# SC1015:
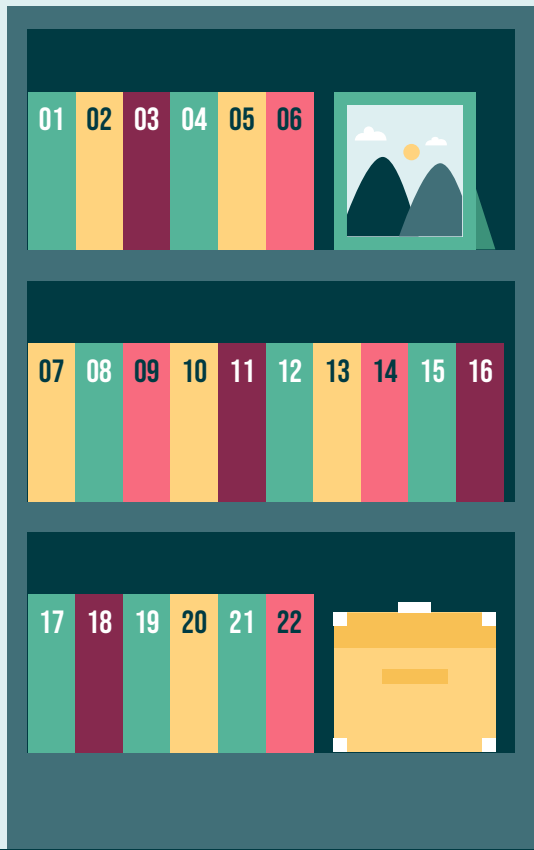# GENRE-RATOR

**SC12 - A Multilabel Classifier of Genres**
- Nathaniel Chin Yi Xuan
- Marcus Soh Yi Qing
- Tan Yan Chi

*https://github.com/natisaver/GoodReads-Multilabel-Genre-Prediction*

# TABLE OF CONTENTS

1. Problem Formulation

2. Data Pre-Processing & Cleaning

3. Exploratory Data Analysis

4. Train Test Split / Standardisation

5. Pipeline (Encoder + Classification Algorithms)

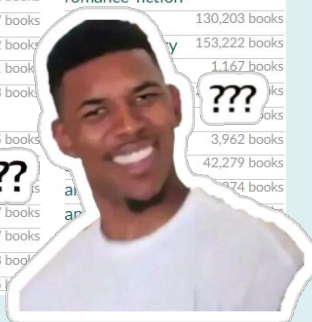6. App Demonstration & Conclusions

7. Future Work

# PROBLEM FORMULATION

Real-Life Problem:

- Over **900** genres on Goodreads website =>
  **Unstandardized** definition of different genres
- Confusing for user navigation

Project Resolution:

- Aims to produce an algorithm that is able to
  predict the combination of genres of a book
  given its **plot description** and **features** to
  create **standardisation**

| | | | | | |
|---|---|---|---|---|---|
| 10th-century | 7,223 books | adult-fiction | 1,772,022 books | american-classics | 117,722 books |
| 11th-century | 8,739 books | adventure | 2,345,871 books | american-fiction | 201,858 books |
| 12th-century | 13,368 books | adventurers | 8,242 books | american-history | 414,817 books |
| 13th-century | 11,838 books | aeroplanes | 906 books | | |
| 14th-century | 19,292 books | africa | 334,114 books | american-novels | 54,505 books |
| 15th-century | 22,449 books | african-american | 199,808 books | | |
| 16th-century | 63,414 books | | | american-revolution | 35,761 books |
| 17th-century | 78,865 books | african-american-literature | 31,547 books | | |
| 1864-shenandoah-campaign | 37 books | | | american-revolutionary-war | 3,533 books |
| 18th-century | 138,519 books | african-american-romance | 14,292 books | | |
| 1917 | 2,272 books | | | americana | 258,767 books |
| 19th-century | 531,321 books | african-literature | 38,781 books | amish | 72,916 books |
| 1st-grade | 86,092 books | | | amish-fiction | 19,365 books |
| 20th-century | 1,045,119 books | agender | 307 books | amish-historical-romance-fiction | 959 books |
| 21st-century | 701,272 books | agriculture | 28,077 books | | |
| 2nd-grade | 86,558 books | ahistory | 52 books | | |
| 40k | 24,297 books | aircraft | 3,747 books | | 130,203 books |
| ableism | 3,666 books | airliners | 22 books | | 153,222 books |
| abuse | 477,476 books | airships | 2,721 books | | 1,167 books |
| academia | 214,259 books | albanian-literature | 2,163 books | | |
| academic | 478,765 books | | | | |
| academics | 84,704 books | alchemy | 35,105 books | | 3,962 books |
| accounting | 7,880 books | alcohol | | | 42,279 books |
| accra | 615 books | alexandria | | | 74 books |
| action | 999,753 books | algebra | 3,207 books | | |
| activism | 101,814 books | algeria | 11,107 books | | |
| adaptations | 121,362 books | algiers | 1,308 books | | |
| addis-ababa | 68 books | algorithms | 7,935 books | | |

# DATASET

The Zenodo Dataset we have chosen contains the following columns.

```
Data columns (total 25 columns):
 #   Column          Non-Null Count    Dtype
---  ------          --------------    -----
 0   bookId          52478 non-null    object
 1   title           52478 non-null    object
 2   series          23470 non-null    object
 3   author          52478 non-null    object
 4   rating          52478 non-null    float64
 5   description     51140 non-null    object
 6   language        48672 non-null    object
 7   isbn            52478 non-null    object
 8   genres          52478 non-null    object
 9   characters      52478 non-null    object
 10  bookFormat      51005 non-null    object
 11  edition         4955 non-null     object
 12  pages           50131 non-null    object
 13  publisher       48782 non-null    object
 14  publishDate     51598 non-null    object
 15  firstPublishDate 31152 non-null   object
 16  awards          52478 non-null    object
 17  numRatings      52478 non-null    int64
 18  ratingsByStars  52478 non-null    object
 19  likedPercent    51856 non-null    float64
 20  setting         52478 non-null    object
 21  coverImg        51873 non-null    object
 22  bbeScore        52478 non-null    int64
 23  bbeVotes        52478 non-null    int64
 24  price           38113 non-null    object
```

# PROJECT PROCESS OVERVIEW

**Preparing Data**
- Setting top 30 genres
- Clean Text Data
- Hot Binary Encoding of Categories
- Extract RGB Features & Relative Brightness of Cover Images

