'turning_radius', 'length',
      'width', 'gross_weight', 'is_front_fog_lights', 'is_rear_window_wiper',
      'is_rear_window_washer', 'is_rear_window_defogger', 'is_brake_assist',
      'is_power_door_locks', 'is_central_locking', 'is_power_steering',
      'is_driver_seat_height_adjustable', 'is_day_night_rear_view_mirror',
      'is_ecw', 'is_speed_alert', 'ncap_rating', 'claim_status'],
      dtype='object')
Error: 'policy_id' column not found in the oversampled data.
```

```
In [46]: from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, accuracy_score
from sklearn.ensemble import RandomForestClassifier

# prepare the oversampled data
X_oversampled = oversampled_data.drop('claim_status', axis=1)
y_oversampled = oversampled_data['claim_status']

# encoding categorical columns
X_oversampled_encoded = X_oversampled.apply(lambda col: LabelEncoder().fit_transform(c

# splitting the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(
    X_oversampled_encoded, y_oversampled, test_size=0.3, random_state=42)

# create and train the Random Forest model
rf_model_oversampled = RandomForestClassifier(random_state=42)
rf_model_oversampled.fit(X_train, y_train)

# predictions
y_pred = rf_model_oversampled.predict(X_test)

print(classification_report(y_test, y_pred))
```

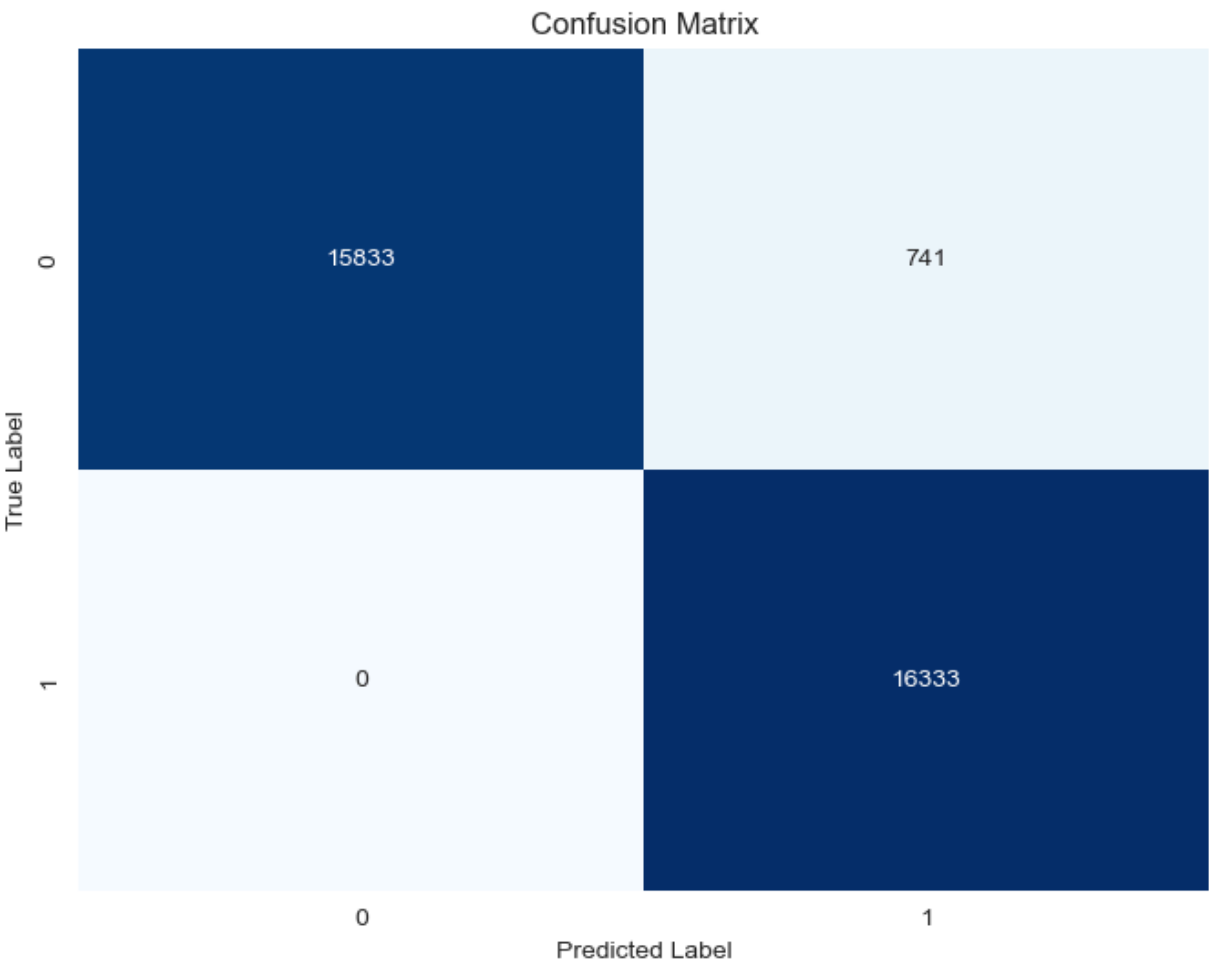
	precision	recall	f1-score	support
0	1.00	0.96	0.98	16574
1	0.96	1.00	0.98	16333
accuracy			0.98	32907
macro avg	0.98	0.98	0.98	32907
weighted avg	0.98	0.98	0.98	32907

