

THIS PAPER IS NOT TO BE REMOVED FROM THE EXAMINATION HALLS

UNIVERSITY OF LONDON

CO3354 ZB

BSc Examination

**COMPUTING AND INFORMATION SYSTEMS, CREATIVE COMPUTING
and COMBINED DEGREE SCHEME**

Introduction to Natural Language Processing

Date and Time: Monday 16 May 2016: 14.30 – 16.45

Duration: 2 hours 15 minutes

There are FIVE questions on this paper. Candidates should answer THREE questions. All questions carry equal marks and full marks can be obtained for complete answers to THREE questions. The marks for each part of a question are indicated at the end of the part in [.] brackets.

Only your first THREE answers, in the order that they appear in your answer book, will be marked.

There are 75 marks available on this paper.

A handheld calculator may be used when answering questions on this paper but it must not be pre-programmed or able to display graphics, text or algebraic equations. The make and type of machine must be stated clearly on the front cover of the answer book.

QUESTION 1

- a) The following sentences are all ambiguous in some way. Express their different meanings using paraphrases and explain the source of the ambiguity in each case:

- i. John punched Bill and he fell over.
- ii. Every man loves a woman.
- iii. Mary watched the cat with one eye.
- iv. Flying planes can be dangerous.

[6]

- b) Explain what is meant by a **corpus** in the context of Natural Language Processing. [2]

What is meant by the following terms in the context of corpus linguistics? (One mark each for i-viii.)

- i. Monitor corpus
- ii. Parallel corpora
- iii. Comparable corpora
- iv. Treebank
- v. Sampling frame
- vi. Supervised learning
- vii. Concordance
- viii. Collocation

[8]

- c) Annotate the text below with POS (part of speech) tags, using the simplified tagset given in Appendix A as in this example:

An/DET Indian/ADJ teenager/N scored/VD more/ADJ than/P
1,000/NUM runs/N in/P a/DET single/ADJ innings/N to/TO set/V a/DET
new/ADJ world/N record/N in/P school/N cricket/N ./.

Where a word has more than one possible POS, explain how you have decided which one to use.

Text (BBC News Website, 5th January 2016):

"I was not thinking of a record," Dhanawade told BBC Hindi. "It was not in my mind at all but as soon as I got close to the feat it was clear to me that I could achieve it."

Dhanawade said his father, an auto-rickshaw driver, had pushed him to play and was partly responsible for his success. He said he was ready for international cricket, but intended to first play in the under-19 state team.

[9]

QUESTION 2

Consider the following grammar:

$S \rightarrow NP VP$

$NP \rightarrow Det N$

$NP \rightarrow NP PP$

$VP \rightarrow V$

$VP \rightarrow V NP$

$VP \rightarrow V NP PP$

$PP \rightarrow P NP$

$Det \rightarrow the \mid a \mid some \mid \epsilon$

$N \rightarrow teacher \mid student \mid classroom \mid books \mid table \mid pen \mid diagram$

$V \rightarrow slept \mid put \mid saw \mid called \mid gave \mid has \mid have \mid wrote \mid drew \mid read$

$P \rightarrow in \mid on \mid at \mid with \mid by \mid from \mid into$

a)

- i. Write out two grammatical sentences of at least 8 words which are generated by this set of rules.
- ii. Draw syntax trees for both of these sentences, according to the grammar rules above. If your example is structurally ambiguous, you should give all applicable trees.

[9]

b) This grammar will generate sequences that are not grammatical sentences. Explain how it can be modified so that it will generate the grammatical examples (i-ii) below but not the ungrammatical (iii-iv).

- i. The student slept in the classroom.
- ii. The teacher put some books on the table.
- iii. *The student slept the classroom.
- iv. *The teacher put on the table.

[8]

c) Explain further how the grammar can be modified to generate examples like i – ii below:

- i. The young student quickly read a little red book.
- ii. The teacher usually sat at a big old wooden desk.