**Train Test Split**
- Iterative Stratification
- Followed by Vectorisation of textual data via TF-IDF encoder

**Scale**
- MinMax or Standard Scaler
- Dependent on classification model used

**All Pipelines**
- Encoder + Classification Method + Model

Input data

Train Data

Test Data

Vectorise & Scale

TF-IDF + Binary Relevance + Logistic Regression

TF-IDF + Label Power Set + Naïve Bayes

TF-IDF + Clustered Label Power Set + Linear SVC
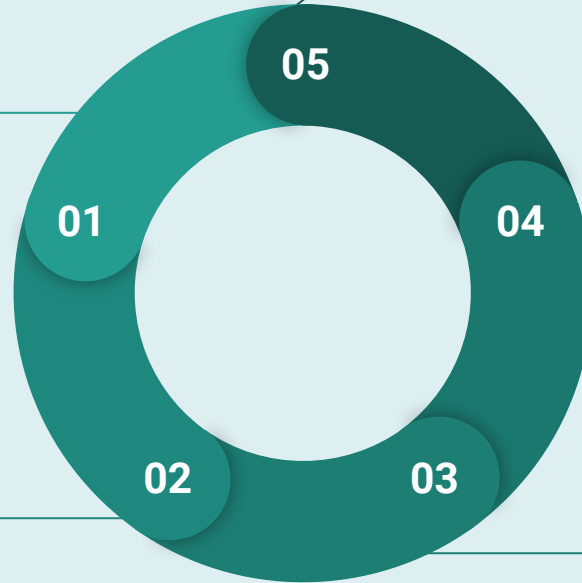
Evaluate Accuracy & F1-Score

# DATA CLEANING

**Removing Non-English Books**

- To ensure that the description of the books are **purely English content**
- Ensuring **consistency** when analysing textual data.

**Dropping columns that are irrelevant**

- Columns like 'price', 'publisher' contain too many **NULL** values to be useful for analysis
- Columns like 'language', 'bookFormat' are **irrelevant** to predicting 'genres'

**Multi-Hot Binary Encoding of Genre**

- Each of the 30 genres were encoded to a **one-hot binary representation**
- If a book belongs to a genre, the value is 1 ("hot") else 0

**Cleaning Genre Column**

- Reducing the initial **967** "genres" to fixed **top 30 genres**

**Cleaning Description Column**

- Cleaning NULL
- Lemmatizing
- Stop Words
- Accented Characters
- Punctuations
- Lower Case

01
02
03
04
05

# PROCESSING COVER IMAGE

- To get additional features, we did image processing to extract the **Red, Green, Blue** Values of the cover images for all the books

- Luminance was then calculated using the formula:

  - ```
    math.sqrt(0.241*(row.r**2) +
    0.691*(row.g**2) + 0.068*(row.b**2))
    ```

# MODIFIED GENRE LIST

Our team manually analysed the top 60 genres and reduced them to a "top 30 genre list"

- **Purpose**
    - Standardisation
    - Combine overlapped genres
- **Example**
    - History, 11th Century, Historical → History
    - Adult Fiction → Adult and Fiction

