

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

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.metrics import classification_report, roc_auc_score
5 import numpy as np
6
7 # Load datasets with consistent dtypes and low_memory=False to avoid DtypeWarning
8 app_df = pd.read_csv('app.csv', dtype={'identity_id_hash': 'object'}, low_memory=False)
9 campa_df = pd.read_csv('campa.csv', dtype={'identity_id_hash': 'object'}, low_memory=False)
10 ib_df = pd.read_csv('ib.csv', dtype={'identity_id_hash': 'object'}, low_memory=False)
11 prodeje_df = pd.read_csv('prodeje.csv', dtype={'identity_id_hash': 'object'}, low_memory=False)
12 web_df = pd.read_csv('web.csv', dtype={'identity_id_hash': 'object'}, low_memory=False)
13
14 # Merge datasets on 'identity_id_hash'
15 merged_df = app_df.merge(campa_df, on='identity_id_hash', how='outer', suffixes=('_app', '_campa')) \
16     .merge(ib_df, on='identity_id_hash', how='outer', suffixes=('', '_ib')) \
17     .merge(prodeje_df, on='identity_id_hash', how='outer', suffixes=('', '_prodeje')) \
18     .merge(web_df, on='identity_id_hash', how='outer', suffixes=('', '_web'))
19
20 # Drop columns with all missing values
21 columns_to_drop = [col for col in merged_df.columns if merged_df[col].isnull().sum() == len(merged_df)]
22 cleaned_df = merged_df.drop(columns=columns_to_drop)
23
24 # Fill missing values for categorical and numerical columns
25 categorical_cols = cleaned_df.select_dtypes(include=['object']).columns
26 numerical_cols = cleaned_df.select_dtypes(include=['float64', 'int64']).columns
27 cleaned_df[categorical_cols] = cleaned_df[categorical_cols].fillna('Unknown')
28 cleaned_df[numerical_cols] = cleaned_df[numerical_cols].fillna(cleaned_df[numerical_cols].median())
29
30 # Feature engineering
31 extra_features = pd.DataFrame({
32     'avg_session_duration': cleaned_df[['session_time', 'session_time_ib']].mean(axis=1),
33     'interaction_count': cleaned_df.apply(lambda row: sum([1 for event in ['page_view', 'custom_event', 'screen_view', 'back_intent_event', 'show_content_event'] if event in row]), axis=1),
34     'campaign_engagement': cleaned_df['campaign_planning_name'].apply(lambda x: 0 if x == 'Unknown' else 1),
35     'unique_product_interest': cleaned_df[['product_l1', 'product_l2', 'product_l2_prodeje']].nunique(axis=1)
36 })
37
38 cleaned_df = pd.concat([cleaned_df, extra_features], axis=1)
39
40 # Define target variable
41 cleaned_df['purchase'] = cleaned_df['agreement_status'].apply(lambda x: 1 if x == 'Y' else 0)
42 # cleaned_df['purchase'] = cleaned_df['nbi_fictive_flag'].apply(lambda x: 1 if x == '>2100' else 0)
43
44 # Select only numeric columns for correlation to avoid ValueError
45 numeric_df = cleaned_df.select_dtypes(include=[np.number])
46
```

