



Travel Insurance Claim Prediction

Asyraf Al Rasyid

Table of Content

**Business
Problem**

**Feature
Engineering**

Conclusion

**Data
Understanding**

**Data
Preprocessing**

Recommendation

**Exploratory Data
Analysis (EDA)**

Modelling

Business Problem

Problem Statement

pengajuan klaim yang tidak terduga dapat memberikan dampak terhadap perusahaan. Model ini diharapkan dapat membantu perusahaan dalam membuat keputusan yang lebih tepat dalam menentukan apakah seseorang akan melakukan klaim atau tidak pada asuransi perjalannya

Goals

Mengembangkan model prediktif yang dapat dengan akurat memprediksi kemungkinan klaim berdasarkan data yang tersedia.

Analytic Approach

menemukan suatu pola yang dapat membedakan pemegang polis yang akan melakukan klaim asuransi dan yang tidak

		Actual Class	
		1	0
Predicted Class	1	True Positive	False Positive
	0	False Negative	True Negative

Seseorang yang dianggap klaim tapi tidak klaim

Seseorang yang dianggap tidak klaim tapi klaim

Data Understanding

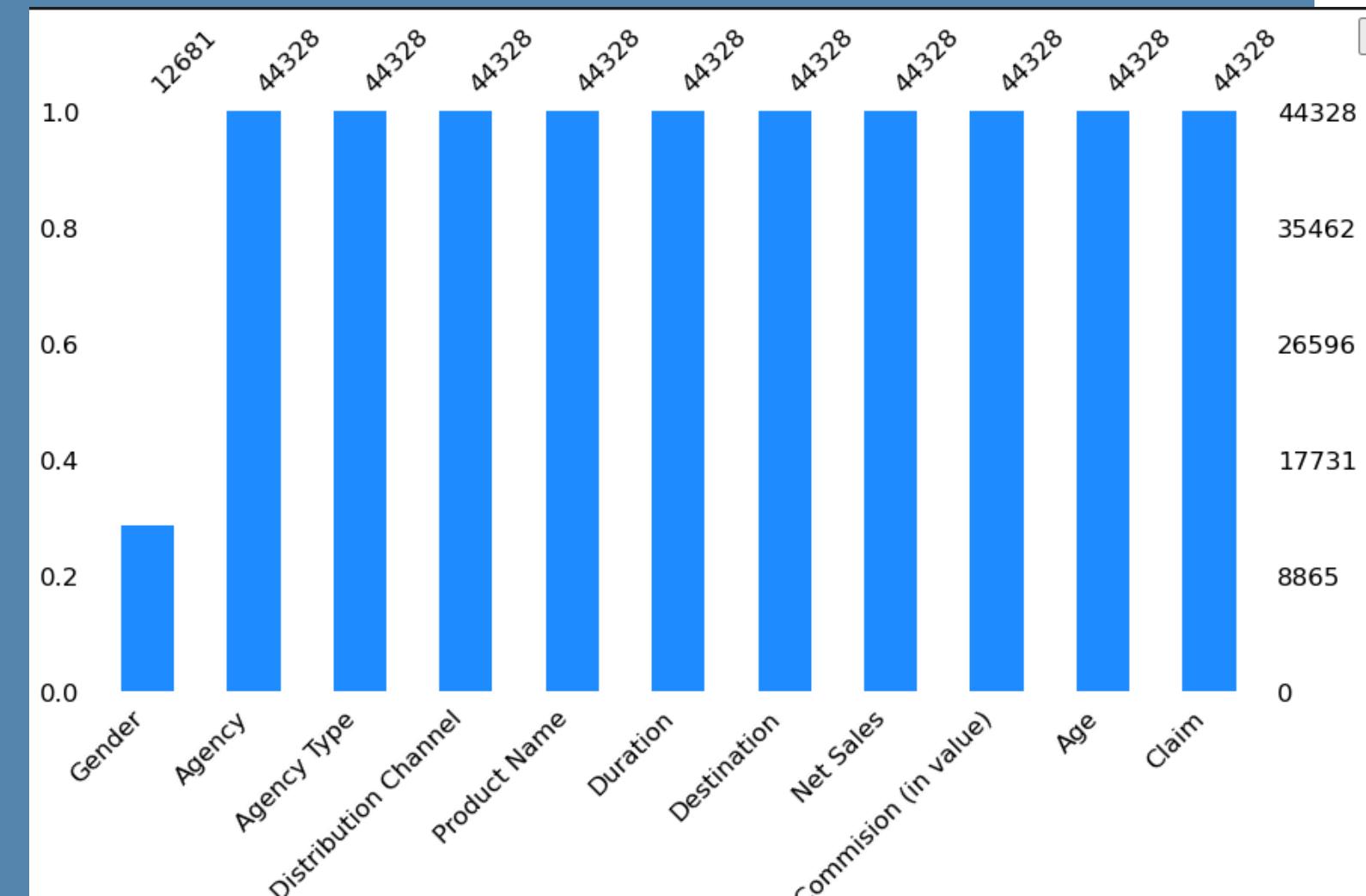
Jumlah baris dan kolom: (44328, 11)

- Agency: Nama agensi asuransi perjalanan
- Agency Type: Tipe asuransi perjalanan
- Distribution Channel: Penyaluran produk kepada konsumen
- Product Name: Jenis produk asuransi yang digunakan turis
- Gender: Jenis kelamin

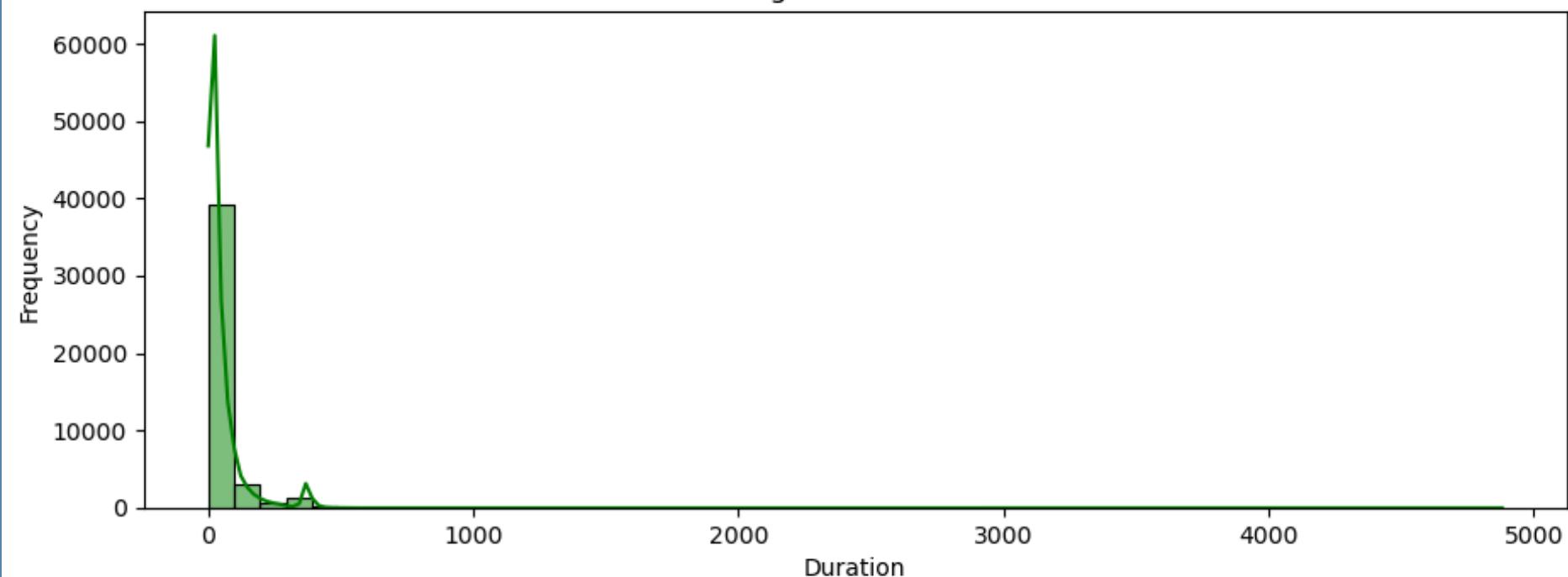
- Duration: Durasi perjalanan dalam hari
- Destination: Tujuan perjalanan
- Net Sales: Jumlah harga penjualan asuransi perjalanan (dalam dollar singapore)
- Commission (in value): Komisi yang didapatkan agensi asuransi perjalanan
- Age: Usia turis
- Claim: Claim status.

	Duration	Net Sales	Commision (in value)	Age
count	44328.000000	44328.000000	44328.000000	44328.000000
mean	49.424292	40.550948	9.707692	39.925600
std	109.153961	48.661970	19.625637	13.954926
min	-1.000000	-357.500000	0.000000	0.000000
25%	9.000000	18.000000	0.000000	35.000000
50%	22.000000	26.500000	0.000000	36.000000
75%	53.000000	48.000000	11.550000	43.000000
max	4881.000000	810.000000	283.500000	118.000000

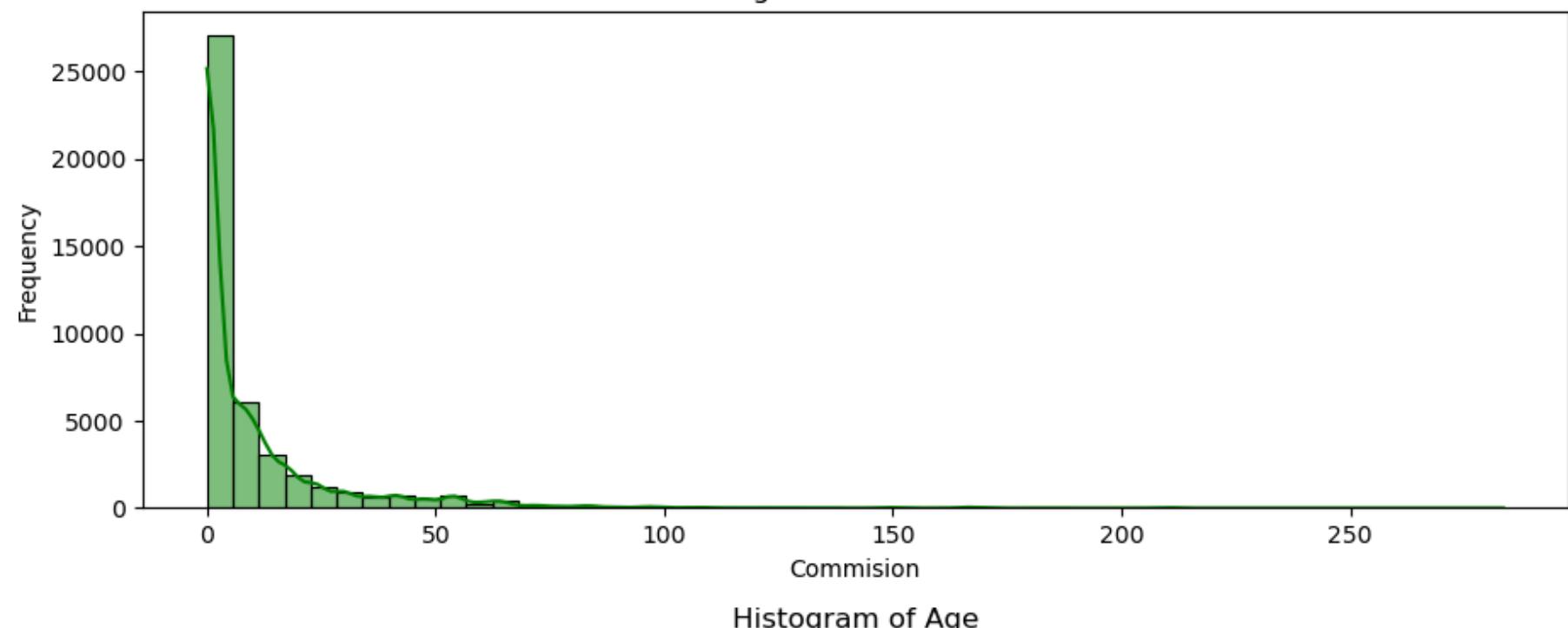
	Agency	Agency Type	Distribution Channel	Product Name	Gender	Destination	Claim
count	44328	44328	44328	44328	12681	44328	44328
unique	16	2	2	26	2	138	2
top	EPX	Travel Agency	Online	Cancellation Plan	M	SINGAPORE	No
freq	24656	32113	43572	12979	6504	9267	43651



