# Introduction

Spam detection is an essential problem in natural language processing (NLP) and machine learning.

## Why is it a big data problem?

**Large Volume**: The volume of emails is huge

**High Velocity**: millions of emails every second.

**Variety of Data**: Spam can take many forms, including email, social media posts, and comments.

**Complexity of Algorithms**: Modern spam detection algorithms are complex and require a lot of computational power to run. This can be a big data problem if the algorithms need to be run on huge data sets.

**False positives**: False positives are a problem in spam detection, as they can result in legitimate messages being classified as spam. To reduce the number of false positives, spam detection algorithms must be very accurate, which requires a lot of data and computational power for training and optimization.

## Spam Detection Applications:

**Email spam detection**: Email service providers use spam detection techniques to filter out unwanted emails from a user's inbox. These techniques may include rule-based filters or machine-learning algorithms.

**SMS spam detection**: SMS service providers also use spam detection techniques to filter out unwanted text messages from a user's inbox. These techniques may include analyzing message content and sender information.

**Comment spam detection**: Websites that allow users to leave comments often use spam detection techniques to filter out unwanted comments that may contain spam or malicious links.

**Social media spam detection**: Social media platforms use detection techniques to identify and remove spam accounts and content.

**Online advertising spam detection**: Advertisers use spam detection techniques to identify and filter out fraudulent ad clicks or impressions.

In this project, we will use PySpark to build a spam detection model that can classify emails as either spam or ham (non-spam).

# Data

The data used in this project is a publicly available dataset of SMS messages, which can be downloaded from here. The dataset contains 5,572 SMS messages, with a class label indicating whether each statement is spam or ham.

# Methodology

The overall approach used in this project is to preprocess the data using PySpark's built-in NLP functions and then build a Naive Bayes classifier to predict whether a given message is spam or ham.

## Preprocessing

The first step in preprocessing the data is to read it into a PySpark DataFrame. This is done using the read.csv function, with the inferSchema parameter set to True to automatically infer the schema from the data. We then drop unnecessary columns and rename the remaining columns for clarity.

We then use PySpark's built-in NLP functions to tokenize the text, remove stop words, and compute the term frequency-inverse document frequency (TF-IDF) features for each message. We also convert the class labels from strings to numeric values.

Finally, we use PySpark's VectorAssembler function to combine the TF-IDF features and message lengths into a single feature vector, and then train a Naive Bayes classifier on the preprocessed data.

# Feature Engineering

After cleaning the text data, we will use PySpark's machine learning libraries to perform feature engineering. In this step, we will transform the text data into numerical features that can be used for modeling. We will use the following feature engineering techniques:

**Tokenization**: Break up the text into individual words, or tokens.
**Stop Word Removal**: Remove common words like "the", "and", and "or" that don't contribute to the meaning of the text.
**Count Vectorization**: Count the number of times each token appears in each text document.
**Term Frequency-Inverse Document Frequency (TF-IDF) Vectorization**: Weight the count vectors by the importance of each token in the corpus.
We will implement these techniques using PySpark's MLlib library.

# Modeling

In this step, we will train a machine learning model to predict whether a message is spam or not. We will use the Naive Bayes algorithm, which is commonly used for text classification tasks. We will train the model using the features we generated in the previous step.

# Evaluation

Finally, we will evaluate the performance of the model using the test data set. We will calculate the accuracy of the model, which is the percentage of test data points that the model correctly classified as spam or not.