

[illegible]

# Muhajir Akbar Hasibuan

## Data Scientist



✉ [muhajirakbarhsb@gmail.com](mailto:muhajirakbarhsb@gmail.com)

☎ +62 85297023234

📍 Jakarta, Indonesia

🌐 [linkedin.com/in/muhajirakbarhsb](https://www.linkedin.com/in/muhajirakbarhsb)

🐙 [github.com/muhajirakbarhsb](https://github.com/muhajirakbarhsb)

📖 [medium.com/@muhajirakbarhsb](https://medium.com/@muhajirakbarhsb)

## Experience

### Telkom Indonesia

Sept 2020 - present

- Develop internal and external use cases
- Provide data understanding in making a model
- Provide Preparation and data engineering according to the use case implemented
- Provide data validation so that the analysis results are as expected
- Building modeling for the development or improvement of internal and external programs
- Provide descriptive and diagnostic insight into data processing
- Recommend and define new growth hacking strategy for digital marketing team

### Project

- Pioneered a Robust Big Data Solution for MyIndihome, Revolutionizing Customer Experiences
- Orchestrated a High-Performing ML Team, Elevating myIndihome TV's Personalized Content Impact by 25%
- Envisioned and Executed a Cutting-Edge Big Data Solution for Elevated Customer Engagement on myIndihomeTV
- Fueled Business Insights via Dynamic Data Profiling, Performance Dashboards, and Insights for Langit Musik, RBT, and Upoint
- Engineered Innovative Big Data Solutions that Propelled Growth for PadiUkm (E-Commerce)
- Powered Success for GameQoo through Strategic Big Data Solutions
- Architected and Established MLOps Framework, Elevating Telkom Indonesia's Digital Business Products
- Crafted Visionary Video Analytics Solutions for Telkom Indonesia's Revolutionary Digital Business IoT Product

MEMBER OF DATA SCIENTIST TASK FORCE | NOVEMBER 2021 - PRESENT

MEMBER OF AI TASK FORCE | NOVEMBER 2022- PRESENT

A Pool of data scientists and AI Engineer Expert in Telkom Indonesia. It was established to leverage data-driven culture for decision-making within the organization. (Applied Research, Standardization, Consultation)

### Achievement

BEST Talent of the Year at Telkom Indonesia, Digital Business and Technology Division - Digital Technology and Platform -2022

### Codex by Telkom Indonesia

Apr 2020 - Sept 2020

- Building Data Pipeline for Langit Musik Recommender System
- Business Analytic for Langit Musik
- Build Recommender System for Langit Musik

### Universitas Syiah Kuala

2015-2020

- Bachelor of Science, Statistics

## Skills

### Hard Skills

- Data Analytics
- Statistics
- Machine Learning
- Deep Learning
- MLOps
- Business Intelligence
- Data Engineering
- Cloud Computing
- Time Series & Demand Forecasting
- Natural Language Processing
- Fraud & Anomaly Detection
- Computer Vision

### Tools Skills

- Python
- Pyspark
- Sql
- Apache Airflow
- GCP
- Docker
- MLflow
- Prometheus
- Evidently
- Grafana
- Redash, Metabase, Superset, Looker Studio
- Pytorch