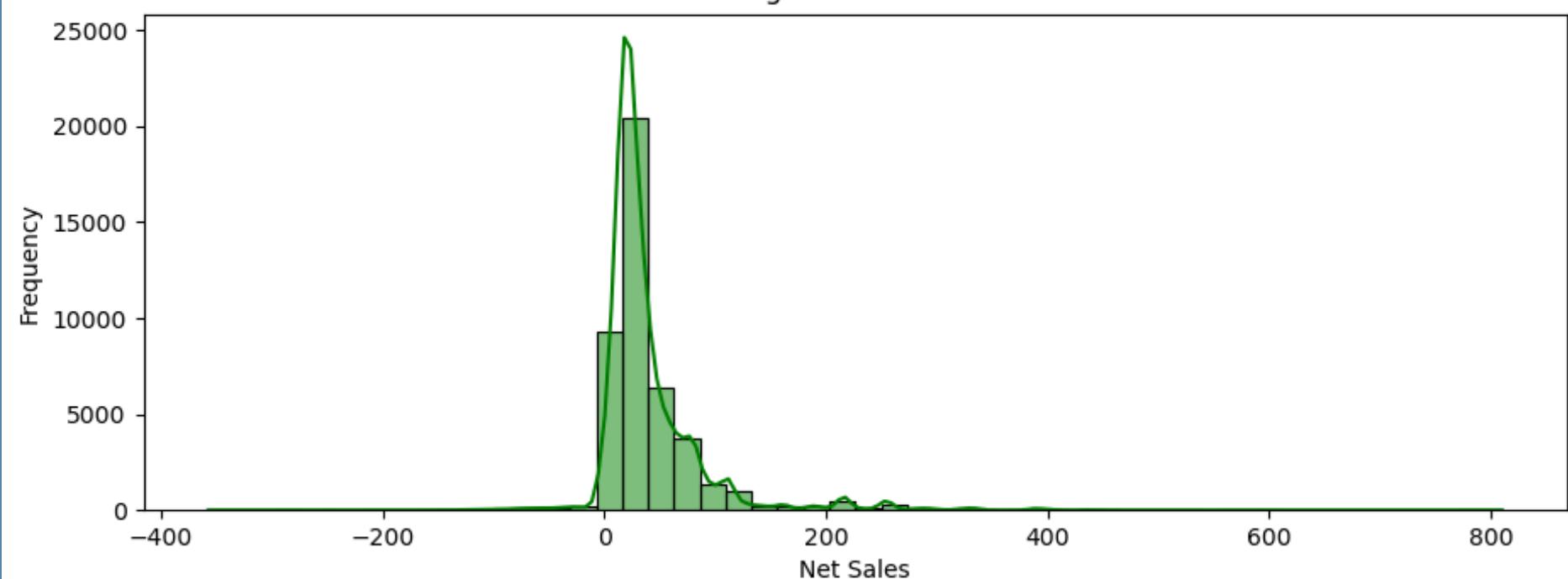
Histogram of Duration



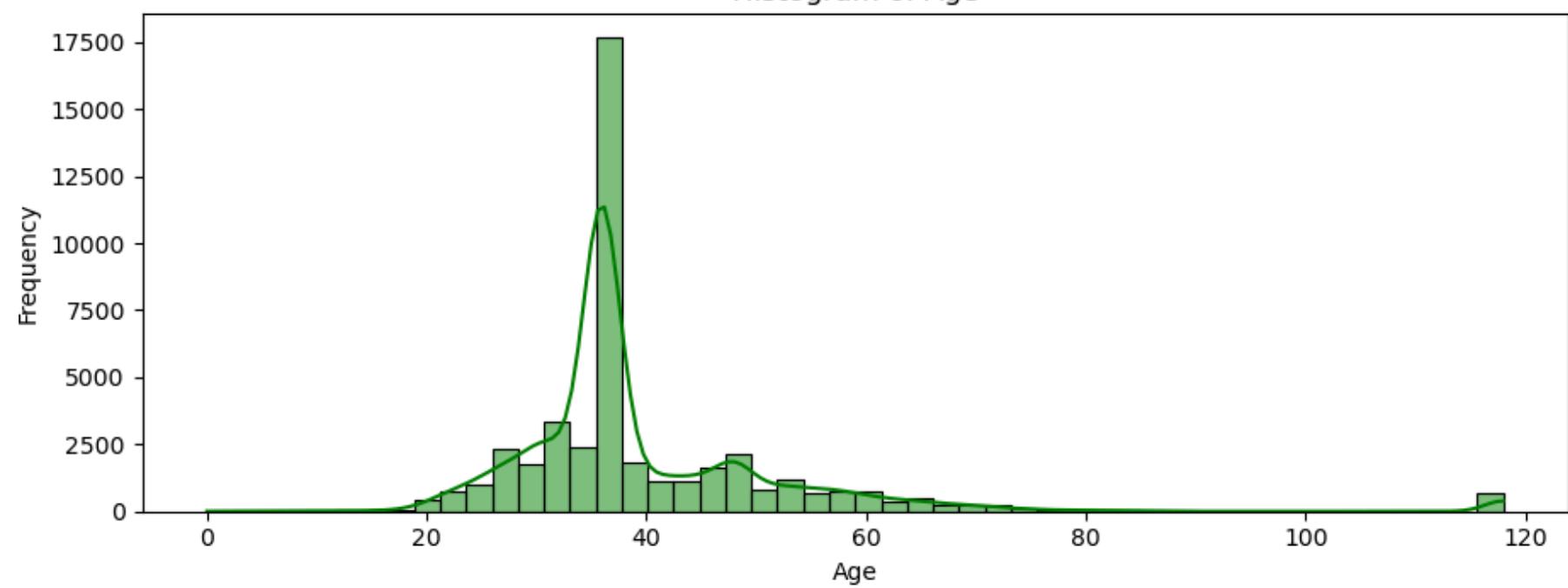
Histogram of Commision



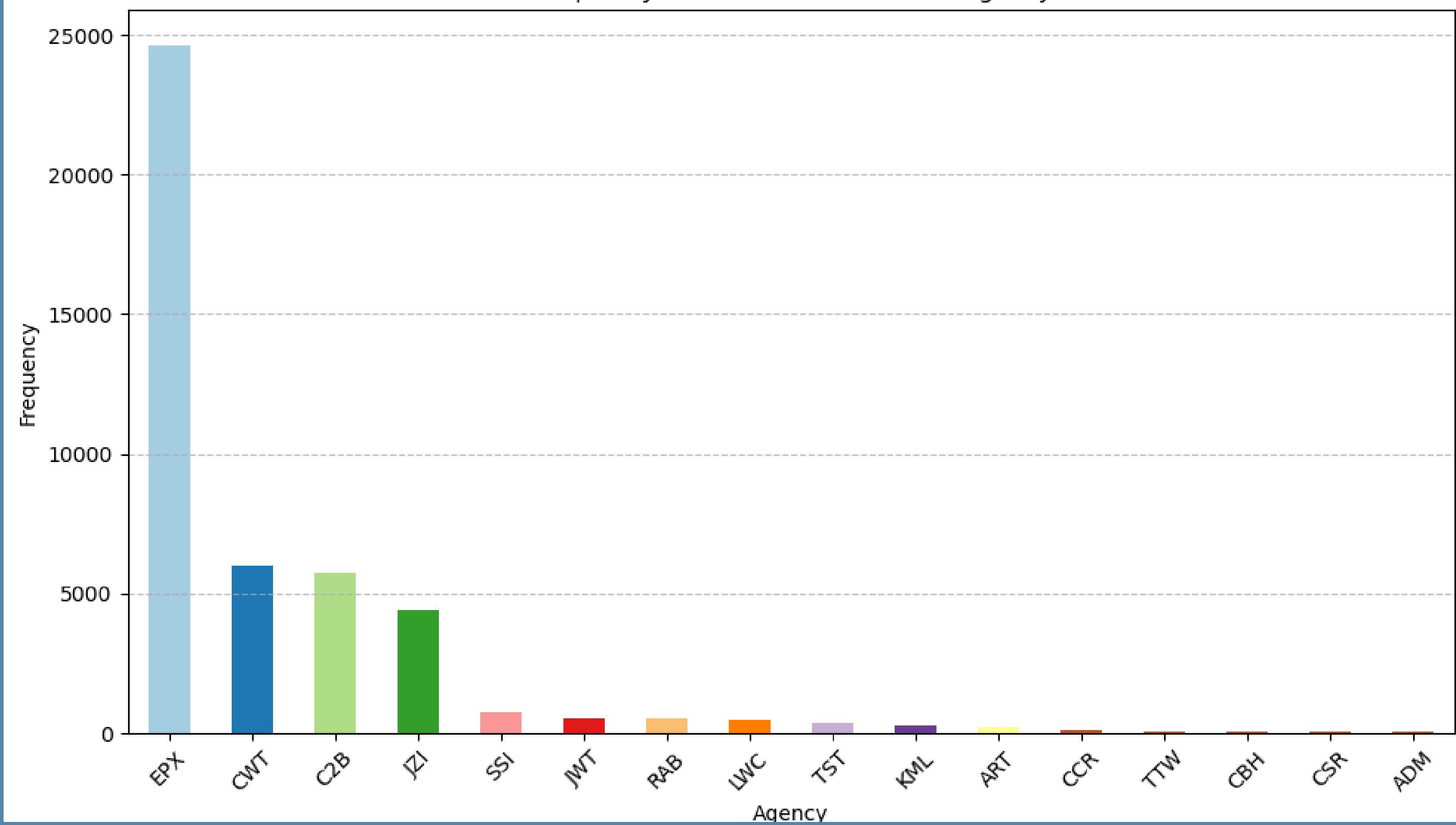
Histogram of Net Sales



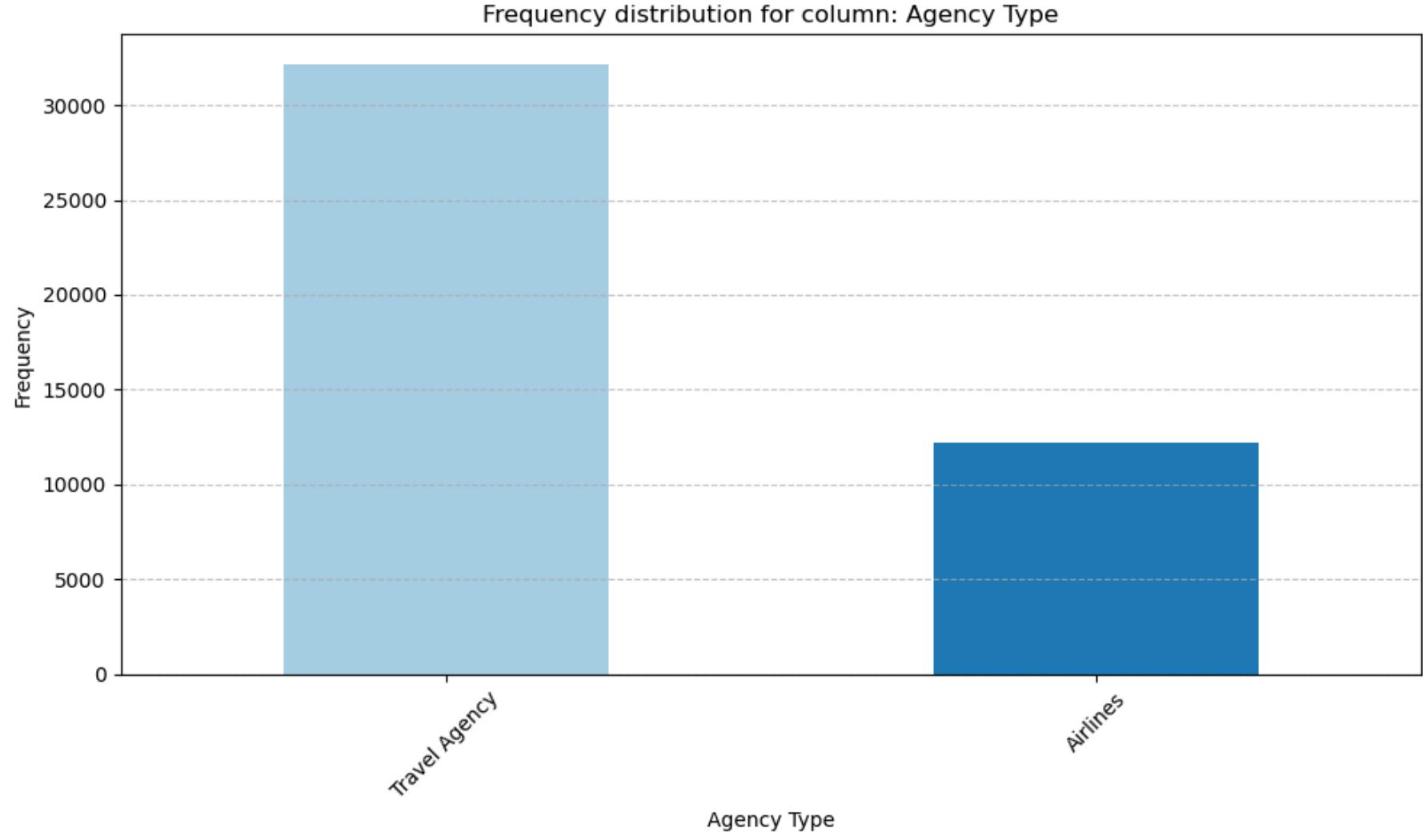
Histogram of Commision



