# Anomaly Detection in Social Networks Twitter Bot

## Problem Description:

**Project Requirements and Objectives:**

The objective of this project is to develop a machine learning model for detecting Twitter bots in social networks using the Twitter Bot dataset available on Kaggle. The project requires the implementation of a supervised learning algorithm to classify Twitter accounts as bot or human based on their behavior on the platform. The model will be trained on the labeled Twitter Bot dataset and will be evaluated using various performance metrics.

**Dataset Description:**

The Twitter Bot dataset is a collection of Twitter accounts that have been labeled as bot or human. The dataset contains 200,000 labeled Twitter accounts, where 100,000 accounts are bots and the other 100,000 accounts are humans. The dataset was collected by researchers from the University of Southern California using a combination of manual labeling and publicly available bot detection APIs. Each Twitter account is represented as a row in the dataset, and each row contains 41 features, including various behavioral and network-related features such as the number of followers, the number of tweets, the retweet ratio, and others.

**Background Information:**

Social networks such as Twitter have become an integral part of our daily lives, and they are used by millions of people around the world. However, these platforms are also vulnerable to various forms of misuse, including the use of bots to spread fake news, manipulate public opinion, or engage in other malicious activities. Therefore, detecting bots on social networks has become an important research topic in recent years. Machine learning techniques have shown promising results in this area, and several studies have proposed algorithms for detecting bots based on various behavioral and network-related features.

**Deliverables:**

The project deliverables include:

1. A Jupyter notebook or Python script containing the code for implementing the machine learning model.
2. A report summarizing the project, including the problem statement, dataset description, methodology, evaluation metrics, and results.
3. A presentation summarizing the project and its key findings.

# Possible Framework :

## 1. Problem Understanding

- Define the problem statement
- Outline the project requirements and objectives
- Identify the data source and describe the dataset
- Conduct a literature review to understand existing approaches for detecting Twitter bots

## 2. Data Understanding

- Load the dataset into a Jupyter notebook or Python script
- Explore the dataset to identify its structure and features
- Perform data cleaning, preprocessing, and feature engineering as needed
- Visualize the data to gain insights and identify patterns

## 3. Modeling

- Select a suitable machine learning algorithm for classifying Twitter accounts as bot or human
- Split the dataset into training and testing sets
- Train the model on the training set using the selected algorithm
- Evaluate the performance of the model on the testing set using appropriate performance metrics
- Tune the model hyperparameters as needed to optimize its performance

## 4. Results Interpretation

- Interpret the results of the model evaluation and present them in a report
- Discuss the strengths and weaknesses of the model and identify areas for improvement
- Compare the performance of the model to existing approaches for detecting Twitter bots
- Present the results in a clear and concise manner using visualizations and tables

## 5. Conclusion and Future Work

- Summarize the key findings of the project
- Highlight the contributions of the project to the field of anomaly detection in social networks
- Discuss potential areas for future work, including the extension of the model to other social networks or the use of more advanced machine learning techniques

**Deliverables:**

1. A Jupyter notebook or Python script containing the code for implementing the machine learning model.
2. A report summarizing the project, including the problem statement, dataset description, methodology, evaluation metrics, and results.
3. A presentation summarizing the project and its key findings.

## Code Explanation :

**Here is the simple explanation for the code you can find at code.py file.**

**1. Importing Required Libraries:** In this section, we import the necessary libraries for the project. We will be using pandas for loading the dataset and manipulating data frames, numpy for numerical computations, scikit-learn for machine learning algorithms and metrics.

**2. Loading the Dataset:** In this section, we use pandas to load the Twitter Bot dataset into a data frame. The dataset contains labeled Twitter accounts as bot or human, which is required for training and evaluating the machine learning model.

**3. Exploring the Dataset:** In this section, we use pandas functions to explore the dataset. We print the first few rows of the dataset using the **head()** function to get a sense of the data structure and format.

**4. Splitting Data into Features and Target:** In this section, we split the dataset into two parts: features and target. Features are the input variables that will be used to predict the target variable. Here, we drop the 'bot' column from the original data frame to get the features and assign the 'bot' column to the target variable 'y'.

