# Microsoft Malware Prediction
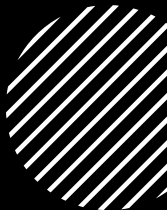
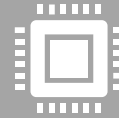Can you predict if a machine will soon be hit with malware?

Wenyu Pan

- Problem & Motivation

- Approach

- Implementation

- Evaluation & Results

- Future Directions

# Problem & Motivation

**Problem:**

Predict a Windows machine's probability of getting infected by various families of malware, based on different properties of that machine.

**Motivation:**

This predictive analysis is crucial as it enables users and system administrators to take proactive measures to enhance their security. A proactive stance in this domain can significantly reduce disruptions, financial losses, and maintain user trust, thereby contributing to a more secure cyber environment.

# Approach

- In my project, I utilized the Microsoft Malware Prediction dataset from Kaggle, applying a unified approach to data processing and modeling.

  DataSet size: DataFrame has 7,853,253 rows and 83 columns.


- I split the data into training and testing sets, then experimented with various models to optimally predict malware occurrences.

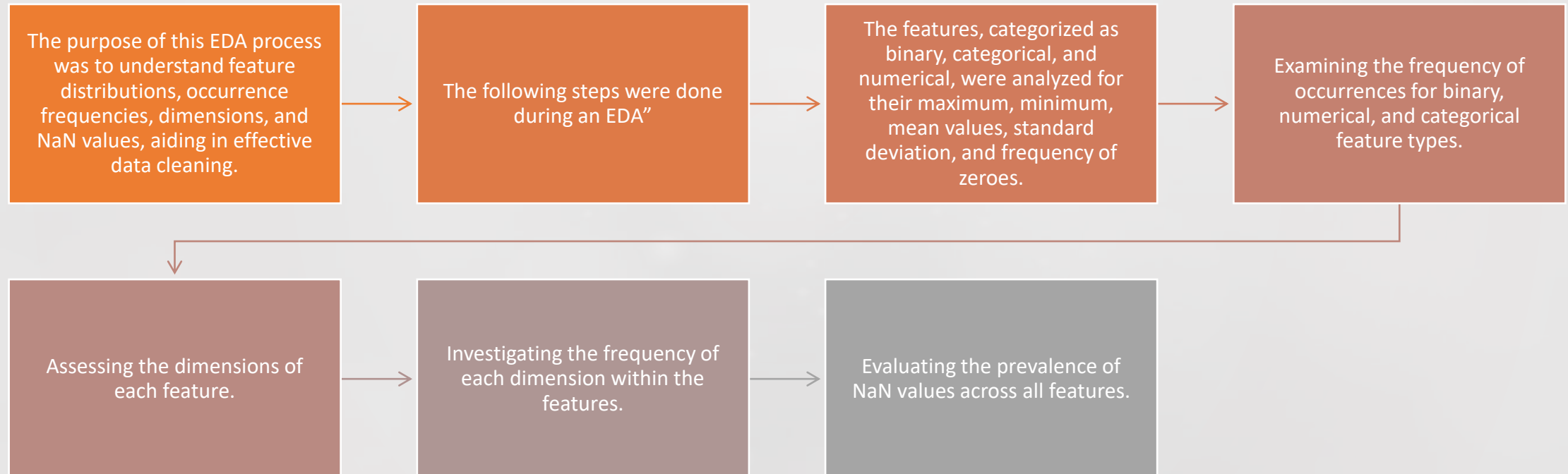# Introduction to Microsoft Malware Prediction dataset

**HasDetections** is the ground truth and indicates that Malware was detected on the machine.

Each row in the dataset corresponds to a machine, uniquely identified by a **MachineIdentifier**

# Before Data Cleaning...EDA

The purpose of this EDA process was to understand feature distributions, occurrence frequencies, dimensions, and NaN values, aiding in effective data cleaning.

The following steps were done during an EDA"

The features, categorized as binary, categorical, and numerical, were analyzed for their maximum, minimum, mean values, standard deviation, and frequency of zeroes.

