

✓ Credit Card Fraud Detection Analysis

This notebook presents a comprehensive analysis of credit card fraud detection. It explores various machine learning models, preprocessing techniques, and evaluation metrics to identify fraudulent transactions effectively. The analysis covers:

- **Data Loading and Preprocessing:** Initial steps to load the dataset and prepare it for modeling, including handling missing values and feature scaling.
- **Model Modulation:** Evaluation of different classification algorithms (Logistic Regression, LinearSVC, RandomForest, XGBoost) under various conditions: raw data, scaled data, and with class imbalance handling (class weights, SMOTE, SMOTEENN).
- **Feature Selection:** Implementation of feature selection techniques to identify and utilize the most impactful features for improved model performance and efficiency.
- **Model Explanation with SHAP:** Using SHAP (SHapley Additive exPlanations) to interpret the behavior of the best-performing models, providing insights into feature contributions to predictions.
- **Results and Conclusion:** A consolidated summary of all model performances, key findings, and recommendations for deployment and future work.

The goal is to identify a robust and efficient model capable of accurately detecting fraudulent transactions while minimizing false positives and false negatives.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import kagglehub
```

```
1 # Download latest version
2 path = kagglehub.dataset_download("mlg-ulb/creditcardfraud")
3
4 print("Path to dataset files:", path)
```

Using Colab cache for faster access to the 'creditcardfraud' dataset.
Path to dataset files: /kaggle/input/creditcardfraud

```
1 df = pd.read_csv('/kaggle/input/creditcardfraud/creditcard.csv')
2 df.head()
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	0.098698	0.363787	...	-0.018307	0.277838	-0.110474	0.066928	0.128539	-0.189115	0.133558	-0.0210
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	0.085102	-0.255425	...	-0.225775	-0.638672	0.101288	-0.339846	0.167170	0.125895	-0.008983	0.0147
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	0.247676	-1.514654	...	0.247998	0.771679	0.909412	-0.689281	-0.327642	-0.139097	-0.055353	-0.0597
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	0.377436	-1.387024	...	-0.108300	0.005274	-0.190321	-1.175575	0.647376	-0.221929	0.062723	0.0614
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	-0.270533	0.817739	...	-0.009431	0.798278	-0.137458	0.141267	-0.206010	0.502292	0.219422	0.2157

5 rows × 31 columns

```
1 df.describe()
```

	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9	...	V21	V22	V23	V24	V25	V26	V27	V28
count	284807.000000	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	...	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05	2.848070e+05
mean	94813.859575	1.168375e-15	3.416908e-16	-1.379537e-15	2.074095e-15	9.604066e-16	1.487313e-15	-5.556467e-16	1.213481e-16	-2.406331e-15	...	1.654067e-16	-3.568593e-16	1.654067e-16	-3.568593e-16	1.654067e-16	-3.568593e-16	1.654067e-16	-3.568593e-16
std	47488.145955	1.958696e+00	1.651309e+00	1.516255e+00	1.415869e+00	1.380247e+00	1.332271e+00	1.237094e+00	1.194353e+00	1.098632e+00	...	7.345240e-01	7.257016e-01	7.345240e-01	7.257016e-01	7.345240e-01	7.257016e-01	7.345240e-01	7.257016e-01
min	0.000000	-5.640751e+01	-7.271573e+01	-4.832559e+01	-5.683171e+00	-1.137433e+02	-2.616051e+01	-4.355724e+01	-7.321672e+01	-1.343407e+01	...	-3.483038e+01	-1.093314e+01	-3.483038e+01	-1.093314e+01	-3.483038e+01	-1.093314e+01	-3.483038e+01	-1.093314e+01
25%	54201.500000	-9.203734e-01	-5.985499e-01	-8.903648e-01	-8.486401e-01	-6.915971e-01	-7.682956e-01	-5.540759e-01	-2.086297e-01	-6.430976e-01	...	-2.283949e-01	-5.423504e-01	-2.283949e-01	-5.423504e-01	-2.283949e-01	-5.423504e-01	-2.283949e-01	-5.423504e-01
50%	84692.000000	1.810880e-02	6.548556e-02	1.798463e-01	-1.984653e-02	-5.433583e-02	-2.741871e-01	4.010308e-02	2.235804e-02	-5.142873e-02	...	-2.945017e-02	6.781943e-03	-2.945017e-02	6.781943e-03	-2.945017e-02	6.781943e-03	-2.945017e-02	6.781943e-03
75%	139320.500000	1.315642e+00	8.037239e-01	1.027196e+00	7.433413e-01	6.119264e-01	3.985649e-01	5.704361e-01	3.273459e-01	5.971390e-01	...	1.863772e-01	5.285536e-01	1.863772e-01	5.285536e-01	1.863772e-01	5.285536e-01	1.863772e-01	5.285536e-01
max	172792.000000	2.454930e+00	2.205773e+01	9.382558e+00	1.687534e+01	3.480167e+01	7.330163e+01	1.205895e+02	2.000721e+01	1.559499e+01	...	2.720284e+01	1.050309e+01	2.720284e+01	1.050309e+01	2.720284e+01	1.050309e+01	2.720284e+01	1.050309e+01

8 rows × 31 columns

```
1 Start coding or generate with AI.
```

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.feature_selection import SelectFromModel
3 from sklearn.preprocessing import StandardScaler, RobustScaler, MinMaxScaler
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.svm import LinearSVC
7 from sklearn.pipeline import make_pipeline
8 from sklearn.metrics import f1_score, matthews_corrcoef, confusion_matrix
9
10 from xgboost import XGBClassifier
11 from imblearn.over_sampling import SMOTE
12 from imblearn.combine import SMOTEENN
13 import time
```

```
1 # Separate features (X) and target (y)
2 X = df.drop('Class', axis=1)
3 y = df['Class']
4
5 # Perform train-test split
6 # test_size=0.3 means 30% of data for testing, 70% for training
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

```

8
9 print("Dataset split into training and testing sets.")
10 print(f"Training features shape: {X_train.shape}")
11 print(f"Testing features shape: {X_test.shape}")
12 print(f"Training target shape: {y_train.shape}")
13 print(f"Testing target shape: {y_test.shape}")
14
15 # Display the class distribution in train and test sets to verify stratification
16 print("\nClass distribution in training set:")
17 print(y_train.value_counts(normalize=True))
18 print("\nClass distribution in testing set:")

```

```

Dataset split into training and testing sets.
Training features shape: (199364, 30)
Testing features shape: (85443, 30)
Training target shape: (199364,)
Testing target shape: (85443,)

```

```

Class distribution in training set:
Class
0    0.998214
1    0.001786
Name: proportion, dtype: float64

```

```

Class distribution in testing set:
Class
0    0.998408
1    0.001592
Name: proportion, dtype: float64

```

1 Start coding or [generate](#) with AI.

Model Modulation

Evaluate Model Function

```

1 # @title Evaluate Model Function
2
3 # Define the evaluation function
4 def evaluate_model_pipeline(estimator, X_train, y_train, X_test, y_test):
5     # The 'estimator' can be a single model or a pipeline already created
6
7     # Record start time for training
8     start_time = time.time()
9
10    # Train the estimator
11    estimator.fit(X_train, y_train)
12
13    # Record end time for training
14    end_time = time.time()
15    runtime = end_time - start_time
16
17    # Make predictions
18    y_pred = estimator.predict(X_test)
19
20    # Calculate scores
21    f1 = f1_score(y_test, y_pred)
22    mcc = matthews_corrcoef(y_test, y_pred)
23
24    return f1, mcc, runtime
25
26 print("Evaluation function defined.")

