

Stance Detection

General Requirements

Students are expected to work on the coursework individually.

The students will work on the task of fake news detection. The task focuses on solving fake news detection by predicting the stance associated to every new article using the FNC dataset (<https://github.com/FakeNewsChallenge/fnc-1>). FNC dataset consists of 49,973 pairs of headlines and article bodies. The body text is annotated by the following classes: agree, disagree, discuss, or unrelated to the headline.

Input: A headline and a body text

Output:

Classify the stance of the body text relative to the claim made in the headline into one of four categories:

- Agrees: The body text agrees with the headline.
- Disagrees: The body text disagrees with the headline.
- Discusses: The body text discusses the same topic as the headline, but does not take a position
- Unrelated: The body text discusses a different topic than the headline

The data suffers from an imbalance problem where around 70% of the articles are unrelated.

Students are expected to:

- 1- Implement word embeddings using standard semantic based techniques and neural techniques.
- 2- Understand the challenges of the provided problem and suggest solutions to handle the imbalance nature of the data.
- 3- Implement natural language processing models and classifiers to predict the right category of a given test example.
- 4- Develop hierarchical deep learning models.

Individual Report [100%]

Each student should separately develop their own NLP models to classify news articles into one of the four categories. Write a report (max 1,500 words) on the **challenges** the dataset present, the **solutions**, and your **findings** which will be assessed as follows:

- 1) Apply the following feature extraction techniques and explain how they work and discuss their advantages and disadvantages
 - a) Term Frequency-Inverse Document Frequency (TF-IDF) [5%]
 - b) A Transformer of your choice (e.g., BERT, GPT) [10%]
- 2) Two steps Classification:
 - a) **Related/Unrelated classification:**
 - i) Use the features extracted using TF-IDF (1.a) and your chosen transformer (1.b) to train a standard Machine Learning method e.g., SVM, Logistic Regression, Random Forest, and discuss its performance on the testing set to classify whether the article body is related or unrelated given the headline. [10%]
 - ii) Train one Deep Learning model (e.g., LSTM, RNN, CNN, hybrid model) using TF-IDF (1.a) and transformer embedding (1.b). Explain the architecture of the deep learning model, the hyper-parameters used, and the loss function. Discuss the performance on the testing set to classify whether the article body is related or unrelated given the headline. [15%]
 - iii) Analyse and compare the performance results for both ML and DL models. [5%]
 - b) **Agree/Disagree/Discuss classification:** Build a new neural model on top of the best performing models you implemented in a)
 - i) Build a deep learning model of your choice to classify articles into the remaining three categories (Agree/Disagree/Discuss) [15%].
 - ii) Analyse the performance of your model and report the results. [10%]
 - c) Combine the two models in a) and b) to **test** your model end to end, report and discuss the overall performance of your solution. [15%]
- 3) What are the ***ethical implications*** of your proposed solutions? What are the potential biases and future misuse cases? [10%]
- 4) Academic English writing, with good use of technical vocabulary, correct grammar, appropriate document structure and referencing where relevant. [5%]

The summative submission deadline is 05/05/2022

The coursework aims at evaluating the students' knowledge and their understanding of the fundamentals and advances in NLP and not their programming skills. Therefore, we ask you to implement the solutions using any Python libraries you are most comfortable with, this includes and is not limited to, PyTorch, Keras, TensorFlow, SpaCy, HuggingFace, Gensim, NLTK...

The report should include the following sections: Introduction, Problem Definition, Proposed Solutions, Analysis of Results, Discussion, Ethical Implications, and Conclusion. The report should use diagrams, figures, and tables to demonstrate the results and analysis.

You should submit your 1,500-word report and the associated Jupyter notebook used to produce your analysis and graphs. **Jupyter notebook file should be saved along with all the produced outputs, results, and figures.**

We understand that not every student has access to the same equipment and therefore this could introduce bias in model performance regardless of the quality of proposed solutions. Therefore, using high spec GPUs that can accelerate the performance with longer runs (e.g., epochs) will not grant the student any extra marks and will not be considered.

The report word count should:

- *Include* all the text, including title, preface, introduction, in-text citations, quotations, footnotes, and any other item not specifically excluded below.
- *Exclude* diagrams, tables (including tables/lists of contents and figures), equations, executive summary/abstract, acknowledgements, declaration, bibliography/list of references and appendices. However, it is not appropriate to use diagrams or tables merely as a way of circumventing the word limit. If a student uses a table or figure as a means of presenting his/her own words, then this is included in the word count.

Examiners will stop reading once the word limit has been reached, and work beyond this point will not be assessed. Checks of word counts will be carried out on submitted work. Checks may take place manually and/or with the aid of the word count provided via an electronic submission. Students are strongly advised to use Arial font size 12 for their assignments.

PLAGIARISM and COLLUSION

Your assignment will be put through the plagiarism detection service on Ultra.

Students suspected of plagiarism, either of published work or work from unpublished sources, including the work of other students, or of collusion will be dealt with according to Computer Science and University guidelines.