Examining the frequency of occurrences for binary, numerical, and categorical feature types.

Assessing the dimensions of each feature.

Investigating the frequency of each dimension within the features.

Evaluating the prevalence of NaN values across all features.

# Data Cleaning

## 1. Deleting features with too much NaN-values

```python
nan_count = df.isnull().sum().to_frame('count')
nan_count['count'] = nan_count['count'].div(8921483).round(2)
irrelevant_feature = nan_count[nan_count['count'] > 0.5]
irrelevant_feature
```

```python
def assess_balance(df, column):
    value_counts = df[column].value_counts()
    max_count = value_counts.max()
    balance_ratio = max_count / len(df)
    return balance_ratio
```

```python
unbalanced_df = balance_ratios_df[balance_ratios_df['balance_ratio'] > 0.98]
```

```python
change_to_irrelevant(output_df, unbalanced_df)
```

## 2. Deleting features that are highly unbalanced

Define a function to calculate if the target feature is balanced. Here, we calculate a balance ratio between max count input and total input count. Ratio close to 1 indicates more imbalance.

# 3.

Features were then separate into three categories:

Numerical (replace NaN values with "-1")

Categorical (rename NaN-Values as '-1' in all features with tpye 'not category')

Binary(Reassign all NaN-Values to the most fequent feature)

This action ensures numerical features have no missing values, enhancing model accuracy and allowing similar preprocessing for test data.

# Now that we are done with data cleaning...

- We move on to data encoding:

1. **Label Encoding**: Transforms categorical values into a numerical range, where each unique label gets a unique integer.

2. **Frequency Encoding**: A variant of Label Encoding where values are encoded based on their frequency, assigning numbers based on the frequency of occurrence.

- This process ensures smooth processing and modeling with machine learning algorithms that typically require numerical input

# Implementation

To predict the data, I implemented three model:

- Logistic Regression Model

- Random Forest

- LightGBM/Keras

# Logistic Regression Model

Split the dataset into training and testing.

```python
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from imblearn.over_sampling import SMOTE

# Define your target and data ID columns
target_column = "HasDetections"
data_id = 'MachineIdentifier'


X = df.drop([target_column, data_id], axis=1)
y = df[target_column]


X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# SMOTE for class imbalance on training data
sm = SMOTE(random_state=42)
X_train_res, y_train_res = sm.fit_resample(X_train, y_train)
```

Delete the loaded whole dataset to save memory.

```python
del df
gc.collect()
```

6775

Train random forest model.

```python
# Initialize the Logistic Regression model
model = LogisticRegression()

# Train the model
model.fit(X_train, y_train)
```

Make prediction on testing dataset and find accuracy.

```python
# Make predictions
predictions = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, predictions)
print(f"Model Accuracy: {accuracy}")
print(classification_report(y_test, predictions))
```

Model Accuracy: 0.5083472510228101

```
              precision    recall  f1-score   support

           0       0.60      0.05      0.09    222980
           1       0.50      0.97      0.66    223095

    accuracy                           0.51    446075
   macro avg       0.55      0.51      0.38    446075
weighted avg       0.55      0.51      0.38    446075
```

# Random Forest

Load encoded dataset

```
[ ] df = pd.read_csv('./drive/MyDrive/train_encoded.csv')
    # df = pd.read_csv('./data/train_encoded.csv')
```

Split the dataset into training and testing.

```
from sklearn.model_selection import train_test_split
from imblearn.over_sampling import SMOTE

# Your original columns and data splitting
target_column = "HasDetections"
data_id = 'MachineIdentifier'
X = df.drop([target_column, data_id], axis=1)
y = df[target_column]

# Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Apply SMOTE only on training data to handle class imbalance
sm = SMOTE(random_state=42)
X_train_res, y_train_res = sm.fit_resample(X_train, y_train)
```

Delete the loaded whole dataset to save memory.

```
[ ] del df
    gc.collect()

    0
```

