```python
In [32]:
import pandas as pd
import numpy as np

import seaborn as sns
import matplotlib.pyplot as plt

from sklearn.svm import LinearSVC
from sklearn.naive_bayes import GaussianNB

from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score, classification_report, confusi

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.utils import resample
from sklearn.preprocessing import StandardScaler

from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score, classification_report


sns.set_style('whitegrid')
sns.set_palette('pastel')
import warnings

warnings.simplefilter("ignore")
```

```python
In [33]:  df = pd.read_csv(r"C:\Users\Simi\Downloads\churn prediction.csv")
```

```python
In [34]:  print(df.head())
```

```
   RowNumber  CustomerId   Surname  CreditScore Geography  Gender  Age  \
0          1    15634602  Hargrave          619    France  Female   42
1          2    15647311      Hill          608     Spain  Female   41
2          3    15619304      Onio          502    France  Female   42
3          4    15701354      Boni          699    France  Female   39
4          5    15737888  Mitchell          850     Spain  Female   43

   Tenure    Balance  NumOfProducts  HasCrCard  IsActiveMember  \
0       2       0.00              1          1               1
1       1   83807.86              1          0               1
2       8  159660.80              3          1               0
3       1       0.00              2          0               0
4       2  125510.82              1          1               1

   EstimatedSalary  Exited
0        101348.88       1
1        112542.58       0
2        113931.57       1
3         93826.63       0
4         79084.10       0
```

```
In [35]: print(df.isnull().sum())
```

```
RowNumber          0
CustomerId         0
Surname            0
CreditScore        0
Geography          0
Gender             0
Age                0
Tenure             0
Balance            0
NumOfProducts      0
HasCrCard          0
IsActiveMember     0
EstimatedSalary    0
Exited             0
dtype: int64
```

```
In [36]: df.columns
```

```
Out[36]: Index(['RowNumber', 'CustomerId', 'Surname', 'CreditScore', 'Geography',
                'Gender', 'Age', 'Tenure', 'Balance', 'NumOfProducts', 'HasCrCar
        d',
                'IsActiveMember', 'EstimatedSalary', 'Exited'],
              dtype='object')
```

```
In [37]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 14 columns):
 #   Column           Non-Null Count  Dtype
---  ------           --------------  -----
 0   RowNumber        10000 non-null  int64
 1   CustomerId       10000 non-null  int64
 2   Surname          10000 non-null  object
 3   CreditScore      10000 non-null  int64
 4   Geography        10000 non-null  object
 5   Gender           10000 non-null  object
 6   Age              10000 non-null  int64
 7   Tenure           10000 non-null  int64
 8   Balance          10000 non-null  float64
 9   NumOfProducts    10000 non-null  int64
 10  HasCrCard        10000 non-null  int64
 11  IsActiveMember   10000 non-null  int64
 12  EstimatedSalary  10000 non-null  float64
 13  Exited           10000 non-null  int64
dtypes: float64(2), int64(9), object(3)
memory usage: 1.1+ MB
```

```
In [38]: df.head(3)
```

Out[38]:

| | RowNumber | CustomerId | Surname | CreditScore | Geography | Gender | Age | Tenure | Ba |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 15634602 | Hargrave | 619 | France | Female | 42 | 2 | |
| 1 | 2 | 15647311 | Hill | 608 | Spain | Female | 41 | 1 | 838 |
| 2 | 3 | 15619304 | Onio | 502 | France | Female | 42 | 8 | 1596 |