<p><b>Class Meeting 1:</b> Introduction to NLP (1 session)</p> <p>Introduction to NLP</p> <p>Basics of text data and its characteristics</p> <p>The importance of NLP in today's world</p> <p>Overview of the course structure and objectives</p> <p><b>Class Meeting 2:</b> Text Preprocessing (1 session)</p> <p>Understanding the text preprocessing pipeline</p> <p>Tokenization, stemming, and lemmatization</p> <p>Stop words removal</p> <p>Hands-on exercises with Python for text preprocessing</p> <p><b>Class Meeting 3:</b> Text Representation (1 session)</p> <p>Introduction to text representation techniques</p> <p>Bag of Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF)</p> <p>n-gram</p> <p>Word embeddings (Word2Vec, GloVe)**just an introduction</p> <p>Practical exercises on text representation</p> <p><b>Class Meeting 4:</b> Introduction to ML, DL Supervised and Unsupervised (1 session)</p> <p>Introduction to deep learning (ML)</p> <p>Introduction to deep learning (DL)</p> <p>DL vs. traditional machine learning (ML)</p> <p>Machine Learning and their applications</p> <p>Neural networks and their applications</p> <p>Basic ML DL concepts and terminology</p> <p><b>Class Meeting 5:</b> Sentiment Analysis (1 session)</p> <p>What is sentiment analysis?</p> <p>Data collection and labeling for sentiment analysis</p> <p>Building a sentiment analysis model</p> <p>Practical sentiment analysis examples</p> <p><b>Class Meeting 6:</b> Text Classification (1 session)</p> <p>Introduction to text classification</p> <p>Binary and multi-class classification</p> <p>Building a text classification model</p> <p>Real-world text classification examples</p>	<p><b>Class Meeting 6:</b> Introduction to Topic Modeling(1 session)</p> <p>Introduction to topic Modeling</p> <p>Topic Modeling Implementation</p> <p>LDA(Latent dirichlet allocation)</p> <p>Hands-on exercises with Topic Modeling</p> <p><b>mid term</b></p> <p><b>Class Meeting 8: Sesi Khusus deep Learning</b></p> <p><b>Class Meeting 9:</b> Introduction to Word Embeddings (1 session)</p> <p>Fundamentals of word embeddings</p> <p>Word2Vec, GloVe, and FastText</p> <p>Word embedding applications</p> <p>Hands-on exercises with word embedding</p> <p><b>Class Meeting 10:</b> Attention Mechanisms (1 session)</p> <p>Introduction to attention mechanisms</p> <p>Self-attention and multi-head attention</p> <p>Transformers architecture</p> <p>Practical examples of attention in NLP</p> <p><b>Class Meeting 11:</b> Transformer Models (1 session)</p> <p>In-depth study of the Transformer model</p> <p>Pre-trained transformer models (BERT, GPT-2)</p> <p>Fine-tuning transformers for NLP tasks</p> <p>Transformer-based applications</p> <p><b>Class Meeting 12:</b> Advanced NLP Topics (1 session)</p> <p>Advanced NLP topics such as BERT, XLNet, and RoBERTa</p> <p>Transfer learning in NLP</p> <p>Ethics in NLP</p> <p>Recent developments and trends in NLP</p> <p><b>Class Meeting 13:</b> Advanced NLP Topics and Deployment Process in Industry (1 session)</p> <p>Introduction to LLM</p> <p>Introduction how industry utilize NLP to generate revenue</p> <p>Introduction MLOps for NLP(bonuses from practitioners)</p>	<p><b>Class Meeting 14-15:</b> Project Work (2 sessions)</p> <p>Dedicated sessions for students to work on NLP projects with guidance and assistance.</p> <p><b>Class Meeting 16:</b> Project Presentations and Conclusion (1 session)</p> <p>Students present their NLP projects</p> <p>Recap of key takeaways from the course</p> <p>Discuss further resources for NLP enthusiasts</p> <p>Course conclusion and feedback</p>
--	---	--

