

✓ Assignment 2 - Naive Bayes and Logistic Regression

Dataset: Network Attack - NB15

The raw network packets of the **UNSW-NB 15 dataset** was created by the IXIA PerfectStorm tool in the Cyber Range Lab of the Australian Centre for Cyber Security (ACCS) for generating a hybrid of real modern normal activities and synthetic contemporary attack behaviours.

Tcpdump tool is used to capture the raw traffic (e.g., Pcap files). This dataset has nine types of attacks, namely, Fuzzers, Analysis, Backdoors, DoS, Exploits, Generic, Reconnaissance, Shellcode and Worms. The Argus, Bro-IDS tools are used and twelve algorithms are developed to generate totally 49 features with the class label.

- [Dataset Link](#)
- [Description of all the 49 features](#)

✓ Install required libraries

```
1
```

```
1 import pandas
2 import numpy
```

✓ Q1. EDA and basic data pre-processing and preparation (4*0.5 = 2M)

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Mounted at /content/drive

```
1 dataset = pandas.read_csv("/content/drive/MyDrive/bits-pillani/PG_AIML/Cassification/assignment2/Network_attacks_UNSW_NI
```

```
1 dataset["attack_cat"].unique()

array(['Normal', 'Backdoor', 'Analysis', 'Fuzzers', 'Shellcode',
       'Reconnaissance', 'Exploits', 'DoS', 'Worms', 'Generic'],
      dtype=object)
```

✓ a. Null/Outlier treatment

```
1 clean_dataset = dataset.dropna()
2
```

✓ b. Remove non-important features

```
1 from sklearn.feature_selection import SelectKBest, mutual_info_classif
```

```
1 clean_dataset = clean_dataset.drop(columns = ["id", "state"])
```

✓ c. Split training and testing data set

```
1 from sklearn.model_selection import train_test_split
2 X = clean_dataset.drop(columns = ["attack_cat"])
3 y = clean_dataset[["attack_cat"]]
4
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
6
7
```

✓ d. Standardise/Normalise the variables whenever required

```

1 from sklearn.preprocessing import StandardScaler
2
3 def std_numeric_columns(df_) :
4     df = df_.copy()
5     numeric_columns = df.select_dtypes(include=['int64', 'float64']).columns
6
7     scaler = StandardScaler()
8     df[numeric_columns] = scaler.fit_transform(df[numeric_columns])
9     return df, scaler
10
11 def std_numeric_cols_test(scaler, df_):
12     numeric_columns = df_.select_dtypes(include=['int64', 'float64']).columns
13     df = df_
14     df[numeric_columns] = std_scaler.transform(df[numeric_columns])
15     return df

```

```

1 def concat(df1_, df2):
2     df1 = df1_.reset_index().drop(columns= ["index"])
3     for col in df2.columns:
4         df1[col] = df2[col]
5     return df1

```

```

1 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
2
3 def label_encoder(df_, column):
4     df = df_.copy()
5     label_encoder = LabelEncoder()
6     df[f"{column}"] = label_encoder.fit_transform(df[column])
7     # df = df.drop(columns = [column])
8     return df, label_encoder
9
10 def encode_text_to_onehot(df_, column):
11     df = df_.copy()
12     hot_encoder = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
13     onehot_encoded = hot_encoder.fit_transform(df[[column]])
14     onehot_encoded_df = pandas.DataFrame(onehot_encoded, columns=hot_encoder.get_feature_names_out([column]))
15     df = concat(df, onehot_encoded_df)
16     df = df.drop(columns = [column])
17     return df, hot_encoder, onehot_encoded_df
18
19 def test_encode_text_to_onehot(df_, column, hot_encoder):
20     df = df_
21     onehot_encoded = hot_encoder.transform(df[[column]])
22     onehot_encoded_df = pandas.DataFrame(onehot_encoded, columns=hot_encoder.get_feature_names_out([column]))
23     df = concat(df, onehot_encoded_df)
24     df = df.drop(columns = [column])
25     return df, onehot_encoded_df
26
27 enc_X_test = X_test
28 enc_X_train = X_train
29
30 # enc_y_train, target_label_encoder = encode_text_to_onehot(y_train, "attack_cat")
31 # enc_y_test = target_label_encoder.transform(y_test)
32 # enc_y_train
33
34 # enc_X_train, proto_hot_encoder, tys = encode_text_to_onehot(X_train, "proto")
35 # enc_X_test["proto"] = proto_hot_encoder.transform(enc_X_test[["proto"]])
36 # enc_X_test, _ = test_encode_text_to_onehot(X_test, "proto", proto_hot_encoder)
37 enc_X_train, proto_label_encoder = label_encoder(enc_X_train, "proto")
38 enc_X_test["proto"] = proto_label_encoder.transform(enc_X_test["proto"])
39
40 # enc_X_train, service_hot_encoder, tys = encode_text_to_onehot(enc_X_train, "service")
41 # enc_X_test["service"] = service_hot_encoder.transform(enc_X_test[["service"]])
42 # enc_X_test, _ = test_encode_text_to_onehot(enc_X_test, "service", service_hot_encoder)
43 enc_X_train, service_label_encoder = label_encoder(enc_X_train, "service")
44 enc_X_test["service"] = service_label_encoder.transform(enc_X_test["service"])
45
46
47 # enc_X_train, state_hot_encoder = encode_text_to_onehot(enc_X_train, "state")
48 # enc_X_test["state"] = state_hot_encoder.transform(enc_X_test[["state"]])

```

```

1 enc_X_train, std_scaler = std_numeric_columns(enc_X_train)
2 enc_X_test = std_numeric_cols_test(std_scaler, enc_X_test)

```

```

1 enc_X_train.dropna()

```

18/02/2024, 07:11

Classification-A2-NBLR-2023aiml066.ipynb - Colaboratory

	dur	proto	service	spkts	dpkts	sbytes	dbytes	rate	sttl	dttl	...	ct_dst_ltm	ct_s
96203	-0.129900	0.153088	1.464739	-0.076656	-0.100649	-0.045975	-0.097123	-0.577575	-1.140489	1.560782	...	-0.395935	
58960	-0.209992	0.465716	-0.701255	-0.137573	-0.174221	-0.049529	-0.105296	0.097298	0.724029	-0.720255	...	-0.395935	
65069	-0.067947	0.153088	-0.701255	-0.015740	-0.063864	-0.014404	-0.100511	-0.577583	0.724029	1.560782	...	-0.644264	
64133	-0.054916	0.153088	1.464739	-0.076656	-0.082256	-0.046166	-0.075515	-0.577660	-1.140489	1.560782	...	-0.644264	
111445	-0.209992	-2.124633	-0.701255	-0.137573	-0.174221	-0.049529	-0.105296	0.637355	0.724029	-0.720255	...	-0.395935	
...
119879	-0.209992	0.421055	0.165143	-0.137573	-0.174221	-0.050044	-0.105296	0.637355	0.724029	-0.720255	...	3.204835	
103694	0.477718	0.153088	0.598341	0.029947	0.028101	-0.043115	-0.093764	-0.577713	-1.140489	1.560782	...	-0.520100	
131932	-0.209992	0.421055	0.165143	-0.137573	-0.174221	-0.050044	-0.105296	0.181682	0.724029	-0.720255	...	1.218203	
146867	-0.209992	0.421055	0.165143	-0.137573	-0.174221	-0.050044	-0.105296	0.941138	0.724029	-0.720255	...	3.701493	
121958	-0.209992	0.421055	0.165143	-0.137573	-0.174221	-0.050044	-0.105296	0.097298	0.724029	-0.720255	...	1.218203	

140272 rows x 41 columns

2. Train a classification model to classify the 9 different attacks (2 + 2 = 4M)

(Use scikit-learn library)