◀ ▶

```
In [39]: df.shape
```

Out[39]: (10000, 14)

```
In [40]: is_Exited = df["Exited"].value_counts()
         print("Yes: ",is_Exited[1])
         print("No: ",is_Exited[0])
```

```
Yes:  2037
No:  7963
```

```
In [41]: print(df.isna().sum().sum())
         print(df.duplicated().sum())
```

```
0
0
```

```
In [42]: fig,axb = plt.subplots(ncols=2,nrows=1,figsize=(15, 8))

         #Gender Distribution
         explode = [0.1, 0.1]
         df.groupby('Gender')['Exited'].count().plot.pie(explode=explode, autopct="

         ax = sns.countplot(x="Gender", hue="Exited", data=df,ax=axb[1])

         # Add values on top of each bar
         for p in ax.patches:
             ax.annotate(f'{p.get_height()}', (p.get_x() + p.get_width() / 2., p.ge
                         ha='center', va='center', xytext=(0, 10), textcoords='offs

         # Set labels and title
         plt.title("Distribution of Gender with Exited Status")
         plt.xlabel("Gender")
         plt.ylabel("Count")

         # Show the plot
         plt.show()
```
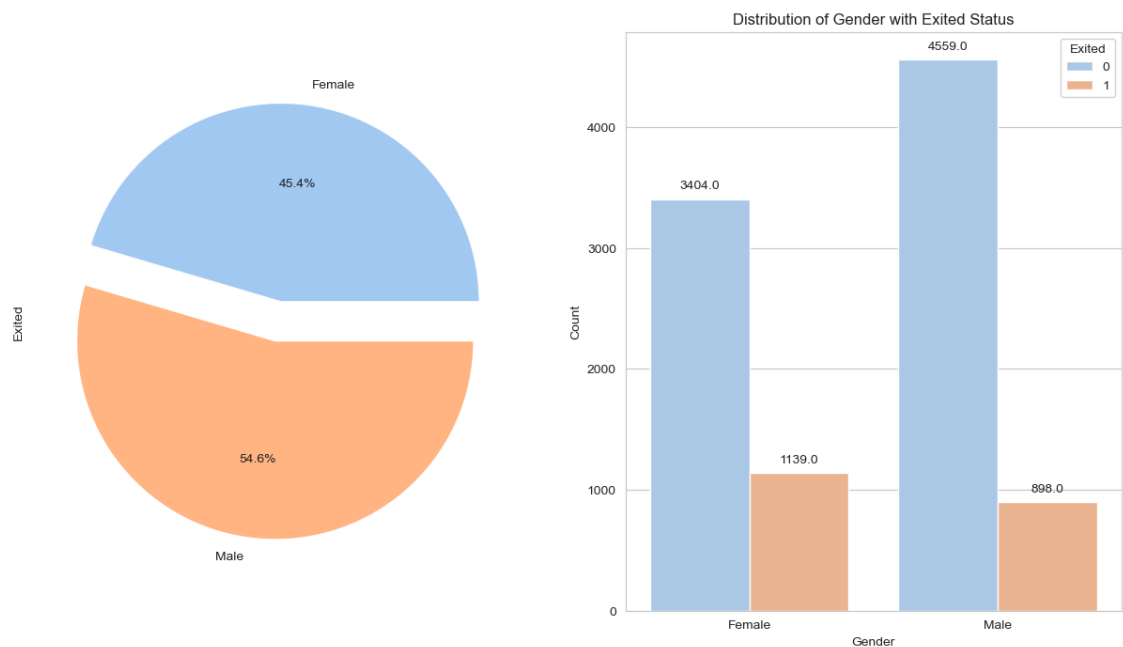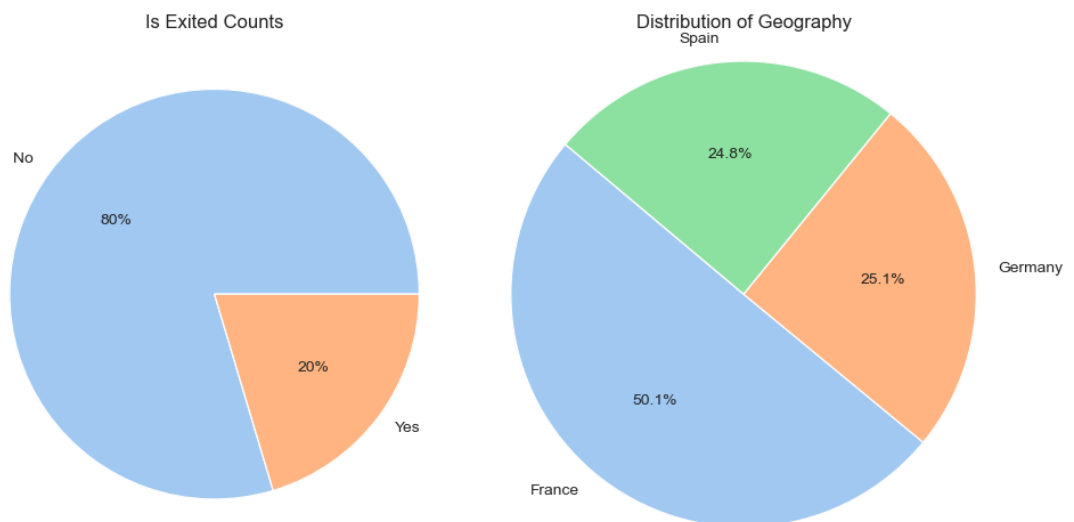
```python
In [43]: is_Exited = df["Exited"].value_counts()
         plt.figure(figsize=(10, 5))  # Set the same figsize for both plots
         plt.subplot(1, 2, 1)  # Subplot for the Exited Counts pie chart
         plt.pie(is_Exited, labels=["No", "Yes"], autopct="%0.0f%%")
         plt.title("Is Exited Counts")

         # Distribution of Geography Pie Chart
         plt.subplot(1, 2, 2)  # Subplot for the Distribution of Geography pie char
         geography_counts = df['Geography'].value_counts()
         plt.pie(geography_counts, labels=geography_counts.index, autopct='%1.1f%%'
         plt.axis('equal')  # Equal aspect ratio ensures that pie is drawn as a cir
         plt.title('Distribution of Geography')

         plt.tight_layout()  # Adjust layout to prevent overlapping
         plt.show()
```
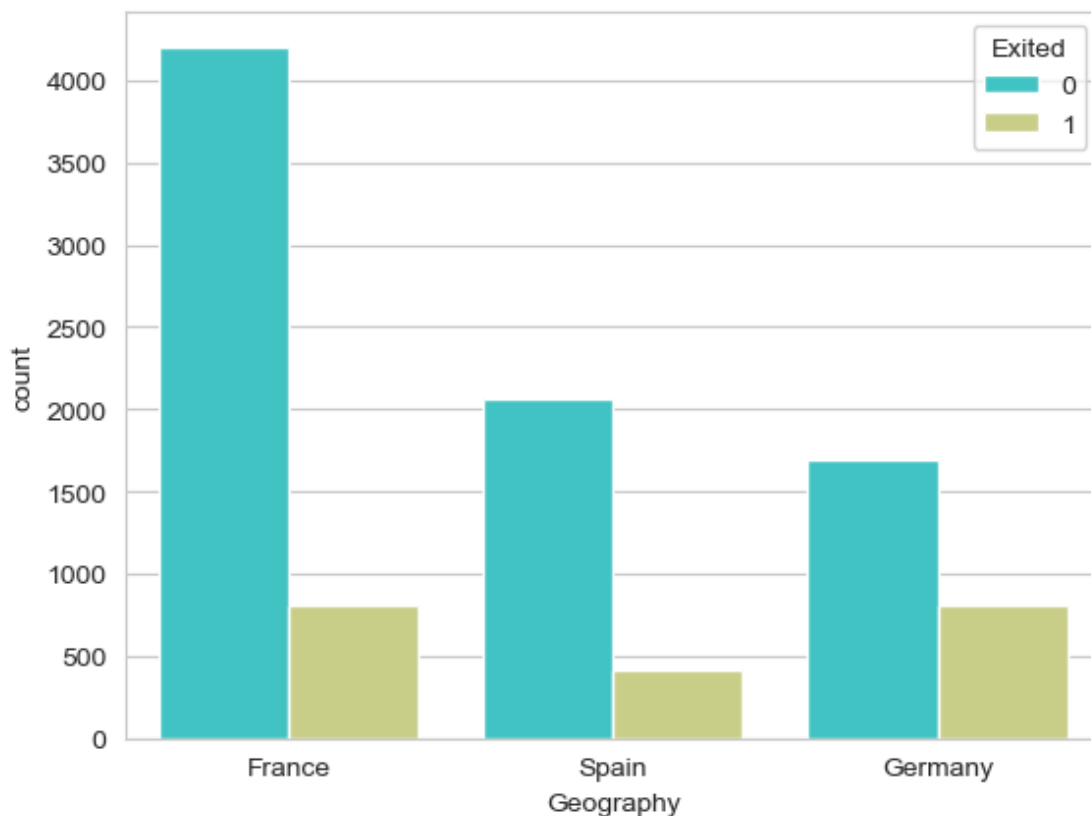
In [44]:
```python
sns.countplot(x='Geography',hue='Exited',data=df, palette='rainbow')
```

