

*What are your  
expectations for the  
course?*



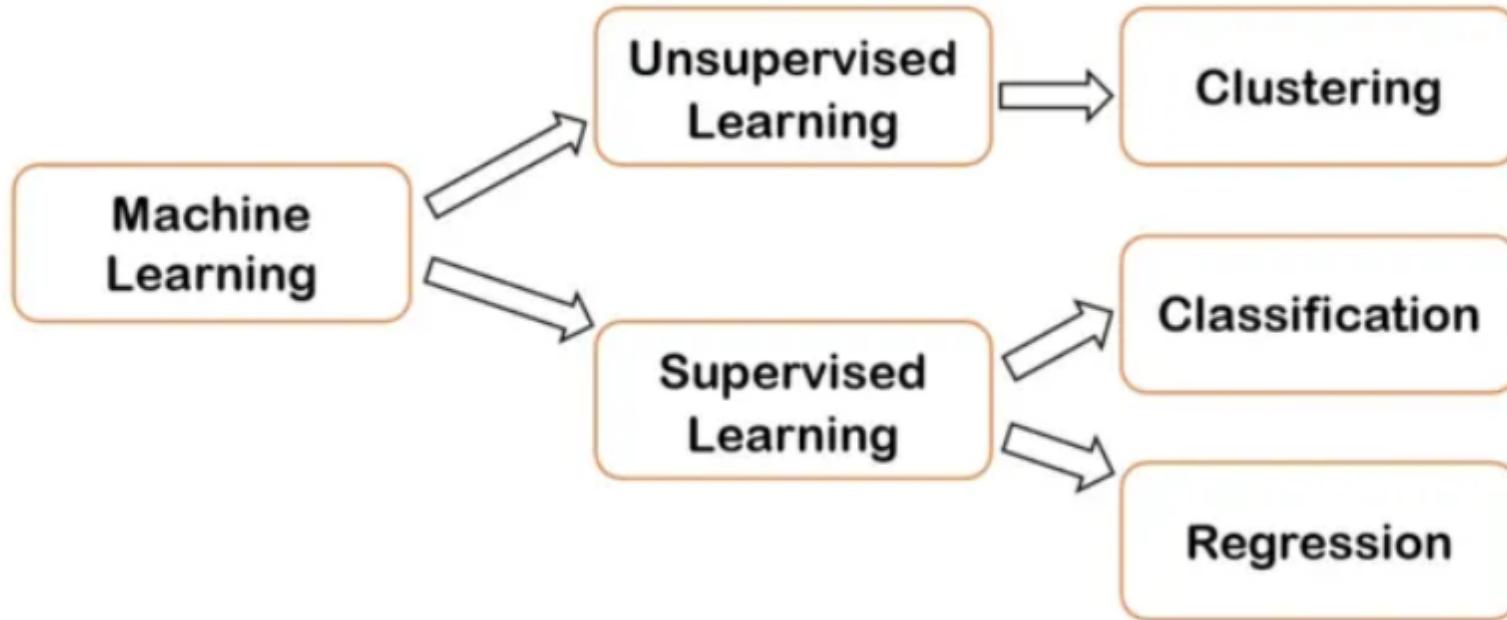
**NLP**

Natural language  
processing



# Agenda

1. Revision on ML
  
1. What is Text Classification?
  - o Binary Classification
  - o Multi class Classification



# Text Classification

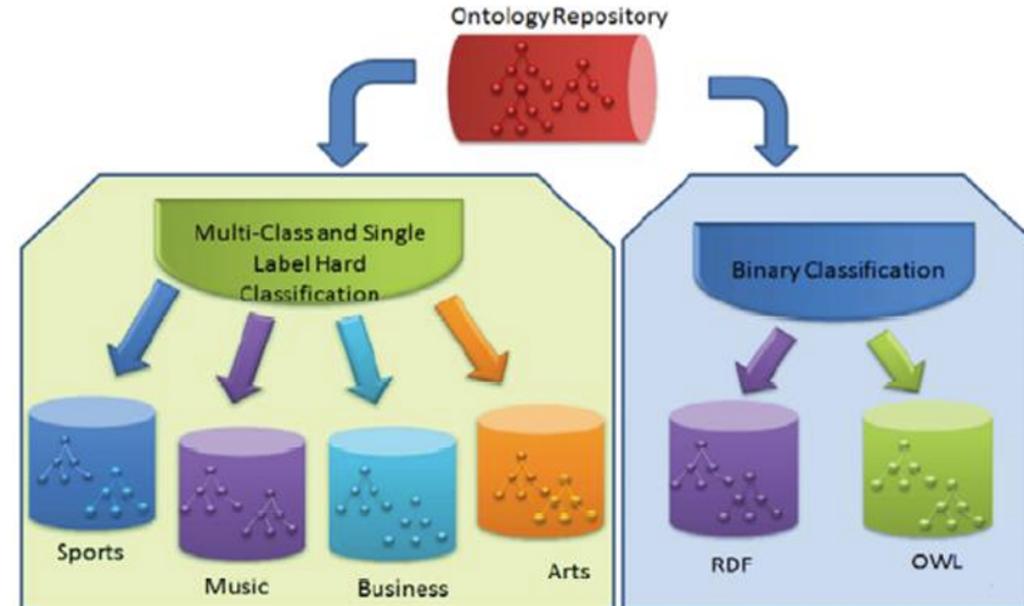
Text classification or Text Categorization is the activity of labeling natural language texts with relevant categories from a predefined set.

- ❑ Assigning subject categories, topics, or genres.
- ❑ Spam detection.
- ❑ Authorship identification.
- ❑ Age/gender identification.
- ❑ Language Identification.
- ❑ Sentiment analysis.

# Text Classification

There are two kinds of classifiers:

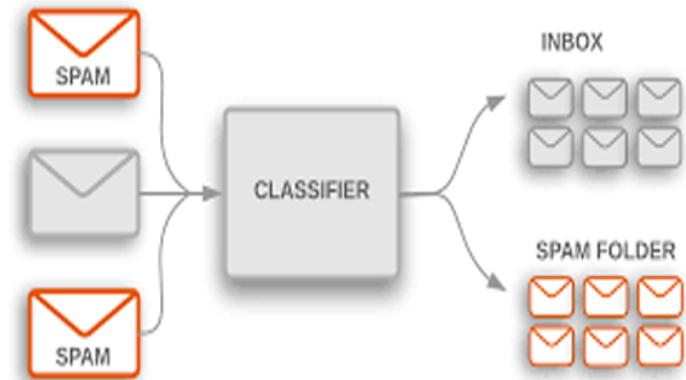
1. Binary Classifiers.
2. Multi-Classifiers