+ Code

+ Text

a. Naive-Bayes

```
1 from sklearn.naive_bayes import GaussianNB
2
3 def train_NB(X_train, y_train):
4     model = GaussianNB()
5     model.fit(X_train, y_train)
6     return model

1 NB_model = train_NB(enc_X_train, y_train)

/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning: A column-vector y was passed as a 1D array, which was interpreted as flattened 2D data if the dimensions are (n_samples, n_features). Please use y = column_or_1d(y, warn=True) instead.
```

b. Logistic Regression

```
1 from sklearn.linear_model import LogisticRegression
2
3 def train_LR(X_train, y_train):
4     model = LogisticRegression()
5     model.fit(X_train, y_train)
6     return model

1
2 LR_model = train_LR(enc_X_train, y_train)

/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning: A column-vector y was passed as a 1D array, which was interpreted as flattened 2D data if the dimensions are (n_samples, n_features). Please use y = column_or_1d(y, warn=True) instead.
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (50 iterations total). Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_i = _check_optimize_result(
```

3. Compute accuracy and F1 scores for both types of models

[1+1 = 2M]

```
1
```

```

1 from sklearn.metrics import accuracy_score, f1_score
2
3 def compute_accuracy_and_f1_score_for_model(model, X_test):
4     test_pred = model.predict(X_test)
5     a_score = accuracy_score(y_test, test_pred)
6     f_score = f1_score(y_test, test_pred, average='weighted')
7     return a_score, f_score
8

```

```

1 nb_accuracy_score, nb_f1_score = compute_accuracy_and_f1_score_for_model(NB_model, enc_X_test)
2 lr_accuracy_score, lr_f1_score = compute_accuracy_and_f1_score_for_model(LR_model, enc_X_test)

```

```

1 print(f"Naive Bayes \nAccuracy Score : {nb_accuracy_score}\nF1 Score : {nb_f1_score}\n\n")
2 print(f"Logistic Regression \nAccuracy Score : {lr_accuracy_score}\nF1 Score : {lr_f1_score}\n\n")

```

```

Naive Bayes
Accuracy Score : 0.4052011748267701
F1 Score : 0.4603245520552212

```

```

Logistic Regression
Accuracy Score : 0.7607858792665887
F1 Score : 0.7418402246262925

```

✓ 4. Apply feature reduction techniques

[1M]

```

1 from sklearn.decomposition import PCA
2 # Performing PCA
3 pca = PCA(n_components=30) # 10 components
4 pca_X_train = pca.fit_transform(enc_X_train)
5 pca_X_test = pca.transform(enc_X_test)
6

```

```
1
```

✓ 5. Train model on the reduced feature subset using NB and Logistic regression

[1+1 = 2M]

✓ a. Naive Bayes

```
1 pca_nb_model = train_NB(pca_X_train, y_train)
```

```

/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning: A column-vector y was passed as a 1D array, which will be deprecated in the future. Use y = column_or_1d(y, warn=True) instead.

```

✓ b. Logistic Regression

```
1 pca_lr_model = train_LR(pca_X_train, y_train)
```

```

/usr/local/lib/python3.10/dist-packages/sklearn/utils/validation.py:1143: DataConversionWarning: A column-vector y was passed as a 1D array, which will be deprecated in the future. Use y = column_or_1d(y, warn=True) instead.
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs failed to converge (50 iterations total). Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
```

```

1 pca_nb_accuracy_score, pca_nb_f1_score = compute_accuracy_and_f1_score_for_model(pca_nb_model, pca_X_test)
2 pca_lr_accuracy_score, pca_lr_f1_score = compute_accuracy_and_f1_score_for_model(pca_lr_model, pca_X_test)

```

```
1
```

```
1 print(f"Naive Bayes \nAccuracy Score : {pca_nb_accuracy_score}\nF1 Score : {pca_nb_f1_score}\n\n")
2 print(f"Logistic Regression \nAccuracy Score : {pca_lr_accuracy_score}\nF1 Score : {pca_lr_f1_score}\n\n")
```

Naive Bayes
Accuracy Score : 0.5066297869913599
F1 Score : 0.5691762439511556

Logistic Regression
Accuracy Score : 0.7500641592289486
F1 Score : 0.7296824554941133

6. Comment on performance of these approaches (NB and Logistic Regression) on reduced dataset (after applying PCA/feature reduction)

[1M]

Read the dataset and removed the null values. Removed the state column. Split the data in train and test.

Trained and tested the data on both Naive Bayes and Logistic regression, Logistic Regression performed very well compare to NB.

After performing PCA, to 30 features there was slight increase in performance of both models, this might due to PCA compacting the important features.

1