Train random forest model.

```
[ ] # Initialize the Random Forest model
    model = RandomForestClassifier(random_state=42)  # Use RandomForestRegressor for regression

    # Train the model
    model.fit(X_train, y_train)
```

Make prediction on testing dataset and find accuracy.

```
[ ] # Make predictions
    predictions = model.predict(X_test)


    # Evaluate the model
    accuracy = accuracy_score(y_test, predictions)
    print(f"Model Accuracy: {accuracy}")


    Model Accuracy: 0.6485097797455585
```

# LightGBM

```python
from sklearn.model_selection import train_test_split
from imblearn.over_sampling import SMOTE

# Your original columns and data splitting
target_column = "HasDetections"
data_id = 'MachineIdentifier'
X = df.drop([target_column, data_id], axis=1)
y = df[target_column]

# Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Apply SMOTE only on training data to handle class imbalance
sm = SMOTE(random_state=42)
X_train_res, y_train_res = sm.fit_resample(X_train, y_train)
```

Delete the loaded whole dataset to save memory.

```python
del df
gc.collect()

30
```

Create the LightGBM data containers

```python
train_data = lgb.Dataset(X_train, label=y_train)
test_data = lgb.Dataset(X_test, label=y_test, reference=train_data)
```

Define the parameters

```python
params = {
    'objective': 'binary',
    'metric': 'binary_logloss',
    'num_leaves': 31,
    'learning_rate': 0.05,
    'num_iterations': 100
}
```

Train the model

```python
gbm = lgb.train(params, train_data, valid_sets=[test_data])
```

```
/usr/local/lib/python3.10/dist-packages/lightgbm/engine.py:172: UserWarning: Found `num_iterations` in params. Will use it instead of argument
  _log_warning(f"Found `{alias}` in params. Will use it instead of argument")
[LightGBM] [Info] Number of positive: 892084, number of negative: 892212
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 1.016717 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 5313
[LightGBM] [Info] Number of data points in the train set: 1784296, number of used features: 66
[LightGBM] [Info] [binary:BoostFromScore]: pavg=0.499964 -> initscore=-0.000143
[LightGBM] [Info] Start training from score -0.000143
```

Predict on the test set

```python
y_pred = gbm.predict(X_test, num_iteration=gbm.best_iteration)
```

Convert probabilities to binary output

```python
y_pred_binary = [1 if x > 0.5 else 0 for x in y_pred]
```

Evaluate the model

```python
accuracy = accuracy_score(y_test, y_pred_binary)
print(f"Accuracy: {accuracy}")

Accuracy: 0.6446673765622373
```

# Evaluation & Results

## Model Evaluation Metric:

- Accuracy was used as the primary metric for classification effectiveness.

## Logistic Regression Model:

- Achieved 50.83% accuracy.
- Struggled with class imbalances and complex features.

## Random Forest Model:

- Reached 64.85% accuracy.
- Benefited from regularization and ensemble methods.

## LightGBM and Keras Models:

- LightGBM achieved 64.4% accuracy.
- Keras achieved 63.61% accuracy.
- Comparable performance to Random Forest.
- Advanced techniques showed potential but didn't significantly outperform traditional methods.

# Future Directions

**Model Evaluation Considerations:**

Accuracy alone may not fully capture model performance.

Imbalanced datasets necessitate a cautious interpretation.

**Recommended Future Metrics:**

Incorporate F1 score, precision, recall, and AUC-ROC.

These metrics provide a more nuanced performance assessment.

**Benefits of a Comprehensive Approach:**

Facilitates deeper insights into each model's strengths and limitations.

# Future Directions for Malware Detection

- **Exploring Additional Data Sources:** Integrating more diverse datasets to cover a wider range of malware signatures and attack vectors.

- **Deep Learning Techniques:** Experimenting with deep neural networks, such as convolutional neural networks (CNNs) or recurrent neural networks (RNNs), to capture complex patterns in data.

- **Ensemble Methods:** Combining predictions from multiple models to improve accuracy and robustness against diverse malware threats.