

```
1 !pip install snowflake-connector-python pandas scikit-learn xgboost
```

Show hidden output

```
1 import snowflake.connector
2 import pandas as pd
3 import numpy as np
```

```
1 conn = snowflake.connector.connect (
2     user = "ARINDHAMFLAKE",
3     password = "Dataanalyst2025",
4     account = "ZIEPNFM-WH94884",
5     warehouse = "GPSFLEET_WH",
6     database = "FLEET_DB",
7     schema = "ANALYTICS"
8 )
```

```
1 df_trip = pd.read_sql("SELECT * FROM trip_summary;", conn)
2 df_driver = pd.read_sql("SELECT * FROM driver_daily;", conn)
3 print(df_trip.shape, df_driver.shape)
```

```
/tmp/ipython-input-4157858477.py:1: UserWarning: pandas only supports SQLAlchemy connectable (engine/connection) or database
df_trip = pd.read_sql("SELECT * FROM trip_summary;", conn)
/tmp/ipython-input-4157858477.py:2: UserWarning: pandas only supports SQLAlchemy connectable (engine/connection) or database
df_driver = pd.read_sql("SELECT * FROM driver_daily;", conn)
(120000, 17) (10, 11)
```

```
1 df_trip.head()
```

	TRIP_ID	DRIVER_ID	VEHICLE_ID	DAY	TRIP_START_AT	TRIP_DURATION_MINUTES	TRIP_DISTANCE_KM	Avg_Speed	Max_Speed	Avg_Accel
0	114387	101	1001	2023-01-02	2023-01-02 07:46:26		66.626622	4.558716	55.215226	55.215226
1	114394	101	3003	2023-01-02	2023-01-02 07:46:33		36.829610	30.989605	51.353421	51.353421
2	114401	101	4004	2023-01-02	2023-01-02 07:46:40		35.398257	25.775307	93.328352	93.328352
3	114408	101	1001	2023-01-02	2023-01-02 07:46:47		53.200739	91.098530	7.532813	7.532813
4	114415	103	3003	2023-01-02	2023-01-02 07:46:54		10.778823	112.873466	2.887324	2.887324

```
1 df_driver.head()
```

	DRIVER_ID	DAY	TRIP_COUNT_DAILY	AVG_TRIP_DISTANCE_DAILY	AVG_SPEED_DAILY	TOTAL_BRAKE_EVENTS_DAILY	ANOMALIES_COUNT_D
0	102	2023-01-02		6740	50.300488	29.966187	33643
1	105	2023-01-02		3274	49.032908	29.668213	16241
2	104	2023-01-02		3468	49.872775	28.997656	17308
3	101	2023-01-01		43094	49.866401	29.951755	215095
4	103	2023-01-01		8621	49.175268	29.836446	43178

```
1 df_trip.head()
```

	TRIP_ID	DRIVER_ID	VEHICLE_ID	DAY	TRIP_START_AT	TRIP_DURATION_MINUTES	TRIP_DISTANCE_KM	Avg_Speed	Max_Speed	Avg_Acceleration
0	114387	101	1001	2023-01-02	2023-01-02 07:46:26	66.626622	4.558716	55.215226	55.215226	1.003904
1	114394	101	3003	2023-01-02	2023-01-02 07:46:33	36.829610	30.989605	51.353421	51.353421	-0.99991
2	114401	101	4004	2023-01-02	2023-01-02 07:46:40	35.398257	25.775307	93.328352	93.328352	29.971987
3	114408	101	1001	2023-01-02	2023-01-02 07:46:47	53.200739	91.098530	7.532813	7.532813	0.000157
4	114415	103	3003	2023-01-02	2023-01-02 07:46:54	10.778823	112.873466	2.887324	2.887324	0.285187

```
1 df_trip.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 120000 entries, 0 to 119999
Data columns (total 17 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   TRIP_ID         120000 non-null   object 
 1   DRIVER_ID       120000 non-null   object 
 2   VEHICLE_ID      120000 non-null   object 
 3   DAY              120000 non-null   datetime64[ns]
 4   TRIP_START_AT   120000 non-null   datetime64[ns]
 5   TRIP_DURATION_MINUTES  120000 non-null   float64
 6   TRIP_DISTANCE_KM  120000 non-null   float64
 7   AVG_SPEED        120000 non-null   float64
 8   MAX_SPEED         120000 non-null   float64
 9   AVG_ACCELERATION 120000 non-null   float64
 10  TOTAL_BRAKE_EVENTS 120000 non-null   int64  
 11  MAX_ROUTE_DEVIATION 120000 non-null   float64
 12  ANY_ANOMALOUS_EVENT 120000 non-null   bool   
 13  ANY_GEOFENCING_VIOLATION 120000 non-null   bool   
 14  ANY_ROUTE_ANOMALY 120000 non-null   bool   
 15  AVG_ROUTE_DEVIATION_SCORE 120000 non-null   float64
 16  BEHAVIORAL_CONSISTENCY_INDEX 120000 non-null   float64
dtypes: bool(3), datetime64[ns](2), float64(8), int64(1), object(3)
memory usage: 13.2+ MB
```

```
1 display(df_trip.describe().T)
```