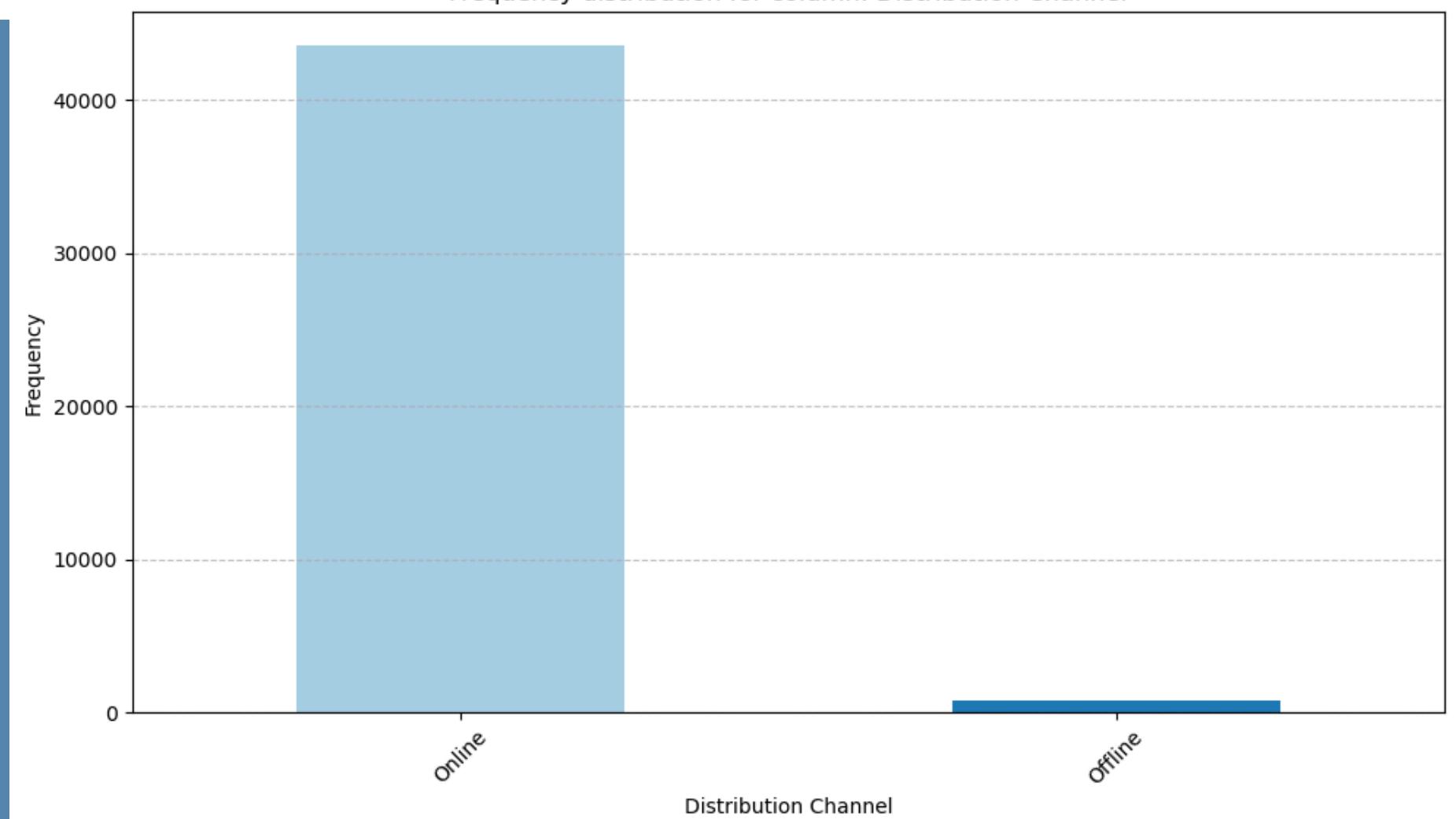
Frequency distribution for column: Agency



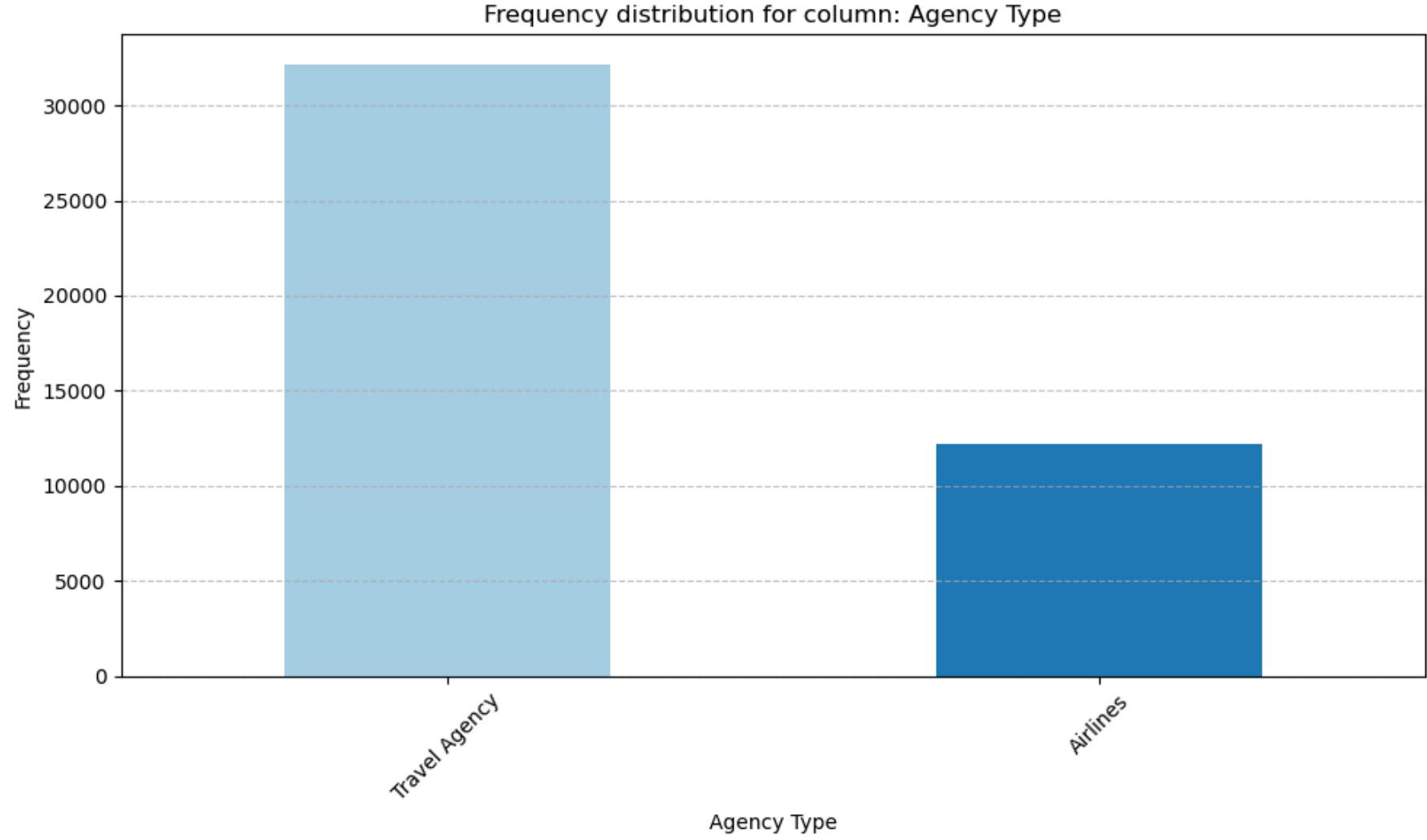
Frequency distribution for column: Agency Type



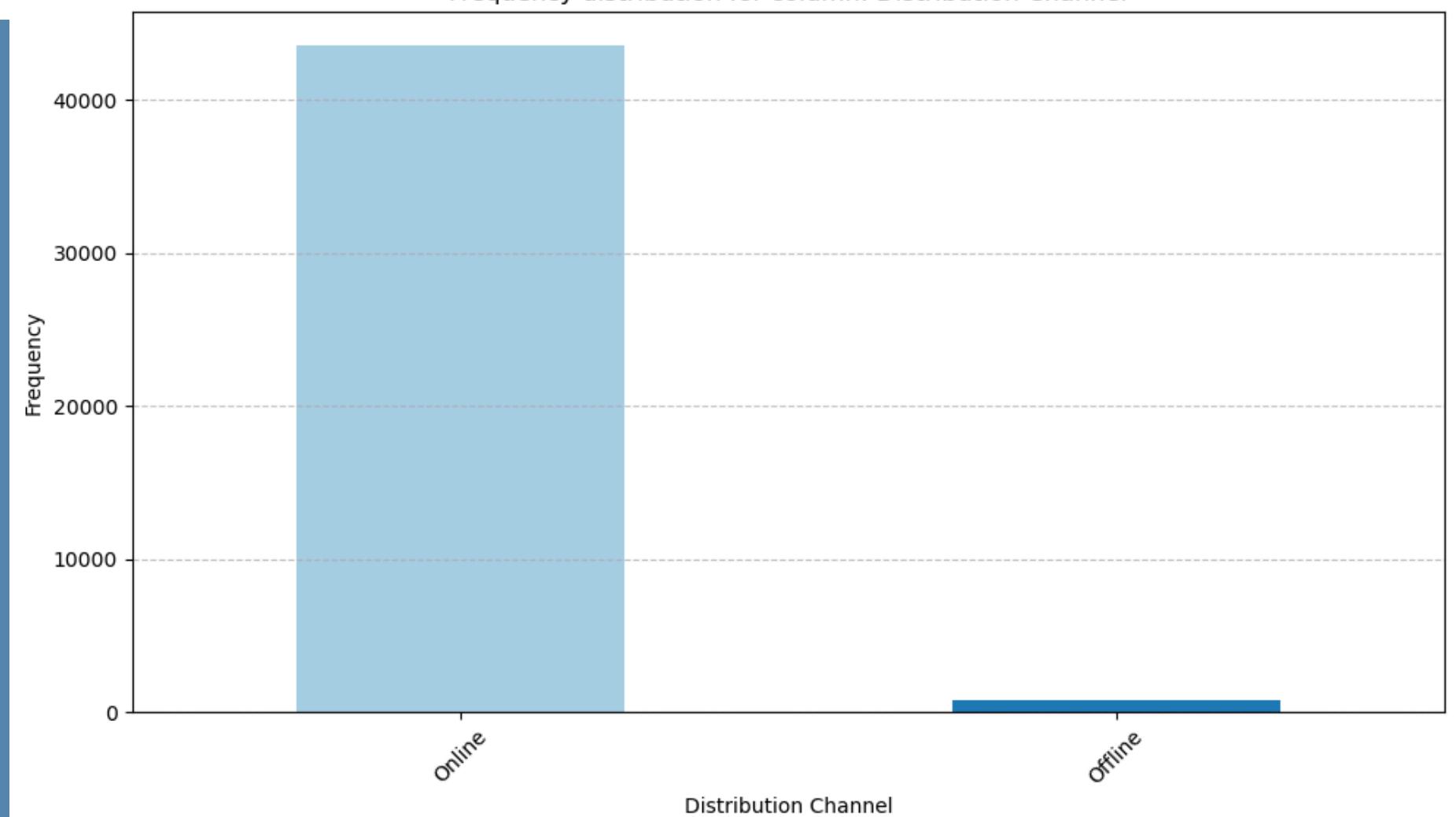
Frequency distribution for column: Distribution Channel