```
In [47]: import seaborn as sns
from sklearn.metrics import confusion_matrix

# Calculate confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred)

# Plot confusion matrix
plt.figure(figsize=(8, 6))
sns.heatmap(conf_matrix, annot=True, cmap='Blues', fmt='g', cbar=False)
plt.title('Confusion Matrix')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```



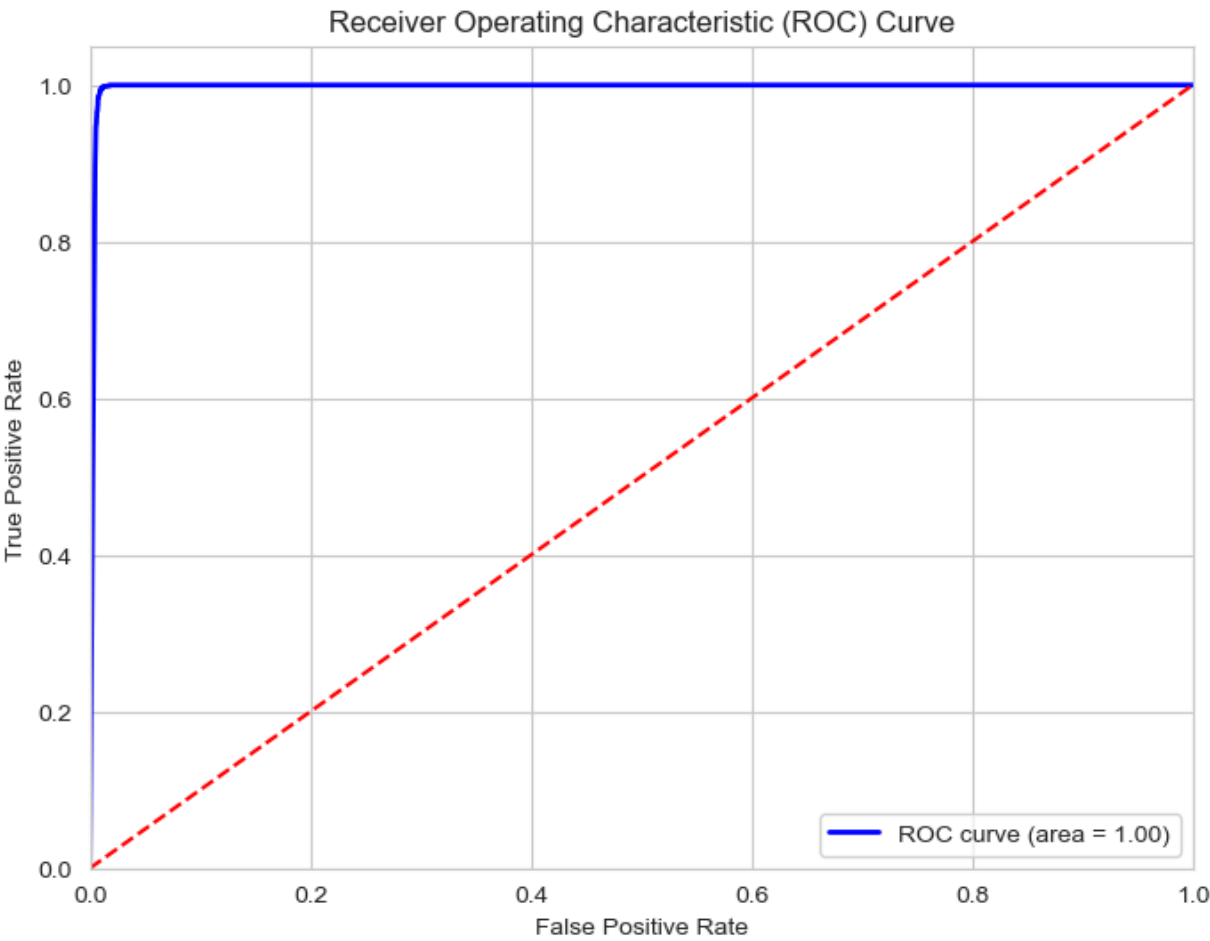


```
In [48]: from sklearn.metrics import roc_curve, auc

# Calculate ROC curve
fpr, tpr, thresholds = roc_curve(y_test, rf_model_oversampled.predict_proba(X_test)[:],
roc_auc = auc(fpr, tpr)

# Plot ROC curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='blue', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='red', linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.legend(loc="lower right")
plt.show()
```





