**Name: Uwase Ketsia Deborah**
**ID: 26244**

*PART 1: PROBLEM DEFINITION & PLANNING*
*I. Sector Selection*
Cybersecurity
*II. Problem Statement*
How can we use flow-based network data to accurately detect and differentiate between
Tor and Non-Tor traffic in darknet environments?
*III. 3. Dataset Identification*
Dataset Title: Darknet
Source Link: DarkNet
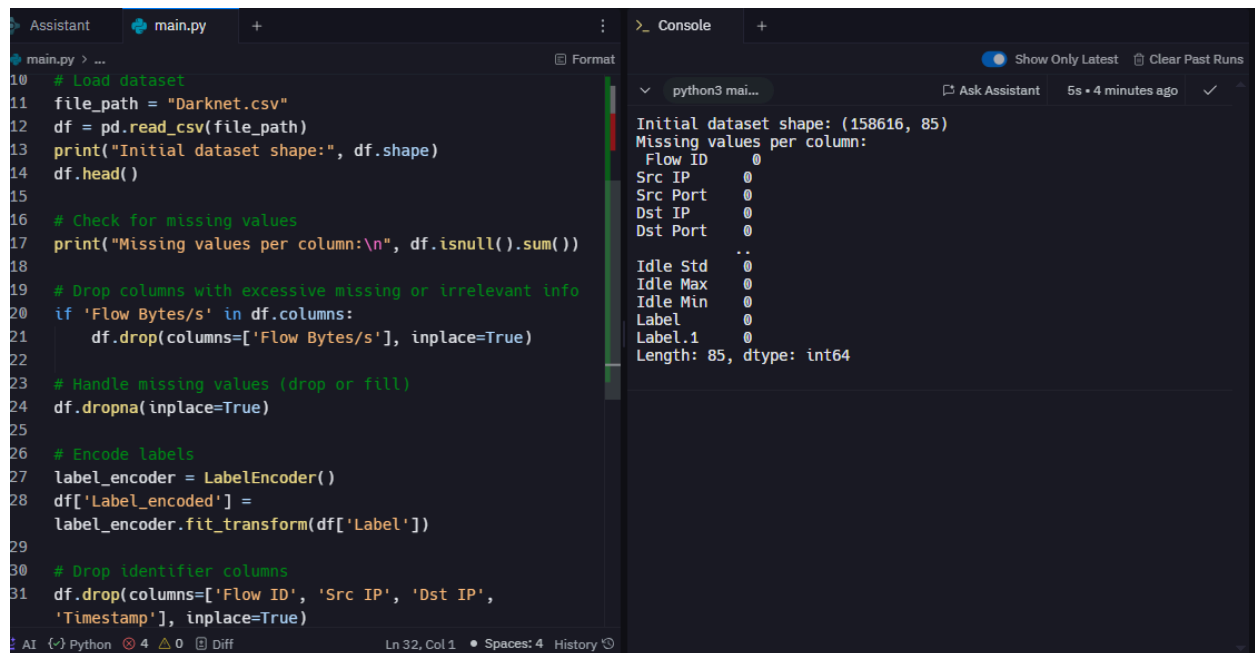Number of Rows and Columns: 158616, 85
Data Structure: csv file
Data Status: needs preprocessing.
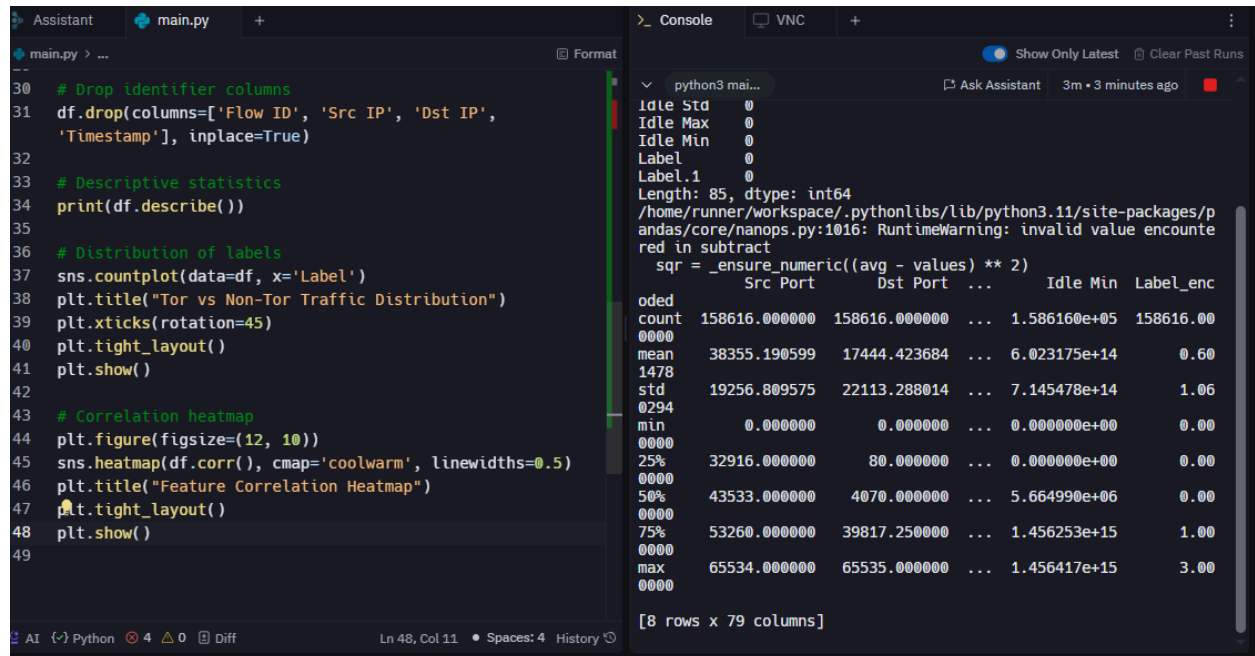
*PART 2: PYTHON ANALYTICS TASKS*
*1.Clean the Dataset*