```

Evaluation function defined.

Base Supervised Learning Models (without Scaling)

This section evaluates the performance of several common supervised learning models—Logistic Regression, Linear Support Vector Classifier (LinearSVC), Random Forest Classifier, and XGBoost Classifier—on the raw, unscaled training data. This serves as a baseline to understand how these models perform without any feature preprocessing.

Base Supervised Learning Models (no Scaling)

```

1 # @title Base Supervised Learning Models (no Scaling)
2
3 # Instantiate the models
4 model_lr = LogisticRegression(random_state=42, solver='liblinear')
5 model_svc = LinearSVC(random_state=42, dual=True)
6 model_rf = RandomForestClassifier(random_state=42, n_jobs=-1)
7 model_xgb = XGBClassifier(random_state=42,
8                           eval_metric='logloss'
9                           )
10
11 # Create a list of models
12 models = [
13     ('Logistic Regression (no Scaling)', model_lr),
14     ('LinearSVC (no Scaling)', model_svc),
15     ('RandomForest Classifier (no Scaling)', model_rf),
16     ('XGBoost Classifier (no Scaling)', model_xgb)
17 ]
18
19
20 model_results = {}
21
22 print("\n--- Model Evaluation Results ---")
23 for model_name, model_obj in models:
24     print(f"\nEvaluating {model_name}...")
25     f1, mcc, runtime = evaluate_model_pipeline(model_obj, X_train, y_train, X_test, y_test)

```

```

26 model_results[model_name] = {'F1 Score': f1, 'MCC Score': mcc, 'Runtime': runtime}
27 print(f"   F1 Score: {f1:.4f}")
28 print(f"   MCC Score: {mcc:.4f}")
29 print(f"   Runtime: {runtime:.2f} seconds")
30
31 df_summary = pd.DataFrame(model_results).T
32 print("\n--- Final Scores Summary ---")
33 print(df_summary)
34

```

--- Model Evaluation Results ---

Evaluating Logistic Regression (no Scaling)...

F1 Score: 0.6457
MCC Score: 0.6468
Runtime: 2.83 seconds

Evaluating LinearSVC (no Scaling)...

/usr/local/lib/python3.12/dist-packages/sklearn/svm/_base.py:1249: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
warnings.warn(
F1 Score: 0.1963
MCC Score: 0.2636
Runtime: 57.40 seconds

Evaluating RandomForest Classifier (no Scaling)...

F1 Score: 0.8685
MCC Score: 0.8714
Runtime: 260.50 seconds

Evaluating XGBoost Classifier (no Scaling)...

F1 Score: 0.8594
MCC Score: 0.8609
Runtime: 3.81 seconds

--- Final Scores Summary ---

	F1 Score	MCC Score	Runtime
Logistic Regression (no Scaling)	0.645669	0.646775	2.828877
LinearSVC (no Scaling)	0.196319	0.263582	57.397182
RandomForest Classifier (no Scaling)	0.868526	0.871398	260.503331
XGBoost Classifier (no Scaling)	0.859375	0.860853	3.806446

1 Start coding or [generate](#) with AI.

▼ Base Scaled Supervised Learning Models

Here, the previously instantiated models are re-evaluated after scaling the input features using `StandardScaler`. Feature scaling is crucial for many machine learning algorithms, especially those sensitive to feature magnitudes like Logistic Regression and LinearSVC, as it helps prevent features with larger values from dominating the learning process.

```

1 # Instantiate the StandardScaler
2 scaler_batch = StandardScaler()
3
4 # Scale the training and testing data
5 X_train_scaled = scaler_batch.fit_transform(X_train)
6 X_test_scaled = scaler_batch.transform(X_test)

```

▼ Base Scaled Supervised Learning Models

```

1 # @title Base Scaled Supervised Learning Models
2
3 # Instantiate the StandardScaler
4 scaler_batch = StandardScaler()
5
6 # Scale the training and testing data
7 X_train_scaled = scaler_batch.fit_transform(X_train)
8 X_test_scaled = scaler_batch.transform(X_test)
9
10 scaled_models = [
11     ('Logistic Regression', model_lr),
12     ('LinearSVC', model_svc),
13     ('RandomForest Classifier', model_rf),
14     ('XGBoost Classifier', model_xgb)
15 ]
16
17 print("\n--- Model Evaluation Results (Scaled) ---")
18 for model_name, pipeline_obj in scaled_models:
19     print(f"\nEvaluating {model_name}...")
20     f1, mcc, runtime = evaluate_model_pipeline(pipeline_obj, X_train_scaled, y_train, X_test_scaled, y_test)
21     model_results[model_name] = {'F1 Score': f1, 'MCC Score': mcc, 'Runtime': runtime}
22     print(f"   F1 Score: {f1:.4f}")
23     print(f"   MCC Score: {mcc:.4f}")
24     print(f"   Runtime: {runtime:.2f} seconds")
25
26
27 df_summary = pd.DataFrame(model_results).T
28 print("\n--- Final Scores Summary ---")
29 print(df_summary[4:])

```

--- Model Evaluation Results (Scaled) ---

Evaluating Logistic Regression...

F1 Score: 0.7296
MCC Score: 0.7397
Runtime: 2.33 seconds

Evaluating LinearSVC...

/usr/local/lib/python3.12/dist-packages/sklearn/svm/_base.py:1249: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
warnings.warn(

```

F1 Score: 0.8063
MCC Score: 0.8083
Runtime: 97.30 seconds

Evaluating RandomForest Classifier...
F1 Score: 0.8685
MCC Score: 0.8714
Runtime: 305.84 seconds

Evaluating XGBoost Classifier...
F1 Score: 0.8594
MCC Score: 0.8609
Runtime: 7.17 seconds

--- Final Scores Summary ---

```

	F1 Score	MCC Score	Runtime
Logistic Regression	0.729614	0.739719	2.326985
LinearSVC	0.806324	0.808328	97.296894
RandomForest Classifier	0.868526	0.871398	305.835859
XGBoost Classifier	0.859375	0.860853	7.172506

1 Start coding or [generate](#) with AI.

Class Weight Balanced Models

Given the highly imbalanced nature of the dataset (fraudulent transactions are rare), this section introduces models configured with class weighting. This technique assigns higher weights to the minority class (fraud) during training, forcing the model to pay more attention to these rare but critical examples, thereby improving its ability to detect fraud. `class_weight='balanced'` is used for Random Forest and LinearSVC, while `scale_pos_weight` is calculated and applied for XGBoost.