```
top30genrelist = ['fiction',
'fantasy',
'romance',
'young adult',
'contemporary',
'adult',
'nonfiction',
'history',
'novels',
'mystery',
'historical fiction',
'audiobook',
'science fiction',
'paranormal',
'literature',
'adventure',
'classics',
'thriller',
'childrens',
'magic',
'humor',
'contemporary romance',
'crime',
'suspense',
'middle grade',
'chick lit',
'biography',
'teen',
'horror',
'philosophy']
```

# EXPLORATORY DATA ANALYSIS

# EDA - NUMBER OF BOOKS PER GENRE



**Key Observations**

- **Fiction** has the highest number of books with **21590**

- Followed by **Romance** which is less than half of fiction at **10862**

- The genre with the least number of books is **philosophy** with **1582** books

# EDA - NUMBER OF GENRE PER BOOK



**Key Observations**

- On average, books have **5.28** genres.

- There are **3 books** that is associated with **11** total different genres!

# EDA - NUMBER OF GENRE PER BOOK

| | fiction | fantasy | romance | young adult | contemporary | adult | nonfiction | history | novels | mystery | historical fiction | audiobook | science fiction | paranormal | literature | adventure | classics | thriller | childrens | magic | humor | contemporary romance | crime | suspense | middle grade | chick lit | biography | teen | horror | philosophy |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| adventure | 1.8 | 2.8 | -1.5 | 2.3 | -1.8 | -1.2 | -1.3 | -0.5 | -0.4 | 0.7 | 0.6 | 0.4 | 2.4 | -0.8 | -0.8 | 10 | | | | | | | | | | | | | | |
| classics | 1 | -1.1 | -1.9 | -1.3 | -1.5 | -1.4 | 0.9 | 2.2 | 2.1 | -0.7 | 1.3 | -0.5 | -0.6 | -1.5 | 4.5 | 0 | 10 | | | | | | | | | | | | | |
| thriller | 1.8 | -1 | -1 | -0.9 | 0.3 | 0 | -1.5 | -1.1 | -0. | 6.3 | -.3 | 1.4 | 0 | -0.2 | -0.2 | 0.3 | -0.6 | 10 | | | | | | | | | | | | |
| childrens | 1.4 | 1 | -2.1 | 3.2 | -0.9 | -1 | -1.4 | -1 | -0.7 | -0.1 | 0.1 | -0.5 | -0.2 | -1.2 | -0.9 | 3 | 0.7 | -0.8 | 10 | | | | | | | | | | | |
| magic | 0.6 | 4.4 | 1.5 | 2.2 | -1.6 | -0.3 | -1.8 | -1.7 | -1.4 | -0.5 | -0.9 | -0.3 | 1 | 3.4 | -1.3 | 1.9 | -1.1 | -1.1 | 0.7 | 10 | | | | | | | | | | |
| humor | 0.5 | -0.5 | -0.2 | 0.1 | 0.9 | 0.7 | -0.4 | -1.1 | 0.5 | -0.5 | -0.8 | 0.5 | -0.4 | -0.7 | 0.3 | -0.1 | 0 | -0.7 | 1 | -0.5 | 10 | | | | | | | | | |
| contemporary romance | 0.8 | -2.1 | 4 | 1.2 | 5.2 | 3 | -1.7 | -1.8 | -0.6 | -0.4 | -1.3 | 0.1 | -1.3 | -0.9 | -1.1 | -1.3 | -1.2 | -0.6 | -1.1 | -1 | 0.8 | 10 | | | | | | | | |
| crime | 1.1 | -1.7 | -1 | -1.2 | 0.2 | 0.1 | -0.8 | -0.9 | 0.1 | 5.1 | -0.3 | 1.5 | -0.7 | -0.9 | 0.2 | -0.3 | -0.3 | 5.7 | -0.8 | -1 | -0.4 | -0.3 | 10 | | | | | | | |
| suspense | 1.2 | -1.3 | 0.4 | -0.9 | 1.5 | 0.9 | -1.4 | -1.2 | -0.4 | 4.7 | -0.7 | 1 | -0.5 | -0.3 | -0.7 | 0 | -1 | 5.7 | -1 | -1 | -0.7 | 1.3 | 4.4 | 10 | | | | | | |
| middle grade | 1.5 | 1.2 | -1.8 | 3.7 | -0.5 | -0.8 | -1.4 | -0.8 | -0.6 | 0.3 | 0.3 | -0.1 | 0 | -0.8 | -0.9 | 3.2 | 0.1 | -0.7 | 7 | .1 | 0.6 | -1 | -0.6 | -0.8 | 10 | | | | | |
| chick lit | 1.4 | -1.7 | 3.2 | 0.9 | 4 | 2.4 | -1.4 | -1.1 | 0.1 | -0.3 | -0.6 | 0.4 | -1.2 | -0.8 | -0.6 | -1.1 | -1 | -0.7 | | -0.8 | 1.5 | 3.8 | -0.6 | -0.3 | -0.6 | 10 | | | | |
| biography | -4.4 | -2.1 | -2.1 | -1.9 | -1.4 | -0.7 | 4.9 | 1.6 | -1.3 | -1.2 | -1.3 | 0.3 | -1.2 | -1.2 | -0.4 | -0.7 | -0.3 | -1 | -0.8 | -1 | 0.1 | -0.9 | -0.4 | -0.9 | -0.8 | -0.7 | 10 | | | |
| teen | 1.8 | 0.7 | 0.1 | 4.5 | 0.8 | -0.9 | -1.6 | -1.3 | -1 | 0.4 | -0.3 | -0.3 | 0.1 | -0.2 | -1.2 | 1.7 | -0.4 | -0.3 | 4 | 0.3 | 0.6 | -0.9 | -0.6 | -0.7 | 4.2 | 0.7 | -0.9 | 10 | | |
| horror | 1 | 2 | -1.1 | 0.1 | -1 | -0.7 | -1.2 | -1.1 | 0.1 | 1.5 | -0.7 | -0.2 | 1.1 | 2.3 | -0.5 | -0.4 | -0.3 | 1.9 | -0.4 | -0.4 | -0.5 | -0.8 | 0.4 | 0.9 | -0.3 | -0.8 | -0.7 | -0.1 | 10 | |
| philosophy | -2.7 | -1.5 | -1.7 | -1.6 | -1.1 | -1.2 | 3.1 | 1.6 | -0.5 | -1.1 | -1 | -0.2 | -0.8 | -1 | 0.4 | -1 | 1.3 | -0.9 | -0.9 | -0.8 | -0.5 | -0.8 | -0.7 | -0.7 | -0.7 | 0.2 | -0.9 | -0.6 | | 10 |

# EDA - NUMBER OF GENRE PER BOOK

# EXPLORATORY DATA ANALYSIS



Book Genre: crime

Book Genre: fiction

**Key Observations**

Certain words such as "**kill**", "**murder**" appears in Crime, Mystery, Suspense and Thriller "**world**" and "**life**" seems to be the most common words that appear in book descriptions

# EDA - BRIGHTNESS



Brightness Distribution

**Key Observations**

- Book covers tend to have **very high brightness** or **very low brightness**