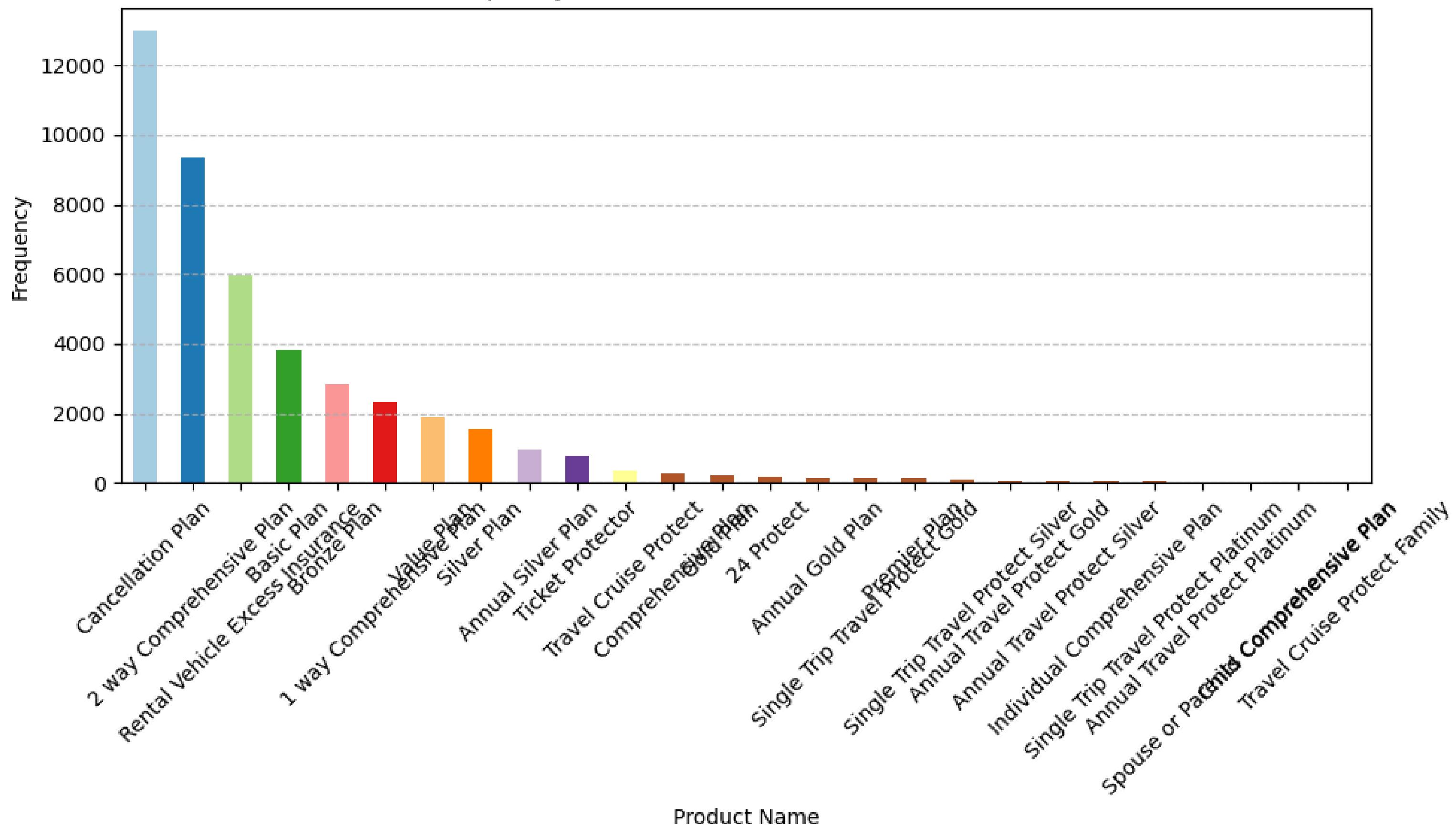
Frequency distribution for column: Agency Type

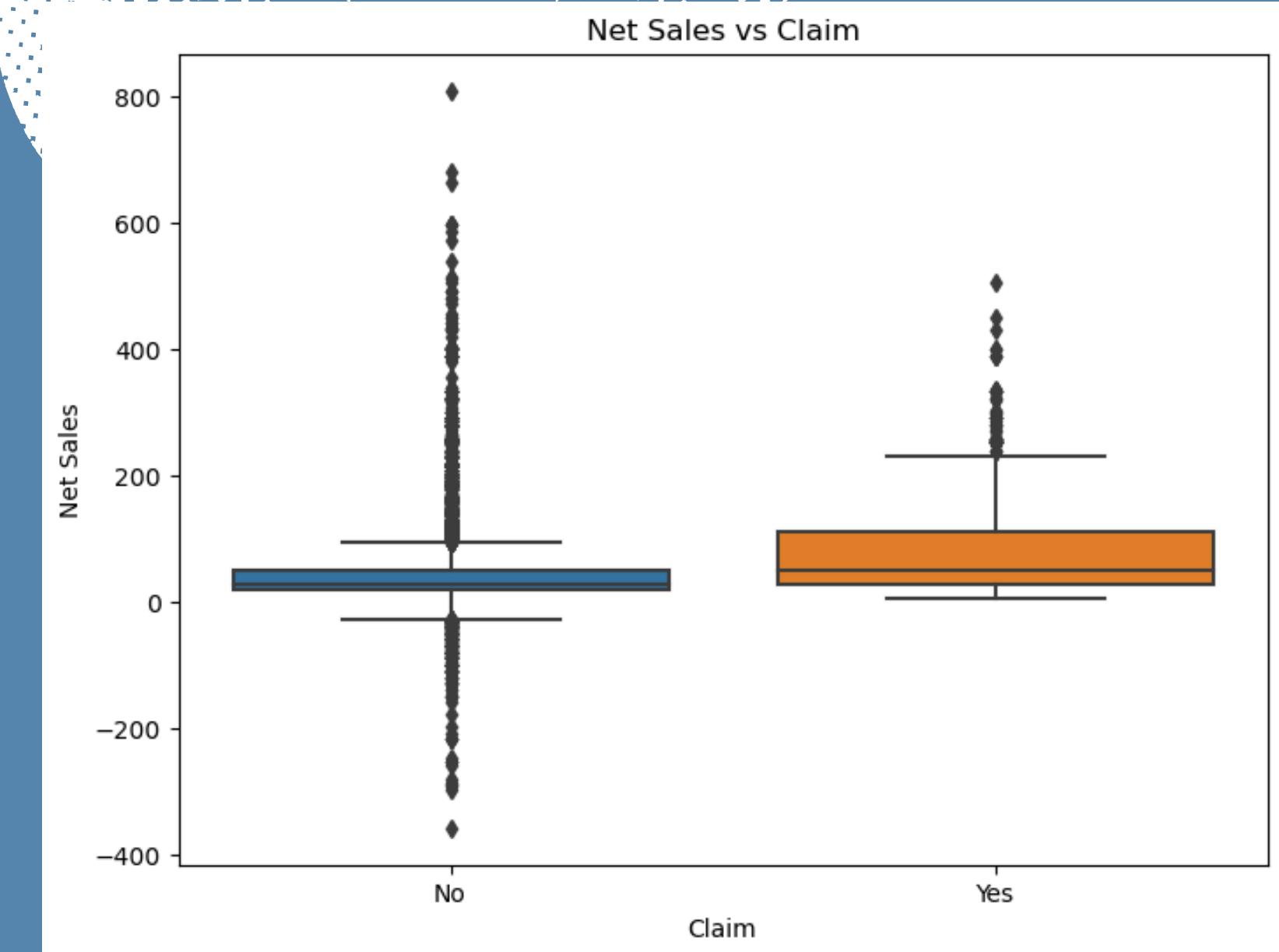
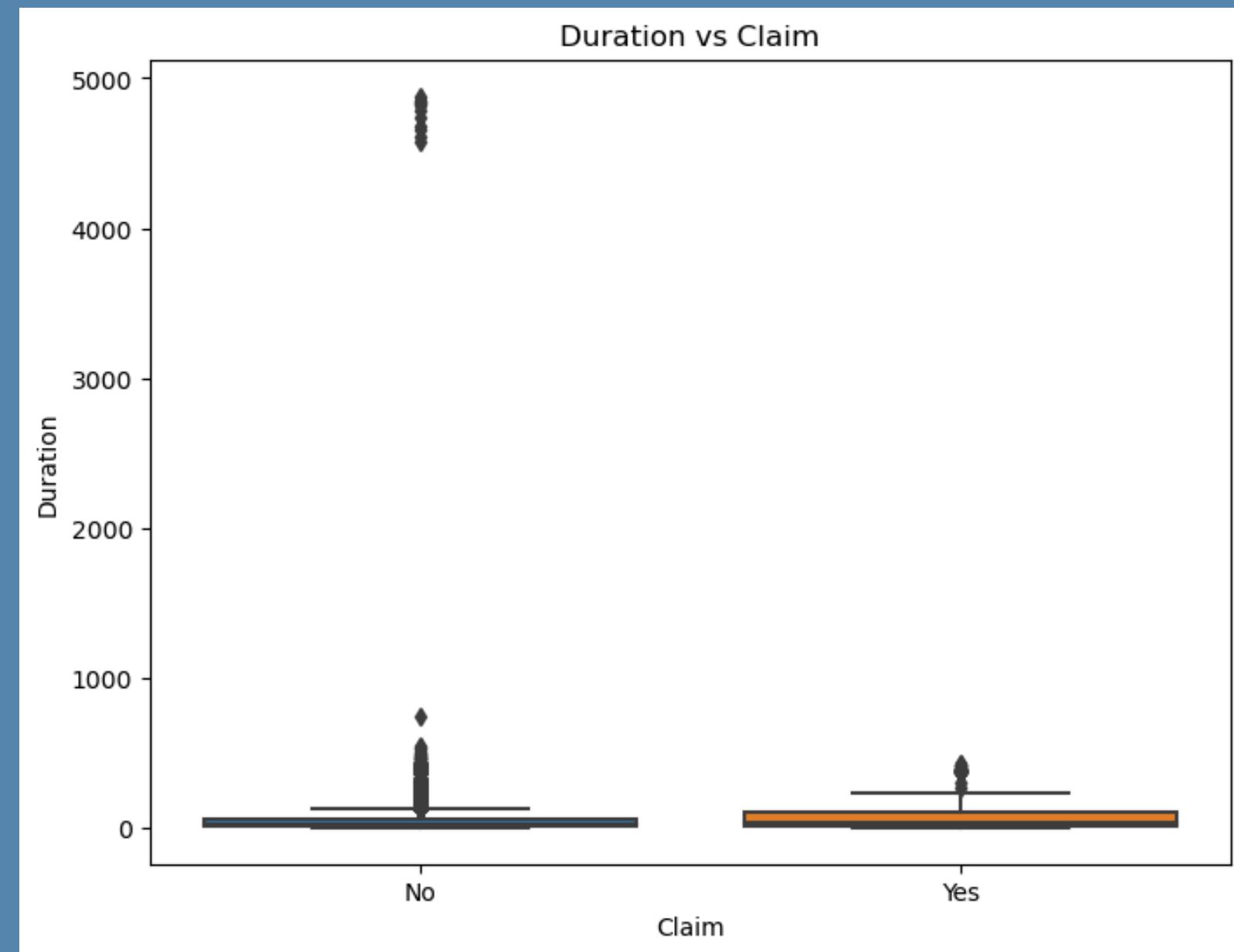