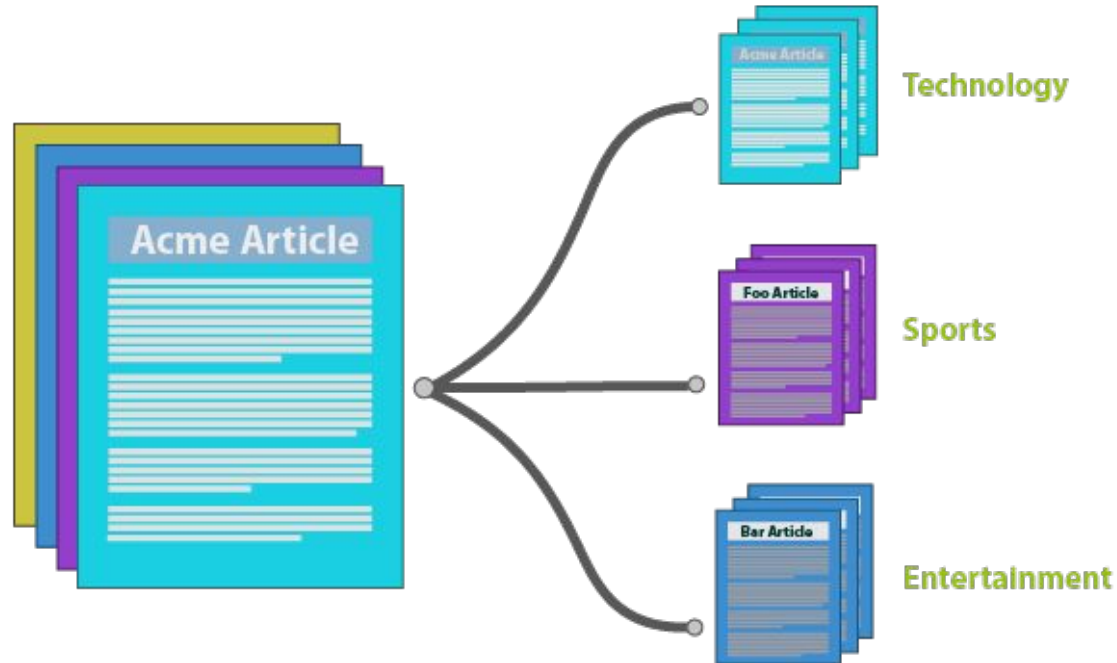
## **Agenda**

- **Introduction to text classification?**
- **Data collection and labeling for Text Classification**
- **Binary and multi-class classification**
- **Building a text classification model**

# **Introduction to Text Classification**

# Text Classification

Text classification is a natural language processing task that involves **categorizing text documents into predefined categories or labels**, such as spam detection or sentiment analysis.



# Machine Learning-Based Text Classification

With advance growth of machine learning, nowadays it is just easier to create model using machine learning and feed data to the model and wait until the model is complete. With the machine learning model, it's much easier and **faster to classify category from input text**. One important step to use machine learning is **feature extraction**. We transform text to **numeral representation in form of vector**, one way of doing it is using **bags of word**, or basically, we **count every words in a text**, or using **tfidf (term frequency inverse document frequency)** etc.



# **Data Collection and Labeling**



# Machine Learning-Based Text Classification

```
[4] df['sentiment_first_scheme'] = df['score'].apply(lambda x: 'negative' if x in [1, 2] else 'positive')
```

```
# Sentiment analysis based on the second scheme
```

```
df['sentiment_second_scheme'] = df['score'].apply(lambda x: 'negative' if x in [1, 2] else ('neutral' if x == 3 else 'positive'))
```

```
[7] df.sentiment_second_scheme.value_counts()
```

```
negative    1717
```

```
positive    1243
```

```
neutral      228
```

```
Name: sentiment_second_scheme, dtype: int64
```

```
df.sentiment_first_scheme.value_counts()
```

```
negative    1717
```

```
positive    1471
```

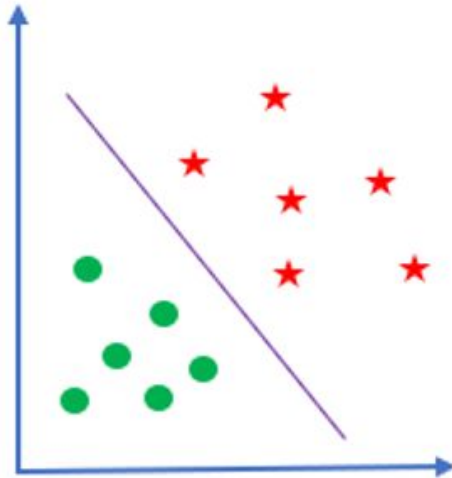
```
Name: sentiment_first_scheme, dtype: int64
```

# **Binary and multi-class classification**

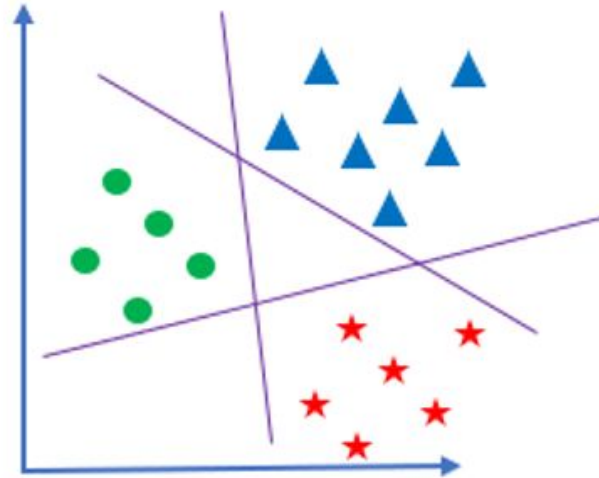
# Binary vs Multiclass

Binary classification involves categorizing data into **two distinct classes**, while multi-class classification extends this task to categorize data into **three or more classes**, making it suitable for a broader range of real-world applications.