The code loads the **Darknet dataset** (158,616 rows × 85 columns), checks for missing
values (none found), removes an irrelevant `"Flow Bytes/s"` column, drops any NaN
rows, encodes the `"Label"` column into numbers, and deletes identifier columns like IP
addresses and timestamps to keep only useful features for analysis or modeling.

## 2. Conduct Exploratory Data Analysis (EDA)



```python
30    # Drop identifier columns
31    df.drop(columns=['Flow ID', 'Src IP', 'Dst IP',
      'Timestamp'], inplace=True)
32
33    # Descriptive statistics
34    print(df.describe())
35
36    # Distribution of labels
37    sns.countplot(data=df, x='Label')
38    plt.title("Tor vs Non-Tor Traffic Distribution")
39    plt.xticks(rotation=45)
40    plt.tight_layout()
41    plt.show()
42
43    # Correlation heatmap
44    plt.figure(figsize=(12, 10))
45    sns.heatmap(df.corr(), cmap='coolwarm', linewidths=0.5)
46    plt.title("Feature Correlation Heatmap")
47    plt.tight_layout()
48    plt.show()
49
```

Console output:
```
Idle Std      0
Idle Max      0
Idle Min      0
Label         0
Label.1       0
Length: 85, dtype: int64
/home/runner/workspace/.pythonlibs/lib/python3.11/site-packages/p
andas/core/nanops.py:1016: RuntimeWarning: invalid value encounte
red in subtract
  sqr = _ensure_numeric((avg - values) ** 2)
              Src Port      Dst Port   ...      Idle Min   Label_enc
oded
count  158616.000000  158616.000000   ...  1.586160e+05  158616.00
0000
mean    38355.190599   17444.423684   ...  6.023175e+14       0.60
1478
std     19256.809575   22113.288014   ...  7.145478e+14       1.06
0294
min         0.000000       0.000000   ...  0.000000e+00       0.00
0000
25%     32916.000000      80.000000   ...  0.000000e+00       0.00
0000
50%     43533.000000    4070.000000   ...  5.664990e+06       0.00
0000
75%     53260.000000   39817.250000   ...  1.456253e+15       1.00
0000
max     65534.000000   65535.000000   ...  1.456417e+15       3.00
0000

[8 rows x 79 columns]
```
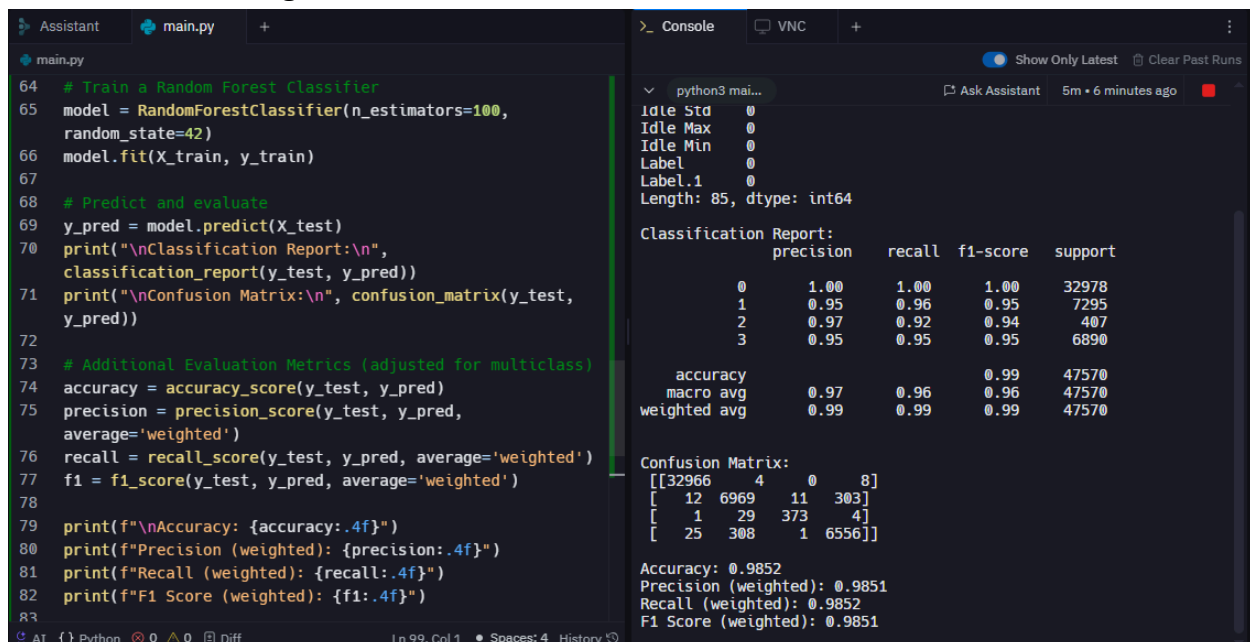
The code shows that after cleaning, the dataset has 79 usable columns and 158,616 records.