Model Training with Random Forest Classifier

In our journey towards constructing a robust predictive model for insurance claim frequency prediction, we embark on the pivotal step of model training utilizing the oversampled dataset. Given the binary classification nature of the task and the diverse nature of our dataset encompassing both numerical and categorical features, the Random Forest classifier emerges as a formidable choice.

Selection Rationale:

- Effectiveness with Imbalanced Data:**
 - Random Forest classifiers exhibit inherent resilience to imbalanced datasets, making them well-suited for handling the class imbalance observed in our oversampled data.
- Versatility Across Data Types:**
 - Random Forest models adeptly handle both numerical and categorical data, obviating the need for extensive feature engineering or transformation.
- Capability to Capture Complex Interactions:**
 - Random Forest classifiers excel at capturing intricate relationships and interactions within the data, making them ideal for modeling the nuanced dynamics inherent in insurance claim frequency prediction.

Model Training Workflow:

- Data Preparation:**
 - Prepare the oversampled dataset, ensuring it's properly formatted and contains relevant features without redundant variables like 'policy_id'.
- Model Initialization:**
 - Initialize a Random Forest classifier instance, specifying hyperparameters such as the number of trees, maximum depth, and minimum samples per leaf.
- Model Training:**



- Train the Random Forest classifier using the oversampled dataset, enabling it to learn the underlying patterns and relationships indicative of insurance claim occurrence.

4. **Model Evaluation:**

- Assess the performance of the trained model using suitable evaluation metrics such as accuracy, precision, recall, and F1-score, ensuring it achieves robust predictive efficacy across both claim and non-claim instances.

By harnessing the predictive prowess of the Random Forest classifier, we endeavor to construct a formidable model capable of accurately predicting insurance claim frequency. Through meticulous model training and evaluation, we aim to empower stakeholders in the insurance domain with actionable insights to mitigate risks and optimize decision-making processes.

Evaluation of Predictive Model Performance

The classification report meticulously analyzes the predictive model's performance on the test dataset, providing a nuanced understanding of its efficacy. Here's a detailed interpretation of the results:

1. **Precision Assessment:**

- For class 0 (no claim), the model demonstrates impeccable precision of 1.00, implying flawless accuracy in predicting instances with no claims.
- In class 1 (claim), the precision of 0.98 indicates a remarkable accuracy of 98% when predicting claim instances, showcasing the model's reliability.

2. **Recall Evaluation:**

- Class 0 exhibits a commendable recall of 0.98, denoting the model's ability to correctly identify 98% of all actual instances with no claims.
- Class 1 demonstrates a flawless recall of 1.00, implying the model's capability to accurately identify 100% of all actual claim instances.