Out[44]: <Axes: xlabel='Geography', ylabel='count'>



In [45]:
```python
df = df.drop(['RowNumber', 'Surname', 'CustomerId'], axis= 1)

df['Balance'] = df['Balance'].astype(int)
df['EstimatedSalary'] = df['EstimatedSalary'].astype(int)
```

In [46]:
```python
df.head(3)
```

Out[46]:

| | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | Is/ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 619 | France | Female | 42 | 2 | 0 | 1 | 1 | |
| 1 | 608 | Spain | Female | 41 | 1 | 83807 | 1 | 0 | |
| 2 | 502 | France | Female | 42 | 8 | 159660 | 3 | 1 | |

```
In [47]:  # Initialize label encoders
          le = LabelEncoder()
          # Fit and transform the data
          df['Gender'] = le.fit_transform(df['Gender'])
          df['Geography'] = le.fit_transform(df['Geography'])

          df.head(4)
```

Out[47]:

| | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | HasCrCard | IsA |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 619 | 0 | 0 | 42 | 2 | 0 | 1 | 1 | |
| 1 | 608 | 2 | 0 | 41 | 1 | 83807 | 1 | 0 | |
| 2 | 502 | 0 | 0 | 42 | 8 | 159660 | 3 | 1 | |
| 3 | 699 | 0 | 0 | 39 | 1 | 0 | 2 | 0 | |

```
In [48]:  No_class = df[df["Exited"]==0]
          yes_class = df[df["Exited"]==1]

          No_class = resample(No_class, replace=False, n_samples=len(yes_class))
          down_samples = pd.concat([yes_class, No_class], axis=0)

          X = down_samples.drop("Exited", axis=1)
          y = down_samples["Exited"]

          X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r

          scaler = StandardScaler()
          X_train = scaler.fit_transform(X_train)
          X_test = scaler.transform(X_test)
```

```python
In [49]:  # Count the occurrences of each class in the original dataset
          original_class_counts = df["Exited"].value_counts()

          # Count the occurrences of each class in the downsampled dataset
          downsampled_class_counts = down_samples["Exited"].value_counts()

          # Calculate the percentage of each class
          original_percentages = original_class_counts / len(df) * 100
          downsampled_percentages = downsampled_class_counts / len(down_samples) * 1

          # Plotting
          plt.figure(figsize=(12, 6))

          # Bar chart for original class distribution
          plt.subplot(1, 2, 1)
          bars_1 = plt.bar(original_class_counts.index, original_class_counts.values
          for bar, label in zip(bars_1, original_percentages):
              plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height() + 5, f'{l
          plt.title('Original Class Distribution')
          plt.xlabel('Class')
          plt.ylabel('Count')
          plt.xticks(original_class_counts.index, ['Not Exited', 'Exited'])

          # Bar chart for downsampled class distribution
          plt.subplot(1, 2, 2)
          bars_2 = plt.bar(downsampled_class_counts.index, downsampled_class_counts.
          for bar, label in zip(bars_2, downsampled_percentages):
              plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height() + 5, f'{l
          plt.title('Downsampled Class Distribution')
          plt.xlabel('Class')
          plt.ylabel('Count')
          plt.xticks(downsampled_class_counts.index, ['Not Exited', 'Exited'])

          plt.tight_layout() # the plots will be automatically adjusted to ensure th
          plt.show()
```
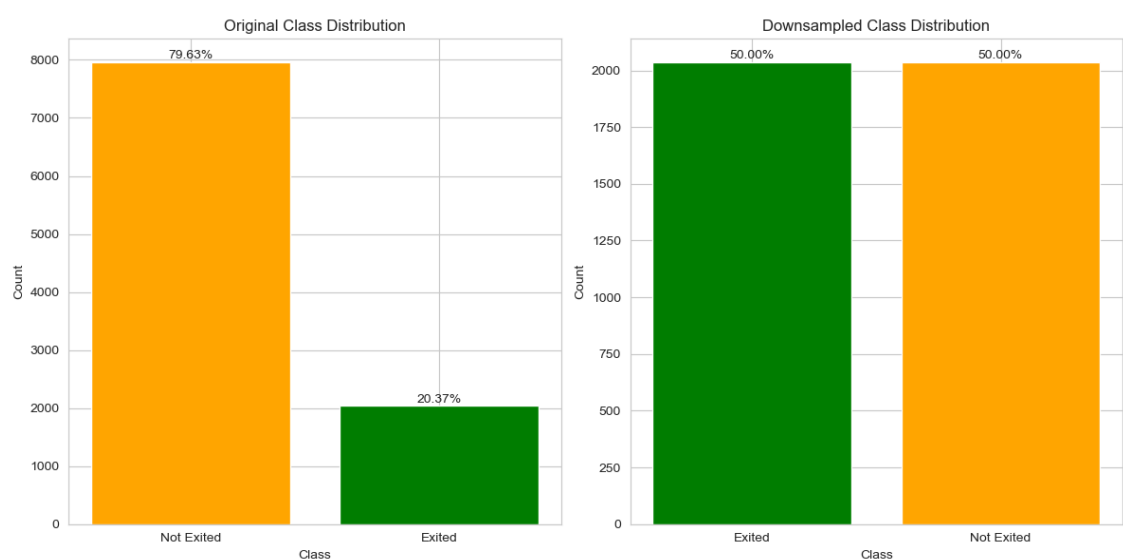
```
In [50]: DT = DecisionTreeClassifier(max_depth=(5), random_state=0)
         DT.fit(X_train, y_train)
         predict_ID3 = DT.predict(X_test)
         print(classification_report(y_test, predict_ID3))