Binary classification



Multi-class classification



# Binary vs Multiclass

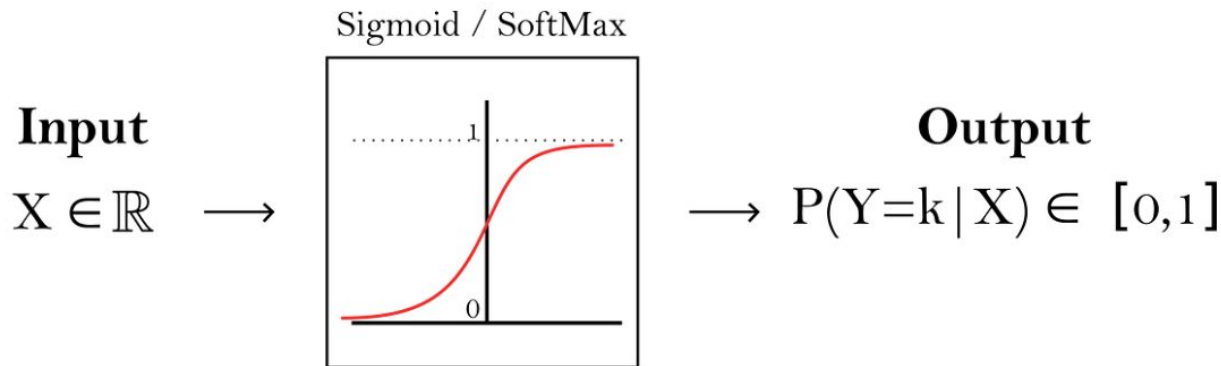
## Binary Classification:

- Binary classification focuses on distinguishing between two exclusive classes or categories, typically referred to as the **positive class and the negative class**.
- Techniques such as logistic regression, support vector machines, and decision trees can be directly applied to binary classification.
- Evaluation metrics often include **accuracy, precision, recall, F1-score, ROC-AUC, and others tailored to two-class scenarios**.

## Multi-class Classification:

- Multi-class classification involves categorizing data into three or more classes or categories, each mutually exclusive.
- Common techniques for multi-class classification include **one-vs-all (OvA or OvR), softmax regression (logistic regression extension)**, and neural network architectures designed for multiple output classes.
- Evaluation metrics differ and can include **macro/micro-average F1-score, cross-entropy loss, confusion matrices, and accuracy calculated across all classes**.

# Binary vs Multiclass



**Figure 1.** Illustration of Sigmoid and SoftMax function. The output is read as “The probability of Y being the class k given the input X”. Image by author

But if both functions map the same transformation (i.e. do the same thing), what is the difference between them?

Sigmoid is used for **binary classification** methods where we only have 2 classes, while SoftMax applies to **multiclass problems**. In fact, the SoftMax function is an extension of the Sigmoid function.

# Binary vs Multiclass

## Sigmoid

2 classes

$$\text{out} = P(Y=\text{class1}|X)$$

## SoftMax

$k > 2$  classes

$$\text{out} = \begin{bmatrix} P(Y=\text{class1}|X) \\ P(Y=\text{class2}|X) \\ P(Y=\text{class3}|X) \\ \vdots \\ P(Y=\text{class}k|X) \end{bmatrix}$$

**Figure 2.** Illustration of the outputs of each function. An important property is the addition of all the probabilities for each Sigmoid class and SoftMax should be equal to 1. In the case of Sigmoid we obtain  $P(Y=\text{class2}|X) = 1 - P(Y=\text{class1}|X)$ . Image by author

# Modeling Process

# Modeling Process

In a machine learning pipeline, what should be the **first step**: splitting the data into training and testing sets or preprocessing steps like imputation, SMOTE (Synthetic Minority Over-sampling Technique), or scaling data?



