

```

1 from google.colab import files
2
3 uploaded = files.upload()

```



Choose Files Churn\_Modelling.csv

• Churn\_Modelling.csv(text/csv) - 684858 bytes, last modified: 3/19/2022 - 100% done  
Saving Churn\_Modelling.csv to Churn\_Modelling.csv

```

1 import pandas as pd
2 data = pd.read_csv('Churn_Modelling.csv')
3
4 # Display the first few rows of the DataFrame to ensure it's loaded correctly
5 data.head()

```



	RowNumber	CustomerId	Surname	CreditScore	Geography	Gender	Age	Tenure	Balance	NumOfProducts	HasCrCard	IsActiveMemb
0	1	15634602	Hargrave	619	France	Female	42	2	0.00	1	1	
1	2	15647311	Hill	608	Spain	Female	41	1	83807.86	1	0	
2	3	15619304	Onio	502	France	Female	42	8	159660.80	3	1	
3	4	15701354	Boni	699	France	Female	39	1	0.00	2	0	
4	5	15737888	Mitchell	850	Spain	Female	43	2	125510.82	1	1	

Next steps:

[Generate code with data](#)

[View recommended plots](#)

[New interactive sheet](#)

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler, OneHotEncoder
7 from sklearn.compose import ColumnTransformer
8 from sklearn.pipeline import Pipeline
9 from sklearn.linear_model import LogisticRegression
10 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
11 from sklearn.metrics import classification_report, roc_auc_score, confusion_matrix, roc_curve

```

```

1 # Drop unnecessary columns
2 data = data.drop(['RowNumber', 'CustomerId', 'Surname'], axis=1)

```

```

1 # Split features and target
2 X = data.drop('Exited', axis=1)
3 y = data['Exited']

```

```

1 # Define categorical and numerical columns
2 categorical_cols = ['Geography', 'Gender']
3 numerical_cols = ['CreditScore', 'Age', 'Tenure', 'Balance', 'NumOfProducts',
4                  'HasCrCard', 'IsActiveMember', 'EstimatedSalary']

```

```

1 # Preprocessing pipeline
2 preprocessor = ColumnTransformer(
3     transformers=[
4         ('num', StandardScaler(), numerical_cols),
5         ('cat', OneHotEncoder(), categorical_cols)]

```