	count	mean	min	25%	50%	75%	
DAY	120000	2023-01-01 06:43:12	2023-01-01 00:00:00	2023-01-01 00:00:00	2023-01-01 00:00:00	2023-01-02 00:00:00	2
TRIP_START_AT	120000	2023-01-01 16:39:59.500000	2023-01-01 00:00:00	2023-01-01 08:19:59.750000128	2023-01-01 16:39:59.500000	2023-01-02 00:59:59.249999872	2
TRIP_DURATION_MINUTES	120000.0	59.891751	0.00003	17.212549	41.402746	82.931203	7
TRIP_DISTANCE_KM	120000.0	49.921316	0.000521	14.3195	34.387936	69.138769	5
AVG_SPEED	120000.0	29.971987	0.000157	8.61525	20.633107	41.439832	3
MAX_SPEED	120000.0	29.971987	0.000157	8.61525	20.633107	41.439832	3
AVG_ACCELERATION	120000.0	1.003904	-0.99991	-0.41903	0.388639	1.777851	1
TOTAL_BRAKE_EVENTS	120000.0	4.995767	0.0	3.0	5.0	6.0	
MAX_ROUTE_DEVIATION	120000.0	0.285187	0.000305	0.160065	0.263214	0.389066	
AVG_ROUTE_DEVIATION_SCORE	120000.0	0.285187	0.000305	0.160065	0.263214	0.389066	

```
1 df_trip.isnull().sum()
```

Show hidden output

```
1 df_trip['ANY_ANOMALOUS_EVENT'].value_counts()
```

```
count
ANY_ANOMALOUS_EVENT
False    107980
True     12020
dtype: int64
```

```
1 print(df_trip['ANY_ROUTE_ANOMALY'].value_counts())
2 print(df_trip['ANY_GEOFENCING_VIOLATION'].value_counts())
```

```
ANY_ROUTE_ANOMALY
False    108035
True     11965
Name: count, dtype: int64
ANY_GEOFENCING_VIOLATION
False    113971
True      6029
Name: count, dtype: int64
```

```
1 df_trip.columns
```

```
Index(['TRIP_ID', 'DRIVER_ID', 'VEHICLE_ID', 'DAY', 'TRIP_START_AT',
       'TRIP_DURATION_MINUTES', 'TRIP_DISTANCE_KM', 'AVG_SPEED', 'MAX_SPEED',
       'AVG_ACCELERATION', 'TOTAL_BRAKE_EVENTS', 'MAX_ROUTE_DEVIATION',
       'ANY_ANOMALOUS_EVENT', 'ANY_GEOFENCING_VIOLATION', 'ANY_ROUTE_ANOMALY',
       'AVG_ROUTE_DEVIATION_SCORE', 'BEHAVIORAL_CONSISTENCY_INDEX',
       'BRAKE_PER_KM', 'ROUTE_DEV_PER_KM', 'HOUR', 'TIME_OF_DAY'],
      dtype='object')
```

```
1 # brake events per km
2 df_trip['BRAKE_PER_KM'] = (df_trip['TOTAL_BRAKE_EVENTS'] / df_trip['TRIP_DISTANCE_KM']).replace(0, np.nan)
```

```
1 # route deviation per km
2 df_trip['ROUTE_DEV_PER_KM'] = (df_trip['MAX_ROUTE_DEVIATION'] / df_trip['TRIP_DISTANCE_KM']).replace(0, np.nan)
```

```
1 # hour of day
2 df_trip['TRIP_START_HOUR'] = df_trip['TRIP_START_AT'].dt.hour
```

```
1 # bucket time of day
2 df_trip['TIME_OF_DAY'] = pd.cut(df_trip['TRIP_START_HOUR'], bins=[-1, 6, 12, 18, 24], labels=['night', 'morning', 'afternoon'])
```

```
1 df_trip.head()
```

