

Data Science Capstone

Mini Project

Ajay Anand and Cantay Caliskan



Contact Information

Instructor

- Cantay Caliskan (pronounced 'Jantai')
- Background in computer science, stats, social sciences
- E-mail: cantay.caliskan@rochester.edu
- Office: Wegmans Hall, 1205
- Office Hours: By appointment.



A Basic Rule

- You can work in a team of maximum two (2) people in this project (including yourself)



A Short Overview

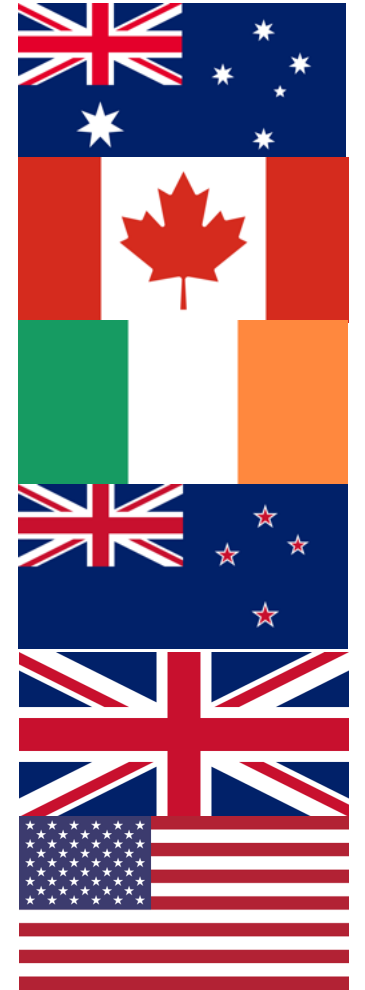
- A **Kaggle** Challenge !
 - *You will be classifying an outcome by using different features.*
 - *Deadline:* September 25, Friday (11:59 PM)



Data

Covid-19 Tweets

- Covid-19 related tweets collected from six different English-speaking countries:
 - *Australia, Canada, Ireland, New Zealand, the United Kingdom, the United States*
- **240,000 tweets** in the training dataset
- **60,000 tweets** in the test dataset
- Outcome category / label is '**country**' (found in the training dataset)



Information about the dataset

text: The text of tweet (including emojis, htmls, hashtags)

reply_to_screen_name: The Twitter screen name of the user the owner of the tweet is replying to (if any)

is_quote: A Boolean variable that indicates if the owner of the tweet is quoting someone else's tweet

is_retweet: A Boolean variable that indicates if the owner of the tweet is retweeting someone else's tweet

hashtags: A list of hashtags included in the tweet

country: The label of the country in which the tweet was posted (found only in the **training** dataset)*

id: An index number associated with tweets (found only on the **test** dataset)

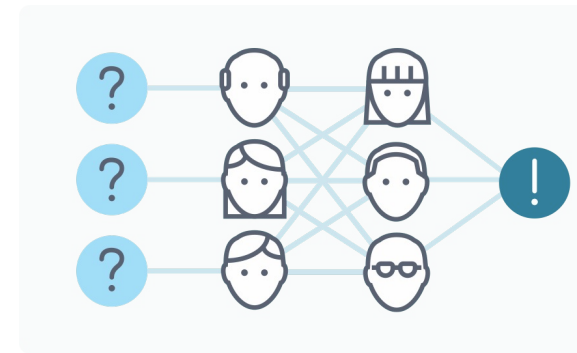
*Countries are labeled as '***us***' (United States), '***uk***' (United Kingdom), '***canada***' (Canada), '***australia***' (Australia), '***ireland***' (Ireland), and '***new_zealand***' (New Zealand).