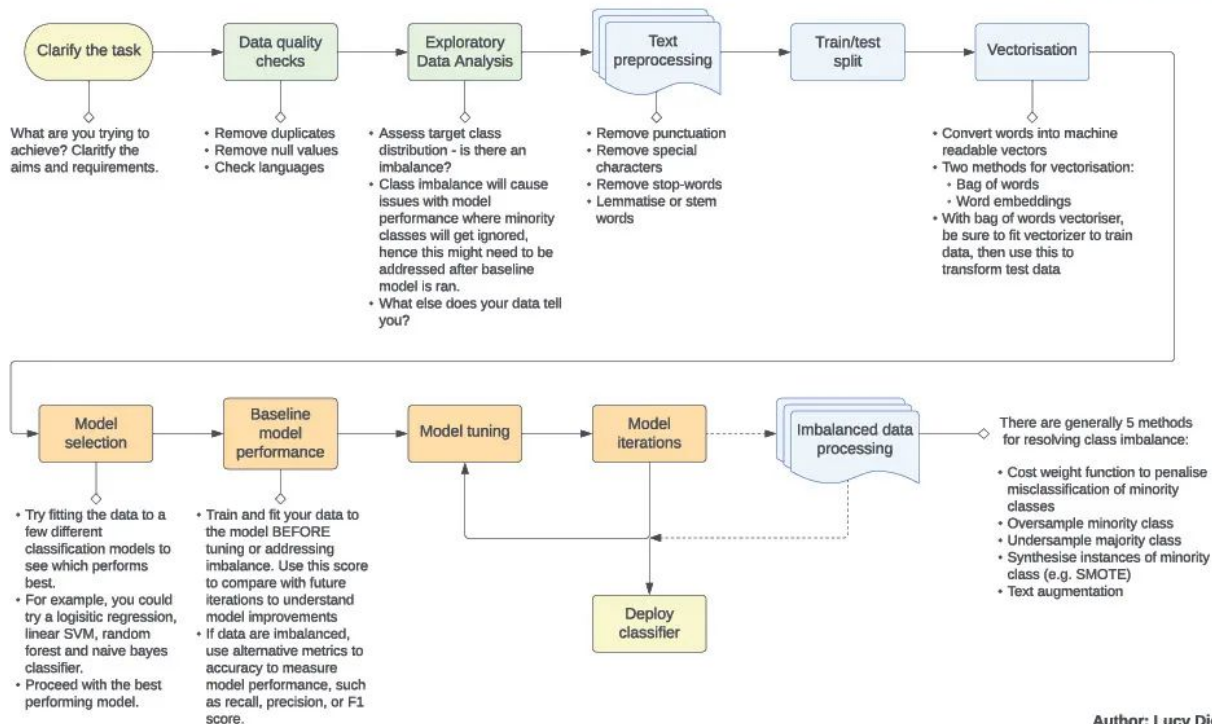
# Modeling Process

## Text Classifier Flow Diagram

{ Supervised machine  
learning approach }

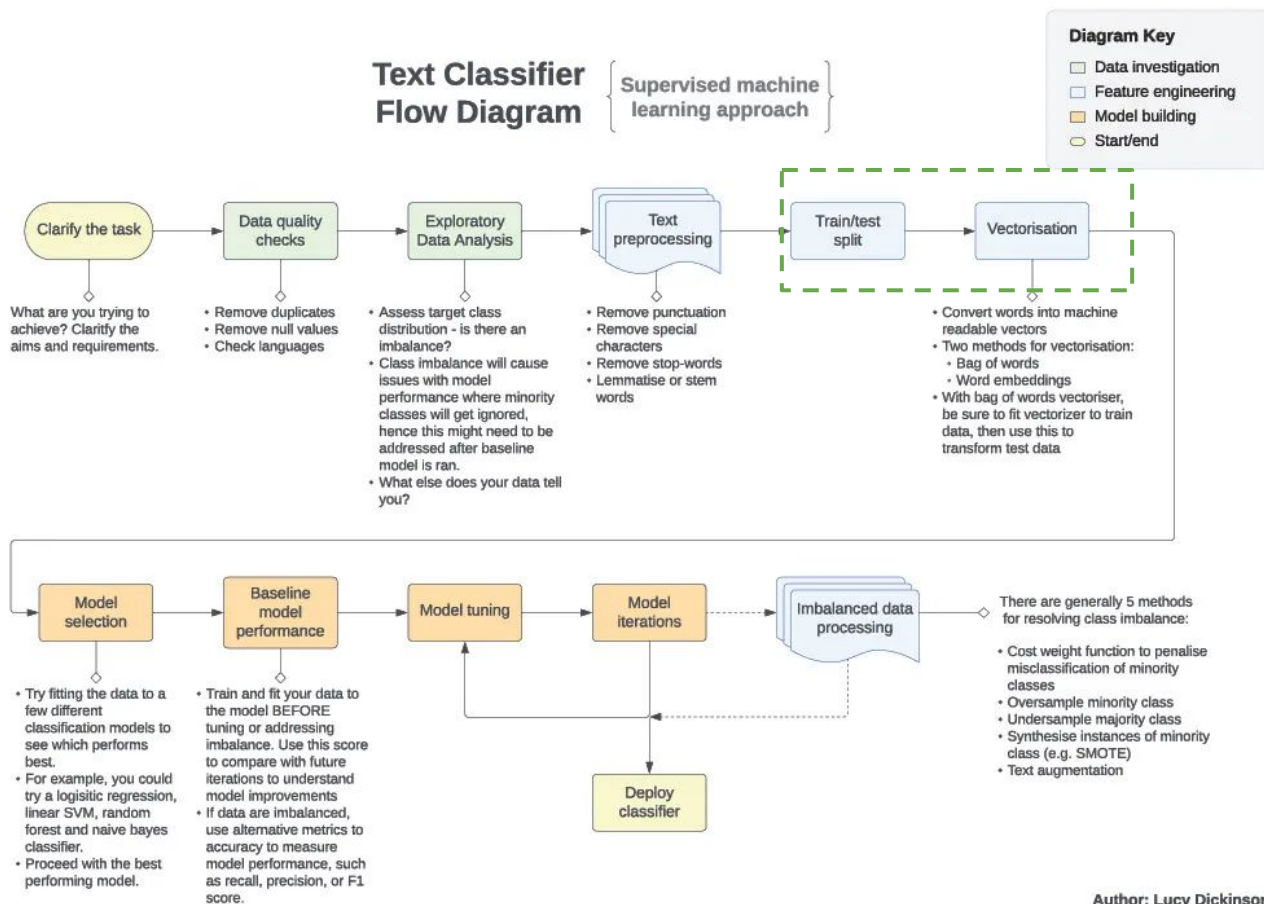
### Diagram Key

- Data investigation
- Feature engineering
- Model building
- Start/end



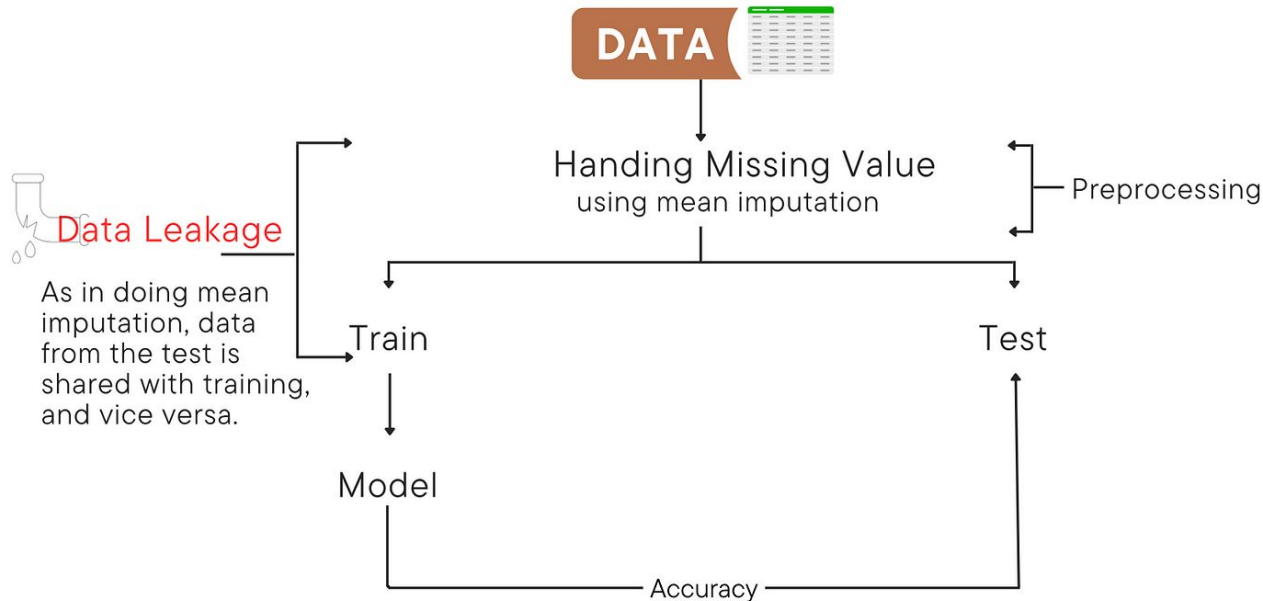
Author: Lucy Dickinson

# Modeling Process

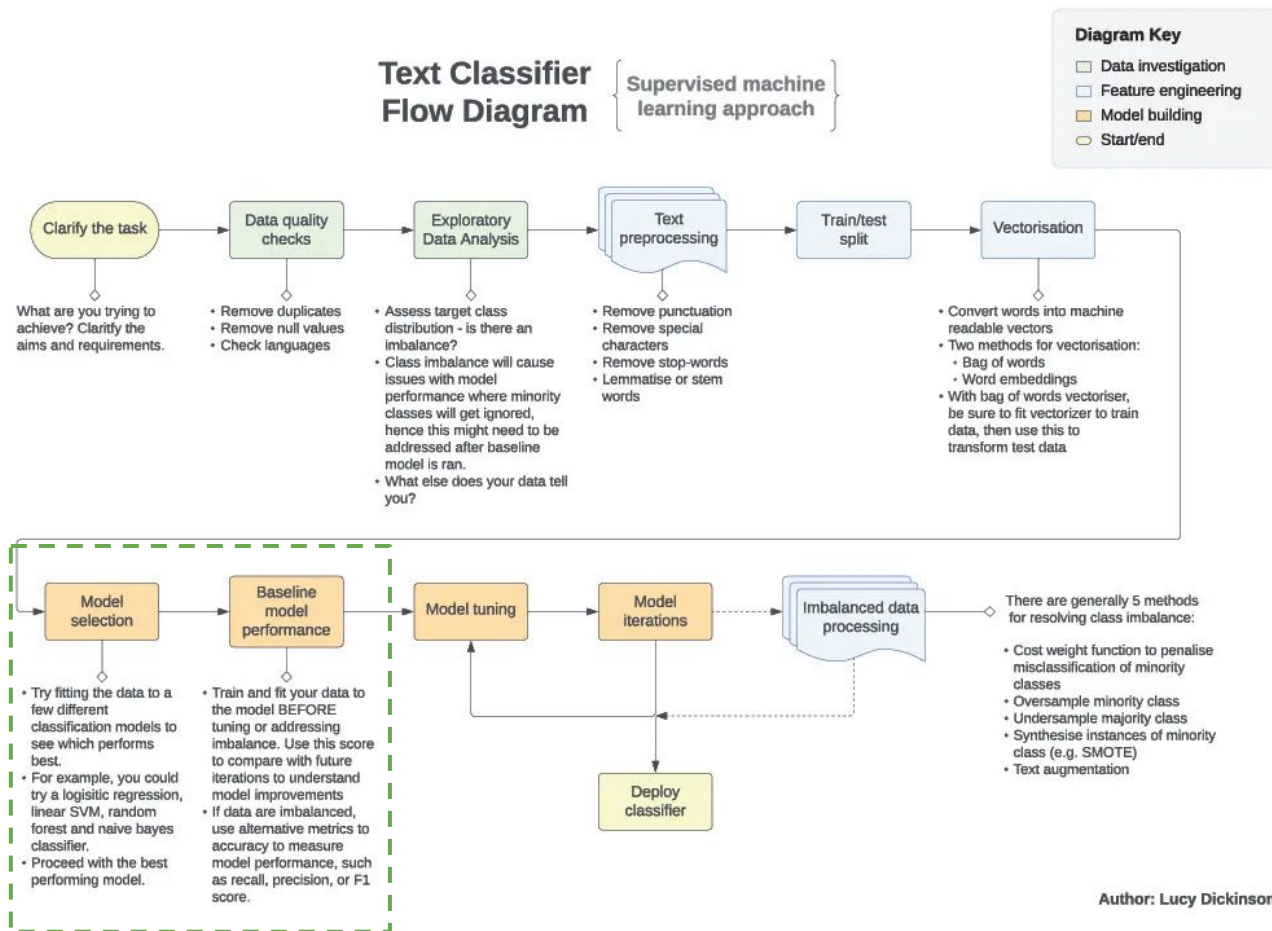


# Data Leakage Problem

This is getting its own sub heading because it is so important to do this step **BEFORE** your start fiddling with the features. Split your data using sklearn's **train\_test\_split()** function and then leave the **test data alone** so there's no risk of **data leakage**.



# Modeling Process



**Author: Lucy Dickinson**

# Classification Evaluation Metrics

## Classification evaluation metrics

In contrast to regression, classification metrics have