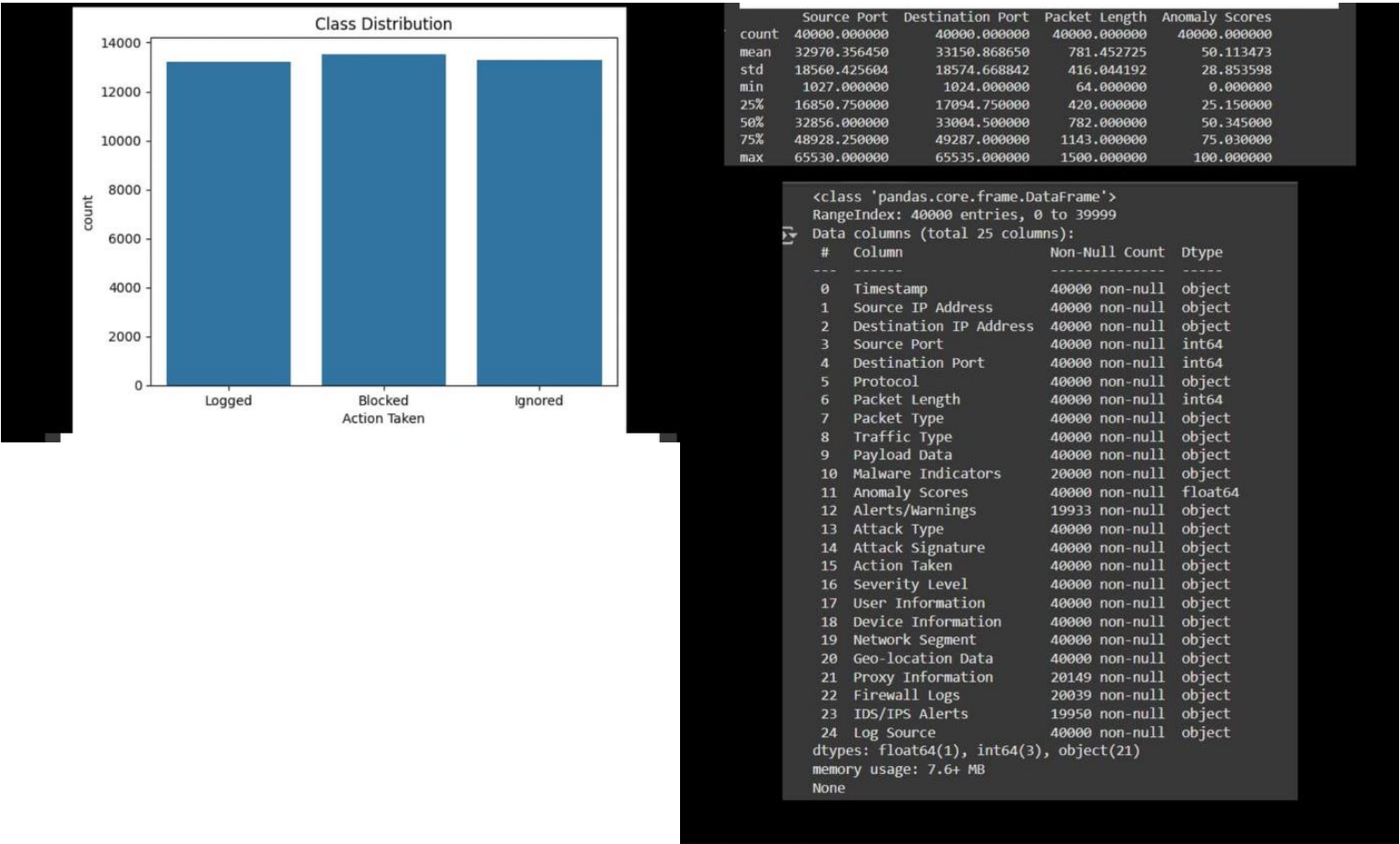


# Exploratory Data Analysis & Hyperparameter Tuning

## 1. Exploratory Data Analysis (EDA)

During the EDA phase, we performed the following steps and observations:

- **Class Distribution:** Plotted counts of each "Action Taken" category (Logged, Blocked, Ignored) to confirm balance and identify any skew.
- **Feature Overview:** Displayed dataset info to verify non-null counts and data types across all 25 original columns.
- **Statistical Summary:**
  - Source/Destination Port: range 1 024–65 535, mean  $\approx 33\,000$ , high variability.
  - Packet Length: mean  $\approx 780$  bytes,  $\sigma \approx 416$ , min 64, max 1 500.
  - Anomaly Scores: uniform distribution on  $[0,100]$ , mean  $\approx 50$ ,  $\sigma \approx 29$ .
- **Missing Values:** Identified five columns ( $\sim 50\%$  missing) and replaced nulls with explicit categories ("None", "No Data", "No Detection") to preserve all samples.
- **Categorical Distributions:** Visualized protocol usage (TCP/UDP/ICMP), traffic types (HTTP, DNS, FTP, etc.), and attack types (DDoS, Intrusion, Malware).



## 2. Hyperparameter Tuning (Random Forest)

To optimize the Random Forest classifier efficiently, we applied successive halving search:

- Subsampling: Extracted 10% of training data to accelerate tuning.
- Halving Strategy: Used HalvingRandomSearchCV with factor=2 and 2-fold CV, progressively narrowing from 529 to 2 candidates over 10 iterations.
- Parameter Space:  $n\_estimators \in [50, 200]$ ,  $max\_depth \in [5, 20]$ .
- Results: Best params =  $\{n\_estimators: 178, max\_depth: 12\}$ , CV accuracy = 62.91%.

This method reduced computational cost compared to a full search while effectively finding strong hyperparameters.

```
n_iterations: 10
n_required_iterations: 10
n_possible_iterations: 10
min_resources_: 8
max_resources_: 4235
aggressive_elimination: False
factor: 2
-----
iter: 0
n_candidates: 529
n_resources: 8
Fitting 2 folds for each of 529 candidates, totalling 1058 fits
-----
iter: 1
n_candidates: 265
n_resources: 16
Fitting 2 folds for each of 265 candidates, totalling 530 fits
-----
iter: 2
n_candidates: 133
n_resources: 32
Fitting 2 folds for each of 133 candidates, totalling 266 fits
-----
iter: 3
n_candidates: 67
n_resources: 64
Fitting 2 folds for each of 67 candidates, totalling 134 fits
```

```
-----
iter: 4
n_candidates: 34
n_resources: 128
Fitting 2 folds for each of 34 candidates, totalling 68 fits
-----
iter: 5
n_candidates: 17
n_resources: 256
Fitting 2 folds for each of 17 candidates, totalling 34 fits
-----
iter: 6
n_candidates: 9
n_resources: 512
Fitting 2 folds for each of 9 candidates, totalling 18 fits
-----
iter: 7
n_candidates: 5
n_resources: 1024
Fitting 2 folds for each of 5 candidates, totalling 10 fits
-----
iter: 8
n_candidates: 3
n_resources: 2048
Fitting 2 folds for each of 3 candidates, totalling 6 fits
-----
iter: 9
n_candidates: 2
n_resources: 4096
Fitting 2 folds for each of 2 candidates, totalling 4 fits
✓ Best params: {'max_depth': 12, 'n_estimators': 178}
✓ Best CV score: 62.91%
```