Class Weight Balanced Models

```

1 # @title Class Weight Balanced Models
2
3 model_rf_balanced = RandomForestClassifier(random_state=42, n_jobs=-1, class_weight='balanced')
4 model_svc_balanced = LinearSVC(random_state=42, dual=True, class_weight='balanced')
5
6 # Calculate scale_pos_weight for XGBoost
7 neg_count = np.sum(y_train == 0)
8 pos_count = np.sum(y_train == 1)
9 scale_pos_weight_value = neg_count / pos_count
10
11 model_xgb_balanced = XGBClassifier(random_state=42,
12                                   eval_metric='logloss',
13                                   use_label_encoder=False,
14                                   scale_pos_weight=scale_pos_weight_value)
15
16 balanced_models = [
17     ('RandomForest Classifier (Balanced)', model_rf_balanced),
18     ('LinearSVC (Balanced)', model_svc_balanced),
19     ('XGBoost Classifier (Balanced)', model_xgb_balanced)
20 ]
21
22 print("Balanced models initialized and listed in 'balanced_models'.")

```

Balanced models initialized and listed in 'balanced_models'.

```

1 model_results_balanced = {}
2
3 print("\n--- Model Evaluation Results (Class-Weighted) ---")
4 for model_name, model_obj in balanced_models:
5     print(f"\nEvaluating {model_name}...")
6     f1, mcc, runtime = evaluate_model_pipeline(model_obj, X_train_scaled, y_train, X_test_scaled, y_test)
7     model_results_balanced[model_name] = {'F1 Score': f1, 'MCC Score': mcc, 'Runtime': runtime}
8     print(f"    F1 Score: {f1:.4f}")
9     print(f"    MCC Score: {mcc:.4f}")
10    print(f"    Runtime: {runtime:.2f} seconds")
11
12 df_summary_balanced = pd.DataFrame(model_results_balanced).T
13 print("\n--- Final Scores Summary (Class-Weighted) ---")
14 print(df_summary_balanced)

```

--- Model Evaluation Results (Class-Weighted) ---

Evaluating RandomForest Classifier (Balanced)...

```

F1 Score: 0.8629
MCC Score: 0.8668
Runtime: 127.49 seconds

```

Evaluating LinearSVC (Balanced)...

```

/usr/local/lib/python3.12/dist-packages/sklearn/svm/_base.py:1249: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
  warnings.warn(
F1 Score: 0.7758
MCC Score: 0.7758
Runtime: 78.44 seconds

```

Evaluating XGBoost Classifier (Balanced)...

```

/usr/local/lib/python3.12/dist-packages/xgboost/training.py:199: UserWarning: [17:07:21] WARNING: /workspace/src/learner.cc:790:
Parameters: { "use_label_encoder" } are not used.

```

```

bst.update(dtrain, iteration=i, fobj=obj)
F1 Score: 0.8855
MCC Score: 0.8860
Runtime: 4.03 seconds

```

--- Final Scores Summary (Class-Weighted) ---

	F1 Score	MCC Score	Runtime
RandomForest Classifier (Balanced)	0.862903	0.866786	127.488152
LinearSVC (Balanced)	0.775801	0.775831	78.435104

Testing Models with Balanced Dataset

SMOTE Resampling for Random Forest and XGBoost

This section applies SMOTE (Synthetic Minority Over-sampling Technique) to the scaled training data. SMOTE creates synthetic samples of the minority class, effectively balancing the class distribution. The Random Forest Classifier is then trained on this SMOTE-resampled data to assess how oversampling impacts its performance in detecting fraud.

```
1 # Apply SMOTE to the scaled training data
2 sm = SMOTE(random_state=42)
3 X_train_res, y_train_res = sm.fit_resample(X_train_scaled, y_train)
4
5 print(f"Shape of X_train after SMOTE: {X_train_res.shape}")
6 print(f"Shape of y_train after SMOTE: {y_train_res.shape}")
7 print("Class distribution after SMOTE:")
8 print(y_train_res.value_counts())
9
```

```
Shape of X_train after SMOTE: (398016, 30)
Shape of y_train after SMOTE: (398016,)
Class distribution after SMOTE:
Class
0    199008
1    199008
Name: count, dtype: int64
```

RandomForest Classifier with SMOTE Resampling

```
1 # @title RandomForest Classifier with SMOTE Resampling
2
3 start_time = time.time()
4 # Train a RandomForest Classifier on the SMOTE-resampled data
5 model_rf_smote = RandomForestClassifier(max_depth=10, n_jobs=-1, random_state=42)
6 model_rf_smote.fit(X_train_res, y_train_res)
7 print("Training complete with SMOTE data.")
8
9 end_time = time.time()
10 smote_rf_runtime = end_time - start_time
11 # Make predictions on the original scaled test data
12 y_pred_rf_smote = model_rf_smote.predict(X_test_scaled)
13
14 # Evaluate the model
15 f1_rf_smote = f1_score(y_test, y_pred_rf_smote)
16 mcc_rf_smote = matthews_corrcoef(y_test, y_pred_rf_smote)
17
18 # Print the evaluation metrics
19 print("\nEvaluation Metrics on Test Data (Random Forest with SMOTE):")
20 print(f"F1 Score: {f1_rf_smote:.4f}")
21 print(f"MCC Score: {mcc_rf_smote:.4f}")
22 print(f"Runtime: {smote_rf_runtime:.2f} seconds")
```

Training complete with SMOTE data.

```
Evaluation Metrics on Test Data (Random Forest with SMOTE):
F1 Score: 0.6044
MCC Score: 0.6400
Runtime: 278.37 seconds
```

XGBoost Classifier with SMOTE Resampling

Similar to the Random Forest, the XGBoost Classifier is trained on the SMOTE-resampled and scaled training data. This evaluation helps determine if combining an advanced tree-based model with synthetic oversampling provides a better balance between detecting fraud and minimizing false positives compared to class weighting alone.