Tasks

Tasks

- You will work on three (3) groups of tasks:
 - ***Descriptive Analysis***
 - Descriptive tables, descriptive graphs, Natural Language Processing (NLP)
 - ***Kaggle Competition***
 - You will create a model that provides the highest ***Accuracy*** value for predicting correct '***country***' labels by using the tweets
 - ***Lab Report***
 - Provide your findings in a traditional IEEE format
 - *Abstract, Introduction, Data, Methods, Results*

kaggle



Descriptive Analysis

3 questions for undergraduate students, 6 questions for graduate students:

- a) Descriptive table that contains a summary of the dataset
- b) Stacked bar chart showing hashtag-related trends
- c) Latent-Dirichlet-Allocation (LDA) analysis to find the ‘topics’ in your dataset
- d) Non-Negative Matrix Factorization (NMF) to find the ‘topics’ in your dataset (one more time)
- e) Text cleaning and text lemmatization, expanding the descriptive analysis
- f) Calculating cosine similarities between the country-specific hashtag count vectors to explore similarities/differences between countries

Online Competition

- Online competition you can enter on **Kaggle**:
 - <https://www.kaggle.com/competitions/classification-of-covid-19-tweets-capstonefall2022/>
- **Goal:** Develop a classification model that classifies the tweets with the highest **Accuracy** possible
 - **No model restrictions!**
 - You can:
 - Use any classification algorithm that you think will give the highest accuracy
 - Perform any type of feature engineering
 - Perform different *NLP* tasks that you can later use as input for your model
 - Perform weighting, dimensionality reduction etc.
 - Use any external dataset to enrich your training and test datasets

Important:

- Use **training_data.csv** to *train your model*
- Use **sample_submission.csv** to submit your answers
- You can send up to 10 submissions every day (competition is currently open!)
- Produce a *confusion matrix*, and your *Accuracy* score in your report

Online Competition: Further Do's and Don'ts

- Code:

- You can use **Python** (only)
- Your code should be *executable*, i.e.:
 - We should be able to run your code by running the cells *consecutively*
 - We should also be able to run your code on a *laptop*
 - We should be able to *understand* what your code is doing. So, please make sure that:
 - you write **a lot of comments** describing your code
 - you only include the code that works
 - you only include your best solution
 - you name your variables mutually intelligibly (i.e. **tweet_data**, not **td123** etc.)

- Model:

- Your classification must give one of the labels provided in the instructions:
 - **'us', 'uk', 'canada', 'australia', 'ireland', 'new_zealand'**

Lab Report

- We will be using a **LATEX** template (*IEEE*) to produce our reports.
 - Link to the template found in the instructions.

Preparation of Papers for IEEE Sponsored Conferences & Symposia*

Huibert Kwakernaak¹ and Pradeep Misra²

Abstract—This electronic document is a “live” template. The various components of your paper [title, text, heads, etc.] are already defined on the style sheet, as illustrated by the portions given in this document.

I. INTRODUCTION

This template, modified in MS Word 2003 and saved as “Word 97-2003 & 6.0/95 – RTF” for the PC, provides authors with most of the formatting specifications needed for preparing electronic versions of their papers. All standard paper components have been specified for three reasons: (1) ease of use when formatting individual papers, (2) automatic compliance to electronic requirements that facilitate the concurrent or later production of electronic products, and (3) conformity of style throughout a conference proceedings. Margins, column widths, line spacing, and type styles are built-in; examples of the type styles are provided throughout this document and are identified in italic type, within parentheses, following the example. Some components, such as multileveled equations, graphics, and tables are not prescribed, although the various table text styles are provided. The formatter will need to create these components, incorporating the applicable criteria that follow.

II. PROCEDURE FOR PAPER SUBMISSION

A. Selecting a Template (Heading 2)

First, confirm that you have the correct template for your paper size. This template has been tailored for output on the US-letter paper size. Please do not use it for A4 paper since the margin requirements for A4 papers may be different from Letter paper size.