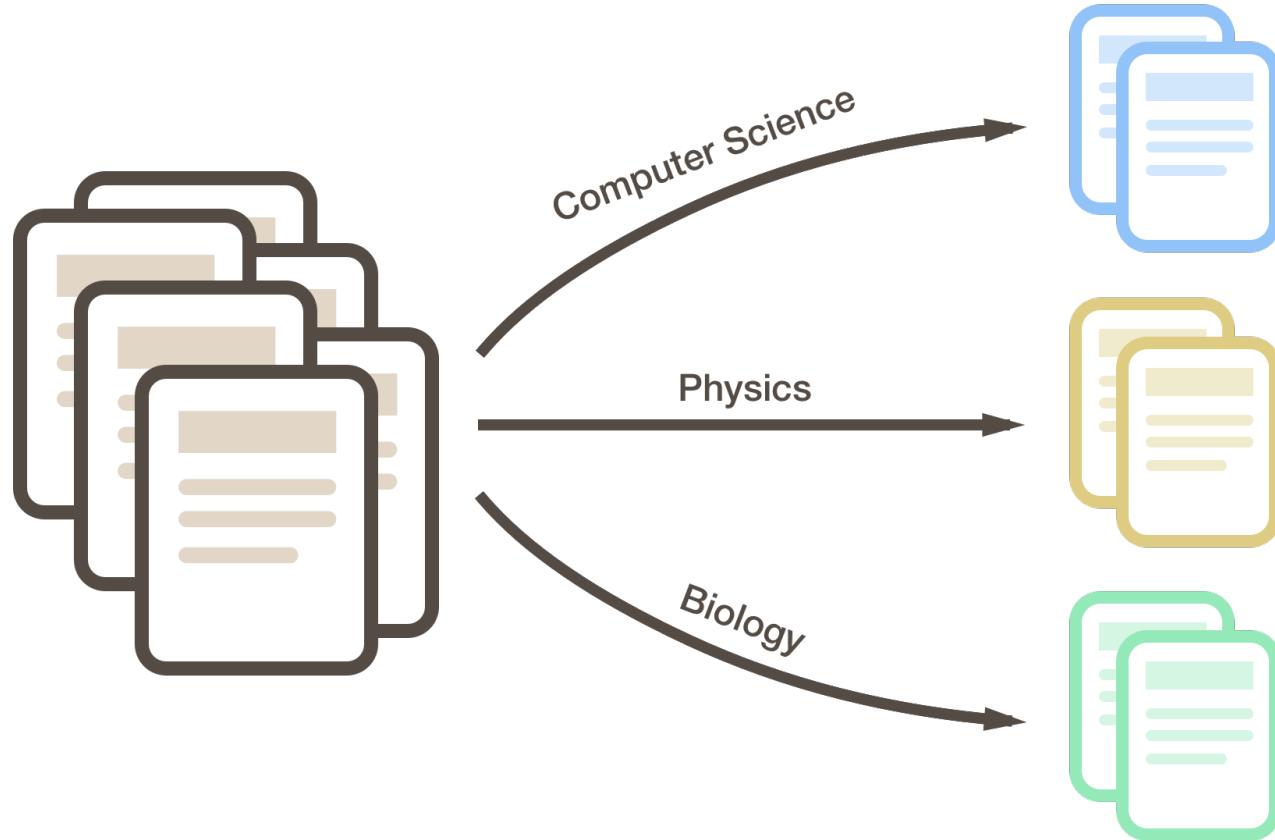
# Binary Classifiers

Suppose we have a dataset of SMS messages, each labeled as 'spam' or 'ham'. Our goal is to train a model that can accurately classify new messages.

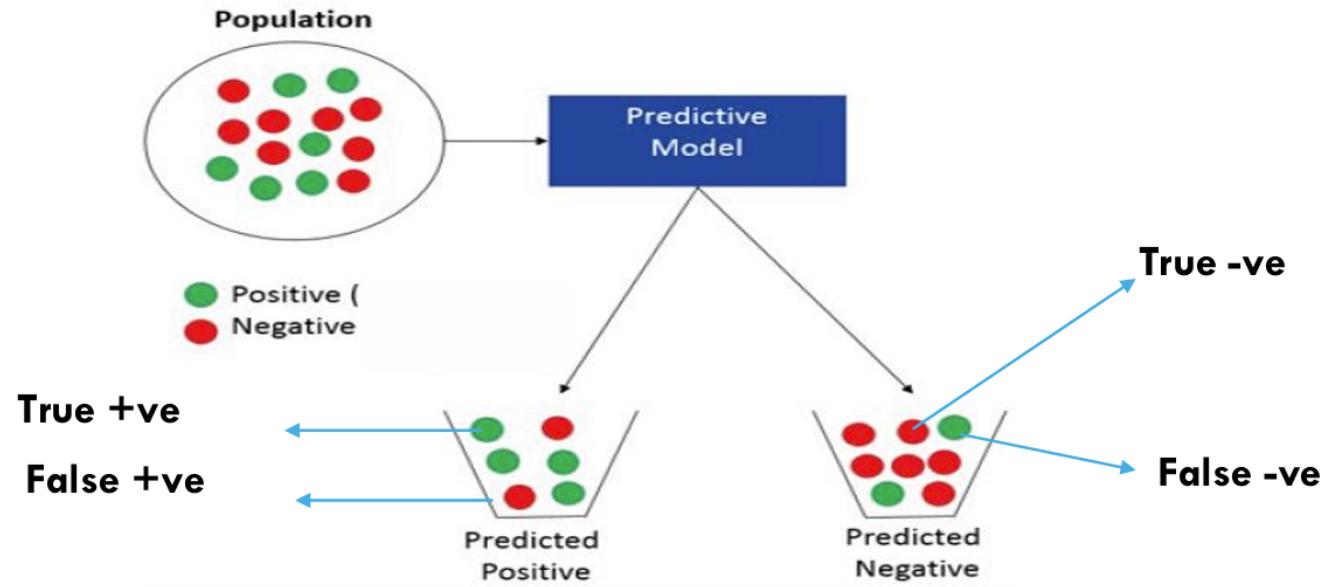
First, we preprocess the messages — removing numbers, special characters, and converting everything to lowercase. Then, we extract features using techniques like Bag-of-Words or TF-IDF ."