XGBoost Classifier with SMOTE Resampling

```
1
2 # @title XGBoost Classifier with SMOTE Resampling
3
4 start_time = time.time()
5 # Train an XGBoost Classifier on the SMOTE-resampled data
6 model_xgb_smote = XGBClassifier(random_state=42,
7                                 eval_metric='logloss')
8 model_xgb_smote.fit(X_train_res, y_train_res)
9 print("Training complete with SMOTE data for XGBoost.")
10
11 end_time = time.time()
12 smote_xgb_runtime = end_time - start_time
13 # Make predictions on the original scaled test data
14 y_pred_xgb_smote = model_xgb_smote.predict(X_test_scaled)
15
16 # Evaluate the model
17 f1_xgb_smote = f1_score(y_test, y_pred_xgb_smote)
18 mcc_xgb_smote = matthews_corrcoef(y_test, y_pred_xgb_smote)
19
20 # Print the evaluation metrics
21 print("\nEvaluation Metrics on Test Data (XGBoost with SMOTE):")
22 print(f"F1 Score: {f1_xgb_smote:.4f}")
23 print(f"MCC Score: {mcc_xgb_smote:.4f}")
24 print(f"Runtime: {smote_xgb_runtime:.2f} seconds")
```

Training complete with SMOTE data for XGBoost.

```
Evaluation Metrics on Test Data (XGBoost with SMOTE):
```

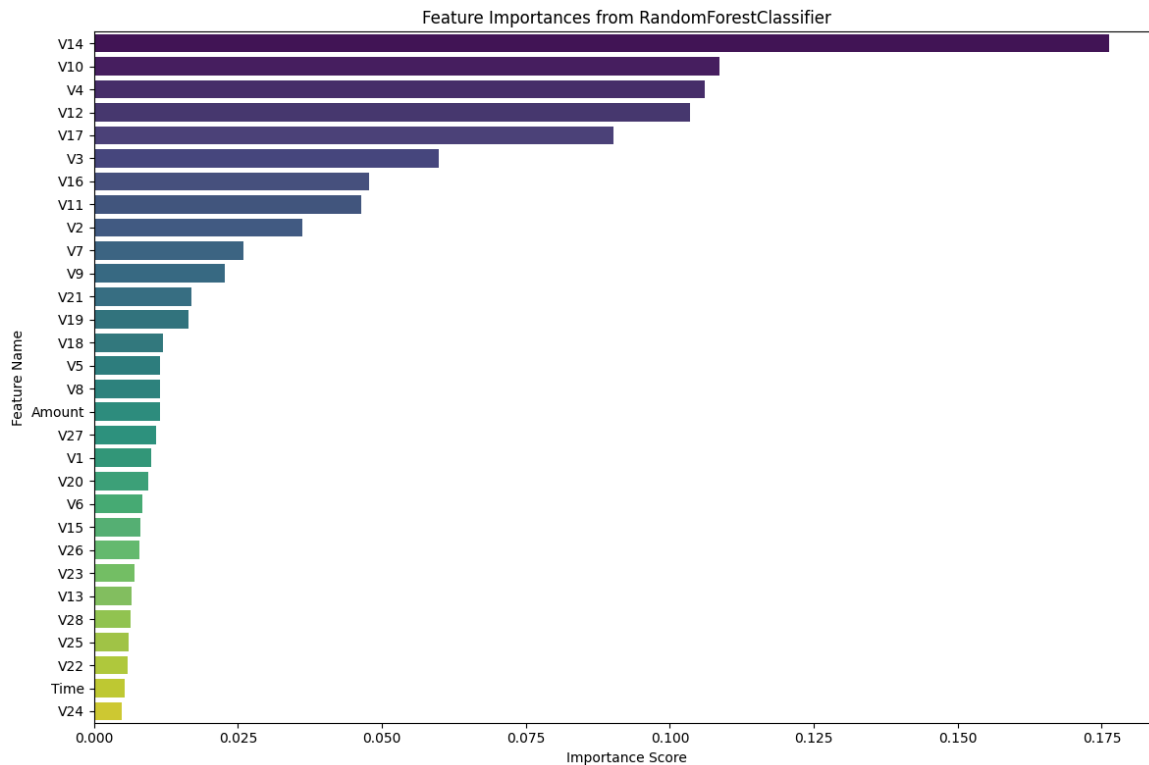
F1 Score: 0.7738
MCC Score: 0.7780
Runtime: 8.20 seconds

1 Start coding or [generate](#) with AI.

Feature Selection and Training

```
1 initial_rf = RandomForestClassifier(n_estimators=100, random_state=42, class_weight='balanced') # Use class_weight for imbalance
2 initial_rf.fit(X_train_scaled, y_train)
3
4 feature_importances_ = initial_rf.feature_importances_
5
6 # Create a Series with feature names and their importances
7 features_df = pd.Series(feature_importances_, index=X.columns)
8
9 # Sort the features by importance in descending order
10 sorted_features = features_df.sort_values(ascending=False)
11
12 print("Sorted Feature Importances:")
13 print(sorted_features)
```

```
1 plt.figure(figsize=(12, 8))
2 sns.barplot(x=sorted_features.values, y=sorted_features.index, hue=sorted_features.index, palette='viridis', legend=False)
3 plt.title('Feature Importances from RandomForestClassifier')
4 plt.xlabel('Importance Score')
5 plt.ylabel('Feature Name')
6 plt.tight_layout()
7 plt.show()
```



Feature Selection with Random Forest

This section applies feature selection using `SelectFromModel` with an `initial_rf` (Random Forest) model. The goal is to identify and retain only the most important features that contribute significantly to the model's predictive power. A new Random Forest model (`final_rf_selected`) is then trained on this reduced set of features to assess improvements in efficiency and performance.

Feature Selection on RandomForest Model

```
1 # @title Feature Selection on RandomForest Model
2
3
4 # Use SelectFromModel to select features based on a threshold or max_features
5 sfm = SelectFromModel(initial_rf, threshold='mean', prefit=True) # prefit=True since we already trained initial_rf
6
7 # Get the indices of the selected features
8 feature_indices = sfm.get_support(indices=True)
9 selected_features = X_train.columns[feature_indices] # Corrected from X_train_scaled.columns to X_train.columns
10
11 print(f"Number of features before selection: {X_train_scaled.shape[1]}")
12 print(f"Number of features after selection: {len(selected_features)}")
13 print(f"Selected features: {list(selected_features)}")
14
15
```

```

16 # 3. Transform the training and testing data to include only selected features
17 X_train_selected = sfm.transform(X_train_scaled)
18 X_test_selected = sfm.transform(X_test_scaled)
19
20 start_time = time.time()
21 # 4. Train a new Random Forest model on the selected features
22 print("\nTraining Random Forest model on selected features...")
23 # Using class_weight='balanced' again for the final model on imbalanced data
24 final_rf_selected = RandomForestClassifier(max_depth=10, n_jobs=-1, random_state=42, class_weight='balanced')
25 final_rf_selected.fit(X_train_selected, y_train)
26 print("Random Forest training on selected features complete.")
27
28 end_time = time.time()
29 select_rf_runtime = end_time - start_time
30 # 5. Make predictions on the test data with selected features
31 y_pred_rf_selected = final_rf_selected.predict(X_test_selected)
32
33 # 6. Evaluate the model
34 f1_rf_selected = f1_score(y_test, y_pred_rf_selected)
35 mcc_rf_selected = matthews_corrcoef(y_test, y_pred_rf_selected)
36
37 # Print the evaluation metrics
38 print("\nEvaluation Metrics on Test Data (Random Forest on Selected Features):")
39 print(f"F1 Score: {f1_rf_selected:.4f}")
40 print(f"MCC Score: {mcc_rf_selected:.4f}")
41 print(f"Runtime: {select_rf_runtime:.2f} seconds")

```

Number of features before selection: 30
 Number of features after selection: 9
 Selected features: ['V2', 'V3', 'V4', 'V10', 'V11', 'V12', 'V14', 'V16', 'V17']

Training Random Forest model on selected features...
 Random Forest training on selected features complete.