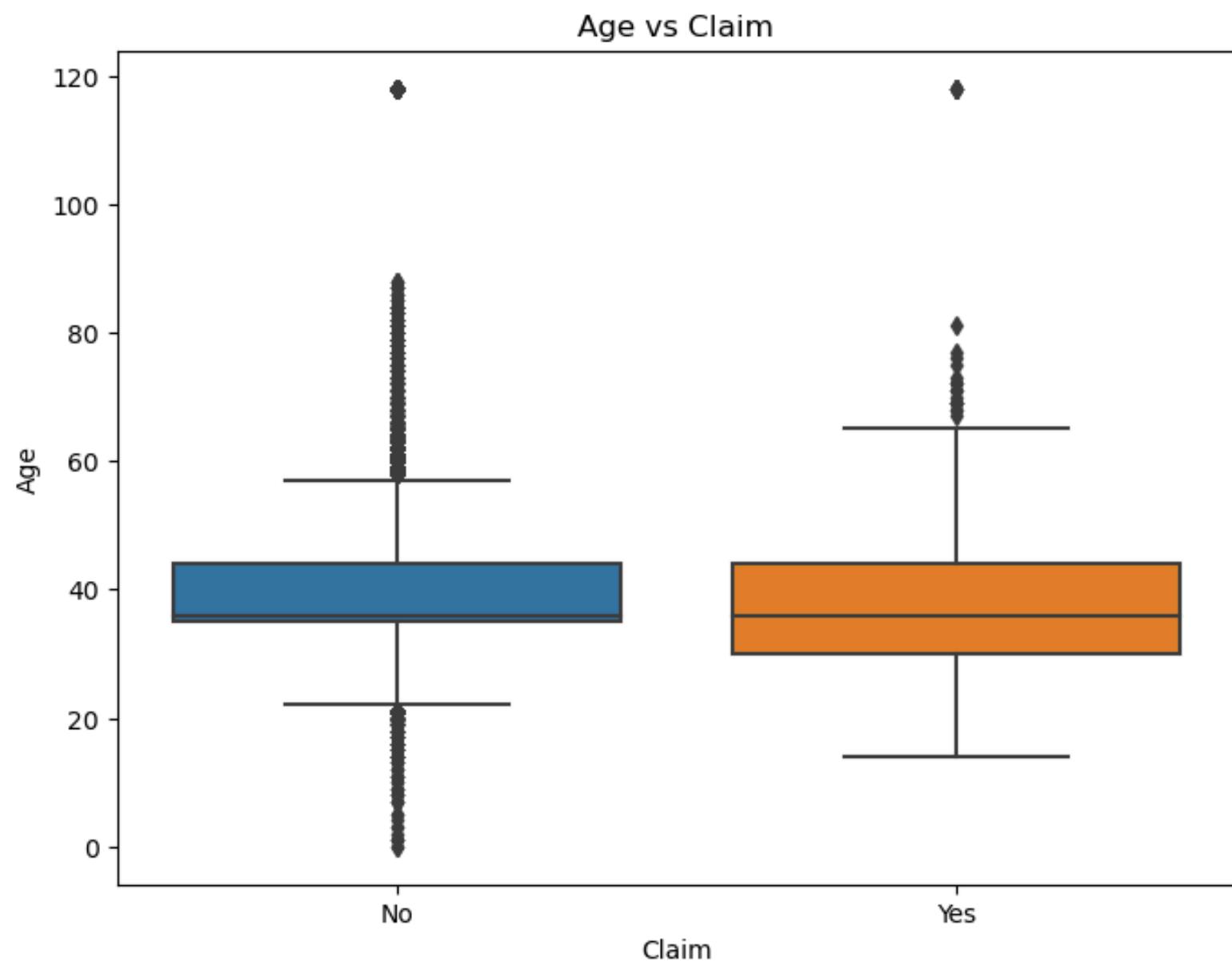
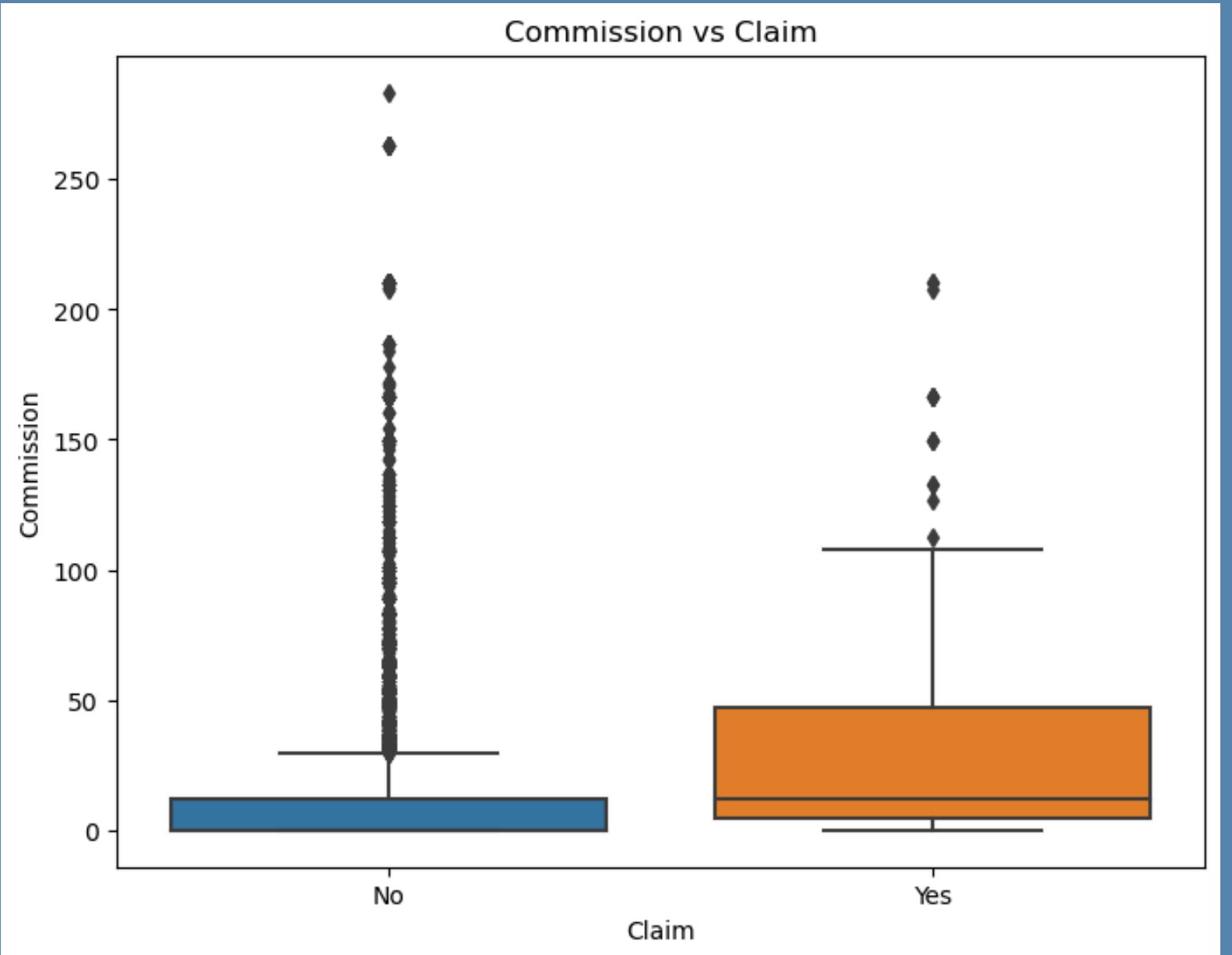
Frequency distribution for column: Distribution Channel



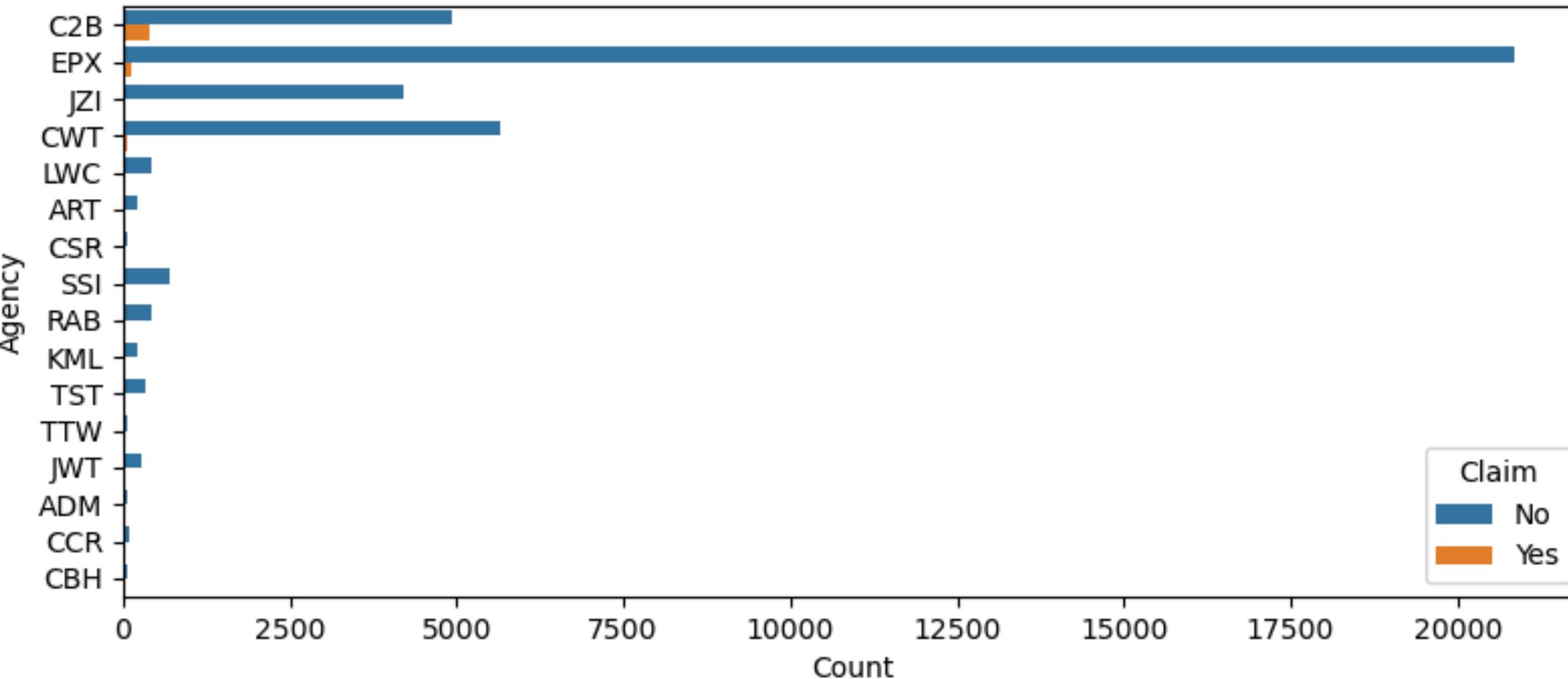
Frequency distribution for column: Product Name