# FEATURE RELEVANCE TO GENRE

# CLASSIFYING & MACHINE LEARNING

# TYPES OF CLASSIFICATION



| Binary Classification | Multiclass Classification | Multi-label Classification |
|---|---|---|
| | Labels (t) | Labels (t) |
| | [0 0 1]  [1 0 0]  [0 1 0] | [1 0 1]  [0 1 0]  [1 1 1] |
| • Spam<br>• Not spam | • Dog<br>• Cat<br>• Horse<br>• Fish<br>• Bird | • Dog<br>• Cat<br>• Horse<br>• Fish<br>• Bird |

# PROJECT PROCESS OVERVIEW

**Preparing Data**
- Setting top 30 genres
- Clean Text Data
- Hot Binary Encoding of Categories
- Extract RGB Features & Relative Brightness of Cover Images

**Train Test Split**
- Iterative Stratification
- Followed by Vectorisation of textual data via TF-IDF encoder

**Scale**
- MinMax or Standard Scaler
- Dependent on classification model used

**All Pipelines**
- Encoder + Classification Method + Model

TF-IDF + Binary Relevance + Logistic Regression

TF-IDF + Label Power Set + Naïve Bayes

TF-IDF + Clustered Label Power Set + Linear SVC

Input data

Train Data

Test Data

Vectorise & Scale

Evaluate Accuracy & F1-Score

# STRATIFICATION



Stratification Based on Labelsets

One labelset

# STRATIFICATION



## Example

**Firstly**
Distribute the positive examples of $\lambda_2$

| Instance | $\lambda_1$ | $\lambda_2$ | $\lambda_3$ |
|----------|-------------|-------------|-------------|
| $i_1$ | 1 | 0 | 1 |
| $i_2$ | 0 | 0 | 1 |
| $i_3$ | 0 | 1 | 0 |
| $i_4$ | 1 | 0 | 0 |
| $i_5$ | 0 | 1 | 1 |
| $i_6$ | 1 | 1 | 0 |
| $i_7$ | 1 | 0 | 1 |
| $i_8$ | 1 | 0 | 1 |
| $i_9$ | 0 | 0 | 1 |
| sum | 5 | 3 | 6 |

**1st Fold**

| | | | |
|---|---|---|---|
| desired | 1.7 | 1 | 2 |

**2nd Fold**

| | | | |
|---|---|---|---|
| desired | 1.7 | 1 | 2 |

**3rd Fold**

| | | | |
|---|---|---|---|
| desired | 1.7 | 1 | 2 |

**Secondly**
Distribute the positive examples of $\lambda_1$

| Instance | $\lambda_1$ | $\lambda_2$ | $\lambda_3$ |
|----------|-------------|-------------|-------------|
| $i_2$ | 0 | 0 | 1 |
| $i_9$ | 0 | 0 | 1 |
| sum | - | - | 2 |

**1st Fold**

| | | | |
|---|---|---|---|
| $i_3$ | 0 | 1 | 0 |
| $i_1$ | 1 | 0 | 1 |
| $i_8$ | 1 | 0 | 1 |
| desired | -0.3 | 0 | 0 |

**2nd Fold**

| | | | |