```

47 # Data leakage check by examining correlation with the target variable
48 correlations = numeric_df.corr()
49 high_corr_features = correlations['purchase'].loc[lamba x: abs(x) > 0.8].drop('purchase').index.tolist()
50
51 if high_corr_features:
52     print("Potential data leakage detected with high correlation features:", high_corr_features)
53 else:
54     print("No high correlation features detected, data leakage is less likely.")
55
56 # Select features and target, excluding highly correlated features if necessary
57 feature_columns = ['avg_session_duration', 'interaction_count', 'campaign_engagement', 'unique_product_interest']
58 if high_corr_features:
59     feature_columns = [col for col in feature_columns if col not in high_corr_features]
60
61 X = cleaned_df[feature_columns]
62 y = cleaned_df['purchase']
63
64 # Split data into train and test sets with an 80-20 split
65 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
66
67 # Set up RandomForest with GridSearchCV for hyperparameter tuning
68 rf_model = RandomForestClassifier(random_state=42)
69 param_grid = {
70     'n_estimators': [50],
71     'max_depth': [5],
72     'min_samples_split': [2]
73 }
74 grid_search = GridSearchCV(rf_model, param_grid, cv=5, scoring='roc_auc', n_jobs=-1)
75
76 # Fit model with GridSearch
77 grid_search.fit(X_train, y_train)
78
79 # Best model and cross-validation scores
80 best_model = grid_search.best_estimator_
81 cross_val_scores = cross_val_score(best_model, X_train, y_train, cv=5, scoring='roc_auc')
82
83 # Make predictions with the best model on the test set
84 y_pred = best_model.predict(X_test)
85 y_prob = best_model.predict_proba(X_test)[:, 1]
86
87 # Evaluation metrics
88 classification_rep = classification_report(y_test, y_pred)
89 roc_auc = roc_auc_score(y_test, y_prob)
90
91 # Output results
92 print("Best Parameters:", grid_search.best_params_)
93 print("Cross-validation AUC scores:", cross_val_scores)
94 print("Average CV AUC:", cross_val_scores.mean())
95 print("Test Set Classification Report:\n", classification_rep)
96 print("Test Set AUC-ROC:", roc_auc)
97

```

```

98 # Segment customers based on predicted purchase probabilities
99 customer_probabilities = best_model.predict_proba(X_test)[:, 1] # Probability of being in class 1 (purchase)
100
101 # Define segments based on probability thresholds
102 segments = []
103 for prob in customer_probabilities:
104     if prob > 0.8:
105         segments.append('High Value Engaged')
106     elif 0.5 < prob <= 0.8:
107         segments.append('Potential Upsell')
108     elif 0.3 < prob <= 0.5:
109         segments.append('Nurture and Educate')
110     else:
111         segments.append('Low Engagement and Awareness')
112
113 # Add segment labels and probability to DataFrame
114 X_test['purchase_probability'] = customer_probabilities
115 X_test['segment'] = segments
116
117 # Product recommendations based on segments
118 recommendations = []
119 for segment in segments:
120     if segment == 'High Value Engaged':
121         recommendations.append('Offer premium services, investment opportunities, and loyalty programs.')
122     elif segment == 'Potential Upsell':
123         recommendations.append('Suggest additional products like credit cards or loans with personalized offers.')
124     elif segment == 'Nurture and Educate':
125         recommendations.append('Provide educational content on products, focusing on convenience and value.')
126     else: # 'Low Engagement and Awareness'
127         recommendations.append('Target with brand awareness campaigns and simple introductory offers.')
128
129 # Add recommendations to the DataFrame
130 X_test['recommendation'] = recommendations
131
132 # Display a sample of the segmented data with recommendations
133 X_test[['purchase_probability', 'segment', 'recommendation']].head(10)

```

Best Parameters: {'max\_depth': 5, 'min\_samples\_split': 2, 'n\_estimators': 50}  
Cross-validation AUC scores: [0.95188881 0.94734026 0.94265207 0.95369071 0.95975599]  
Average CV AUC: 0.951065568125418  
Test Set Classification Report:

	precision	recall	f1-score	support
0	0.95	0.88	0.91	695
1	0.81	0.92	0.86	391
accuracy			0.89	1086
macro avg	0.88	0.90	0.89	1086
weighted avg	0.90	0.89	0.89	1086

Test Set AUC-ROC: 0.9418498960422454

	purchase_probability	segment	recommendation	
79	0.777550	Potential Upsell	Suggest additional products like credit cards ...	
4139	0.909583	High Value Engaged	Offer premium services, investment opportuniti...	
1643	0.005575	Low Engagement and Awareness	Target with brand awareness campaigns and simp...	
167	0.838086	High Value Engaged	Offer premium services, investment opportuniti...	
439	0.777550	Potential Upsell	Suggest additional products like credit cards ...	
2531	0.065151	Low Engagement and Awareness	Target with brand awareness campaigns and simp...	
1421	0.005575	Low Engagement and Awareness	Target with brand awareness campaigns and simp...	
2834	0.065151	Low Engagement and Awareness	Target with brand awareness campaigns and simp...	
4356	0.777550	Potential Upsell	Suggest additional products like credit cards ...	
168	0.838086	High Value Engaged	Offer premium services, investment opportuniti...	