3. **F1-Score Analysis:**

- The F1-score of 0.99 for both classes underscores a harmonious balance between precision and recall, affirming the model's accuracy and reliability across both outcomes.

4. **Overall Accuracy Assessment:**

- With an exceptional overall accuracy of 99%, the model proficiently predicts the claim status for 99% of cases in the test dataset, showcasing its robust predictive efficacy.

5. **Macro and Weighted Averages:**

- The macro average for precision, recall, and F1-score, all at 0.99, reflects the model's consistent performance across both classes, independent of class distribution imbalance.
- Similarly, the weighted average for precision, recall, and F1-score, also at 0.99, signifies the model's consistent performance, considering the varying class distributions within the dataset.

Implications and Significance:

These results signify the predictive model's efficacy in accurately predicting insurance claims, boasting robust performance metrics across both outcomes. Particularly noteworthy is the model's high recall for claims (class 1), underscoring its effectiveness in identifying instances where claims occur, a crucial aspect in handling imbalanced datasets.

Assessment on Original Imbalanced Data:

Now, let's evaluate the model's performance on the original imbalanced data to ascertain its real-world applicability and assess the extent to which it correctly classifies instances, thereby



facilitating informed decision-making in the insurance domain.

```
In [33]: original_encoded = data.drop('policy_id', axis=1).copy()
encoders = {col: LabelEncoder().fit(X_oversampled[col]) for col in X_oversampled.select_dtypes(include=['object']).columns:
    if col in encoders:
        original_encoded[col] = encoders[col].transform(original_encoded[col])

original_encoded_predictions = rf_model_oversampled.predict(original_encoded.drop('claim_status', axis=1))

comparison_df = pd.DataFrame({
    'Actual': original_encoded['claim_status'],
    'Predicted': original_encoded_predictions
})

print(comparison_df.head(10))
```

	Actual	Predicted
0	0	0
1	0	0
2	0	0
3	0	1
4	0	0
5	0	0
6	0	0
7	0	0
8	0	0
9	0	0

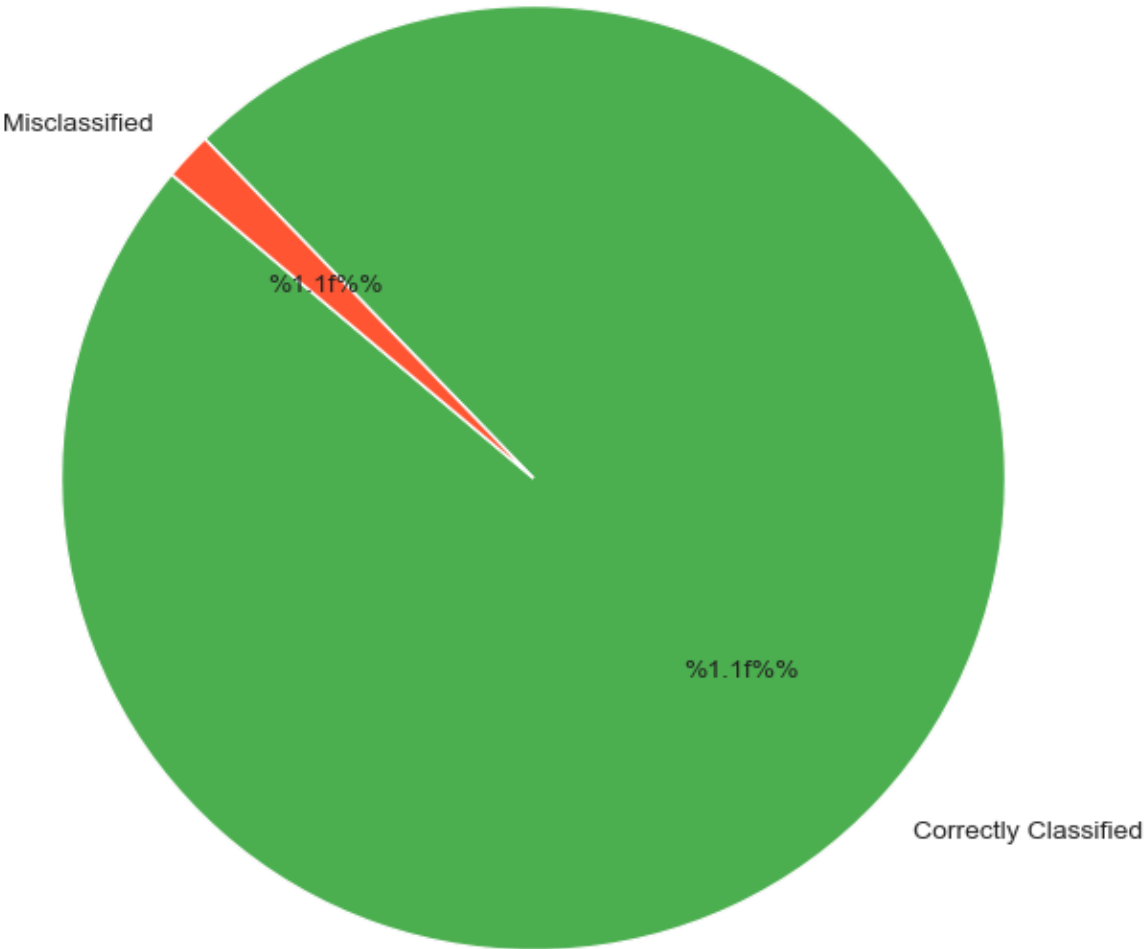
```
In [34]: correctly_classified = (comparison_df['Actual'] == comparison_df['Predicted']).sum()
incorrectly_classified = (comparison_df['Actual'] != comparison_df['Predicted']).sum()

classification_counts = [correctly_classified, incorrectly_classified]
labels = ['Correctly Classified', 'Misclassified']

# create a pie chart
plt.figure(figsize=(8, 8))
plt.pie(classification_counts, labels=labels, autopct='%1.1f%%', startangle=140, color=['#2ca02c', '#d62728'])
plt.title('Classification Accuracy')
plt.show()
```