	TRIP_ID	DRIVER_ID	VEHICLE_ID	DAY	TRIP_START_AT	TRIP_DURATION_MINUTES	TRIP_DISTANCE_KM	AVG_SPEED	MAX_SPEED	AVG_ACCELERATION
0	114387	101	1001	2023-01-02	2023-01-02 07:46:26		66.626622	4.558716	55.215226	55.215226
1	114394	101	3003	2023-01-02	2023-01-02 07:46:33		36.829610	30.989605	51.353421	51.353421
2	114401	101	4004	2023-01-02	2023-01-02 07:46:40		35.398257	25.775307	93.328352	93.328352
3	114408	101	1001	2023-01-02	2023-01-02 07:46:47		53.200739	91.098530	7.532813	7.532813
4	114415	103	3003	2023-01-02	2023-01-02 07:46:54		10.778823	112.873466	2.887324	2.887324

5 rows × 22 columns

```
1 df_trip.describe().T
```

		count	mean	min	25%	50%	75%
DAY		120000	2023-01-01 06:43:12	2023-01-01 00:00:00	2023-01-01 00:00:00	2023-01-01 00:00:00	2023-01-02 00:00:00
TRIP_START_AT		120000	2023-01-16:39:59.500000	2023-01-01 00:00:00	2023-01-01 08:19:59.750000128	2023-01-01 16:39:59.500000	2023-01-02 00:59:59.249999872
TRIP_DURATION_MINUTES		120000.0	59.891751	0.00003	17.212549	41.402746	82.931203
TRIP_DISTANCE_KM		120000.0	49.921316	0.000521	14.3195	34.387936	69.138769
Avg_Speed		120000.0	29.971987	0.000157	8.61525	20.633107	41.439832
Max_Speed		120000.0	29.971987	0.000157	8.61525	20.633107	41.439832
Avg_Accelaration		120000.0	1.003904	-0.99991	-0.41903	0.388639	1.777851
Total_Brake_Events		120000.0	4.995767	0.0	3.0	5.0	6.0
Max_Route_Deviation		120000.0	0.285187	0.000305	0.160065	0.263214	0.389066
Avg_Route_Deviation_Score		120000.0	0.285187	0.000305	0.160065	0.263214	0.389066
Behavioral_Conistency_Index		120000.0	0.285855	0.000295	0.160986	0.264427	0.389694
Brake_Per_Km		120000.0	1.452215	0.0	0.061677	0.134322	0.337682 15.0
Route_Dev_Per_Km		120000.0	0.093352	0.000001	0.003194	0.007321	0.01897
Hour		120000.0	9.45	0.0	4.0	8.0	15.0

```
1 print(df_trip.columns)
```

```
Index(['TRIP_ID', 'DRIVER_ID', 'VEHICLE_ID', 'DAY', 'TRIP_START_AT',
       'TRIP_DURATION_MINUTES', 'TRIP_DISTANCE_KM', 'AVG_SPEED', 'MAX_SPEED',
       'AVG_ACCELARATION', 'TOTAL_BRAKE_EVENTS', 'MAX_ROUTE_DEVIATION',
       'ANY_ANOMALOUS_EVENT', 'ANY_GEOFENCING_VIOLATION', 'ANY_ROUTE_ANOMALY',
       'AVG_ROUTE_DEVIATION_SCORE', 'BEHAVIORAL_CONSISTENCY_INDEX',
       'BRAKE_PER_KM', 'ROUTE_DEV_PER_KM', 'HOUR', 'TIME_OF_DAY',
       'TRIP_START_HOUR'],
      dtype='object')
```

```
1 feature_cols = [
2     'TRIP_DURATION_MINUTES',
3     'TRIP_DISTANCE_KM',
4     'AVG_SPEED',
5     'MAX_SPEED',
6     'AVG_ACCELARATION',
7     'TOTAL_BRAKE_EVENTS',
8     'MAX_ROUTE_DEVIATION',
9     'AVG_ROUTE_DEVIATION_SCORE',
10    'BEHAVIORAL_CONSISTENCY_INDEX',
11    'BRAKE_PER_KM',
12    'ROUTE_DEV_PER_KM',
13    'TRIP_START_HOUR'
14 ]
15
16 X = df_trip[feature_cols]
```

```
1 y = df_trip['ANY_ANOMALOUS_EVENT'].astype(int)
2
```

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, stratify = y, random_state = 42)
```