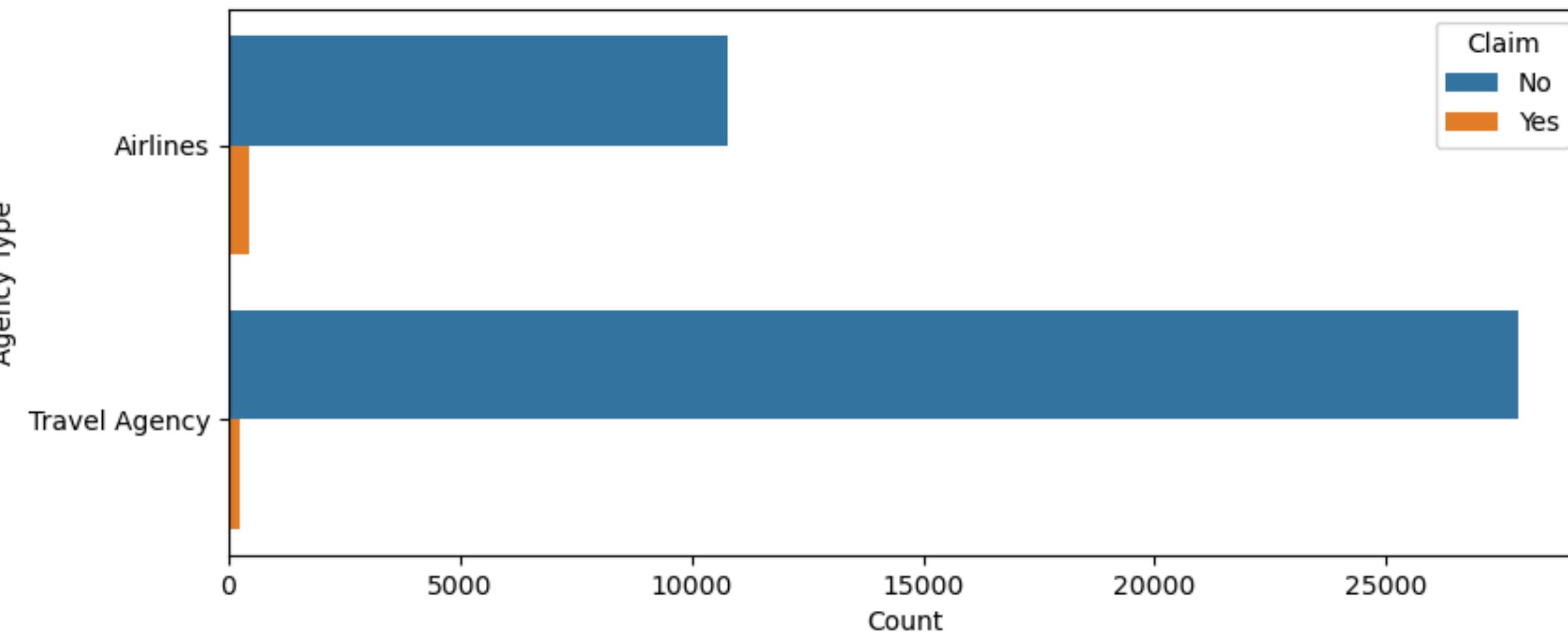




Agency Distribution with Claim Status

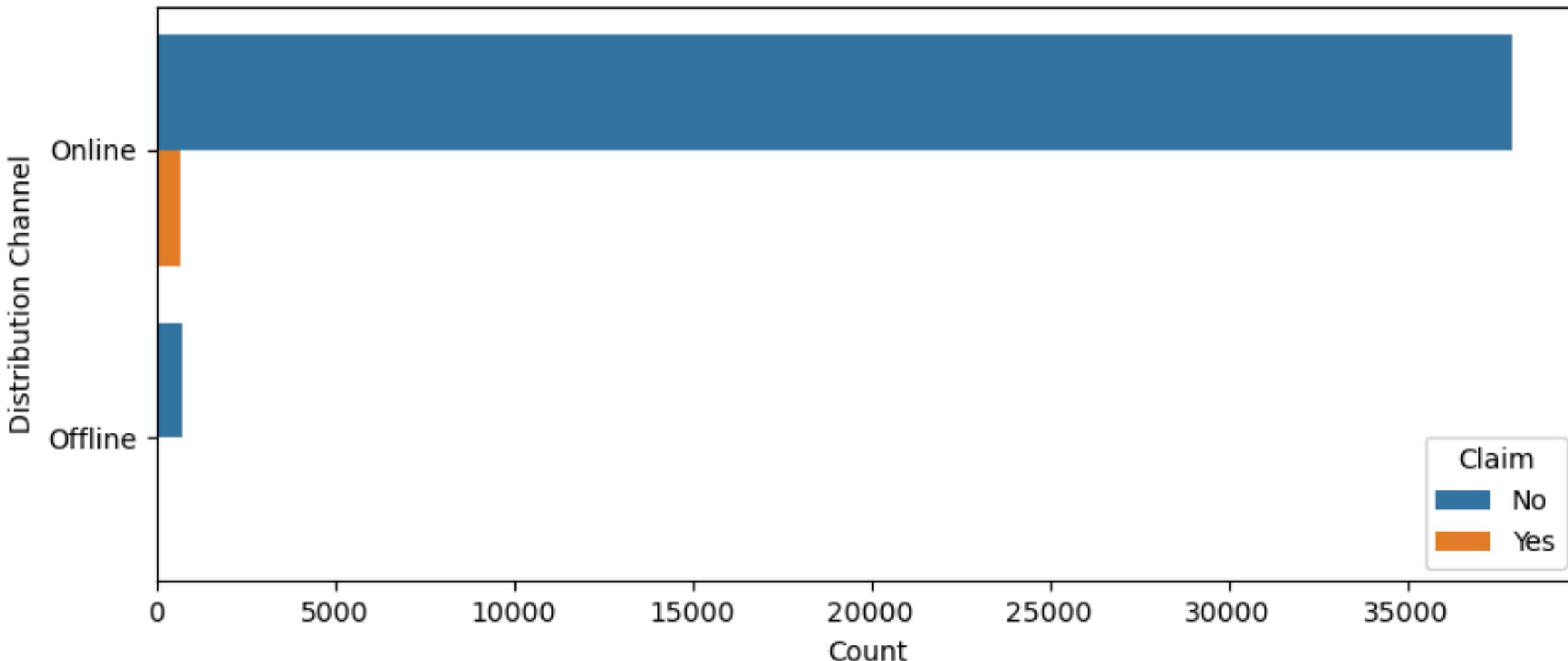


Agency Type Distribution with Claim Status

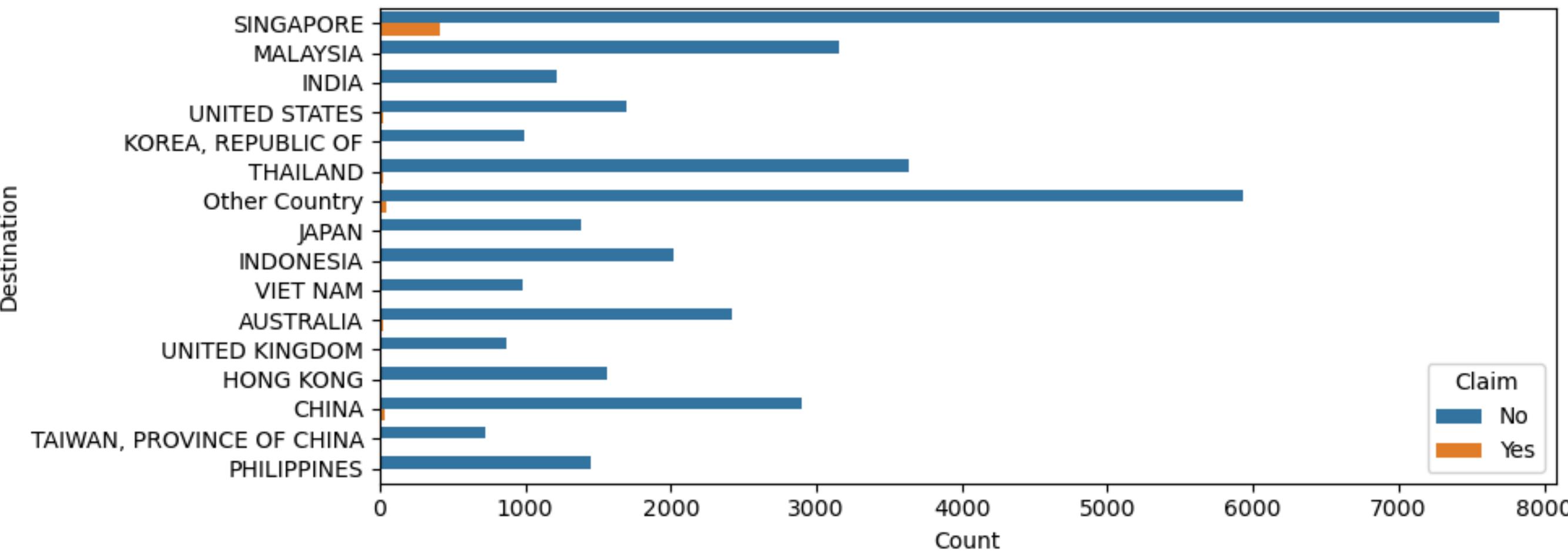


Distribution Channel Distribution with Claim Status

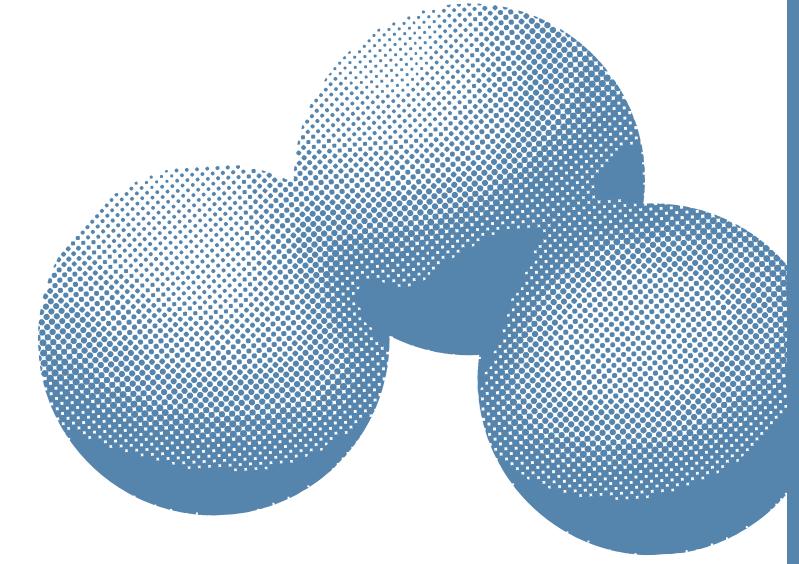
Distribution Channel Distribution with Claim Status



Destination Distribution with Claim Status



Feature Engineering



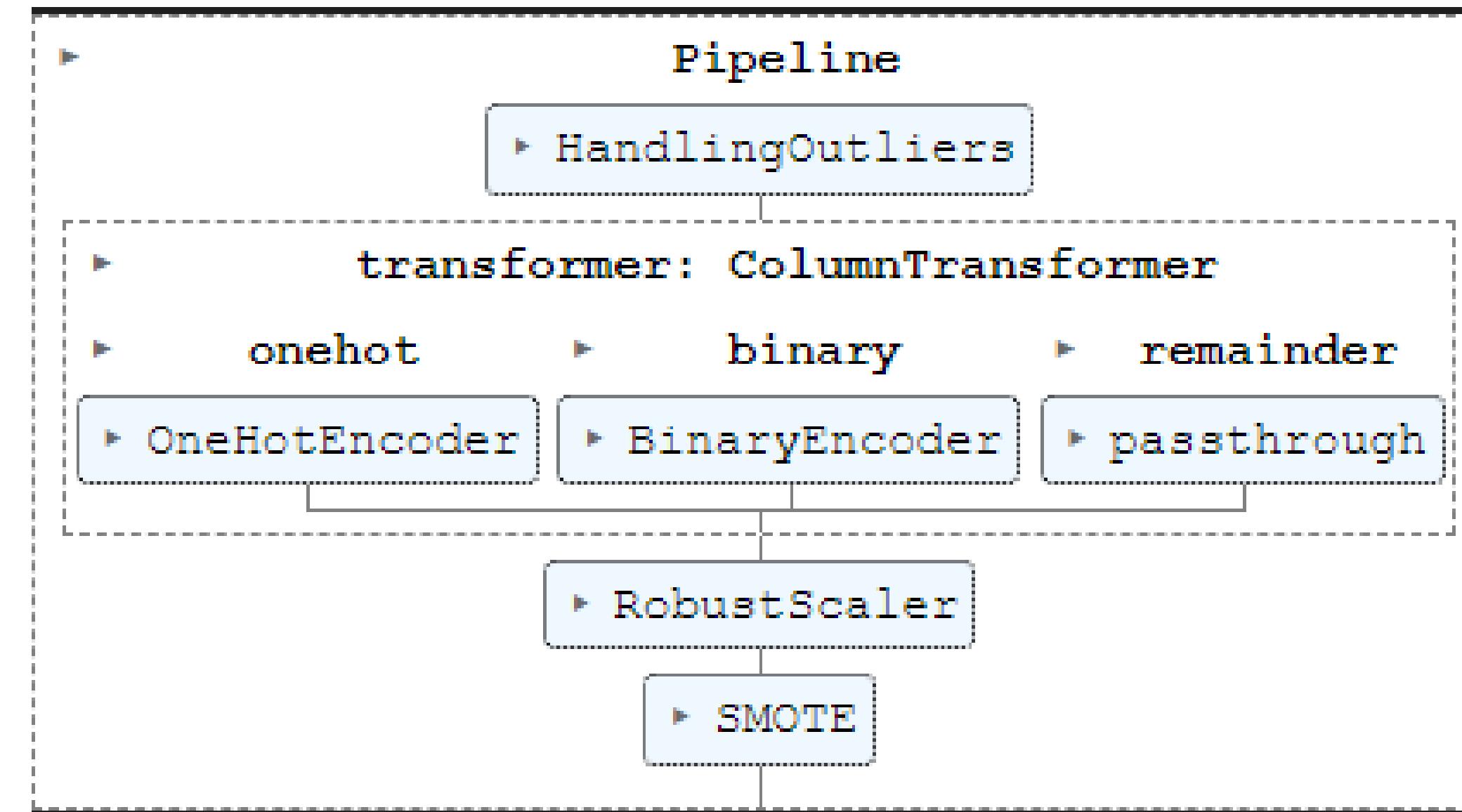
- Kolom Gender di Drop karena Missing Value nya 71%
- Data Duplicate di drop sebanyak 5000 Data
- Negative Value Duration di drop
- Destination menggunakan 15 data terbanyak, sisanya dikelompokkan
- Cardinality pada Agency, Destination, Product Name
- Imbalance Data pada tingkat mild pada target

Data Preprocessing

- Encoding pada feature kategorikal dengan One Hot Encoding dan Binary Encoding
- Membuat Feature dan Target
- Handling Outlier dengan Windsorizing
- Scaling dengan Robust Scaler
- Resampling Menggunakan SMOTE
- Dimasukkan kedalam Pipeline



Data Preprocessing



Modelling

Logistic Regression, KNN Classifier,
Decision Tree Classifier, Random Forest
Classifier, Adaboost Classifier, Gradient
Boost Classifier, dan XGBoost Classifier