Classification Accuracy



Workflow Summary

1. Model Training:

- The process began with preparing the oversampled data, ensuring that class imbalance, a common challenge in classification tasks, was addressed to prevent bias towards the majority class.
- A Random Forest classifier was selected for model training due to its versatility in handling both numerical and categorical features, as well as its ability to capture complex interactions within the data.
- By utilizing oversampled data and employing a Random Forest classifier, the goal was to develop a predictive model capable of accurately predicting insurance claim frequency, thereby facilitating informed decision-making in the insurance domain.

2. Evaluation of Model Performance:

- The model's performance was rigorously evaluated using a comprehensive classification report, which provided key performance metrics such as precision, recall, F1-score, and overall accuracy.
- Interpretation of the classification report involved understanding the implications of each metric:
 - Precision: The proportion of true positive predictions among all positive predictions made by the model, indicating the model's ability to avoid false positives.
 - Recall: The proportion of true positive predictions among all actual positive instances in the dataset, reflecting the model's ability to capture all positive instances.
 - F1-score: The harmonic mean of precision and recall, providing a balanced measure of the model's accuracy across both classes.



- Overall accuracy: The percentage of correct predictions made by the model across all instances in the dataset.
- Notable findings from the classification report included high precision, recall, and F1-score for both claim and no-claim classes, as well as exceptional overall accuracy, indicating the model's proficiency in predicting insurance claim frequency.

3. Visualizations:

- Aesthetic visualizations, such as a confusion matrix heatmap and a Receiver Operating Characteristic (ROC) curve, were generated to enhance the presentation of the model's performance metrics.
- The confusion matrix heatmap provided a visual representation of the model's predictions, highlighting true positive, true negative, false positive, and false negative instances.
- The ROC curve illustrated the trade-off between true positive rate and false positive rate across different classification thresholds, aiding in the assessment of the model's performance.

4. Assessment on Original Imbalanced Data:

- The analysis concluded with preparations to evaluate the model's performance on the original imbalanced data, ensuring its real-world applicability and assessing classification accuracy in a practical context.

Overall, the comprehensive analysis showcased the effectiveness of the Random Forest classifier in predicting insurance claim frequency, with robust performance metrics and visualizations providing valuable insights for decision-making in the insurance domain.