# Multi class Classification



# Text Classification Model



# Text Classification Model

## Key Components of a Confusion Matrix

- **True Positive (TP)**: The model correctly predicts a positive outcome.
- **True Negative (TN)**: The model correctly predicts a negative outcome.
- **False Positive (FP)**: The model incorrectly predicts a positive outcome (Type I error).
- **False Negative (FN)**: The model incorrectly predicts a negative outcome (Type II error)

## Confusion Matrix

	Actually Positive (1)	Actually Negative (0)
Predicted Positive (1)	True Positives (TPs)	False Positives (FPs)
Predicted Negative (0)	False Negatives (FNs)	True Negatives (TNs)

# Text Classification Model

		Actual	
		Spam (+ve)	Not Spam (-ve)
Predictions	Spam (+ve)	TP	FP
	Not Spam (-ve)	FN	TN

Number of **Positive (P)** predictions that are correct or **True (T)**

Number of **Positive (P)** predictions that are wrong or **False (F)**

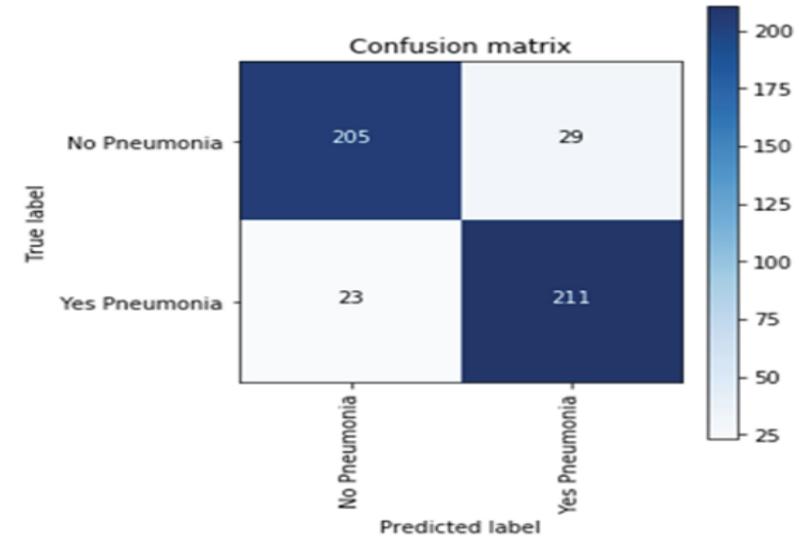
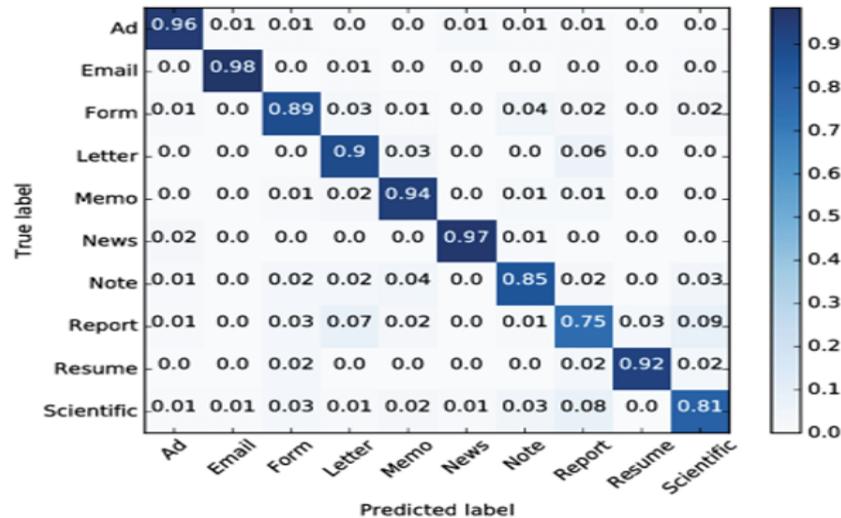
Number of **Negative (N)** predictions that are wrong or **False (F)**