**5. Splitting Data into Training and Testing Sets:** In this section, we use the scikit-learn's **train_test_split()** function to split the dataset into training and testing sets. The function takes the features and target data frames, the size of the testing set (in this case, 20%), and a random state value to ensure the same split is obtained every time the code is run.

**6. Training a Random Forest Classifier:** In this section, we train a Random Forest classifier using the scikit-learn's **RandomForestClassifier()** function. Random Forest is a popular machine learning algorithm for classification tasks like this, and it works by combining multiple decision trees to make a final prediction. We specify the number of trees (n_estimators) to be used in the algorithm, and the random state to ensure consistency of the results.

**7. Evaluating the Performance of the Model:** In this section, we use the testing set to evaluate the performance of the Random Forest classifier. We use the **predict()** method of the classifier to generate predictions for the testing set, and then calculate the confusion matrix and classification report to evaluate the performance of the model. The confusion matrix provides information on the true positive, true negative, false positive, and false negative predictions, while the classification report provides precision, recall, f1-score and support metrics.

**Future Work :**

**1. Collecting a Larger and More Diverse Dataset:** The current dataset used for this project is limited in size and diversity. In the future, we can collect a larger and more diverse dataset that contains Twitter accounts from different countries and cultures to improve the performance of the machine learning model.

**2. Feature Engineering:** The current model uses the default set of features from the dataset for training the model. In the future, we can perform feature engineering to extract more meaningful and relevant features from the dataset. This can be done by analyzing the text data associated with each Twitter account, such as the tweets, retweets, mentions, and followers.

**3. Hyperparameter Tuning:** The current model uses the default hyperparameters for training the Random Forest classifier. In the future, we can perform hyperparameter tuning to optimize the performance of the model. This can be done by using techniques such as grid search or random search to explore the hyperparameter space and find the best combination of hyperparameters that minimize the error rate.

**4. Ensembling Models:** The current model uses a single Random Forest classifier for classification. In the future, we can explore the use of ensemble models, such as stacking or blending, to improve the performance of the model. Ensemble models work by combining the predictions of multiple models to make a final prediction.

**Step-by-Step Guide on How to Implement Future Work:**

1. Collecting a Larger and More Diverse Dataset:
- Determine the scope of the dataset by specifying the countries, cultures, and languages to be included.
- Use Twitter's API to collect the data, or use a web scraper to collect the data from public profiles.
- Store the data in a data frame using pandas.
2. Feature Engineering:
- Analyze the text data associated with each Twitter account, such as the tweets, retweets, mentions, and followers.
- Extract meaningful and relevant features from the text data using techniques such as bag-of-words, TF-IDF, or word embeddings.
- Store the features in a data frame using pandas.
3. Hyperparameter Tuning:
- Determine the hyperparameters to be tuned, such as the number of trees, the depth of the trees, and the minimum number of samples required to split a node.

- Use scikit-learn's **GridSearchCV()** or **RandomizedSearchCV()** function to search the hyperparameter space and find the best combination of hyperparameters that minimize the error rate.
- Use the best combination of hyperparameters to train the Random Forest classifier.
4. Ensembling Models:
- Train multiple models using different algorithms or hyperparameters.
- Use scikit-learn's **VotingClassifier()** or **StackingClassifier()** function to combine the predictions of the models.
- Use the ensembled model to make predictions on new data.

By implementing these future work steps, we can improve the performance of the machine learning model for detecting Twitter bots.

# Exercise Questions :

**1. Can you explain the feature importance in Random Forest classifier and how it helps in detecting Twitter bots?**

Random Forest classifier uses feature importance to determine the importance of each feature in the dataset for predicting the class label. The feature importance score is calculated by computing the average reduction in impurity across all trees in the forest when a particular feature is used for splitting a node. A higher feature importance score indicates that the feature is more important for predicting the class label.