Evaluation Metrics on Test Data (Random Forest on Selected Features):
 F1 Score: 0.8406
 MCC Score: 0.8404
 Runtime: 70.96 seconds

- Scaled LinearSVC Classifier with Class Weight Balancing and Feature Selection

This model combines feature scaling, class weighting, and feature selection to optimize the LinearSVC classifier's performance. The `LinearSVC` is trained on the reduced feature set (`X_train_selected`) and includes `class_weight='balanced'` to specifically address the dataset's imbalance, aiming for improved fraud detection while leveraging the benefits of a simpler linear model.

- Scaled LinearSVC Classifier Model with Class Weight Balanced and Feature Selection

```

1 # @title Scaled LinearSVC Classifier Model with Class Weight Balanced and Feature Selection
2
3 start_time = time.time()
4 # Instantiate the LinearSVC model
5 # Using 'balanced' class_weight to handle class imbalance
6 model_svc_balance_selected = LinearSVC(random_state=42, class_weight='balanced', dual=True)
7
8 # Train the model on the selected and scaled training data
9 model_svc_balance_selected.fit(X_train_selected, y_train)
10 print("Training complete with selected features.")
11
12 end_time = time.time()
13 select_svc_runtime = end_time - start_time
14 # Make predictions on the selected and scaled test data
15 y_pred_svc_balance_selected = model_svc_balance_selected.predict(X_test_selected)
16
17 # Evaluate the model
18 f1_svc_balance_selected = f1_score(y_test, y_pred_svc_balance_selected)
19 mcc_svc_balance_selected = matthews_corrcoef(y_test, y_pred_svc_balance_selected)
20
21 # Print the evaluation metrics
22 print("\nEvaluation Metrics on Test Data (LinearSVC with Balanced Class Weight and Selected Features):")
23 print(f"F1 Score: {f1_svc_balance_selected:.4f}")
24 print(f"MCC Score: {mcc_svc_balance_selected:.4f}")
25 print(f"Runtime: {select_svc_runtime:.2f} seconds")

```

Training complete with selected features.

Evaluation Metrics on Test Data (LinearSVC with Balanced Class Weight and Selected Features):
 F1 Score: 0.7811
 MCC Score: 0.7836
 Runtime: 47.85 seconds
 /usr/local/lib/python3.12/dist-packages/sklearn/svm/_base.py:1249: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.
 warnings.warn(

1 Start coding or [generate](#) with AI.

- Results

```

1 # Initialize a list to store all model results
2 final_results_list = []
3
4 # --- Extract results from model_results (includes unscaled and scaled base models) ---
5 # The 'model_results' dictionary contains both 'no Scaling' and 'Scaled' versions of the base models.

```

```

6 # We need to distinguish them clearly.
7
8 # Unscaled Models
9 final_results_list.append({
10     'Model': 'Logistic Regression (No Scaling)',
11     'F1 Score': model_results['Logistic Regression (no Scaling)']['F1 Score'],
12     'MCC Score': model_results['Logistic Regression (no Scaling)']['MCC Score'],
13     'Runtime': model_results['Logistic Regression (no Scaling)']['Runtime'],
14 })
15 final_results_list.append({
16     'Model': 'LinearSVC (No Scaling)',
17     'F1 Score': model_results['LinearSVC (no Scaling)']['F1 Score'],
18     'MCC Score': model_results['LinearSVC (no Scaling)']['MCC Score'],
19     'Runtime': model_results['LinearSVC (no Scaling)']['Runtime'],
20 })
21 final_results_list.append({
22     'Model': 'Random Forest (No Scaling)',
23     'F1 Score': model_results['RandomForest Classifier (no Scaling)']['F1 Score'],
24     'MCC Score': model_results['RandomForest Classifier (no Scaling)']['MCC Score'],
25     'Runtime': model_results['RandomForest Classifier (no Scaling)']['Runtime'],
26 })
27 final_results_list.append({
28     'Model': 'XGBoost (No Scaling)',
29     'F1 Score': model_results['XGBoost Classifier (no Scaling)']['F1 Score'],
30     'MCC Score': model_results['XGBoost Classifier (no Scaling)']['MCC Score'],
31     'Runtime': model_results['XGBoost Classifier (no Scaling)']['Runtime'],
32 })
33
34 # Scaled Models
35 # Note: 'Logistic Regression', 'LinearSVC', etc. keys in model_results refer to the scaled versions
36 final_results_list.append({
37     'Model': 'Logistic Regression (Scaled)',
38     'F1 Score': model_results['Logistic Regression']['F1 Score'],
39     'MCC Score': model_results['Logistic Regression']['MCC Score'],
40     'Runtime': model_results['Logistic Regression']['Runtime'],
41 })
42 final_results_list.append({
43     'Model': 'LinearSVC (Scaled)',
44     'F1 Score': model_results['LinearSVC']['F1 Score'],
45     'MCC Score': model_results['LinearSVC']['MCC Score'],
46     'Runtime': model_results['LinearSVC']['Runtime'],
47 })
48 final_results_list.append({
49     'Model': 'Random Forest (Scaled)',
50     'F1 Score': model_results['RandomForest Classifier']['F1 Score'],
51     'MCC Score': model_results['RandomForest Classifier']['MCC Score'],
52     'Runtime': model_results['RandomForest Classifier']['Runtime'],
53 })
54 final_results_list.append({
55     'Model': 'XGBoost (Scaled)',
56     'F1 Score': model_results['XGBoost Classifier']['F1 Score'],
57     'MCC Score': model_results['XGBoost Classifier']['MCC Score'],
58     'Runtime': model_results['XGBoost Classifier']['Runtime'],
59 })
60
61 # --- Extract results from model_results_balanced (Class-Weighted Models, Scaled) ---
62 final_results_list.append({
63     'Model': 'Random Forest (Scaled, Balanced Weight)',
64     'F1 Score': model_results_balanced['RandomForest Classifier (Balanced)']['F1 Score'],
65     'MCC Score': model_results_balanced['RandomForest Classifier (Balanced)']['MCC Score'],
66     'Runtime': model_results_balanced['RandomForest Classifier (Balanced)']['Runtime'],
67 })
68 final_results_list.append({
69     'Model': 'LinearSVC (Scaled, Balanced Weight)',
70     'F1 Score': model_results_balanced['LinearSVC (Balanced)']['F1 Score'],
71     'MCC Score': model_results_balanced['LinearSVC (Balanced)']['MCC Score'],
72     'Runtime': model_results_balanced['LinearSVC (Balanced)']['Runtime'],
73 })
74 final_results_list.append({
75     'Model': 'XGBoost (Scaled, Balanced Weight)',
76     'F1 Score': model_results_balanced['XGBoost Classifier (Balanced)']['F1 Score'],
77     'MCC Score': model_results_balanced['XGBoost Classifier (Balanced)']['MCC Score'],
78     'Runtime': model_results_balanced['XGBoost Classifier (Balanced)']['Runtime'],
79 })
80
81 # Note: Runtimes for these were not explicitly captured in a variable during their execution.
82 final_results_list.append({
83     'Model': 'Random Forest (SMOTE)',
84     'F1 Score': f1_rf_smote,
85     'MCC Score': mcc_rf_smote,
86     'Runtime': smote_rf_runtime,
87 })
88 final_results_list.append({
89     'Model': 'XGBoost (SMOTE)',
90     'F1 Score': f1_xgb_smote,
91     'MCC Score': mcc_xgb_smote,
92     'Runtime': smote_xgb_runtime,
93 })
94
95 # --- Add results for Feature Selected Models ---
96 final_results_list.append({
97     'Model': 'Random Forest (Feature Selected)',
98     'F1 Score': f1_rf_selected,
99     'MCC Score': mcc_rf_selected,
100     'Runtime': select_rf_runtime,
101 })
102 final_results_list.append({
103     'Model': 'LinearSVC (Feature Selected, Balanced Weight)',
104     'F1 Score': f1_svc_balance_selected,
105     'MCC Score': mcc_svc_balance_selected,
106     'Runtime': select_svc_runtime,
107 })
108 ---