         ID3_accuracy = accuracy_score(predict_ID3,y_test)
         print('ID3 model accuracy is: {:.2f}%'.format(ID3_accuracy*100))
```

```
              precision    recall  f1-score   support

           0       0.76      0.82      0.79       410
           1       0.80      0.73      0.76       405

    accuracy                           0.78       815
   macro avg       0.78      0.78      0.78       815
weighted avg       0.78      0.78      0.78       815

ID3 model accuracy is: 77.67%
```

```
In [51]: LR_model = LogisticRegression()
         LR_model.fit(X_train, y_train)
         predict_LR = LR_model.predict(X_test)
         print(classification_report(y_test, predict_LR))
         LR_accuracy = accuracy_score(predict_LR,y_test)
         print('Logistic Regression accuracy is: {:.2f}%'.format(LR_accuracy*100))
```

```
              precision    recall  f1-score   support

           0       0.72      0.71      0.71       410
           1       0.71      0.71      0.71       405

    accuracy                           0.71       815
   macro avg       0.71      0.71      0.71       815
weighted avg       0.71      0.71      0.71       815

Logistic Regression accuracy is: 71.29%
```

```
In [52]: svm_model = LinearSVC()
         svm_model.fit(X_train, y_train)
         predict = svm_model.predict(X_test)

         print(classification_report(y_test, predict))
         svm_accuracy = accuracy_score(predict,y_test)
         print('SVC model accuracy is: {:.2f}%'.format(svm_accuracy*100))
```

```
              precision    recall  f1-score   support

           0       0.71      0.72      0.72       410
           1       0.71      0.71      0.71       405

    accuracy                           0.71       815
   macro avg       0.71      0.71      0.71       815
weighted avg       0.71      0.71      0.71       815

SVC model accuracy is: 71.29%
```

```
In [53]: RF = RandomForestClassifier(n_estimators=60, random_state=0)
         RF.fit(X_train, y_train)

         predict_RF = RF.predict(X_test)

         # Evaluate the model
         print(classification_report(y_test, predict_RF))
         RF_accuracy = accuracy_score(predict_RF, y_test)
         print('Random Forest model accuracy is: {:.2f}%'.format(RF_accuracy * 100)
```

```
              precision    recall  f1-score   support

           0       0.77      0.82      0.79       410
           1       0.81      0.75      0.77       405

    accuracy                           0.78       815
   macro avg       0.79      0.78      0.78       815
weighted avg       0.79      0.78      0.78       815


Random Forest model accuracy is: 78.40%
```

```
In [54]: gb_classifier = GradientBoostingClassifier(n_estimators=100, learning_rate

         gb_classifier.fit(X_train, y_train)
         y_pred = gb_classifier.predict(X_test)

         # Generate classification report
         report = classification_report(y_test, y_pred)
         print("Classification Report:\n", report)

         # Calculate accuracy
         gb_accuracy = accuracy_score(y_test, y_pred)
         print('XGBoost model accuracy is: {:.2f}%'.format(gb_accuracy * 100))
```

```
Classification Report:
              precision    recall  f1-score   support

           0       0.79      0.83      0.81       410
           1       0.82      0.78      0.80       405

    accuracy                           0.80       815
   macro avg       0.80      0.80      0.80       815
weighted avg       0.80      0.80      0.80       815


XGBoost model accuracy is: 80.37%
```

```
In [ ]:
```