B. Maintaining the Integrity of the Specifications

The template is used to format your paper and style the text. All margins, column widths, line spaces, and text fonts are prescribed; please do not alter them. You may note peculiarities. For example, the head margin in this template measures proportionately more than is customary. This measurement and others are deliberate, using specifications that anticipate your paper as one part of the entire proceedings, and not as an independent document. Please do not revise any of the current designations

*This work was not supported by any organization

¹H. Kwakernaak is with Faculty of Electrical Engineering, Mathematics and Computer Science, University of Twente, 7500 AE Enschede, The Netherlands h.kwakernaak at papercept.net papercept.net

²P. Misra is with the Department of Electrical Engineering, Wright State University, Dayton, OH 45435, USA p.misra at ieee.org

III. MATH

Before you begin to format your paper, first write and save the content as a separate text file. Keep your text and graphic files separate until after the text has been formatted and styled. Do not use hard tabs, and limit use of hard returns to only one return at the end of a paragraph. Do not add any kind of pagination anywhere in the paper. Do not number text heads—the template will do that for you.

Finally, complete content and organizational editing before formatting. Please take note of the following items when proofreading spelling and grammar:

A. Abbreviations and Acronyms

Define abbreviations and acronyms the first time they are used in the text, even after they have been defined in the abstract. Abbreviations such as IEEE, SI, MKS, CGS, sc, dc, and rms do not have to be defined. Do not use abbreviations in the title or heads unless they are unavoidable.

B. Units

- Use either SI (MKS) or CGS as primary units. (SI units are encouraged.) English units may be used as secondary units (in parentheses). An exception would be the use of English units as identifiers in trade, such as “3.5-inch disk drive”.
- Avoid combining SI and CGS units, such as current in amperes and magnetic field in oersteds. This often leads to confusion because equations do not balance dimensionally. If you must use mixed units, clearly state the units for each quantity that you use in an equation.
- Do not mix complete spellings and abbreviations of units: “Wb/m²” or “webers per square meter”, not “webers/m²”. Spell out units when they appear in text: “...a few henries”, not “...a few H”.
- Use a zero before decimal points: “0.25”, not “.25”. Use “cm³”, not “cc”. (bullet list)

C. Equations

The equations are an exception to the prescribed specifications of this template. You will need to determine whether or not your equation should be typed using either the Times New Roman or the Symbol font (please no other font). To create multileveled equations, it may be necessary to treat the equation as a graphic and insert it into the text after your paper is styled. Number equations consecutively. Equation numbers, within parentheses, are to position flush right, as in (1), using a right tab stop. To make your equations more compact, you may use the solidus (/), the exp function, or appropriate exponents. Italicize Roman symbols

1) Abstract

2) Introduction

3) Data

4) Methods

5) Results

Lab Report: Do's and Don'ts

- A few important notes:
 - Think of the report as an **essay**!
 - The report should have the following sections:
 - *Abstract, Introduction, Data, Methods, Results*
 - Try to create a good flow, and a 'story-like' report
 - Spend enough time on explaining:
 - Your data
 - Your goals
 - What you did to achieve these goals
 - What you think you could have done to achieve better results
 - Criticize yourself!

Deliverables

Deliverables

- Your code in **.ipynb** format
 - Add a lot of comments to your code!
- Your ranking in **Kaggle** system
- The lab report in **IEEE** format published as a **.pdf** file
 - Lab report should include all of the visuals, tables, and the confusion matrix.
- Submit everything through **BlackBoard**



Grading

Mini Project: Grading

- You will be graded based on the following criteria:
 - Code
 - Cleanliness/understandability (i), executability (ii), format (iii)
 - Ranking
 - Ranking in the **Kaggle** competition
 - Lab Report
 - *Introduction* (i), *Data* (ii), *Methods* (iii), *Results* (iv)
 - Flow, readability, level of detail, quality of visuals/tables, adherence to the guidelines

More about Grading

- Other important information about Kaggle competition:
 - The lowest grade you can get from the **ranking** component will be **75/100**.
 - The highest ranked project will get **100/100** for the **ranking** component.
 - However:
 - If your accuracy is close to the accuracy achieved by randomization of the outcome variable, your grade may be lower (and it may be zero, as well).
- **Graduate** and **undergraduate** students will be ranked separately.

kaggle

And one last point...

- Let's say you have achieved a really good (or maybe a really bad) accuracy and you are done with model training:
 - **Please do not post the solutions online!**
 - **Or, simply said, please do not post any related code online 😊**

kaggle