# Assignment 5

# Assignment on Azure Cloud Platform

PartB

1.

We aim to develop a robust machine learning model to identify and classify spam comments on YouTube—a growing concern for content creators and viewers alike. Utilizing the YouTube Spam Collection dataset, which consists of 1,956 comments from five popular videos, our goal is to accurately distinguish between class 1 (irrelevant or inappropriate messages) and class 0 (relevant, appropriate comments).

This classification is crucial for maintaining a healthy digital environment on YouTube, as spam comments can be disruptive, misleading, or even harmful. By addressing this issue, we not only seek to enhance the user experience on YouTube but also to gain valuable insights into the nature of spam comments on social media platforms. Ultimately, our model will contribute to creating a safer and more enjoyable space for users to engage with content online.

Solution:
We propose a two-fold approach utilizing Logistic Regression and Random Forest classifiers:

**Logistic Regression:**

Why: Logistic Regression is a simple yet effective linear model for binary classification problems. It provides a clear interpretation of feature importance and performs well on linearly separable data.
How: We will use Logistic Regression to create a baseline model, focusing on its ability to identify relationships between features and the binary outcome (spam vs. non-spam).

**Random Forest:**

Why: Random Forest is an ensemble learning method that builds multiple decision trees and merges them to achieve higher accuracy and robustness. It excels at handling complex, non-linear relationships in data and is less prone to overfitting compared to individual decision trees.
How: We will train a Random Forest model to capture the intricate patterns in the data, leveraging its ability to manage feature interactions and provide feature importance insights.

I have written the answers for Q2,Q3,Q4 in the Markdown block of the notebook and exported the codes along with the explanations. See A5_Notebook.pdf for details.

Q5 is answered on the following pages

## Q5:

I will use Automated ML for my data set and explain the best model results.

Configure Automated ML job seetings



Choose my compute instance

Start Running:



It will run mutilple models to find the best one

Choose one we can see the training details of the model:
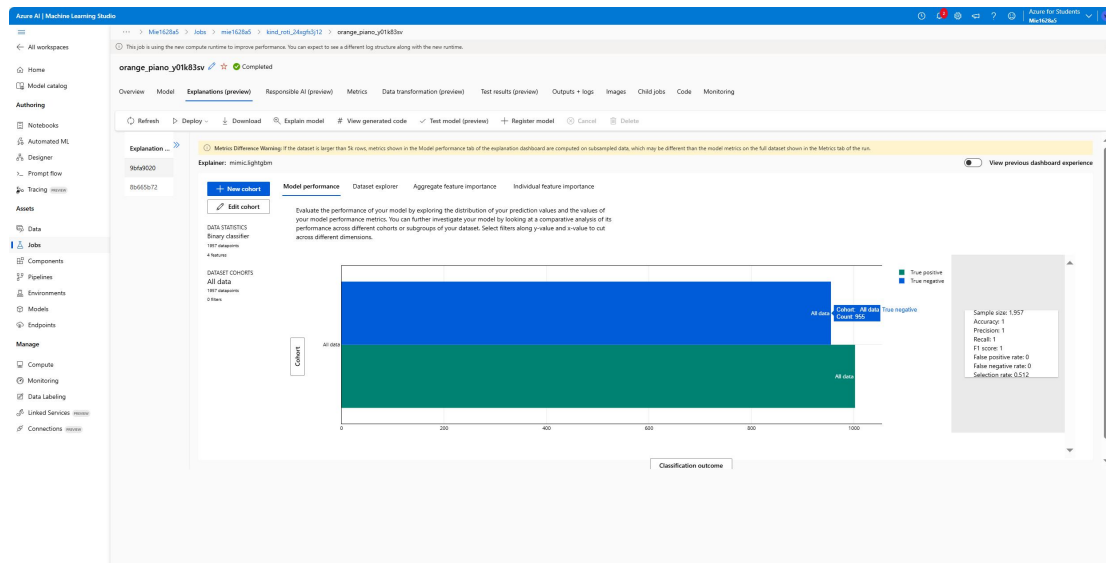


Job completed

We can expolre the best model



Choose the best model "VotingEnsemble" which achieved the highest accuracy

From the **Explanations** dashboard, we can get insight into this trained model:

Model performance:



Explore the top-k important features that impact the overall model predictions: