

University of the Witwatersrand  
School of Computer Science and Applied Mathematics  
COMS4054A& COMS7062A & COMS7066A:  
Natural Language Processing  
Lab 1 (Word2Vec)

## 1 Instructions

For this lab you will be implementing the skip-gram version of Word2Vec discussed in the lecture (mapping from the focal word to context words). You may work in groups of two or three people. To complete the base version of the lab you must:

1. Read in a paragraph of text from a text file from one of the Harry Potter books provided (“HP1.txt” for example is the first Harry Potter book). You **do not** need to use all of the dataset. Only as much as you feel you need to complete the lab. These text files are also not completely cleaned or homogeneous so keep an eye out for that.
2. Split the text by whitespace and remove punctuation so that only the words are left.
3. Extract the unique words from the text and collect them in an array (in order of the word’s first appearance in the text).
4. There are two options here:
  - Map each word to a unique *1-hot representation* where the dictionary size is the number of words in the original text.
  - Map each word to a unique *index in your dictionary* where the dictionary size is the number of words in the original text.
5. For every word in the original paragraph (not unique word - so the same word can be used as input multiple times) create a dataset where the words are inputs and the corresponding 2-word context (on either side) is the labels. One context word is used as a label at a time.
6. Train a neural network with a linear hidden layer (you may use PyTorch, Jax or Tensorflow/Keras) on this data. Note that the first layer of the network depends on how you represent the input words:
  - If you used a 1-hot representation then you must use a linear layer.
  - If you use the index of the word then you must use an embedding layer.
7. Implement an inference function which receives a 1-hot vector and provides the corresponding embedding.

Hint: I strongly recommend you initialize your networks with small random (gaussian) weights (small being around less than 0.1) and a large learning rate. This will put your networks in the so called “feature learning regime” which is what we want when learning embeddings.

For this lab we will be discussing the broad topic of interpreting the embedding space. Some example sub-topics may be:

1. How does the embedding space change when we use data from different books?
2. How useful is the embedding space for auxiliary tasks like part-of-speech prediction?
3. Can we map from one embedding space to another where the embeddings are trained on different data?

Essentially any meaningful deviation from the base implementation which may lead to some insight into the working of the model and the topic of embeddings. You will need to explore one such sub-topics. Please ensure that you have added your group to the google sheet and you are welcome to propose your own sub-topic on this sheet. Anyone with a blank sub-topic after one week from the lab being released will be allocated one (there's no penalty for not suggesting your own one). Suggested sub-topics are not final and I may still change your sub-topic but I will do my best to not be prescriptive.

You will also be expected to write a one-page (double column) report on the topic you covered (a second page may be used for figures, tables, your contribution statement and references). Please use the IEEE format for this and respect the page limit. Please see the rubric below for more details on what should be included in the write-up. This write-up should focus on your methodology, results and insight for your particular topic. Naturally there will be overlap in what is written by each group, however your aim should be to spend as much time discussing decisions or results for your particular setting.

## 2 Submission

**Due Date:** 12 August 2025 at 10:00.

For the submission you may work in pairs or groups of 3. You will be required to:

1. Submit your full code implementation (in a single file called `word2vec.py`).
2. Submit a one-page report which includes a small statement on the contribution of each person if you did not work alone. This must be completed in a new paragraph and does not contribute to the one-page limit.

Table 1 shows the rubric which will be used to assess your write-up.

## 3 On the Use of ChatGPT

These labs do not count much individually for the course, and yet they are crucial to understand the material and prepare for other assessments. Thus, I urge you to take them seriously and engage with the material. I am aware though that there is an incentive to just complete these labs as quickly as possible which would result in the use of generative tools to speed up the process. This is an NLP course though and the use of ChatGPT-like software is at least a learnable experience for us. Thus, you are welcome to use generative models for the written portion of these assessments. However, greater emphasis will now be placed on the factual correctness of what is said and the degree of insight which is shown in the write-ups. If I detect something resembling a clear “hallucination” (a generative model making up facts) you will receive 0. Negative marking will also be implemented for poor writing or formatting and incorrect information or incoherent reasoning. You will also receive less marks for extremely generic facts or information. The strategy then is to use these tools to get started but then add insight in afterwards - particularly insight directed towards the sub-topic that you are investigating. Equally, clever prompt-engineering will likely go a long way here. Appropriate use of both strategies will be rewarded. Similarly, using generative models to help you code is fine but the usual plagiarism rules for the school remain (it is not tolerated).

Table 1: Rubric for Lab 1

Mode	0% to 20%	20% to 40%	40% to 60%	60% to 80%	80% to 100%
Write-up Structure (10%) [Negative Marking]	Adequate use of language and structure. Text in paragraphs and does not go over page limit.	Fair use of language and structure. Paragraphs, appropriate tone and within page limit.	Well written and clear. It is easy to understand the writing and one section leads naturally to the next.	Good use of language and structure. Good linking between sections, easily understood. Figures where appropriate with helpful captions.	Excellent use of language, very well written and structured. Helpful and clear figures with legible captions, labelling and legends.
Method (30%)	No clear description of the architecture, data or approach to training	Some description of high-level details	Fair description on the details of the base model but little elaboration of the approach to answering the sub-topic	Base approach is described in detail and steps are well justified. Adequate discussion on the approach to answering the sub-topic	Excellent description and motivation for all parts of the method including on how the sub-topic was answered
Results (30%)	Results are not given or irrelevant	Some results given with inappropriate metrics	General results presented with appropriate metrics	General and sub-topic results presented with appropriate metrics	Thorough results which are appropriate for answering the sub-topic and display the general correctness of the model
Discussion Section (30%)	No interpretation of results or incorrect interpretation. Little knowledge of the topic displayed	Some general interpretation of results broadly	Results are interpreted which display some understanding of the mechanics of the model. Displays a basic understanding of the topic	Results are interpreted and contextualized to begin to answer the sub-topic. Clear demonstration of knowledge of the general topic	Results are interpreted and contextualized to answer the sub-topic and display insight into the working of the model or original thought on the concepts