



## **Autumn Examinations 2021/2022**

<b>Course Code(s) Exam(s)</b>	Instance 1CSD1, 1CSD2, 1SPE1, 1MAO2, 1MAI1 MSc in Computer Science (Data Analytics), MSc in Computer Science (Artificial Intelligence), MSc in Computer Science (Artificial Intelligence) - Online
<b>Module Code(s) Module(s)</b>	CT5120, CT5146 Introduction to Natural Language Processing, Introduction to Natural Language Processing - Online
Paper No.	1
External Examiner(s)	Dr John Woodward
Internal Examiner(s)	Dr. Michael Madden *Dr. John McCrae Dr Bharathi Raja Chakravarthi Dr Omnia Zayed

**Instructions:** Answer 4 sections out of 5; each section is worth 25 marks (100 marks total). **Use a separate answer book for each section answered.**

<b>Duration</b>	2 hours
<b>No. of Pages</b>	6
<b>Discipline(s)</b>	Computer Science
<b>Course Co-ordinator(s)</b>	Dr. Frank Glavin Dr. Matthias Nickles Dr. James McDermott

**Requirements:**

Release in Exam Venue	Yes	
MCQ		No
Handout	None	
Statistical/ Log Tables	None	
Cambridge Tables	None	
Graph Paper	None	
Log Graph Paper	None	
Other Materials	None	
Graphic material in colour		No

# Introduction to Natural Language Processing

Exam Duration: 2 Hours

**You must answer 4 of the following sections**

## Section 1: Text Classification

**Question 1A:**

**5 Marks**

What does it mean when a word is **out of vocabulary**? List two solutions to this problem and describe how they solve this issue.

**Question 1B:**

**10 Marks**

State the formula for TF-IDF (Term Frequency-Inverse Document Frequency). How can TF-IDF be used with Naive Bayes to perform text classification?

**Question 1C:**

**10 Marks**

Consider the following sentences with sentiment labels.

- This house room was good [POS]
- The roof was not as good as expected [NEG]
- The garden was good for the kids [POS]
- The kids loved the flowers [POS]

Using Bayes' Law, calculate the probability of the labels POS and NEG given a single feature that considers whether the word 'good' occurs in the text.

**PTO**

## Section 2: Linguistic Concept and Parsing

### Question 2A:

5 Marks

Explain what task a tokenizer performs and how it deals with issues such as punctuation. How many types and tokens in the following sentence?

'the cat sat on the mat'

### Question 2B:

10 Marks

Consider the probabilistic context-free grammar below. Draw **one** parse tree and calculate the probability of that parse for the following sentence: "Connor Murphy and Ciara Byrne dance"

Rule	Probability	Rule	Probability
$S \rightarrow NP V$	0.9	$NN \rightarrow \text{Connor}$	0.2
$S \rightarrow CL CONJ CL$	0.1	$NN \rightarrow \text{Murphy}$	0.2
$CL \rightarrow NP V$	1.0	$NN \rightarrow \text{Ciara}$	0.2
$NP \rightarrow NP CONJ NP$	0.2	$NN \rightarrow \text{Byrne}$	0.2
$NP \rightarrow NN NP$	0.3	$NN \rightarrow \text{dance}$	0.2
$NP \rightarrow NN$	0.5	$V \rightarrow \text{Murphy}$	0.1
$CONJ \rightarrow \text{and}$	1.0	$V \rightarrow \text{dance}$	0.9

### Question 2C:

10 Marks

*it's snowing it's falling  
the old lady is snoring  
she went to roof  
and she bumped her head  
and she couldn't get up in the morning.*

For the above calculate all unigram and bigram probabilities. You should treat "it's" and "couldn't" as single tokens. Treat the whole corpus as a single sentence.

PTO

### Section 3: Vector Space Models

Question 3A:

10 Marks

Give **two** reasons that we may create a vector representation of a word. Explain how a vector representation solves these problems.

Question 3B:

10 Marks

*...its name stands for “**language** model for dialogue applications”...  
...an example of a very large **language** model, or a computer program...  
...I wrote recently, using **language** models in place of search engines ...  
...with the help of a **language** model. It has...  
...human-like interfaces such as **language**. For any automated system...*

Source: The Guardian

Create a context vector for the word ‘language’ from the text above using a context window of two words either side (ignore all punctuation and case).

Question 3C:

5 Marks

Suggest **one** change you could make in how the context vector is constructed that may improve performance.

PTO

## Section 4: Information Extraction

*Meat, cereals, dairy, fruit and vegetables are likely to be the worst affected as the war in Ukraine combines with production lockdowns in China and export bans on key food stuffs such as palm oil from Indonesia and wheat from India, the grocery trade body IGD warns.*

*Products that rely on wheat, such as chicken, pork and bakery items, are likely to face the most rapid price rises as problems with exports and production from Ukraine, a big producer of grain, combine with sanctions on Russia, another key producer.*

*Source: The Guardian*

### Question 3A:

**10 Marks**

From the text above extract at least two hypernym relations by means of Hearst patterns

### Question 3B:

**5 Marks**

Indicate two named entities within the text above and provide the appropriate entity class.

### Question 3C:

**10 Marks**

What is a tagging scheme that would allow a tagging model such as a hidden Markov model in order to tag named entities within a text? Give an example of this tagging on one sentence from the text above.

**PTO**

## Section 5: Semantic Analysis

### Question 5A

10 Marks

**List** and **explain** the three levels of semantic analysis and **give** an example of a task for each level.

### Question 5B

10 Marks

**Define** Word Sense Disambiguation (WSD) and **list** three NLP applications that can benefit from it.

### Question 5C

5 Marks

Consider the following sentence:

*Bob hit Scott with the bottle.*

**Identify** the "semantic roles" for each entity of the event expressed by the given sentence.

**END**