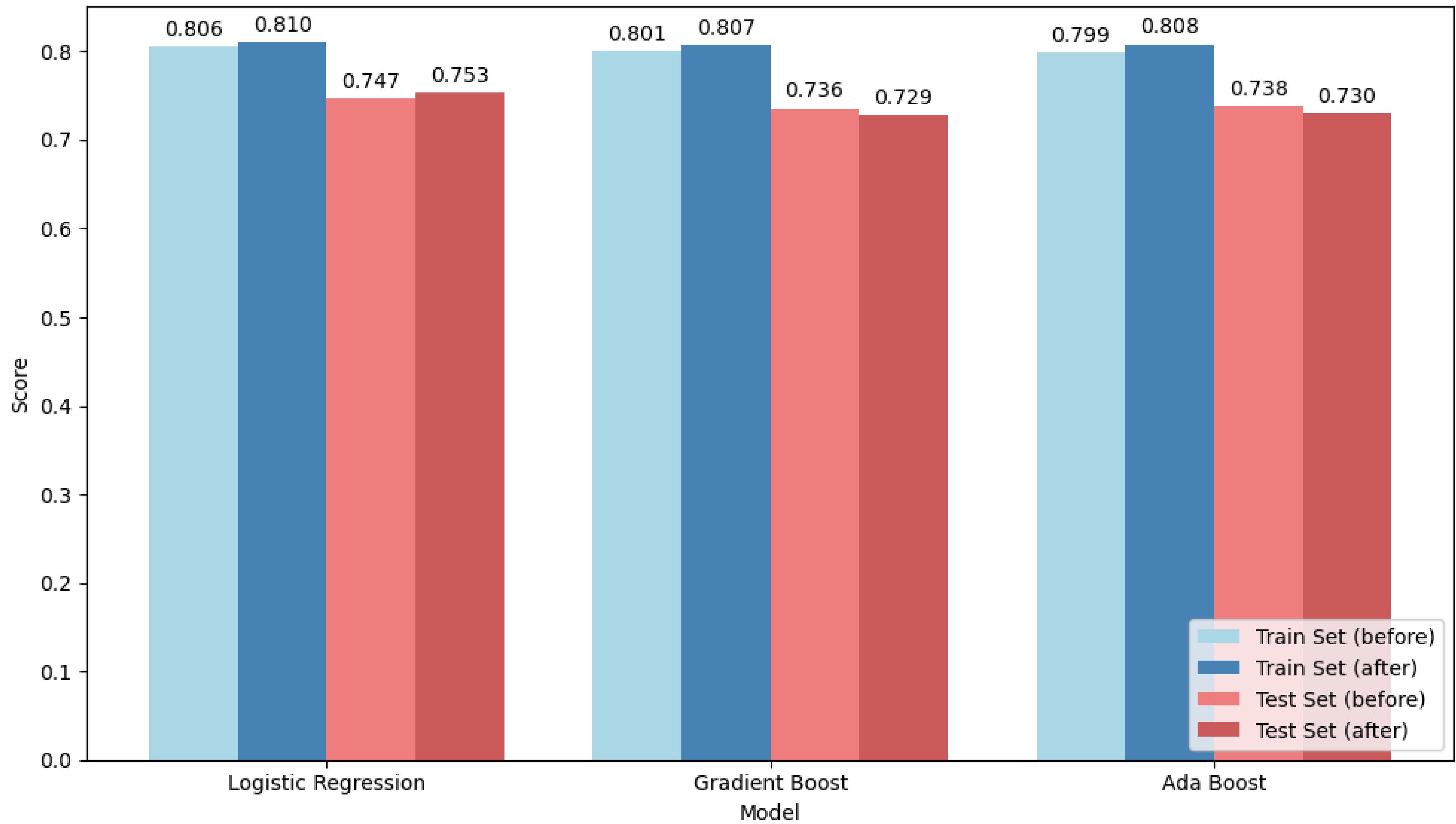
Cross Validation

```
for algoritma in models:  
  
    model_pipe = Pipeline([  
        ('outlier', HandlingOutliers()),  
        ('preprocessing', transformer),  
        ('scaler', scaler),  
        ('resampler', smote),  
        ('model', algoritma)  
    ])  
  
    skfold = StratifiedKFold(n_splits = 5)  
  
    model_cv = cross_val_score(  
        model_pipe,  
        X_train,  
        y_train,  
        cv = skfold,  
        scoring = 'roc_auc',  
        error_score='raise'  
    )  
  
    # model yang sudah dimasukkan ke dalam pipeline  
    # data sebelum di preprocessing
```

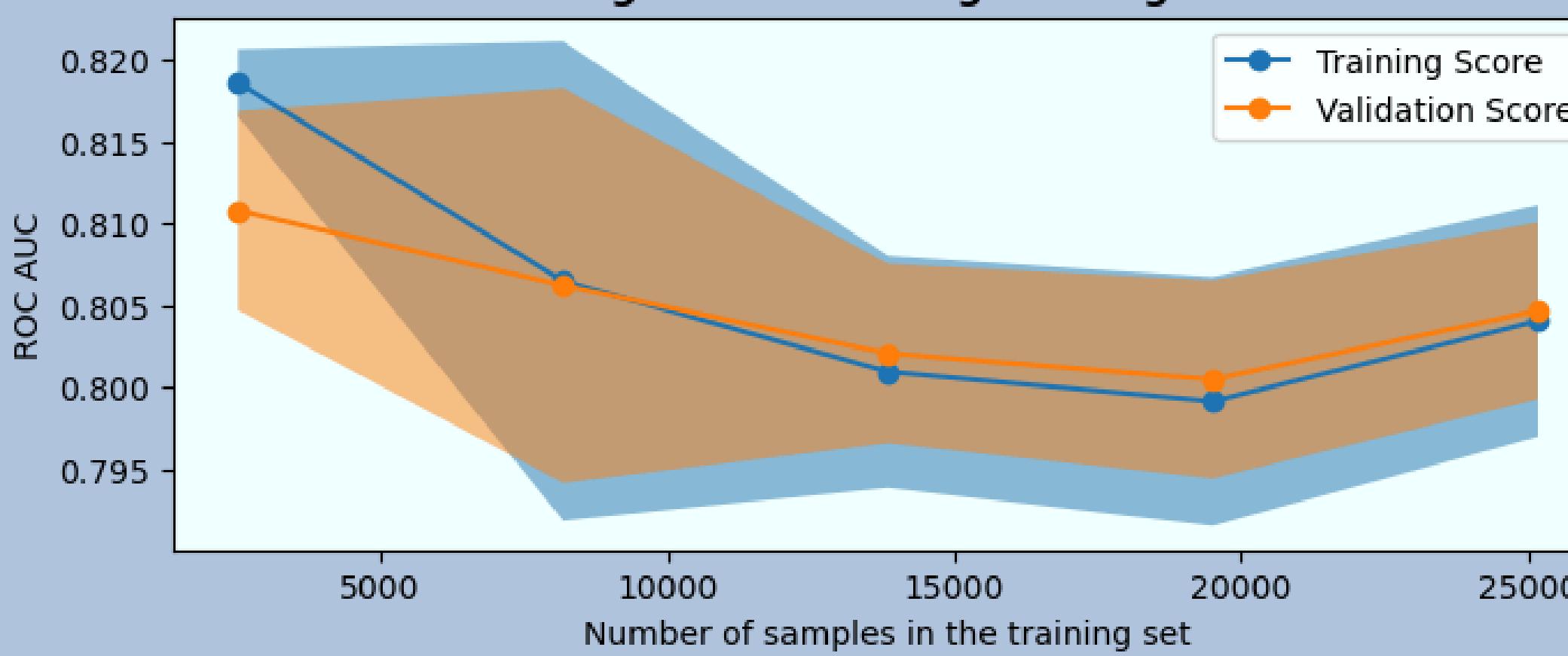
	model	mean	std	all_score
0	LogisticRegression	0.806732	0.023814	[0.8074, 0.8196, 0.7732, 0.8426, 0.7907]
5	GradientBoostingClassifier	0.801139	0.013131	[0.7953, 0.8109, 0.7856, 0.8215, 0.7924]
4	AdaBoostClassifier	0.799833	0.025325	[0.7995, 0.8211, 0.7607, 0.8319, 0.7859]
6	XGBClassifier	0.759585	0.011800	[0.755, 0.763, 0.7631, 0.7764, 0.7404]
3	RandomForestClassifier	0.735154	0.020823	[0.7663, 0.7296, 0.7474, 0.704, 0.7284]
1	KNeighborsClassifier	0.683993	0.011265	[0.7029, 0.6825, 0.6695, 0.6881, 0.677]
2	DecisionTreeClassifier	0.546512	0.006774	[0.5496, 0.5412, 0.5581, 0.5446, 0.5391]

	model	accuracy (test_set)
0	LogisticRegression	0.747032
4	AdaBoostClassifier	0.738925
5	GradientBoostingClassifier	0.736076
6	XGBClassifier	0.618830
1	KNeighborsClassifier	0.595237
3	RandomForestClassifier	0.579673
2	DecisionTreeClassifier	0.527045

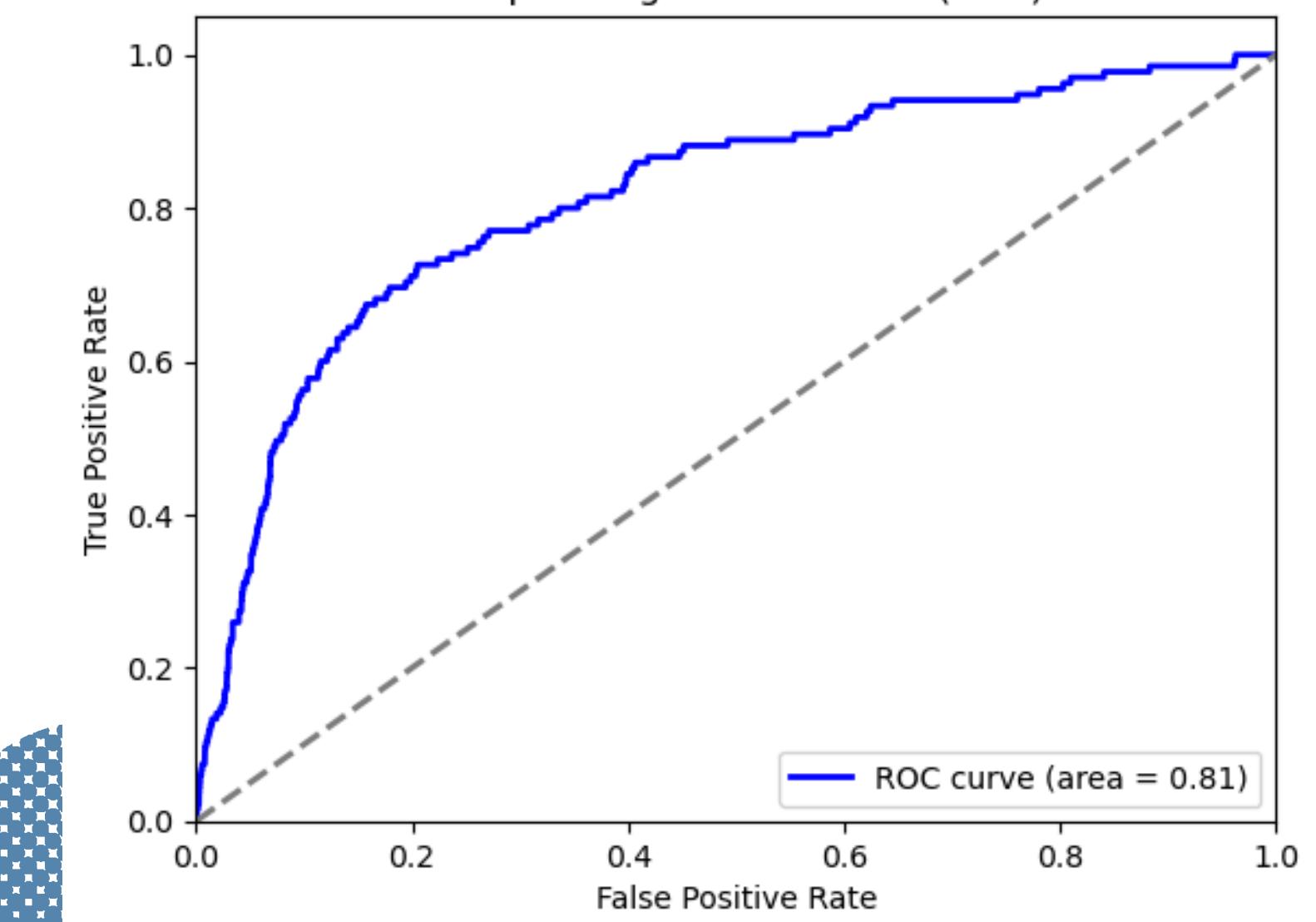
Comparison of Scores for Different Models



