

```

1
2 import pandas as pd
3 from datetime import datetime
4 from sklearn.model_selection import train_test_split
5 from xgboost import XGBRegressor
6 from sklearn.metrics import mean_absolute_error, mean_squared_error
7 import numpy as np
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10
11
12

```

```

1 # Step 2: Load Dataset
2 df = pd.read_excel("/content/Online Retail.xlsx")
3 df.head()
4
5
6


```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom
1	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom
3	536365	810000	KNITTED UNION FLAG HOT WATER	6	2010-12-01	2.80	17850.0	United Kingdom

```

1 df.dropna(subset=['CustomerID'], inplace=True)
2 df = df[~df['InvoiceNo'].astype(str).str.startswith('C')]
3 df['TotalPrice'] = df['Quantity'] * df['UnitPrice']

```

 <ipython-input-5-4b292717a725>:3: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)  
df['TotalPrice'] = df['Quantity'] \* df['UnitPrice']

```


1 snapshot = df['InvoiceDate'].max() + pd.Timedelta(days=1)
2 rfm = df.groupby('CustomerID').agg({
3     'InvoiceDate': lambda x: (snapshot - x.max()).days,
4     'InvoiceNo': 'nunique',
5     'TotalPrice': 'sum'
6 }).rename(columns={'InvoiceDate': 'Recency', 'InvoiceNo': 'Frequency', 'TotalPrice': 'Monetary'})
7
8 rfm['AOV'] = rfm['Monetary'] / rfm['Frequency']
9

```

```

1 X = rfm[['Recency', 'Frequency', 'AOV']]
2 y = rfm['Monetary']
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
5 model = XGBRegressor()
6 model.fit(X_train, y_train)
7
8 # Step 6: Evaluation
9 y_pred = model.predict(X_test)
10 mae = mean_absolute_error(y_test, y_pred)
11 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
12 print("MAE:", mae)
13 print("RMSE:", rmse)

```

 MAE: 208.95508759955652  
RMSE: 1863.0781089577395

```

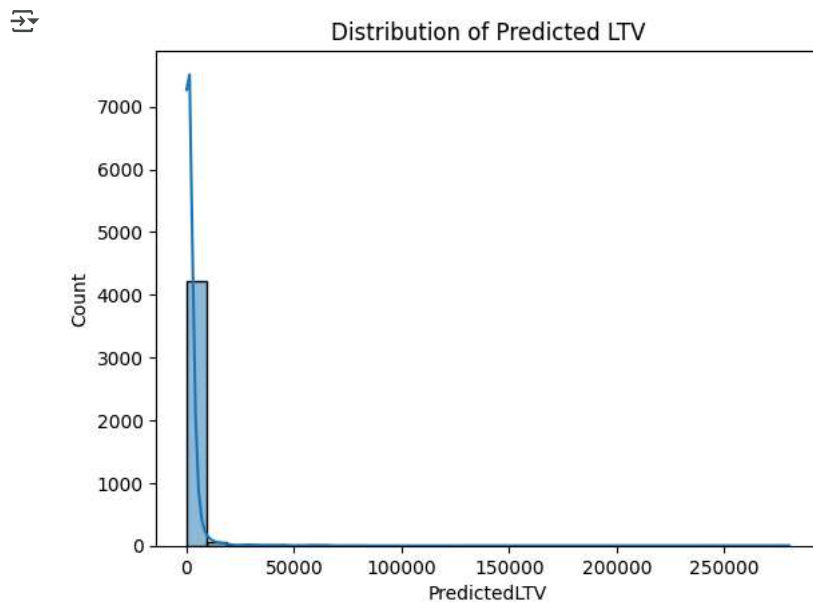
1 rfm['PredictedLTV'] = model.predict(X)
2 rfm['Segment'] = pd.qcut(rfm['PredictedLTV'], 4, labels=['Low', 'Medium', 'High', 'Very High'])
3

```

```

1 sns.histplot(rfm['PredictedLTV'], bins=30, kde=True)
2 plt.title("Distribution of Predicted LTV")
3 plt.show()

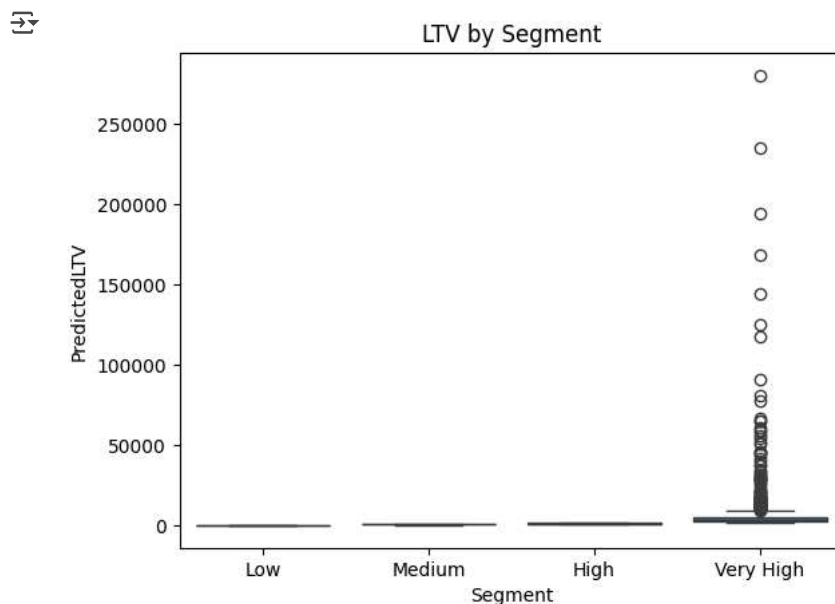
```



```

1 sns.boxplot(x='Segment', y='PredictedLTV', data=rfm)
2 plt.title("LTV by Segment")
3 plt.show()
4

```



```

1 rfm.reset_index()[['CustomerID', 'PredictedLTV', 'Segment']].to_csv("ltv_predictions.csv", index=False)

```