The numeric summary (`df.describe()`) gives ranges, averages, and variation for each feature, showing wide value differences between network traffic attributes.

The label distribution plot reveals how balanced Tor vs Non-Tor classes are.

```python
30   # Drop identifier columns
31   df.drop(columns=['Flow ID', 'Src IP', 'Dst IP',
     'Timestamp'], inplace=True)
32
33   # Descriptive statistics
34   print(df.describe())
35
36   # Distribution of labels
37   sns.countplot(data=df, x='Label')
38   plt.title("Tor vs Non-Tor Traffic Distribution")
39   plt.xticks(rotation=45)
40   plt.tight_layout()
41   plt.show()
42
43   # Correlation heatmap
44   plt.figure(figsize=(12, 10))
45   sns.heatmap(df.corr(), cmap='coolwarm', linewidths=0.5)
46   plt.title("Feature Correlation Heatmap")
47   plt.tight_layout()
48   plt.show()
49
```

## 3. Machine learning model



```python
64   # Train a Random Forest Classifier
65   model = RandomForestClassifier(n_estimators=100,
     random_state=42)
66   model.fit(X_train, y_train)
67
68   # Predict and evaluate
69   y_pred = model.predict(X_test)
70   print("\nClassification Report:\n",
     classification_report(y_test, y_pred))
71   print("\nConfusion Matrix:\n", confusion_matrix(y_test,
     y_pred))
72
73   # Additional Evaluation Metrics (adjusted for multiclass)
74   accuracy = accuracy_score(y_test, y_pred)
75   precision = precision_score(y_test, y_pred,
     average='weighted')
76   recall = recall_score(y_test, y_pred, average='weighted')
77   f1 = f1_score(y_test, y_pred, average='weighted')
78
79   print(f"\nAccuracy: {accuracy:.4f}")
80   print(f"Precision (weighted): {precision:.4f}")
81   print(f"Recall (weighted): {recall:.4f}")
82   print(f"F1 Score (weighted): {f1:.4f}")
83
```

The code trains a Random Forest Classifier on the darknet dataset to classify network traffic and evaluates its performance using a classification report, confusion matrix, and weighted metrics. The model achieved about **98.5% accuracy**, with high precision, recall, and F1-scores across all classes, showing it reliably distinguishes normal and potentially malicious traffic. The confusion matrix indicates very few misclassifications, confirming the model is highly effective for darknet traffic detection.