

The University of Melbourne

School of Computing and Information Systems

COMP90042

Natural Language Processing

Final Exam

Semester 1 2023

Exam duration: 165 minutes (150 minutes + 15 minutes upload time)

Length: This paper has 4 pages (including this page). You should attempt all questions.

Instructions to students

- This exam is worth a total of 100 marks.
- You can read the questions and answers during the exam.
- You are recommended to use a pen or pencil for some answers requiring drawing diagrams or writing.
- You will need to scan the questions and answers using Gradescope. Be sure to label the scans/photographs with the question number (e.g. 1, 2, 3, 4). If 4 questions are unlabeled, a penalty of -10% will be applied.
- Please answer all questions.

Format: Open Book

- While you are undertaking this assessment you are permitted to:
 - make use of the textbooks, lecture slides and workshop materials.
- While you are undertaking this assessment you must **not**:
 - make use of any messaging or communications technology;
 - make use of any world-wide web or internet-based resources such as Wikipedia, Stackoverflow, Google, AI services (e.g. ChatGPT) or other web services;
 - act in any manner that could be regarded as providing assistance to another student who is undertaking this assessment, or will in the future be undertaking this assessment.
- The work you submit must be based on your own knowledge and skills, without assistance from any other person.



COMP90042 Natural Language Processing

Semester 1, 2023

Total marks: 120 (40% of subject)

Students must attempt all questions

Section A: Short Answer Questions [33 marks]

Answer each of the questions in this section briefly. Each answer should be no longer than several sentences.

Question 1: General Concepts [24 marks]

- a) Compare and contrast “lemmatisation” and “stemming”. With the aid of an example word, show how lemmatisation and stemming differ. [6 marks]
- b) What is “Relation Extraction”? What are the assumptions made about the set of relations. [6 marks]
- c) Why is “topic model” effective? What are its strengths and shortcomings (one for each). [6 marks]
- d) “Copy mechanism” is introduced in neural machine translation. To handle out-of-vocabulary words in the source document, justify whether the copy mechanism is still effective. [6 marks]

Question 2: Distributed Representations [9 marks]

- a) Compare and contrast “word embeddings” and “word vectors”. [6 marks]
- b) Explain one limitation of word embeddings when learning “word vectors”. [3 marks]



Section B: Method Questions [46 marks]

In this section you are asked to demonstrate your conceptual understanding of the methods that we have studied in this subject.

Question 3: Hidden Markov Models [20 marks]

“Text chunking” refers to the preprocessing step where we divide a text into syntactically related non-overlapping groups of words. The following is an example of a chunked sentence:

[money] [could be added] [to] [a spending bill] [covering] [the federal emergency management agency] , [which] [coordinates] [federal disaster relief] .

- Show how a “Hidden Markov model” can be used to perform text chunking. In your answer, you should: (i) explain the state inventory; (ii) show how the given sentence would be labelled for training the HMM; (iii) describe how the model parameters can be learned (assuming we have a labelled corpus) and illustrate with examples from the given sentence. [12 marks]
- Explain two drawbacks of using HMM for this task, and suggest a solution for each of the drawbacks [8 marks]

Question 4: Discourse Segmentation

Consider the following paragraph and the cosine similarity algorithm described in the lecture to partition it into segments. \cos_i denotes the cosine similarity between the two segments on either side of gap i :

$\text{sim}_0 = 0.9$	Saturn is the largest planet in the Solar System.
$\text{sim}_1 = 0.3$	It is the second largest planet after Jupiter.
$\text{sim}_2 = 0.1$	Saturn is the only planet in the Solar System with word Saturnus, which means “Saturn” in Latin.
$\text{sim}_3 = 0.2$	The Roman name for Saturn is Kronos.
$\text{sim}_4 = 0.7$	Saturn has a ring system like the planet Mercury.
	The large ring system is made of ice and rock.

- Perform “discourse segmentation” on the paragraph, assuming $t = 0.2$. Your answer should show which gap a boundary will be inserted and the computation involved in producing the decision. [3 marks]
- Discuss one drawback of using bag-of-words vectors for computing similarity, and illustrate how it created an erroneous boundary based on your answer in (a). You may assume that the bag-of-words vectors are created using “TF-IDF” weights with the following preprocessing steps: (1) symbols and numbers are removed; and (2) words are lowercased and tokenised using white space. [5 marks]
- Propose 2 solutions that alleviate the drawback. With the aid of examples, show how the solutions will help improve discourse segmentation for the given paragraph. [6 marks]

Question 5: Ethics [12 marks]

Given the first application described in the guest lecture (automatic triaging of legal requests), discuss **three** ethical implications of this application.

Section C: Algorithmic Questions [41 marks]

In this section you are asked to demonstrate your understanding of the methods that we have studied in this subject, in being able to perform algorithmic calculations.

Question 6: Context-free Grammar [23 marks]

Consider a “context-free grammar” with the following production rules:

$S \rightarrow NP \ VP \ NP$	$S \rightarrow NP \ VP \ PP$	
$NP \rightarrow NP \ PP$	$NP \rightarrow N$	
$VP \rightarrow VP \ NP$	$VP \rightarrow V$	
$PP \rightarrow IN \ NP$		
$N \rightarrow \text{he}$	$N \rightarrow \text{elephants}$	$N \rightarrow \text{trucks}$
$V \rightarrow \text{loads}$		
$IN \rightarrow \text{onto}$		

- a) Convert the grammar into “Chomsky Normal Form”. What rules are modified, removed, and added. [4 marks]
- b) Given the sentence: `he loads elephants onto trucks`, parse it using the converted grammar. Your solution should include the parse tree, which includes the edges. Your solution should also include the probability of the parse tree. [6 marks]
- c) One of the parses is more probable than the other. (The unconverted) grammar above has now been modified to include the probability of the parses. Propose a change to the probabilistic grammar to produce a higher probability for the more probable parse. To illustrate how the change affects the probability, use parse trees to show the parser producing a higher probability for the more probable parse. [6 marks]

Question 7: N-gram language model

This question asks you to compute the probability of a sentence under a “unsmoothed bigram language model”. You should leave your answers as fractions.

i wish to wish the
if you won't wish
i won't wish the wish you wish to wish

- a) Compute the probability of all bigrams given the context word `wish` under a “unsmoothed bigram language model”. [4 marks]
- b) Compute the probability of all bigrams given the context word `wish` under a “bigram language model with absolute discounting” where the discount factor $d = 0.2$. [8 marks]
- c) Compute the probability of **unseen bigrams** given the context word `wish` under a “bigram language model with Kneser-ney smoothing” where the discount factor $d = 0.2$. [6 marks]

— End of Exam —