

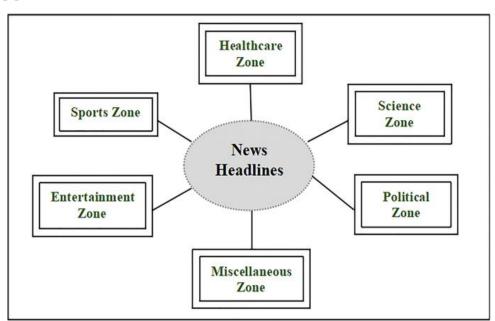
Presented By: Basmah AlQahtani Mariam Alrashdi

### We will discuss ...

- 1. Introduction
- 2. Project Aim
- 3. Description
- 4. Models Used
- 5. Dataset Used
- 6. Workflow Explanation
- 7. Results & Discussion
- 8. Conclusion



In news, categorization is a multi-label text classification issue. The purpose is to allocate a news story to one or more categories.





# Project Aim ...

In this project, we aim to train a model that can properly categorize previously unknown news articles into some classes such as business, science and technology, entertainment, and health.



# **Description** ...

The purpose of dividing news into multiple categories is to make the news reader's experience easier.

We construct a model that automatically classifies these news headlines based on their category.

# **Description** ...



Preprocessing the dataset



Machine Learning



## Model used ...

In this project, two models were used in a dataset for news classification.

Naïve Bayes	Logistic Regression
Collection of classification algorithms based on Bayes' Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle.	A Classification algorithm is used to forecast a result using a set of independent variables. To describe data and explain the relationship between one dependent binary variable and one or more nominal, ordinal.

# **Data Used**

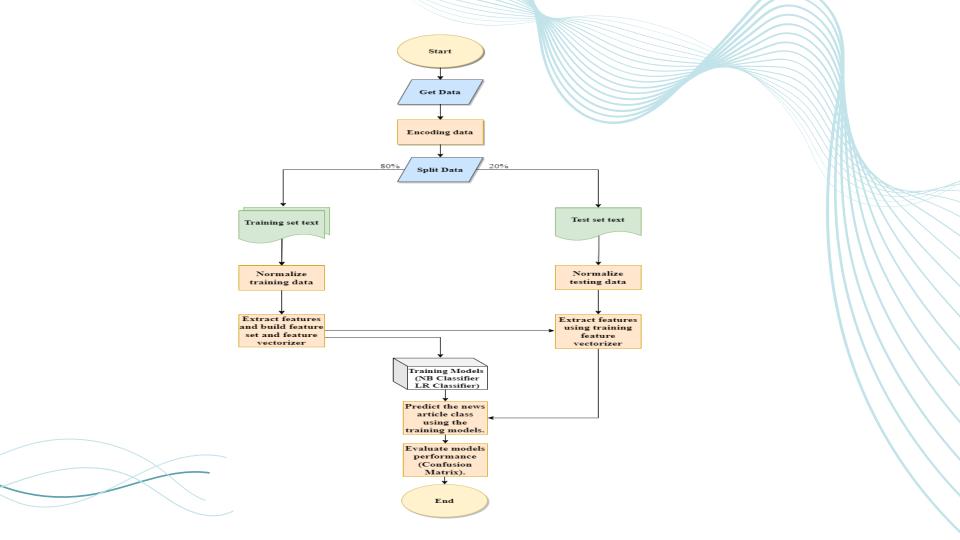
# Data used...

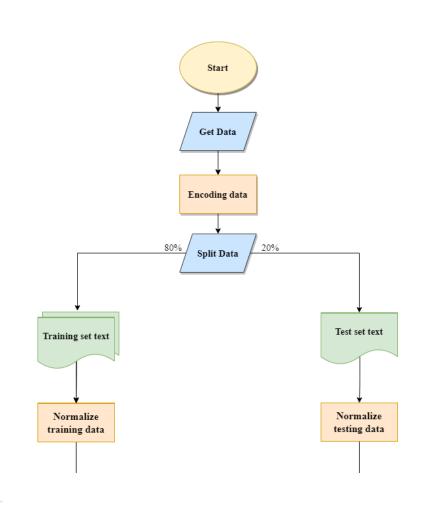
we used a database from Kaggle. This dataset contains headlines, URLs, and categories for 422,937 news stories collected by a web aggregator between March 10th, 2014, and August 10th, 2014.