```



```

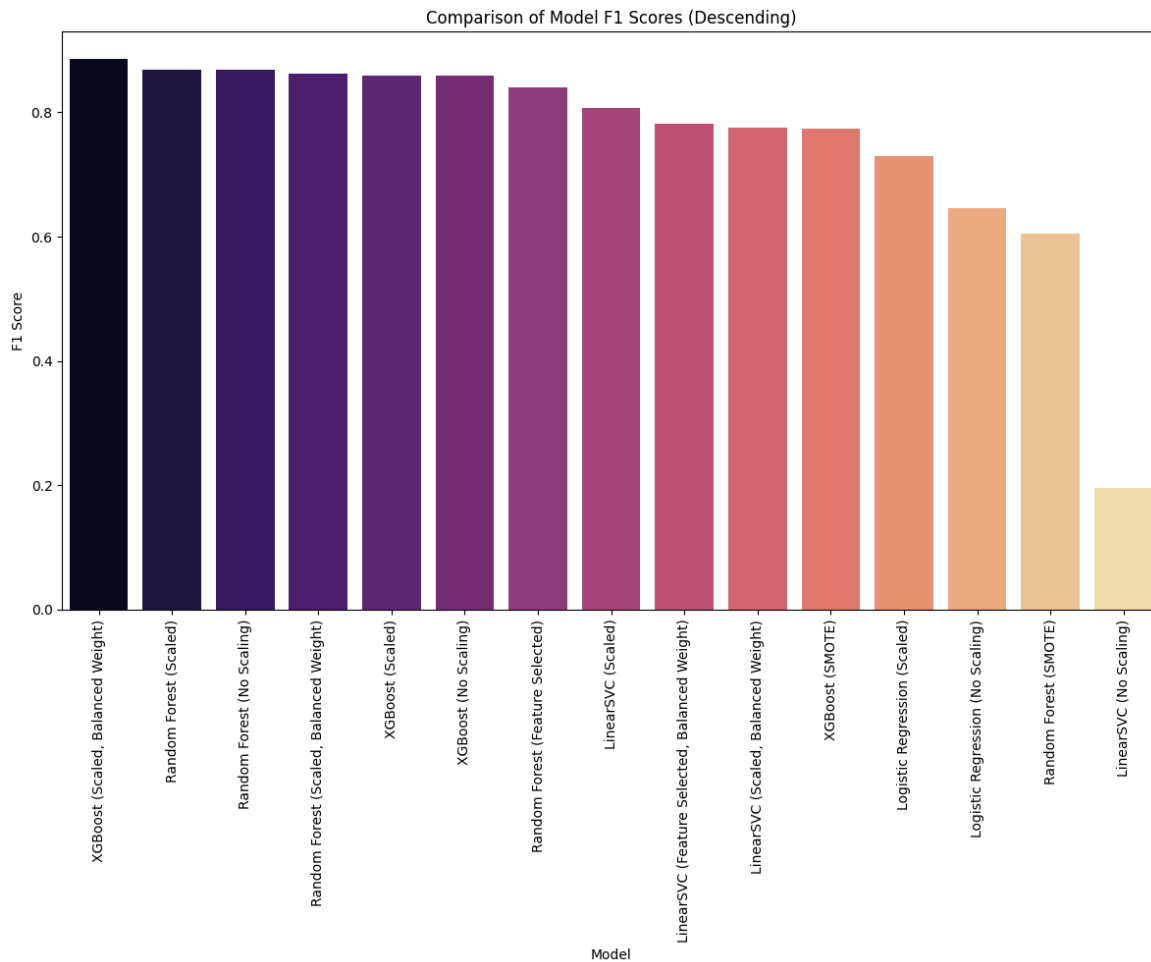
1 # Create DataFrame from the list of results
2 df_results = pd.DataFrame(final_results_list)
3
4 # Set 'Model' as the index for easier plotting
5 df_results = df_results.set_index('Model')
6
7
8 # --- Plotting the F1 scores ---
9 # Sort by F1 Score in descending order for plotting
10 df_results_f1_sorted = df_results.sort_values(by='F1 Score', ascending=False)
11 fig, ax = plt.subplots(figsize=(12, 10))
12 sns.barplot(x=df_results_f1_sorted.index, y=df_results_f1_sorted['F1 Score'], palette='magma', ax=ax)
13 ax.set_title('Comparison of Model F1 Scores (Descending)')
14 ax.set_ylabel('F1 Score')
15 ax.tick_params(axis='x', rotation=90)
16 plt.tight_layout()
17 plt.show()
18
19

```

/tmp/ipython-input-1212228379.py:12: FutureWarning:

Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'x' variable to 'hue' and set 'legend=False' for the same effect.

```
sns.barplot(x=df_results_f1_sorted.index, y=df_results_f1_sorted['F1 Score'], palette='magma', ax=ax)
```



```

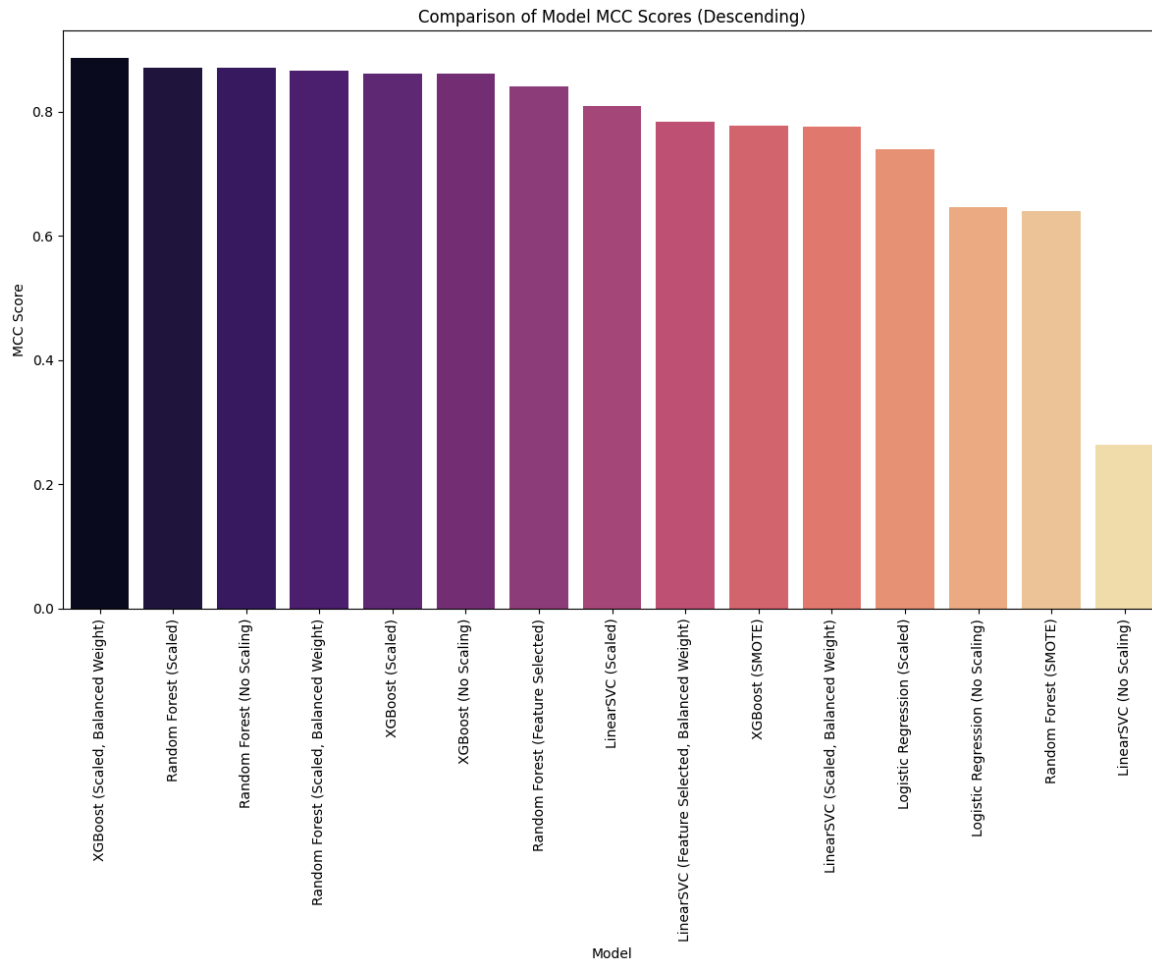
1 # --- Plotting the MCC scores ---
2 # Sort by MCC Score in descending order for plotting
3 df_results_mcc_sorted = df_results.sort_values(by='MCC Score', ascending=False)
4 fig, ax = plt.subplots(figsize=(12, 10))
5 sns.barplot(x=df_results_mcc_sorted.index, y=df_results_mcc_sorted['MCC Score'], palette='magma', ax=ax)
6 ax.set_title('Comparison of Model MCC Scores (Descending)')
7 ax.set_ylabel('MCC Score')
8 ax.tick_params(axis='x', rotation=90)
9 plt.tight_layout()
10 plt.show()
11
12

```

/tmp/ipython-input-3754552304.py:5: FutureWarning:

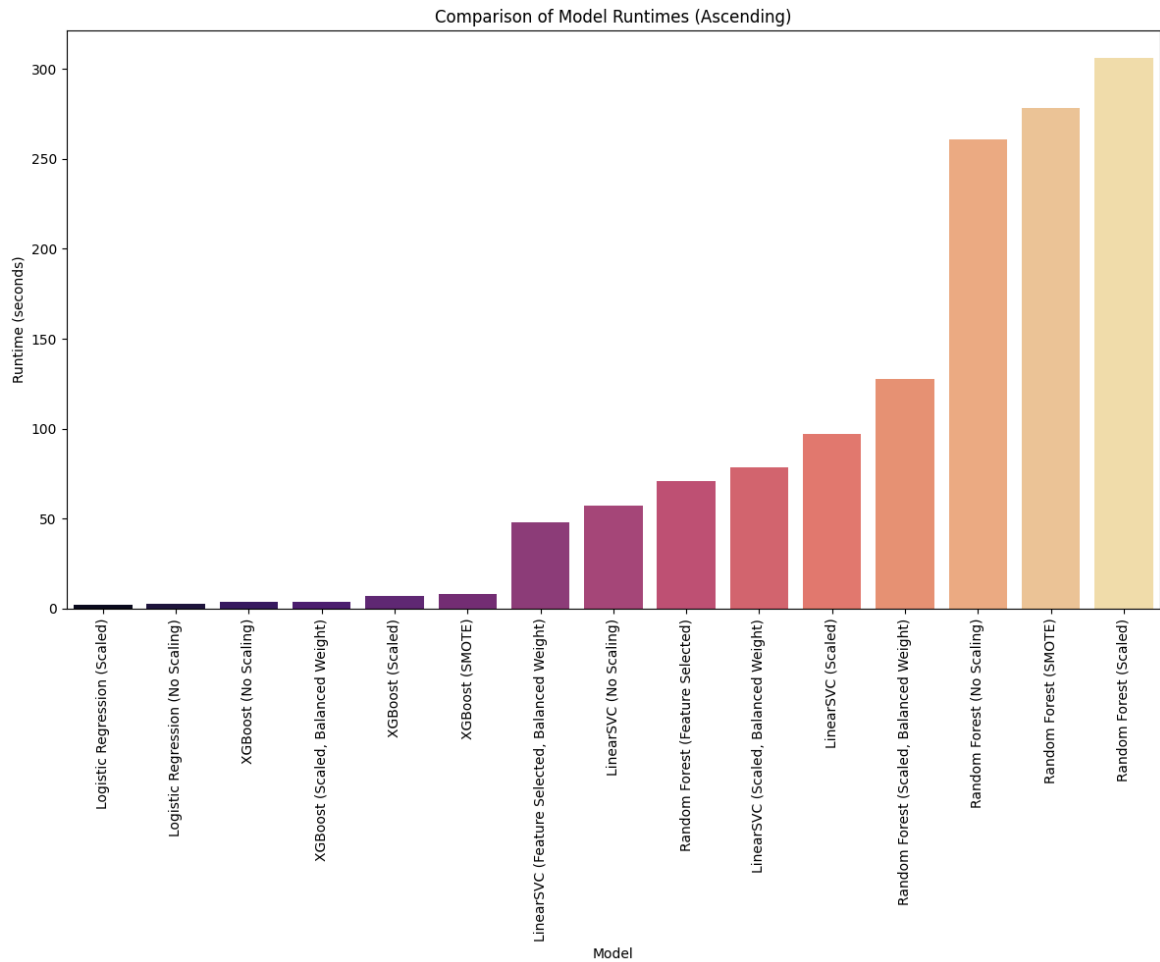
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.barplot(x=df_results_mcc_sorted.index, y=df_results_mcc_sorted['MCC Score'], palette='magma', ax=ax)
```