more variations and must be evaluated carefully. Some metrics are useful for binary and multiclass classification. Some metrics may be chosen when the proportion among classes severely differs.

Often looking at accuracy can be misleading, basically when the data is unbalanced (do you know why? Hint try to calculate the formula when you have 1:90 proportion). There are examples of when to choose other evaluation metrics when accuracy alone is not enough.

Metric	Formula	Interpretation
Accuracy	$\frac{TP + TN}{TP + TN + FP + FN}$	Overall performance of model
Precision	$\frac{TP}{TP + FP}$	How accurate the positive predictions are
Recall Sensitivity	$\frac{TP}{TP + FN}$	Coverage of actual positive sample
Specificity	$\frac{TN}{TN + FP}$	Coverage of actual negative sample
F1 score	$\frac{2TP}{2TP + FP + FN}$	Hybrid metric useful for unbalanced classes

**Choose Recall** if the idea of false positives is far better than false negatives, in other words, **if the occurrence of false negatives is unacceptable/intolerable**, that you'd rather get some extra false positives(false alarms) over saving some false negatives, like in our diabetes example.

You'd rather get some healthy people labeled diabetic over leaving a diabetic person labeled healthy.

**Choose precision** if you want to be more confident of your **true positives**. for example, Spam emails. You'd rather have some spam emails in your inbox rather than some regular emails in your spam box. So, the email company wants to be extra sure that email Y is spam before they put it in the spam box and you never get to see it.

**Choose Specificity** if you want to **cover all true negatives**, meaning you don't want any false alarms, you don't want any false positives. for example, you're running a drug test in which all people who test positive will immediately go to jail, you don't want anyone drug-free going to jail. False positives here are intolerable.

		Predicted class	
		+	-
Actual class	+	<b>TP</b> True Positives	<b>FN</b> False Negatives Type II error
	-	<b>FP</b> False Positives Type I error	<b>TN</b> True Negatives

# Confusion Matrix Multiclass

Actual	12	0	0	0	0	0	1	False Positives	0	0	0			
	0	10	0	0	0	0	0		0	2	0			
	0	0	10	0	0	0	7		0	0	0			
	0	0	1	18	0	0	0		0	TN	0			
	0	0	0	0	18	0	0		0	0	0			
	0	3	2	0	0	15	0		0	0	0			
	0	0	0	0	0	0	17		0	0	0			
	1	0	False Negatives						2	3	TP	0	FN	0
	0	3	1	1	0	0	1		FP	13	TN	0		
	1	1	0	3	0	0	5			0	0	8		
								Predicted						

# Python Time

[https://github.com/muhajirakbarhsb/NLP\\_class\\_2023/blob/main/Week\\_Class\\_5.ipynb](https://github.com/muhajirakbarhsb/NLP_class_2023/blob/main/Week_Class_5.ipynb)