```

1 # Logistic Regression pipeline
2 log_reg_pipeline = Pipeline(steps=[('preprocessor', preprocessor),
3                                   ('classifier', LogisticRegression())])
4
5 # Random Forest pipeline
6 rf_pipeline = Pipeline(steps=[('preprocessor', preprocessor),
7                               ('classifier', RandomForestClassifier(random_state=42))])
8
9 # Gradient Boosting pipeline
10 gb_pipeline = Pipeline(steps=[('preprocessor', preprocessor),
11                              ('classifier', GradientBoostingClassifier(random_state=42))])

```

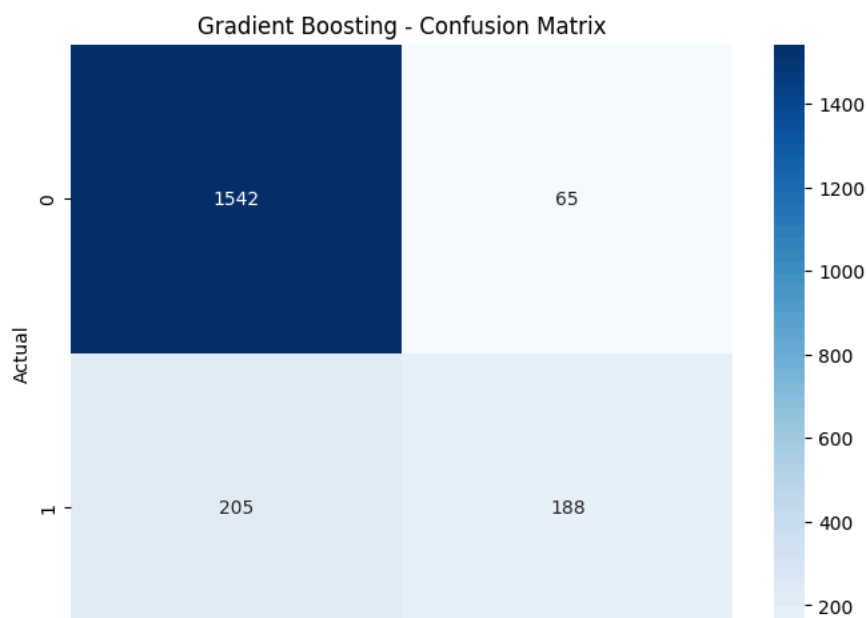
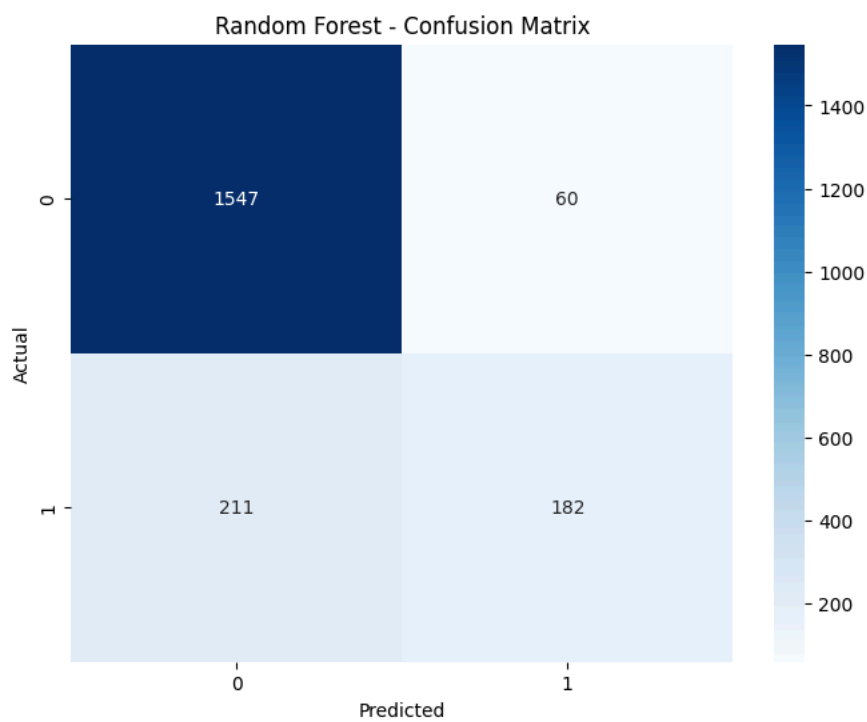
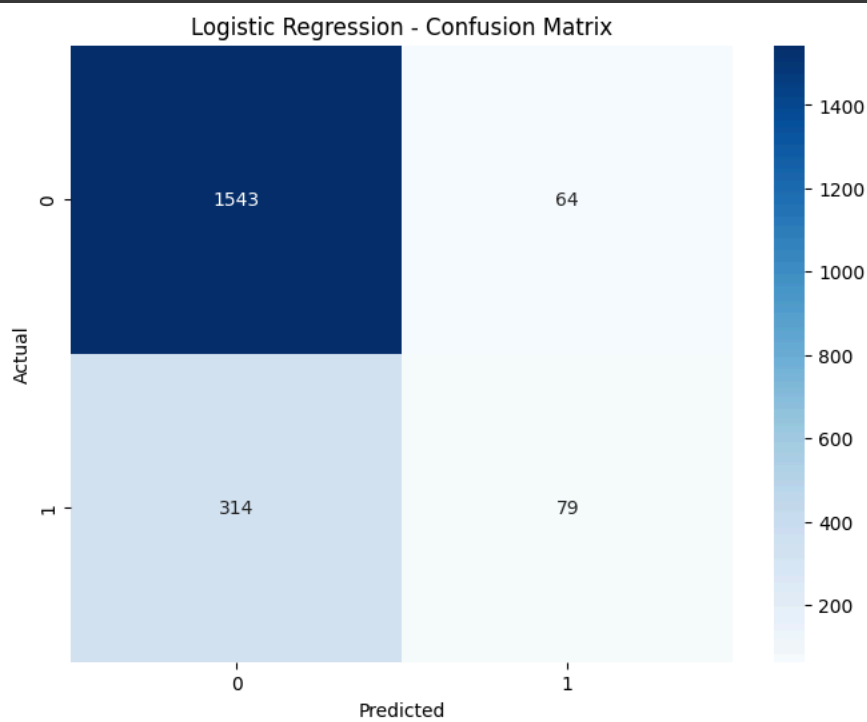
```

1 # Split the data
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

```
1 # Train models
2 pipelines = {}
3 for model_name, model in models.items():
4     pipeline = Pipeline(steps=[('preprocessor', preprocessor), ('classifier', model)])
5     pipeline.fit(X_train, y_train)
6     pipelines[model_name] = pipeline
```