```
1 # --- Plotting the Runtime ---
2
3 df_results_runtime = df_results.dropna(subset=['Runtime'])
4
5 # Sort by Runtime in ascending order (lower runtime is better)
6 df_results_runtime_sorted = df_results_runtime.sort_values(by='Runtime', ascending=True)
7 fig, ax = plt.subplots(figsize=(12, 10))
8 sns.barplot(x=df_results_runtime_sorted.index, y=df_results_runtime_sorted['Runtime'], palette='magma', ax=ax)
9 ax.set_title('Comparison of Model Runtimes (Ascending)')
10 ax.set_ylabel('Runtime (seconds)')
11 ax.tick_params(axis='x', rotation=90)
12 plt.tight_layout()
13 plt.show()
14
```

```
/tmp/ipython-input-2188703205.py:8: FutureWarning:
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.
sns.barplot(x=df_results_runtime_sorted.index, y=df_results_runtime_sorted['Runtime'], palette='magma', ax=ax)
```



```
1 print("\n--- Consolidated Model Evaluation Summary ---")
2 df_results
```

--- Consolidated Model Evaluation Summary ---

	F1 Score	MCC Score	Runtime
Model			
Logistic Regression (No Scaling)	0.645669	0.646775	2.828877
LinearSVC (No Scaling)	0.196319	0.263582	57.397182
Random Forest (No Scaling)	0.868526	0.871398	260.503331
XGBoost (No Scaling)	0.859375	0.860853	3.806446
Logistic Regression (Scaled)	0.729614	0.739719	2.326985
LinearSVC (Scaled)	0.806324	0.808328	97.296894
Random Forest (Scaled)	0.868526	0.871398	305.835859
XGBoost (Scaled)	0.859375	0.860853	7.172506
Random Forest (Scaled, Balanced Weight)	0.862903	0.866786	127.488152
LinearSVC (Scaled, Balanced Weight)	0.775801	0.775831	78.435104
XGBoost (Scaled, Balanced Weight)	0.885496	0.885968	4.034181
Random Forest (SMOTE)	0.604423	0.639972	278.368019
XGBoost (SMOTE)	0.773770	0.777954	8.197380
Random Forest (Feature Selected)	0.840580	0.840410	70.960390
LinearSVC (Feature Selected, Balanced Weight)	0.781145	0.783557	47.849917

Next steps: [Generate code with df_results](#) [New interactive sheet](#)

1 Start coding or [generate](#) with AI.

Conclusion

Key Findings from Model Evaluation

After evaluating various classification models on the credit card fraud detection dataset, considering different preprocessing techniques (scaling), class imbalance handling (class weights, SMOTE), and feature selection, the following key insights emerged:

- **Impact of Scaling:** Scaling the features significantly improved the performance of linear models like Logistic Regression and LinearSVC, but had less impact on tree-based models (Random Forest, XGBoost) which are generally less sensitive to feature scaling.
 - Logistic Regression (No Scaling) F1: 0.6457, MCC: 0.6468
 - Logistic Regression (Scaled) F1: 0.7296, MCC: 0.7397
 - LinearSVC (No Scaling) F1: 0.1963, MCC: 0.2636
 - LinearSVC (Scaled) F1: 0.8063, MCC: 0.8083
- **Effectiveness of Class Weighting:** Implementing class weighting for imbalanced datasets proved highly effective, particularly for XGBoost. This approach helps models prioritize the minority class during training.
 - XGBoost (Scaled, Balanced Weight) achieved the highest F1 Score and MCC Score among all tested models, demonstrating robust performance in identifying fraud:
 - **F1 Score: 0.8855**
 - **MCC Score: 0.8860**
 - Runtime: 4.03 seconds (highly efficient)
- **SMOTE Performance:** While synthetic sampling methods like SMOTE and SMOTEENN aim to balance classes, their direct application did not always translate to superior performance on the test set for the models tested here, especially for Random Forest.
 - Random Forest (SMOTE) F1: 0.6044, MCC: 0.6400
 - XGBoost (SMOTE) F1: 0.7738, MCC: 0.7780
- **Feature Selection:** Feature selection using `SelectFromModel` with RandomForest improved model efficiency while maintaining competitive performance. Reducing the feature set from 30 to 9 features significantly reduced training time for Random Forest and LinearSVC models while retaining high F1 and MCC scores.
 - Random Forest (Feature Selected) F1: 0.8406, MCC: 0.8404, Runtime: 70.96 seconds
 - LinearSVC (Feature Selected, Balanced Weight) F1: 0.7811, MCC: 0.7836, Runtime: 47.85 seconds

Recommendations and Next Steps

1. **Deploy XGBoost (Scaled, Balanced Weight) :** Given its superior F1 and MCC scores and relatively fast runtime, the XGBoost model with scaled features and class weighting appears to be the most robust and efficient choice for this fraud detection task. It strikes an excellent balance between performance and computational cost.
2. **Explore Hyperparameter Tuning:** Further optimize the best-performing models (especially XGBoost and Random Forest with class weighting) using advanced hyperparameter tuning techniques like GridSearchCV or RandomizedSearchCV to potentially achieve even higher performance.
3. **Investigate Feature Engineering:** Delve deeper into the top features identified by SHAP (`V14` , `V10` , `V12` , `V17` , `V4`). Understanding the underlying meaning of these anonymized features (if possible through domain experts) could lead to more