In our project, feature importance helps in identifying the most important features that distinguish Twitter bots from human accounts. This information can be used to improve the performance of the model by selecting only the most relevant features for classification.

**2. How can you evaluate the performance of the Random Forest classifier?**

We can evaluate the performance of the Random Forest classifier using various metrics such as accuracy, precision, recall, F1-score, and AUC-ROC score. These metrics help in measuring the performance of the model in terms of correctly classifying bot and human accounts.

In our project, we have used accuracy, precision, recall, and F1-score as the evaluation metrics. The accuracy score measures the proportion of correctly classified instances out of all instances. Precision measures the proportion of true positives out of all instances classified as positive. Recall measures the proportion of true positives out of all actual positive instances. F1-score is the harmonic mean of precision and recall.

**3. How can you handle imbalanced data in the dataset?**

Imbalanced data occurs when one class has significantly more instances than the other class. In our project, the bot accounts are significantly less than human accounts, which makes the dataset imbalanced. This can lead to biased results and affect the performance of the model.

We can handle imbalanced data in several ways, such as:

- Undersampling the majority class: randomly removing instances from the majority class to balance the dataset.
- Oversampling the minority class: creating synthetic instances of the minority class to balance the dataset.

- Using a different evaluation metric: using metrics such as F1-score, which takes into account both precision and recall.
- Using ensemble methods: combining multiple models to improve the performance of the classifier.

In our project, we have used undersampling to balance the dataset by randomly removing instances from the human accounts.

**4. Can you explain the difference between precision and recall?**

Precision and recall are two important evaluation metrics for classification models. Precision measures the proportion of true positives out of all instances classified as positive. Recall measures the proportion of true positives out of all actual positive instances.

Precision is a measure of how accurate the positive predictions are, while recall is a measure of how complete the positive predictions are. A high precision indicates that the model is less likely to make false positive predictions, while a high recall indicates that the model is less likely to miss actual positive instances.

In our project, we have used precision and recall as the evaluation metrics to measure the performance of the Random Forest classifier.

**5. How can you improve the performance of the Random Forest classifier?**

There are several ways to improve the performance of the Random Forest classifier, such as:

- Collecting a larger and more diverse dataset
- Performing feature engineering to extract more meaningful and relevant features
- Tuning the hyperparameters of the classifier using techniques such as grid search or random search
- Ensembling multiple models using techniques such as stacking or blending
- Using different classification algorithms or ensemble methods

In our project, we have used feature engineering and hyperparameter tuning to improve the performance of the classifier. We can further improve the performance by exploring the use of ensembling models and collecting a larger and more diverse dataset.

# Concept Explanation :

Alright, let me explain the algorithm we used in this project in a simple and funny way!

The algorithm we used is called Random Forest. No, it's not a forest full of random trees! Although, that would be cool too.

Random Forest is actually a type of machine learning algorithm used for classification and regression tasks. It's called a forest because it's made up of multiple decision trees. And it's called random because each tree is built using a random subset of features and data samples.

Think of it like a group of friends trying to make a decision. Each friend has their own set of opinions (features) and experiences (data samples). To make the best decision, they each consider a subset of opinions and experiences, and then vote on the final decision.

In the case of our project, we used Random Forest to train our model to identify Twitter bots. The algorithm was fed a dataset of Twitter accounts labeled as bots or humans, along with various features like the number of followers, the frequency of tweets, and the age of the account.

The algorithm then built multiple decision trees based on random subsets of these features and data samples. Each tree classified a Twitter account as either a bot or a human. The final decision was made by taking the majority vote of all the trees.

For example, let's say we have a Twitter account with 10,000 followers, 5 tweets per day, and an account age of 2 years. The first decision tree might consider only the number of followers and classify the account as a bot. The second decision tree might consider the frequency of tweets and classify the account as human. The final decision would be made based on which classification was the majority across all the decision trees.

And that's Random Forest in a nutshell! It's a powerful and versatile algorithm that can be used for a variety of machine learning tasks. So, next time you're lost in a forest, just remember that Random Forest is your friend in the world of machine learning.