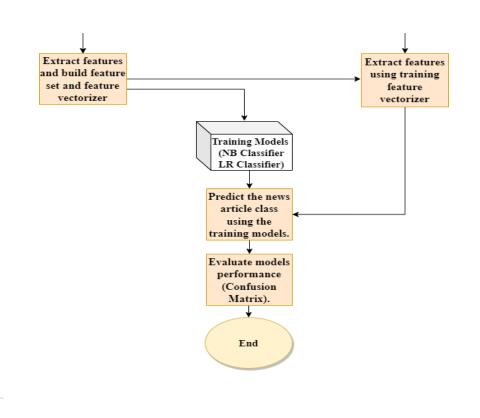
### The columns included in this dataset are:

- **ID:** the numeric ID of the article
- **TITLE:** the headline of the article
- **URL:** the URL of the article
- **PUBLISHER:** the publisher of the article
- **CATEGORY:** the category of the news item; one of:
  - b: business
  - t: science and technology
  - e: entertainment
  - m: health
- **STORY:** alphanumeric ID of the news story that the article discusses
- **HOSTNAME:** hostname where the article was posted
- TIMESTAMP: approximate timestamp of the article's publication, given in Unix time (seconds since midnight on Jan 1, 1970)









# Results & Discussions

### 1- Naïve Bayes:

Naïve Bayes model performance on the training dataset:

accuracy 0.92	2870226522852	203		
	precision	recall	f1-score	support
b	0.90	0.92	0.91	92774
е	0.95	0.97	0.96	121975
m	0.97	0.86	0.91	36511
t	0.91	0.91	0.91	86675
accuracy			0.93	337935
macro avg	0.93	0.91	0.92	337935
weighted avg	0.93	0.93	0.93	337935

### 1- Naïve Bayes:

Naïve Bayes model performance on the test dataset:

Naive Bayes Model Performance Analysis						
accuracy	accuracy 0.9208844278206524					
		precision	recall	f1-score	support	
	ь	0.89	0.91	0.90	23193	
	e	0.94	0.97	0.96	30494	
	m	0.97	0.84	0.90	9128	
	t	0.90	0.90	0.90	21669	
accur	асу			0.92	84484	
macro	avg	0.93	0.90	0.91	84484	
weighted	avg	0.92	0.92	0.92	84484	

### **2- Logistic Regression:**

Logistic Regression model performance on the training dataset:

accuracy 0.95	342595469542	96		
	precision	recall	f1-score	support
b	0.93	0.94	0.94	92774
е	0.97	0.98	0.98	121975
m	0.97	0.93	0.95	36511
t	0.94	0.94	0.94	86675
accuracy			0.95	337935
macro avg	0.95	0.95	0.95	337935
weighted avg	0.95	0.95	0.95	337935

### **2- Logistic Regression:**

Logistic Regression model performance on test dataset:

Logistic Regression Model Performance Analysis					
accuracy	0.94	0793523033947	73		
		precision	recall	f1-score	support
	ь	0.92	0.92	0.92	23193
	е	0.96	0.98	0.97	30494
	m	0.96	0.91	0.93	9128
	t	0.93	0.92	0.92	21669
accui	racy			0.94	84484
macro	avg	0.94	0.93	0.94	84484
weighted	avg	0.94	0.94	0.94	84484





Basmah AlQahtani 443800986@kku.edu.sa Mariam AlRashdi 443800993@kku.edu.sa 10 May 2022