Number of **Negative (N)** predictions that are correct or **True (T)**

T or F : if prediction = actual or not

P or N : based on the prediction only

# Text Classification Model



Whenever the diagonal of the propagation matrix has higher values= better accuracy

# TEXT CLASSIFICATION Steps:

- ① 1. Data Collection: (Labelled Data)
- ② 2. Text Pre-processing (Stop Words removal, normalization,...etc)
- ③ 3. Feature Extraction : BOW, TF-IDF...etc(specify which features are important for your classification, for example chemical labels refer to scientific articles....etc)
- ④ 4. ML Algorithm Selection (Logistic Regression, Naïve Bayes, decision tree...etc)
- ⑤ 5. Classification Evaluation (Precision, recall, F1).

# Challenges in Text Classification

**?**Imbalanced Data : When one class (e.g., ham) dominates the other (spam), the model might become biased.

**?**Ambiguity in Language : Words can have multiple meanings depending on context.

**?**Preprocessing Decisions : Choosing what to clean and what to keep is critical.

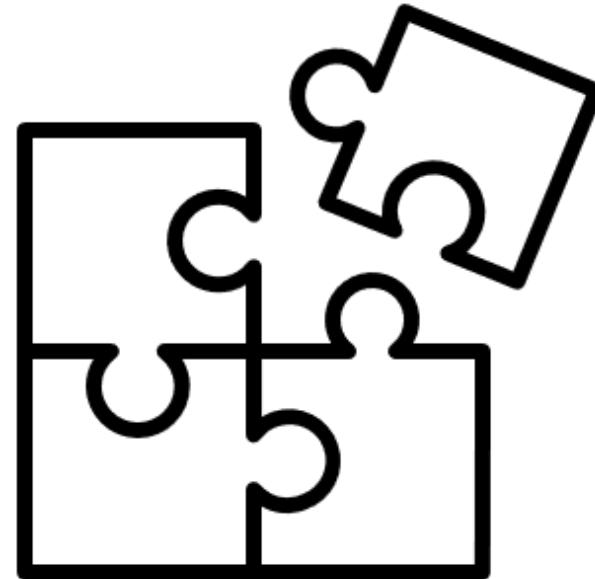
# Hands-On



## Other Useful Datasets for Text Classification:

Use Case	Dataset	Link
Sentiment Analysis	IMDB Movie Reviews	<a href="#">Kaggle IMDB Dataset</a>
Topic Classification	BBC News Dataset	<a href="#">Kaggle BBC News Dataset</a>
General Text Classification	AG News Dataset	<a href="#">FastText Dataset Page</a>

# Quiz





Thanks