|---|---|---|---|
| $i_6$ | 1 | 1 | 0 |
| $i_7$ | 1 | 0 | 1 |
| desired | -0.3 | 0 | 1 |

**3rd Fold**

| | | | |
|---|---|---|---|
| $i_5$ | 0 | 1 | 1 |
| $i_4$ | 1 | 0 | 0 |
| desired | 0.7 | 0 | 1 |

**Λ2 is distributed first**

Fraction of the number of Books of each Genre in Train and Test Sets

**Distribution of Train/Test Split**
- Train - 70%
- Test - 30%

# TRAIN/TEST SPLIT

| Y | |
|---|---|
| **Top 30 Genre** | Standardised 30 Genres we decided on previously |

| X | |
|---|---|
| **Brightness** | Use of R,G,B values to calculate perceived luminance |
| **Number of Ratings** | Number of ratings the book has on GoodReads |
| **Description** | The book description given by the author |

# PROJECT PROCESS OVERVIEW

**Preparing Data**
- Setting top 30 genres
- Clean Text Data
- Hot Binary Encoding of Categories
- Extract RGB Features & Relative Brightness of Cover Images

**Train Test Split**
- Iterative Stratification
- Followed by Vectorisation of textual data via TF-IDF encoder

**Scale**
- MinMax or Standard Scaler
- Dependent on classification model used

**All Pipelines**
- Encoder + Classification Method + Model

Input data

Train Data

Test Data

Vectorise & Scale

TF-IDF + Binary Relevance + Logistic Regression

TF-IDF + Label Power Set + Naïve Bayes

TF-IDF + Clustered Label Power Set + Linear SVC

Evaluate Accuracy & F1-Score

# VECTORIZATION

**Vectorization**

- Similar concept to one hot
  encoding

- Convert text to numerical
  representation

$$IDF_i = \log\left(1 + \frac{N_D}{f_i}\right)$$

**Inverse Document Frequency** for the search term $i$ within the corpus of documents

**The number of documents** in the corpus of documents that contain the term D

**The number of documents** that contain the search term

**TF-IDF**

- Term frequency-inverse document frequency

- Used to quantify the importance or relevance of string
  representation (words, phrases, etc)

# MACHINE LEARNING

**Scaling**
- Our group decided to scale our data because ML algorithms are sensitive to data scales
- If the data is not scaled, the features with a higher value range starts dominating when calculating distances
