

Model Evaluation for Election Campaign Data

I built and compared Logistic Regression, Decision Tree, and Random Forest models on the election campaign data to predict the outcome of a general election (W - Won, L - Lost).

I loaded the dataset and dropped the following columns: `cand_id`, `last_name`, `first_name`, `twitterbirth`, `facebookdate`, `facebookjan`, `youtubebirth`. The target variable is `gen_election`.

Preprocessing Steps

1. Handling Missing Values

- I input numerical columns with the median.
- I imputed categorical columns with the most frequent value.

2. Encoding Categorical Variables

- I applied encoding to categorical variables.

3. Feature Scaling

- I scaled the features using Standard Scaler.

Evaluation Metrics

Model	Accuracy	Precision	Recall	F1 Score
Logistic Regression	0.919492	0.96875	0.853211	0.907317
Decision Tree	0.936441	0.951923	0.912281	0.931507
Random Forest	0.949153	0.980198	0.928571	0.953125

Conclusion

I chose the **Random Forest** model is the best model based on its highest accuracy, precision, and F1 score. These metrics indicate that the Random Forest model has a good balance between precision and recall, making it a reliable choice for predicting the election outcomes.

Assignment4 > Assignment4_Seerat > ...

```
1  import pandas as pd
2  from sklearn.model_selection import train_test_split
3  from sklearn.preprocessing import StandardScaler, OneHotEncoder
4  from sklearn.compose import ColumnTransformer
5  from sklearn.pipeline import Pipeline
6  from sklearn.impute import SimpleImputer
7  from sklearn.linear_model import LogisticRegression
8  from sklearn.tree import DecisionTreeClassifier
9  from sklearn.ensemble import RandomForestClassifier
10 from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score
11 # Correct file path by using a raw string or double backslashes
12 file_path = r'C:\Assignment 4\R\Classification\election_campaign_data.csv'
13 # or
14 # file_path = 'C:\\Assignment 4\\R\\Classification\\election_campaign_data.csv'
15
16 # Load the dataset
17 dataset = pd.read_csv(file_path)
18
19 # Drop the specified variables
20 columns_to_drop = ['cand_id', 'last_name', 'first_name', 'twitterbirth', 'facebookdate', 'facebookjan', 'youtubebirth']
21 dataset = dataset.drop(columns=columns_to_drop)
22
23 # Define the output variable
24 dataset['gen_election'] = dataset['gen_election'].map({'W': 1, 'L': 0})
25
26 # Separate features and target variable
27 X = dataset.drop(columns=['gen_election'])
28 y = dataset['gen_election']
29
30 # Handle missing values and encode categorical variables using pipelines
31 numeric_features = X.select_dtypes(include=['int64', 'float64']).columns
32 categorical_features = X.select_dtypes(include=['object']).columns
33
34 numeric_transformer = Pipeline(steps=[
35     ('imputer', SimpleImputer(strategy='median')),
36     ('scaler', StandardScaler())
37 ])
38
39 categorical_transformer = Pipeline(steps=[
40     ('imputer', SimpleImputer(strategy='most_frequent')),
41     ('onehot', OneHotEncoder(handle_unknown='ignore'))
42 ])
43
44 preprocessor = ColumnTransformer(
45     transformers=[
46         ('num', numeric_transformer, numeric_features),
47         ('cat', categorical_transformer, categorical_features)
48     ])
49
50 # Split the dataset into the Training set and Test set
```

```
50 # Split the dataset into the Training set and Test set
51 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=123)
52
53 # Apply the preprocessor to the training and test data
54 X_train = preprocessor.fit_transform(X_train)
55 X_test = preprocessor.transform(X_test)
56
57 # Build and evaluate models
58 results = []
59
60 # Logistic Regression
61 logistic_model = LogisticRegression()
62 logistic_model.fit(X_train, y_train)
63 y_pred_logistic = logistic_model.predict(X_test)
64 results.append(('Logistic Regression', y_test, y_pred_logistic))
65
66 # Decision Tree
67 tree_model = DecisionTreeClassifier()
68 tree_model.fit(X_train, y_train)
69 y_pred_tree = tree_model.predict(X_test)
70 results.append(('Decision Tree', y_test, y_pred_tree))
71
72 # Random Forest
73 forest_model = RandomForestClassifier()
74 forest_model.fit(X_train, y_train)
75 y_pred_forest = forest_model.predict(X_test)
76 results.append(('Random Forest', y_test, y_pred_forest))
77
78 # Evaluate models
79 evaluation_metrics = []
80
81 for model_name, y_true, y_pred in results:
82     accuracy = accuracy_score(y_true, y_pred)
83     precision = precision_score(y_true, y_pred)
84     recall = recall_score(y_true, y_pred)
85     f1 = f1_score(y_true, y_pred)
86     evaluation_metrics.append((model_name, accuracy, precision, recall, f1))
87
88 evaluation_df = pd.DataFrame(evaluation_metrics, columns=['Model', 'Accuracy', 'Precision', 'Recall', 'F1 Score'])
89
90 print(evaluation_df)
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

2 Random Forest 0.953390 0.990000 0.908257 0.947368
PS C:\Users\8noor\AI-101-Bootcamp-June-2024-1>

& C:/Users/8noor/AppData/Local/Programs/Python/Python312/python.exe c:/Users/8noor/AI-101-Bootcamp-June-2024-1/Assignment4/Assignment4_Seerat

	Model	Accuracy	Precision	Recall	F1 Score
0	Logistic Regression	0.919492	0.968750	0.853211	0.907317
1	Decision Tree	0.915254	0.932039	0.880734	0.905660
2	Random Forest	0.953390	0.990000	0.908257	0.947368

PS C:\Users\8noor\AI-101-Bootcamp-June-2024-1>