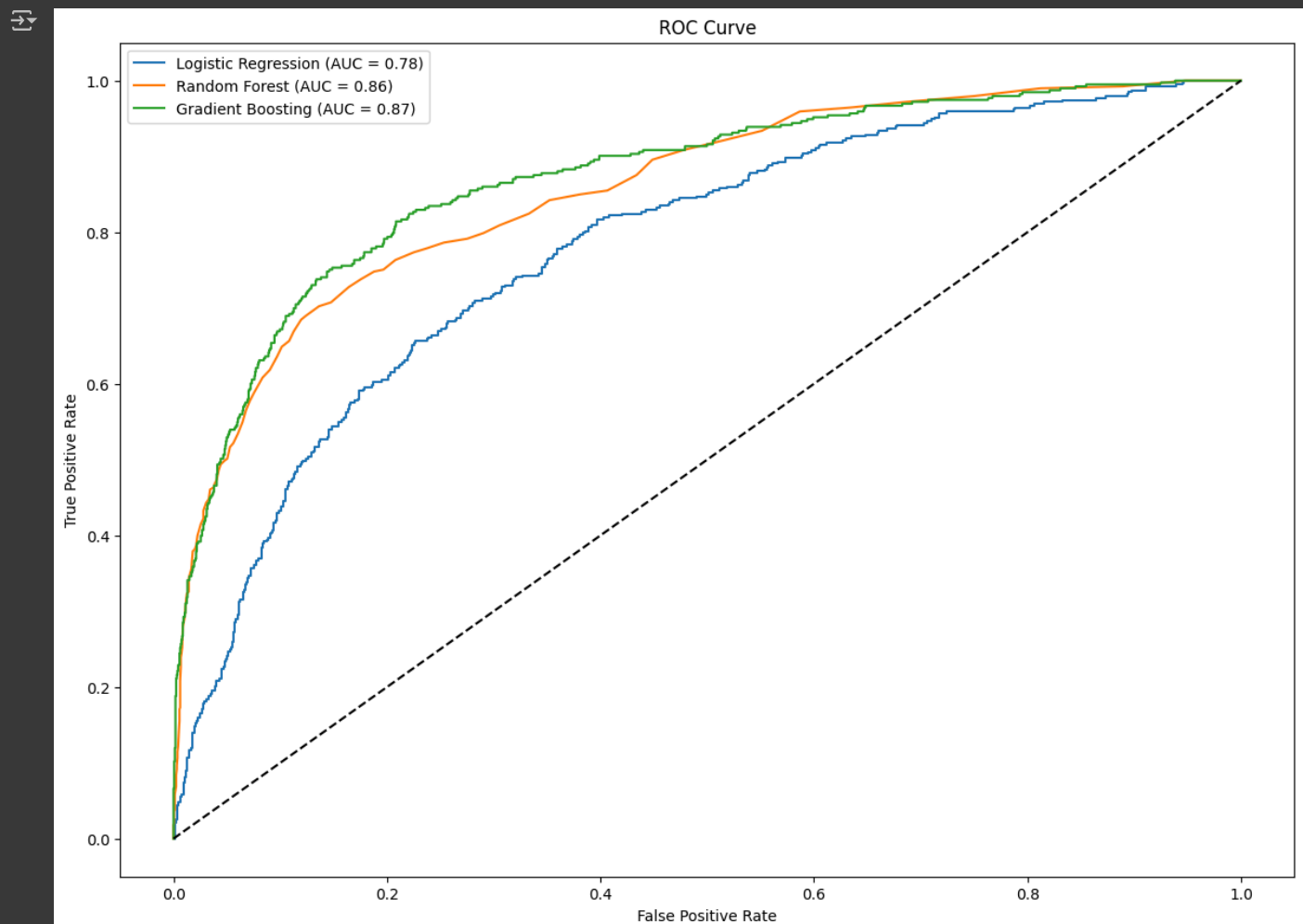
```
1 for model_name, pipeline in pipelines.items():
2     y_pred = pipeline.predict(X_test)
3     cm = confusion_matrix(y_test, y_pred)
4     plt.figure(figsize=(8, 6))
5     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
6     plt.title(f'{model_name} - Confusion Matrix')
7     plt.xlabel('Predicted')
8     plt.ylabel('Actual')
9     plt.show()
```



```

1 plt.figure(figsize=(14, 10))
2
3 for model_name, pipeline in pipelines.items():
4     y_pred_prob = pipeline.predict_proba(X_test)[: , 1]
5     fpr, tpr, _ = roc_curve(y_test, y_pred_prob)
6     plt.plot(fpr, tpr, label=f'{model_name} (AUC = {roc_auc_score(y_test, y_pred_prob):.2f})')
7
8 plt.plot([0, 1], [0, 1], 'k--') # Diagonal line
9 plt.title('ROC Curve')
10 plt.xlabel('False Positive Rate')
11 plt.ylabel('True Positive Rate')
12 plt.legend(loc='best')
13 plt.show()

```



```

1 for model_name, pipeline in pipelines.items():
2     if model_name in ['Random Forest', 'Gradient Boosting']:
3         model = pipeline.named_steps['classifier']
4         feature_importance = model.feature_importances_
5         feature_names = numerical_cols + list(pipeline.named_steps['preprocessor'].transformers_[1][1].get_feature_names_out(ca
6         importance_df = pd.DataFrame({'Feature': feature_names, 'Importance': feature_importance})
7         importance_df = importance_df.sort_values(by='Importance', ascending=False)
8
9         plt.figure(figsize=(10, 8))
10        sns.barplot(x='Importance', y='Feature', data=importance_df)
11        plt.title(f'{model_name} - Feature Importance')
12        plt.show()

```