- Our team choose **MinMax Scalar** because it scales from [0 to 1] which ensures that no negative values are returned when passing it into the different models

**Machine Learning Algorithms**
- Logistic Regression
- Naive Bayes
- Linear Support Vector Machine

# PROJECT PROCESS OVERVIEW

**Preparing Data**
- Setting top 30 genres
- Clean Text Data
- Hot Binary Encoding of Categories
- Extract RGB Features & Relative Brightness of Cover Images

**Train Test Split**
- Iterative Stratification
- Followed by Vectorisation of textual data via TF-IDF encoder

**Scale**
- MinMax or Standard Scaler
- Dependent on classification model used

**All Pipelines**
- Encoder + Classification Method + Model

Input data

Train Data

Test Data

Vectorise & Scale

TF-IDF + Binary Relevance + Logistic Regression

TF-IDF + Label Power Set + Naïve Bayes

TF-IDF + Clustered Label Power Set + Linear SVC

Evaluate Accuracy & F1-Score

# MULTI-LABEL CLASSIFICATION ALGORITHMS

**Binary Relevance**
- Treat each label as a separate class classification
- 30 genres => 30 binary classifiers



**Label Powerset**
- Treat each unique genre combination as a class
- 6436 unique combinations => 6436 classes



**Label Powerset with Clustering**
- Reduce the number of genre combinations by clustering from 6436 to 100
- k=100 (number of clusters) gave us the highest average F1-score

# PIPELINE ANALYSIS

# PROJECT PROCESS OVERVIEW

**Preparing Data**
- Setting top 30 genres
- Clean Text Data
- Hot Binary Encoding of Categories
- Extract RGB Features & Relative Brightness of Cover Images

**Train Test Split**
- Iterative Stratification
- Followed by Vectorisation of textual data via TF-IDF encoder

**Scale**
- MinMax or Standard Scaler
- Dependent on classification model used

**All Pipelines**
- Encoder + Classification Method + Model

Input data

Train Data

Test Data

Vectorise & Scale

TF-IDF + Binary Relevance + Logistic Regression

TF-IDF + Label Power Set + Naïve Bayes

TF-IDF + Clustered Label Power Set + Linear SVC

Evaluate Accuracy & F1-Score

# EVALUATION OF MODELS



relevant elements

false negatives

true negatives

true positives

false positives

retrieved elements

How many retrieved items are relevant?