```
1 from xgboost import XGBClassifier
2 xg_model = XGBClassifier()
```

```
1 xg_model.fit(X_train, y_train)
```

```
XGBClassifier(base_score=None, booster=None, callbacks=None,
              colsample_bylevel=None, colsample_bynode=None,
              colsample_bytree=None, device=None, early_stopping_rounds=None,
              enable_categorical=False, eval_metric=None, feature_types=None,
              feature_weights=None, gamma=None, grow_policy=None,
              importance_type=None, interaction_constraints=None,
              learning_rate=None, max_bin=None, max_cat_threshold=None,
              max_cat_to_onehot=None, max_delta_step=None, max_depth=None,
              max_leaves=None, min_child_weight=None, missing=nan,
              monotone_constraints=None, multi_strategy=None, n_estimators=None,
              n_jobs=None, num_parallel_tree=None, ...)
```

```
1 y_pred = xg_model.predict(X_test)
```

```
1 from sklearn.metrics import *
2 acc_score = accuracy_score(y_test , y_pred)
3 auc_score = roc_auc_score(y_test, y_pred)
4 print("Accuracy score of XGBClassifier model is :", acc_score*100)
5 print("AUC Score of XGBClassifier model is :", auc_score*100)
```

Accuracy score of XGBClassifier model is : 89.98333333333333
AUC Score of XGBClassifier model is : 50.0

```
1 # probabilities (positive class)
2 y_proba = xg_model.predict_proba(X_test)[:,1]
3
4 from sklearn.metrics import roc_auc_score
5 auc_score = roc_auc_score(y_test, y_proba)
6 print("Correct AUC (using probabilities):", auc_score*100)
```

Correct AUC (using probabilities): 50.60528102904919

```
1 import numpy as np
2 k = int(0.05 * len(y_test))    # top 5% of test set
3 top_idx = np.argsort(y_proba)[-k:]    # indices of top-k highest probs
4 precision_at_k = y_test.iloc[top_idx].sum() / len(top_idx)
5 print("Precision@Top5%:", precision_at_k)
6
```

Precision@Top5%: 0.1105555555555556

```
1 model = XGBClassifier(n_estimators=200, max_depth=6, learning_rate=0.08, use_label_encoder=False, eval_metric='logloss'
2 model.fit(X_train, y_train, eval_set=[(X_test,y_test)], verbose=False)
```

/usr/local/lib/python3.12/dist-packages/xgboost/training.py:199: UserWarning: [15:51:16] WARNING: /workspace/src/learner.cc:
Parameters: { "use_label_encoder" } are not used.

```
bst.update(dtrain, iteration=i, fobj=obj)
```

```
XGBClassifier(base_score=None, booster=None, callbacks=None,
              colsample_bylevel=None, colsample_bynode=None,
              colsample_bytree=None, device=None, early_stopping_rounds=None,
              enable_categorical=False, eval_metric='logloss',
              feature_types=None, feature_weights=None, gamma=None,
              grow_policy=None, importance_type=None,
              interaction_constraints=None, learning_rate=0.08, max_bin=None,
              max_cat_threshold=None, max_cat_to_onehot=None,
              max_delta_step=None, max_depth=6, max_leaves=None,
              min_child_weight=None, missing=nan, monotone_constraints=None,
              multi_strategy=None, n_estimators=200, n_jobs=None,
```

```
1 y_proba = model.predict_proba(X_test)[:,1]
2
3
```

```
1 from sklearn.metrics import roc_auc_score, precision_recall_curve
2 auc = roc_auc_score(y_test, y_proba)
3 print("Auc_Score",auc*100)
4 # precision@k: take top K% (or top K trips)
5 k = int(0.05 * len(y_test))    # top 5% as example
6 top_idx = np.argsort(y_proba)[-k:]
7 precision_at_k = y_test.iloc[top_idx].sum() / len(top_idx)
8 print("Precision@Top5%:", precision_at_k*100)
9
```

```
Auc score 50.446081388845485
```

```
1 df_preds = pd.DataFrame({  
2     'TRIP_ID': df_trip.iloc[X_test.index]['TRIP_ID'].values,  
3     'DRIVER_ID': df_trip.iloc[X_test.index]['DRIVER_ID'].values,  
4     'ACTUAL': y_test.values,  
5     'PRED_PROB': y_proba  
6 })  
7 })
```

```
1 from snowflake.connector.pandas_tools import write_pandas  
2  
3 success, nchunks, nrows, _ = write_pandas(conn, df_preds, 'PREDICTIONS', schema='ANALYTICS')  
4 print(success, nrows)  
5
```

```
True 36000
```

```
1 df_preds
```

TRIP_ID	DRIVER_ID	ACTUAL	PRED_PROB	SCORE_AT
0	0A001	101	0	0 105029 2025 12 02 16:55:20 880277