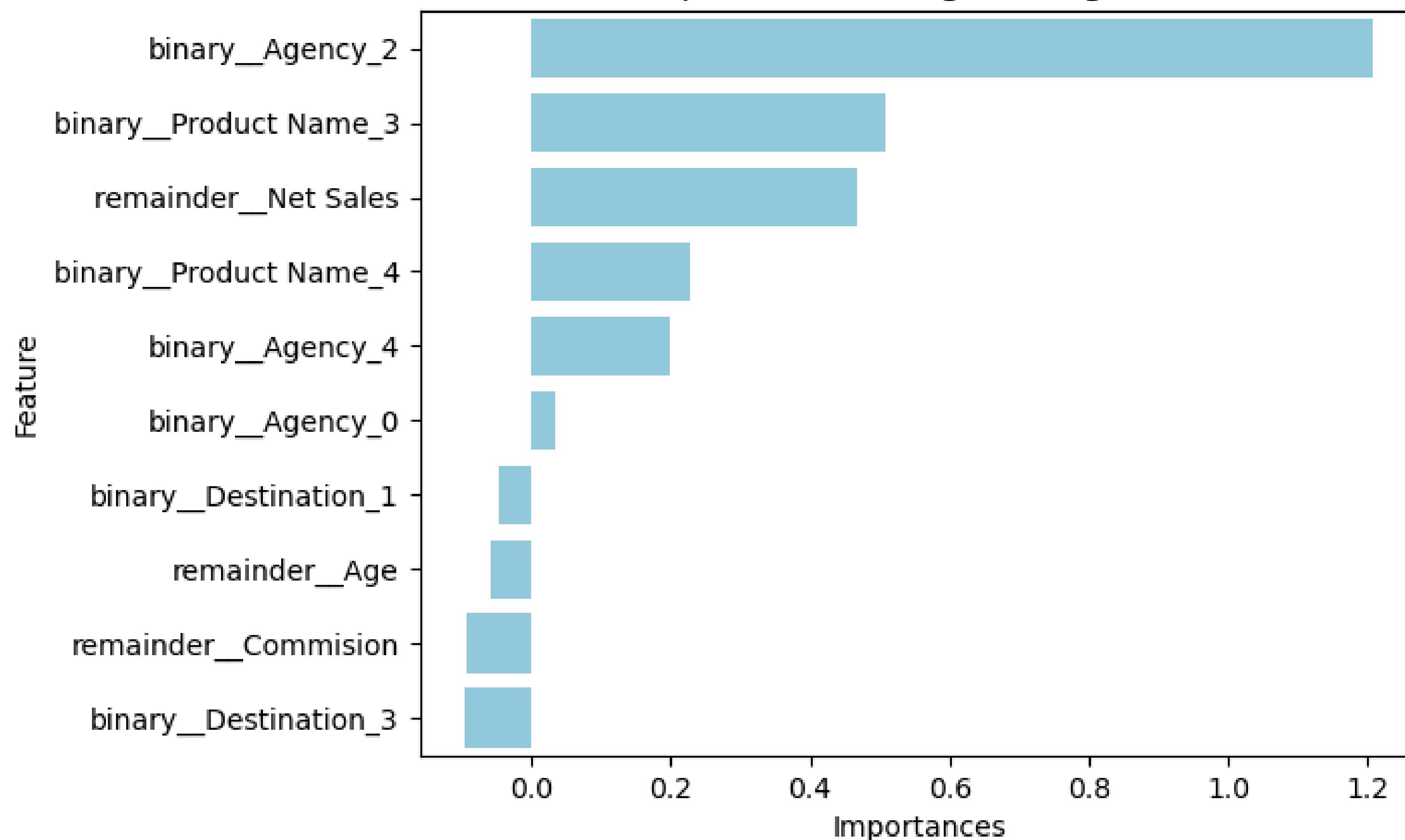
Learning Curve from Logistic Regression



Receiver Operating Characteristic (ROC) Curve

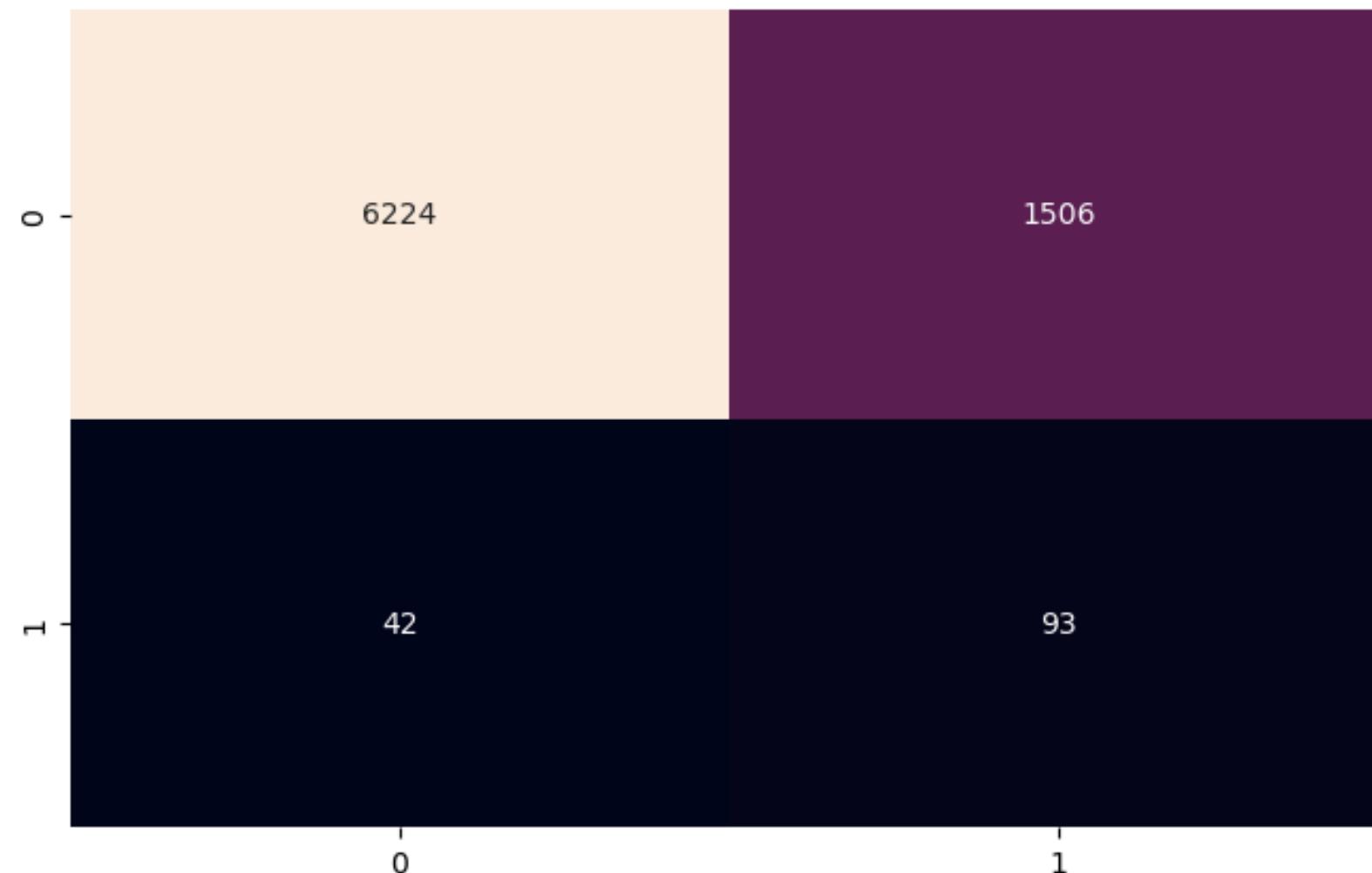


Feature Importances of Logistic Regression model

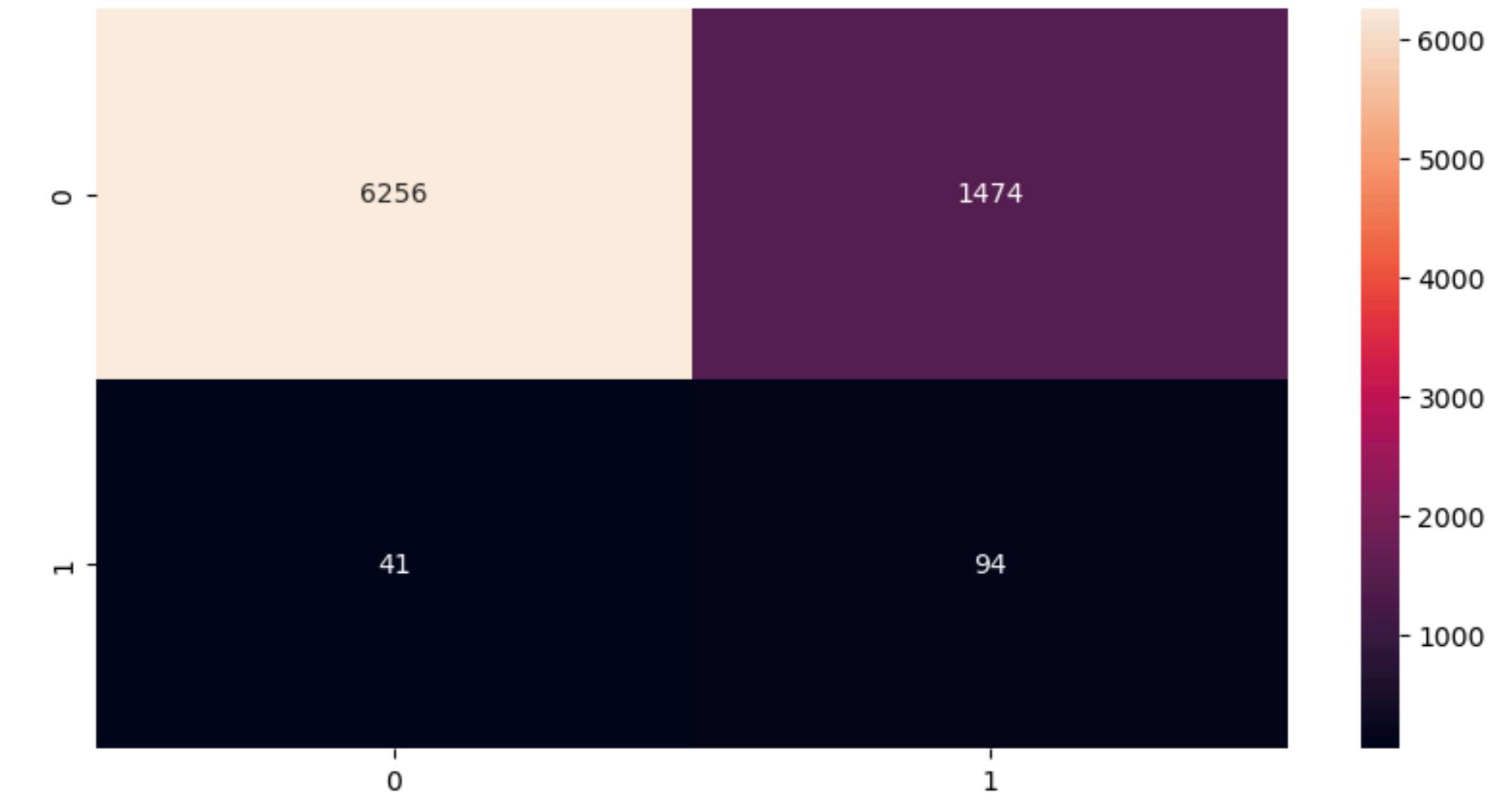


Confusion Matrix

Sebelum



Sesudah



Cost Benefit Analysis

Faktor	Tanpa ML	Dengan ML
Biaya Awal	Rendah (terutama gaji)	Tinggi (infrastruktur, gaji, pelatihan)
Biaya Berkelanjutan	Tinggi (gaji, biaya kesalahan, pelatihan)	Sedang (pemeliharaan, pembaruan)
Akurasi	Lebih rendah (kesalahan manusia)	Lebih tinggi (presisi berbasis data)
Kecepatan Pemrosesan	Lebih lambat (proses manual)	Lebih cepat (proses otomatis)
Skalabilitas	Terbatas oleh sumber daya manusia	Skalabilitas tinggi
Kepuasan Pelanggan	Potensial lebih rendah (waktu pemrosesan lama)	Lebih tinggi (klaim cepat dan akurat)
Deteksi Penipuan	Kurang efektif	Lebih efektif

Conclusion

1. **Logistic Regression** adalah model **terbaik** untuk dataset ini karena performa yang cukup stabil dan baik pada dataset yang besar.
2. Hasil ROC AUC
 - **Train score sebelum tuning:** 75.66%
 - **Test score sebelum tuning:** 74.70%
 - **Train score setelah tuning:** 75.91%
 - **Test score setelah tuning:** 75.28%

3. Setelah tuning, model menunjukkan **peningkatan dalam recall, precision, dan accuracy.** Ini menunjukkan bahwa **tuning membantu model menjadi lebih baik** dalam mengidentifikasi klaim yang benar sambil mengurangi kesalahan dalam prediksi klaim.
4. Berdasarkan **feature importance** didapati bahwa seseorang akan klaim atau tidak klaim dapat dipengaruhi dimana pemegang polis tersebut membeli asuransi perjalanan dari suatu **agen asuransi.**

Recommendation Model

1. penggunaan hyperparameter tuning yang lebih luas akan membantu
2. Feature Engineering yang lebih dalam
3. Penggunaan metriks lain dengan penyesuaian problem statement
4. Mencoba menggunakan resampler lain seperti RandomOverSampling, SMOTENC, SMOTE-ENN

Recommendation Business

1. Akurasi yang Lebih Tinggi
2. Kecepatan Pemrosesan yang Lebih Tinggi
3. Penghematan Biaya
4. Keputusan yang Lebih Tepat
5. Peningkatan Kepuasan Pelanggan

Terima
Kasih