$$Precision = \frac{\phantom{xxxxx}}{\phantom{xxxxx}}$$

How many relevant items are retrieved?

$$Recall = \frac{\phantom{xxxxx}}{\phantom{xxxxx}}$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1\text{-}score = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

# EVALUATION OF MODELS

| | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| fiction | 0.82 | 0.62 | 0.71 | 6477.0 |
| fantasy | 0.83 | 0.42 | 0.55 | 3208.0 |
| romance | 0.72 | 0.55 | 0.62 | 3298.0 |
| young adult | 0.70 | 0.43 | 0.53 | 2774.0 |
| contemporary | 0.67 | 0.27 | 0.39 | 2157.0 |
| adult | 0.32 | 0.47 | 0.38 | 2039.0 |
| nonfiction | 0.47 | 0.79 | 0.59 | 2176.0 |
| biography | 0.20 | 0.59 | 0.30 | 681.0 |
| teen | 0.56 | 0.11 | 0.18 | 1034.0 |
| horror | 0.44 | 0.08 | 0.14 | 575.0 |
| philosophy | 0.64 | 0.31 | 0.42 | 475.0 |
| Avg/Total | 0.58 | 0.42 | 0.43 | 48132.0 |

# EVALUATION OF MODELS

| Machine Learning Model | Precision | Recall | F1-Score |
|---|---|---|---|
| Binary Relevance + Logistic Regression | 0.65 | 0.26 | 0.29 |
| **Label Powerset + Naive Bayes** | **0.58** | **0.42** | **0.43** |
| Label Powerset Clustering + Linear Support Vector Machine | 0.29 | 0.26 | 0.27 |

# FUTURE WORK

**Expand the list of genres**
Increase from size 30 to cover a wider range of genres or make it more in depth

**1**

**Utilise different Text Encoder**
Google Universal Sentence Encoder

**2**

**Utilising different Scaling techniques**
Robust Scalar

**3**

**Utilising other Machine Learning Models**
- Neural Network
- Cosine Similarity
- Adjusting hyperparameters

**4**

# NEW TECHNOLOGIES USED

# THANK YOU!

### References
- https://scikit-learn.org/stable/modules/multiclass.html
- https://realpython.com/image-processing-with-the-python-pillow-library/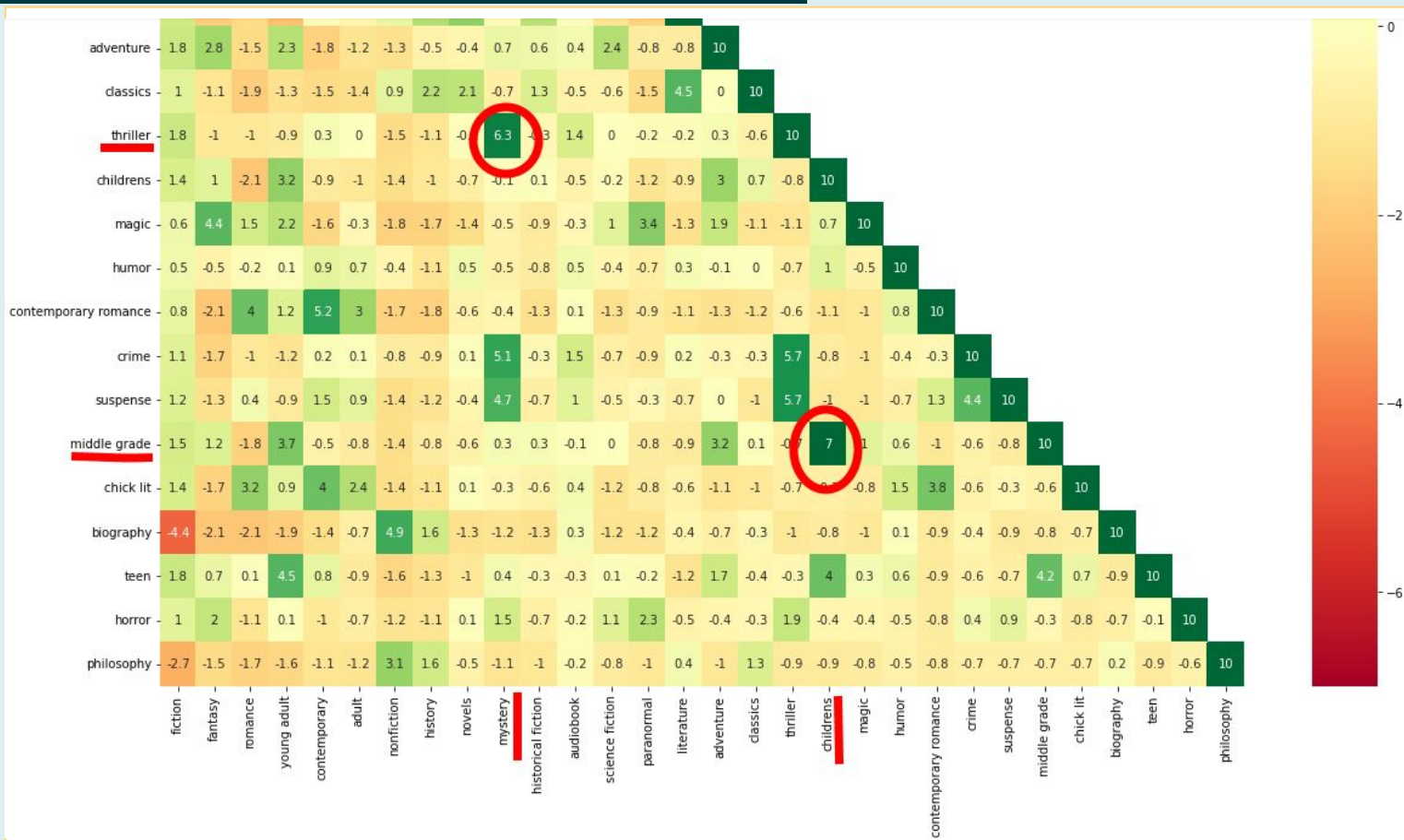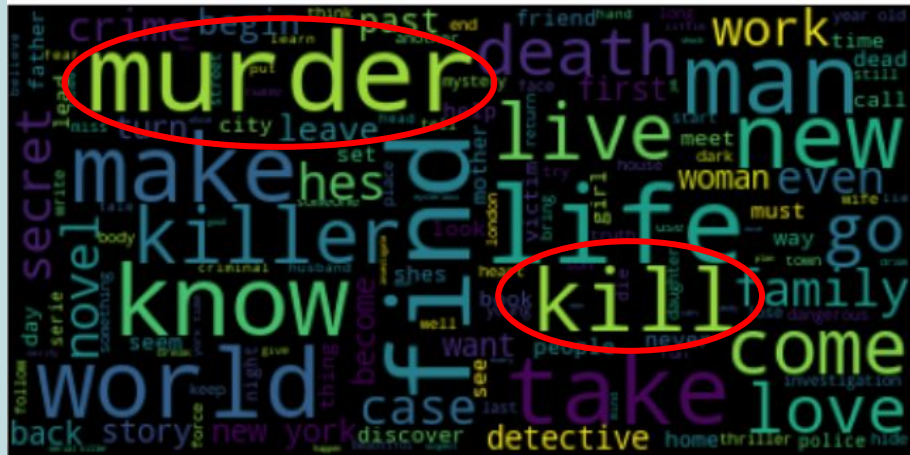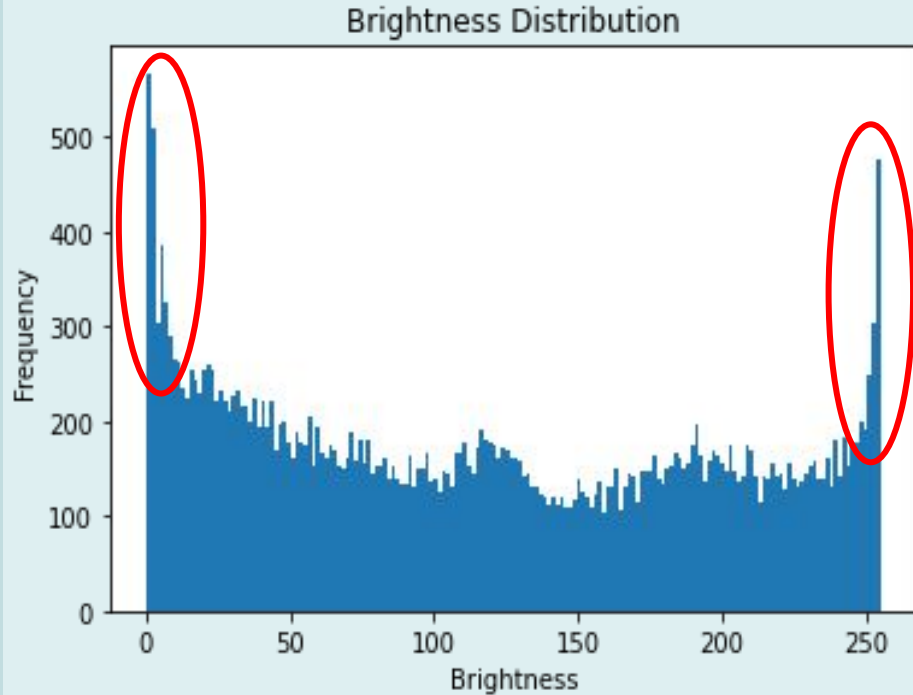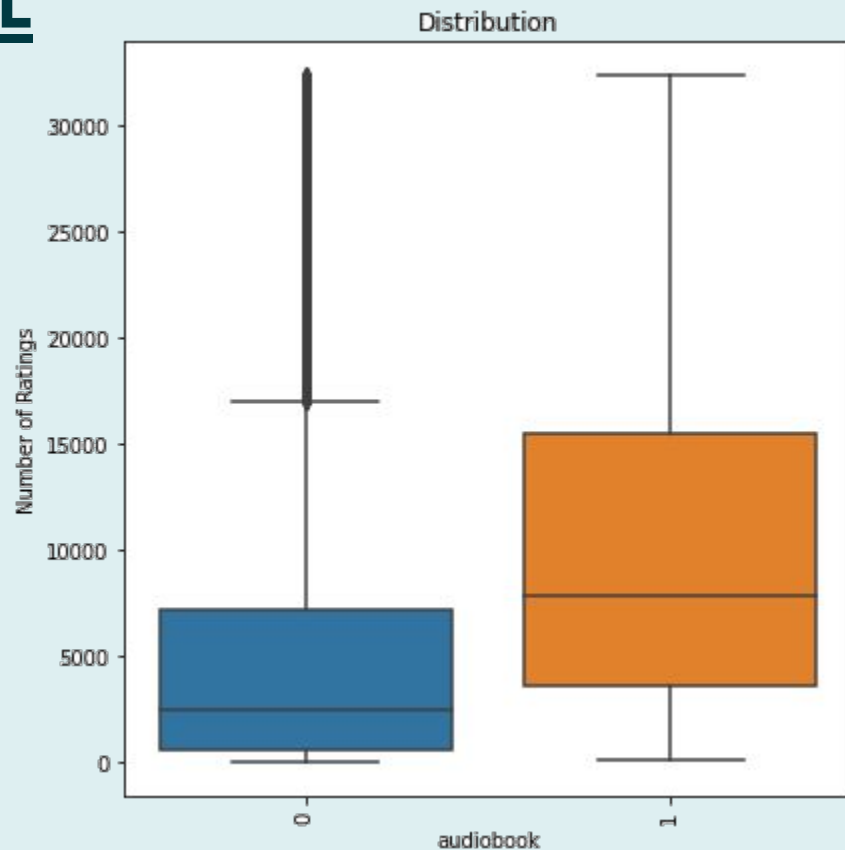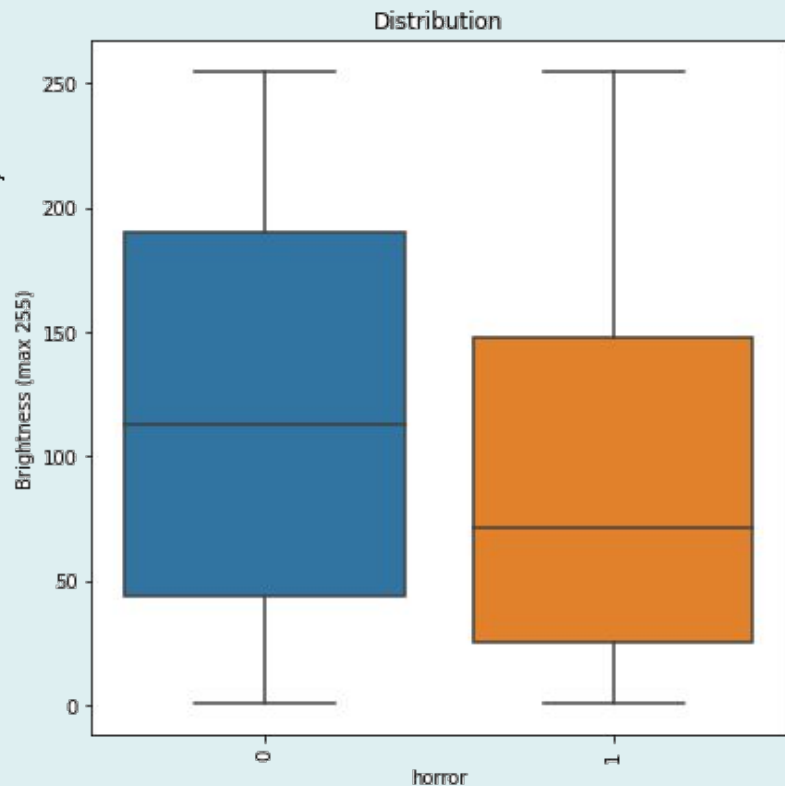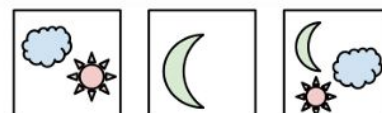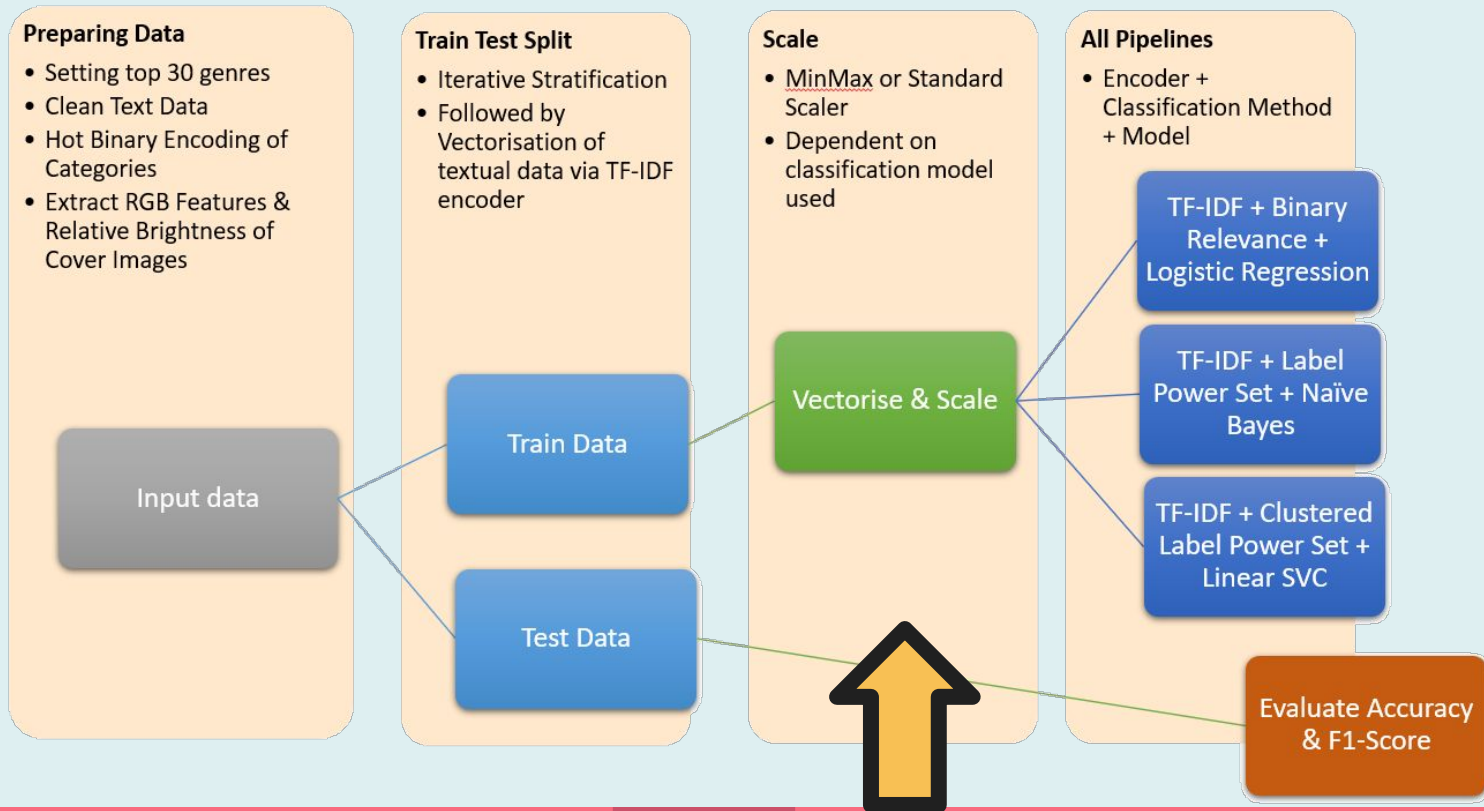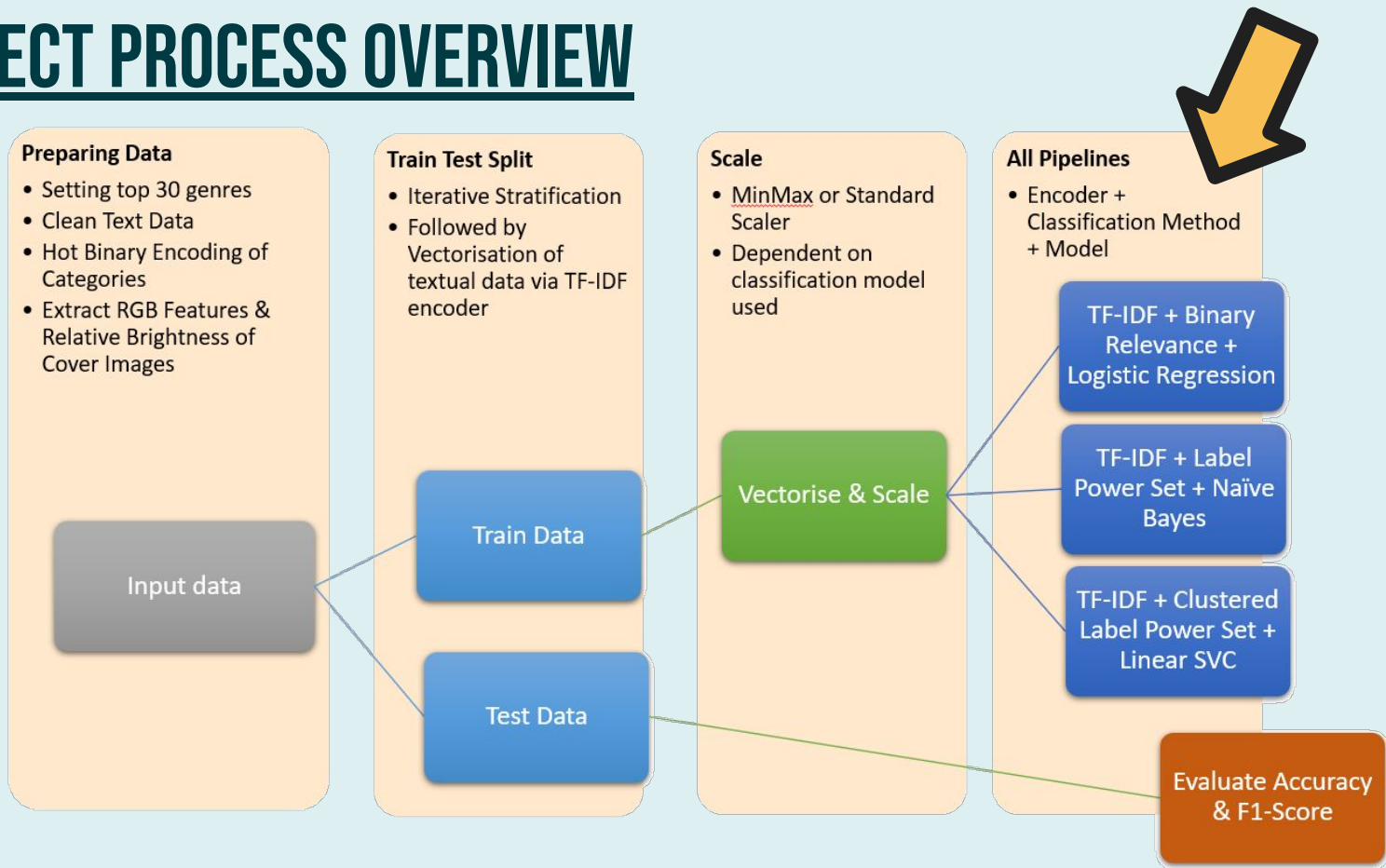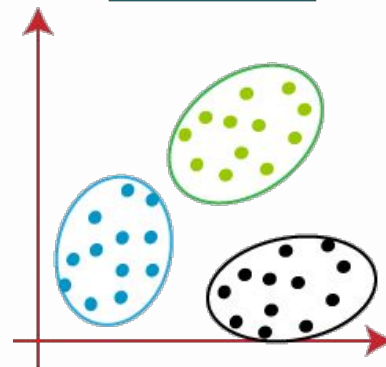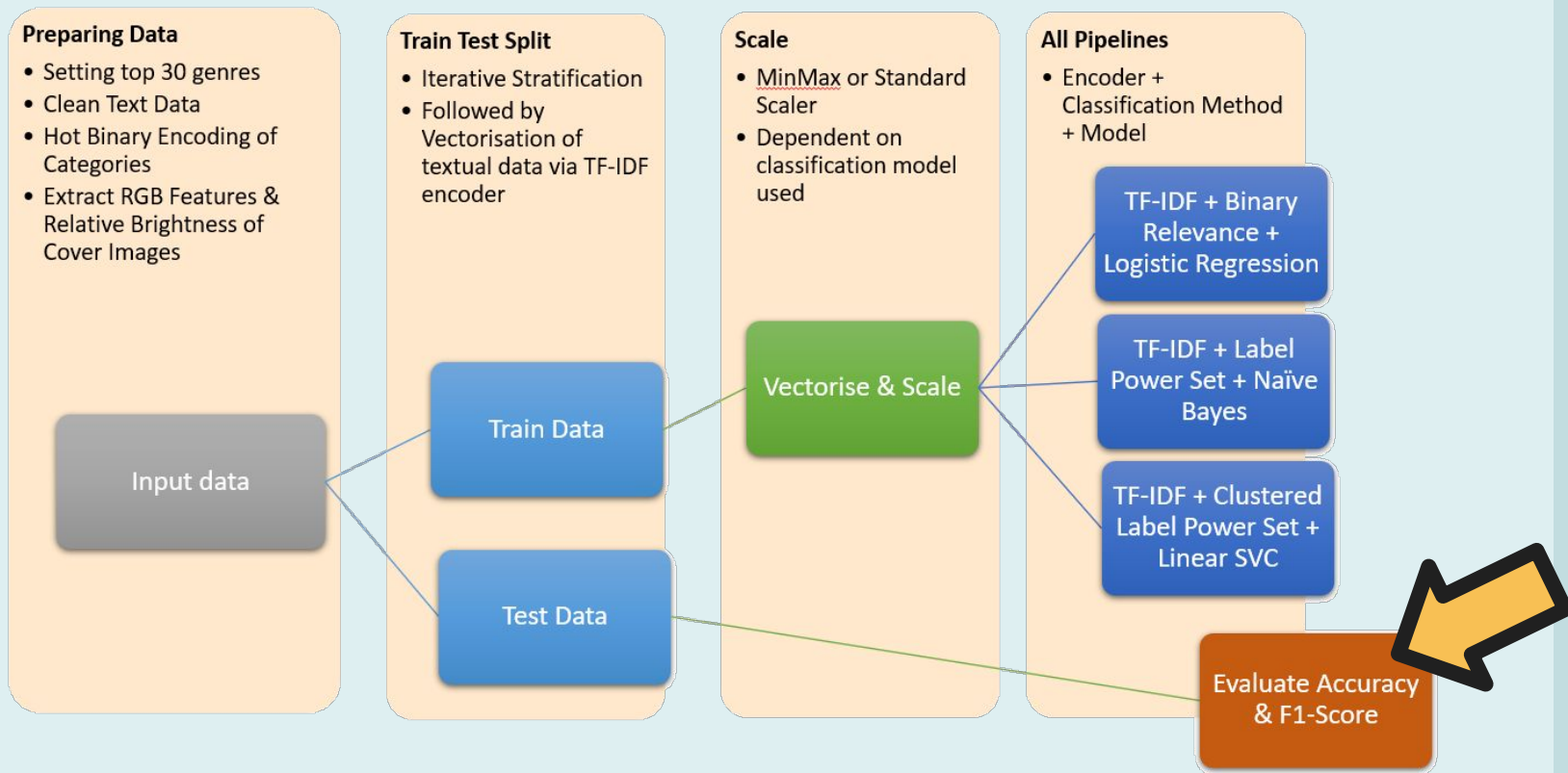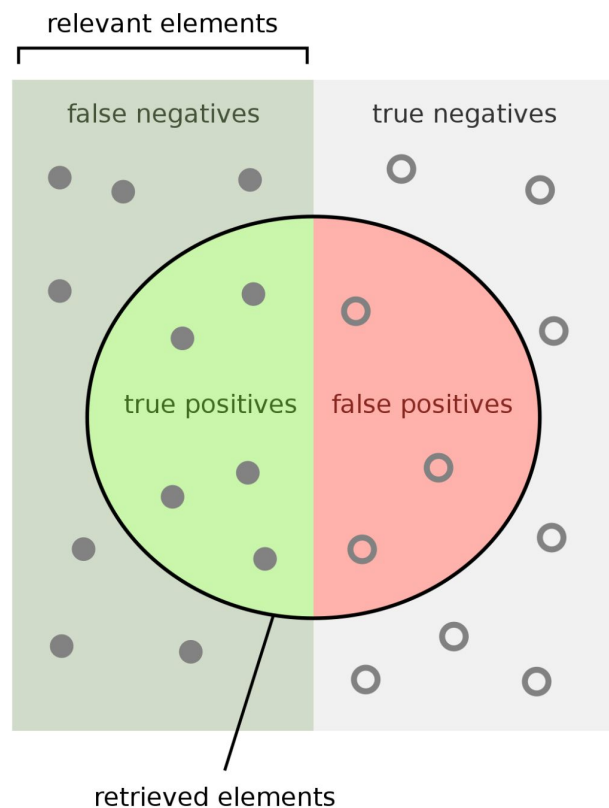