[8]

NB: you are not required to provide complete grammars in answer to (b) and (c), only to explain how the grammars need to be modified, giving appropriate worked examples. Marks will be awarded for compactness and generality of your solutions, not just whether they match the specific examples provided.

QUESTION 3

a) Briefly explain what is meant by each of the following in the context of Natural Language Processing:

- i. N-gram tagging
- ii. Conditional frequency distribution
- iii. Stopwords
- iv. Hidden Markov Model

[8]

b) Using the regular expressions supplied in Appendix B, describe and give examples to illustrate the classes of strings matched by the following regular expressions (for example, $(ac)^*$ matches ϵ , ac , $acac$...):

- i. a^*b^+
- ii. $(a^*b)^+$
- iii. $[a-z]^*[0-9]^*$
- iv. $[a-z0-9]^*$

[9]

c) Suppose a corpus contains 300,000 word-tokens, and 70,000 of these are tagged as N (common noun). The word-form *house* occurs 1,000 times in the corpus, tagged either as N or V. Analysis shows that *house* accounts for 0.3% of all common noun tokens in the corpus. Use Bayes' formula to calculate the probability that a given occurrence of *house* is tagged as N. Explain your answer and show your working.

[8]

QUESTION 4

The following paragraph is taken from a Goldsmiths press release dated 8th December 2015:

In May 2014 Goldsmiths held a nationwide competition to design an art gallery for our south east London campus, a venue that will run exhibitions, projects and residencies by leading artists and curators from the UK and abroad. Assemble's innovative design respectfully incorporates the black steel water tanks of the historic Laurie Grove baths, capitalising on the 'raw and robust' construction of the existing structure while building two additional steel frame structures which allow for a varied range of gallery space.

This is the result of running the above text through the Lancaster Stemmer:

```
['in', 'may', '2014', 'goldsmith', 'held', 'a', 'nationwid',  
'competit', 'to', 'design', 'an', 'art', 'gallery', 'for', 'our',  
'sou', 'east', 'london', 'camp', ' ', 'a', 'venu', 'that', 'wil',  
'run', 'exhibit', ' ', 'project', 'and', 'resid', 'by', 'lead',  
'art', 'and', 'cur', 'from', 'the', 'uk', 'and', 'abroad', ' ',  
'assemble's', 'innov', 'design', 'respect', 'incorp', 'the', 'black',  
'steel', 'wat', 'tank', 'of', 'the', 'hist', 'laury', 'grov', 'bath',  
' ', 'capit', 'on', 'the', '"raw", 'and', 'robust', '"', 'construct',  
'of', 'the', 'ex', 'structure', 'whil', 'build', 'two', 'addit',  
'steel', 'fram', 'structures', 'which', 'allow', 'for', 'a', 'vary',  
'rang', 'of', 'gallery', 'spac', '.']
```

And this is the result of running it through the Snowball Stemmer:

```
['in', 'may', '2014', 'goldsmith', 'held', 'a', 'nationwid',  
'competit', 'to', 'design', 'an', 'art', 'galleri', 'for', 'our',  
'south', 'east', 'london', 'campus', ' ', 'a', 'venu', 'that',  
'will', 'run', 'exhibit', ' ', 'project', 'and', 'resid', 'by',  
'lead', 'artist', 'and', 'curat', 'from', 'the', 'uk', 'and',  
'abroad', ' ', 'assembl', 'innov', 'design', 'respect', 'incorpor',  
'the', 'black', 'steel', 'water', 'tank', 'of', 'the', 'histor',  
'lauri', 'grove', 'bath', ' ', 'capitalis', 'on', 'the', 'raw',  
'and', 'robust', '"', 'construct', 'of', 'the', 'exist', 'structur',  
'while', 'build', 'two', 'addit', 'steel', 'frame', 'structur',  
'which', 'allow', 'for', 'a', 'vari', 'rang', 'of', 'galleri',  
'space', '.']
```

a)

- i. Explain what is meant by **text normalisation** in the context of natural language processing. What benefits derive from normalising text?
- ii. Why is it recommended to use an "off-the-shelf" stemmer rather than coding your own using regular expressions?

