# SICSS-Edinburgh

## 2022

# Text Classification - Practical

**Björn Ross**

17 June 2022

# Lecture Objectives

- <u>Implement</u> your first text classifier easy steps

- This is a practical lecture
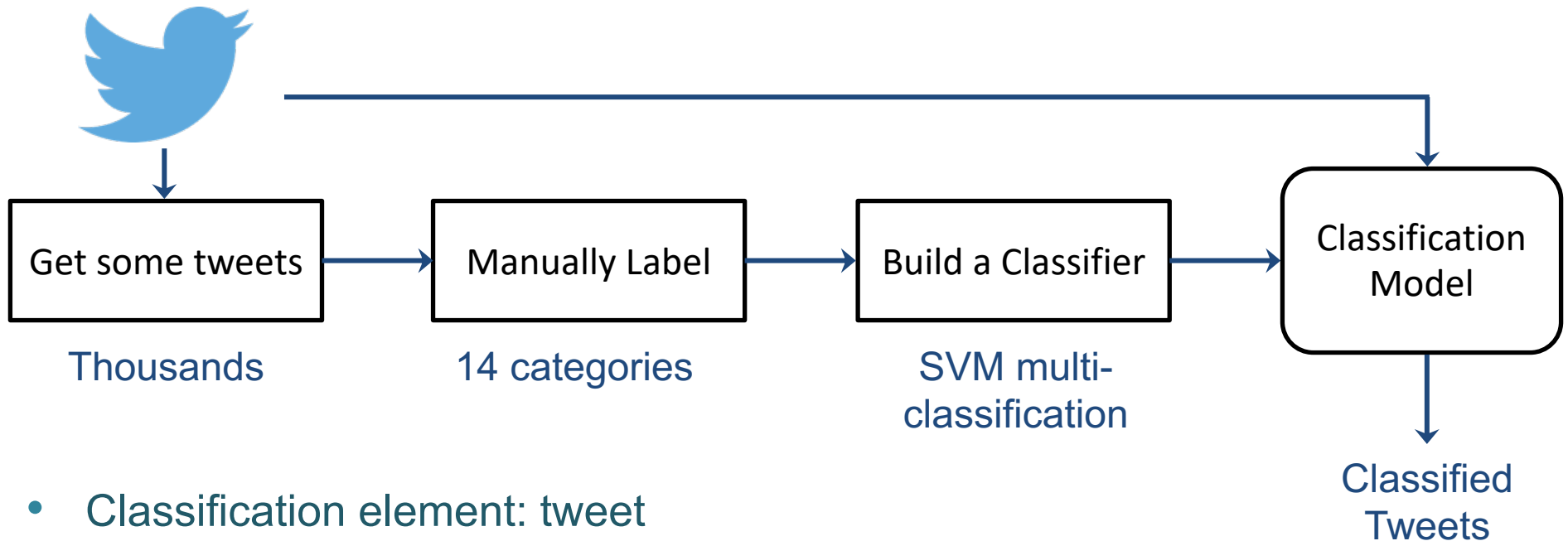  No equations this time ☺

THE UNIVERSITY of EDINBURGH

# My first text classifier: Ingredients

- Text elements to be classified
  - Document, paragraph, sentence

- Set of predefined classes (classification task)
  - At least two (binary)
  - Topical, spam, relevance, sentiment, …

- Training set
  - Enough samples of text elements for each class

- Test set (+ possible validation set)
  - Some samples of each class that not used in training

- Features set
  - A set of features extracted from the text to train the classifier

- Classifier
  - The ML module that learns a classification model

# My first text classifier: Application

- Classifying tweets into general-purpose categories

| Get some tweets | Manually Label | Build a Classifier | Classification Model |
|---|---|---|---|
| Thousands | 14 categories | SVM multi-classification | |

Classified Tweets

- Classification element: tweet
- Classes: 14 categories: sports, politics, comedy, …
- Training/test set: 3129 tweets → 80/20% for train/test
- Features: BOW
- Classifier: SVM multiclass classifier

# My first text classifier: Steps

1. Prepare training data
   required: piece of text (tweet) + label to class

2. Extract features
   1. Pre-process text: lowercase, tokenise, remove useless strings
   2. Create a list of all unique terms in the training data. Give each term a unique ID
   3. Convert the text into features, by replacing each term with its corresponding feature ID. Add value to the feature (simplest: value "1" if exists, or count of occurrences)

3. Prepare test file
   Convert test file text into features using the same mapping from the training data. For terms that are not in the features list, it could be neglected, or assigned to an ID representing OOV.

4. Run the learning process on the training data features to create a model

5. Run the classification on the features of the test data and get predictions

6. Evaluate performance

THE UNIVERSITY
*of* EDINBURGH

# Examples

- Tweet + Label

| Kobe passes Wilt for 4th on all-time scoring list | Sports |
|---|---|

- Learned features (BOW) from training data

- After converting text to feature vectors

Feature ID → | Corresponding word

|   | 0 | | 2943 | 2944 | 2945 | 2946 | | 8330 | 8331 | | 10000 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | … | 0 | 1 | 0 | 0 | | 1 | 0 | … | 0 |
| 1 | 0 | … | 1 | 0 | 0 | 1 | … | 0 | 0 | … | 0 |
| … | | | | | | | … | | | … | |

| 2944 | kobe |
|---|---|
| 2945 | rapping |
| .. | |
| 4525 | 4th |
| 4526 | trevi |
| .. | |
| 8330 | passes |
| 8331 | ducks |
| .. | |
| 9929 | 17 |
| 9930 | wilt |
| ... | |

- SVM prediction output

| 7 |
|---|

Predicted Class ID

THE UNIVERSITY of EDINBURGH

# Practical

THE UNIVERSITY
of EDINBURGH

# Other things to try?

- Try a different classifier

  - Suitable for people with little or no previous Python experience

  - Go to https://scikit-learn.org/stable/modules/multiclass.html

- Try a different dataset

  - Suitable for people with some knowledge of Python

- Try extracting different features

  - Suitable for people familiar with Python

Pair up and experiment, then tell us about your results!

THE UNIVERSITY of EDINBURGH