[9]

- b)
- i. Make a list of rules which the Lancaster stemmer seems to have applied in this example and discuss the motivations for the rules.
 - ii. Make a similar list for the Snowball stemmer and note any cases where the two stemmers have different results.
- [12]**
- c) Are there any cases where you think the rules were applied incorrectly? If so, give up to four examples. Justify your answer.
- [4]**

QUESTION 5

A probabilistic phrase structure grammar

Phrasal rules

$S \rightarrow NP VP$ [1.0]
 $NP \rightarrow Det N$ [0.8]
 $NP \rightarrow NP PP$ [0.2]
 $VP \rightarrow V NP$ [0.9]
 $VP \rightarrow VP PP$ [0.1]
 $PP \rightarrow P NP$ [1.0]

Lexical rules

$Det \rightarrow a$ [0.5] | the [0.5]
 $N \rightarrow equation$ [0.2] | $formula$ [0.2] | $teacher$ [0.2] | $blackboard$ [0.2] | $students$ [0.2]
 $V \rightarrow copied$ [0.5] | $explained$ [0.5]
 $P \rightarrow on$ [1.0]

- a) The rules shown above make up an example of a probabilistic or weighted grammar. What advantages can such grammars have over conventional phrase-structure grammars?

[6]

- b) Using the probabilistic grammar rules and lexical rules given above, draw as many syntax trees as you can for the sentence:

“The students copied the equation on the blackboard”.

Explain any ambiguities it may have by giving paraphrases for the meanings corresponding to the different syntactic analyses.

[9]

- c)
- Calculate the relative probabilities assigned to different analyses of the sentences by the grammar rules. You may omit the lexical probabilities as these make no difference to the outcome.
 - Discuss whether the results agree with your intuitive understanding of the sentence.

[10]

APPENDIX A: NLTK SIMPLIFIED PART-OF-SPEECH TAGSET

ADJ	adjective	<i>new, good, high, special</i>
ADV	adverb	<i>really, already, still, early, now</i>
CNJ	conjunction	<i>and, or, but, if, while, although</i>
DET	determiner	<i>the, a, some, most, every, no</i>
EX	existential	<i>there, there's</i>
FW	foreign word	<i>dolce, ersatz, esprit, quo</i>
MOD	modal verb	<i>will, can, may, must</i>
N	noun	<i>year, home, costs, time</i>
NP	proper noun	<i>Adam, Paris</i>
NUM	number	<i>twenty-four, fourth, 1991, 14:24</i>
PRO	pronoun	<i>he, their, her, its, my, I, us</i>
P	preposition	<i>on, of, at, with, by, into, under</i>
TO	the word to	<i>to be or not to be</i>
UH	interjection	<i>ah, huh, oops,</i>
V	verb	<i>is, has, get, do, make, see, run</i>
VD	past tense	<i>said, took, told,</i>
VG	present participle	<i>making, going, playing</i>
VN	past participle	<i>given, taken, begun, sung</i>
WH	wh-determiner	<i>who, which, when, what, where, how</i>

APPENDIX B: REGULAR EXPRESSIONS

.	Wildcard, matches any character
^abc	Matches some pattern abc at the start of a string
abc\$	Matches some pattern abc at the end of a string
[abc]	Matches one of a set of characters
[A-Z0-9]	Matches one of a range of characters
ed ing s	Matches one of the specified strings (disjunction)
*	Zero or more of previous item, e.g. a*, [a-z]*
+	One or more of previous item, e.g. a+, [a-z]+
?	Zero or one of the previous item (i.e. optional), e.g. a?, [a-z]?
a(b c)+	Parentheses that indicate the scope of the operators

END OF PAPER