

UNIVERSITY OF EDINBURGH
COLLEGE OF SCIENCE AND ENGINEERING
SCHOOL OF INFORMATICS

FOUNDATIONS OF NATURAL LANGUAGE PROCESSING

Monday 17 August 2009

14:30 to 16:30

Year 3 Courses

Convener: K Kalorkoti

External Examiners: A Frisch, J Gurd

INSTRUCTIONS TO CANDIDATES

Answer QUESTION 1 and TWO other questions.

Question 1 is COMPULSORY.

The short answer parts of Question 1 are each worth 1 mark, 20 marks in total. Each of the remaining three questions is worth 15 marks — answer any TWO of these.

Use one script book for each question, that is, three books in all.

CALCULATORS MAY NOT BE USED IN THIS EXAMINATION

Part A

Answer ALL questions in Part A.

Your answers should be more than a single sentence—take the opportunity to show what you know, as opposed to being as brief as possible—but they need not be lengthy. When two terms are contrasted, make sure your short definitions of each make clear where the contrast lies.

1. What is **Zipf's law**? [1 mark]
2. What's the difference between a **selective** account and an **instructional** account of, for example, speech recognition? [1 mark]
3. What is a **gold standard** in data-intensive linguistics? [1 mark]
4. What's the difference between a **representative** corpus and an **opportunistic** one? [1 mark]
5. What is corpus **metadata**? Give four examples. [1 mark]
6. What are two of the benefits of using a markup language (SGML, XML) to deliver corpora? [1 mark]
7. What is the difference between **well-formed** XML and **valid** XML? [1 mark]
8. What are **joint** and **conditional** probabilities? Give an equation relating the two. [1 mark]
9. What is an **N-gram language model**? [1 mark]
10. What is **backoff** and why is it necessary? [1 mark]
11. What is **smoothing** in the context of language models? Name two examples. [1 mark]
12. Give the formula for **Bayes' rule**. [1 mark]
13. In the noisy channel model, what are the **likelihood** and the **prior**? [1 mark]
14. What is the difference between a **Markov chain** and a **Hidden Markov Model**? [1 mark]
15. What are the two main sources of **ambiguity** in the parsing of natural language text? [1 mark]
16. What is a **well-formed substring table** and what use is it? [1 mark]

17. What are the two main weaknesses of **probabilistic context-free phrase-structure grammar**? [1 mark]
18. What does it mean to **lexicalise a context-free phrase-structure grammar**? [1 mark]
19. How do the **categories** of a **categorial grammar** differ from those of a **context-free phrase-structure grammar**? [1 mark]
20. What are **precision** and **recall** and how are they used? [1 mark]

Part B

ANSWER TWO QUESTIONS IN PART B.

1. Exploring text corpora

Use the two corpus extracts in **Appendix A** to illustrate your points in answering this question.

Describe the operations you would perform in analysing a corpus such as the Brown corpus which is composed of material from qualitatively different sources. Frame your description in terms of exploring ways in which you might look for contrasts between the different samples in the corpus. Try to cover all of the following topics, recalling the exercises we have done in tutorials and assessed coursework:

- letter N-grams
- tokenisation
- word length
- word normalisation
- what constitutes a word
- frequency distribution
- types vs. tokens
- probability distribution
- plots
- sentence length
- word N-grams

Wherever possible, use either pseudocode or Python/NLTK to make your description precise. Minor errors in Python syntax or NLTK class names will *not* be penalised. [15 marks]

2. Dynamic programming and POS-tagging

- (a) Both Jurafsky & Martin and the course notes gave worked examples of using dynamic programming to find the minimum edit distance between two strings. They both ignored **transpositions**

Consider the following (genuine) extracts from the log of a meeting:

This is the memo we are talking about whcih i promised to add a use case to.
... to infromation about things, etc.

came rfom discussions of 303 responses

The directiosn coem from ... it comes from teh [XXX] and [YYY] technology

What would the lowest cost of correcting 'teh' to 'the' be, given the standard costs of 1 for both deletion and insertion and 2 for substitution?

[2 marks]

- (b) Using those costs again, draw a dynamic programming lattice to find the lowest-cost transformation of **rfom** into **from**. Be sure to include backpointers.

Show your work. That is, make clear *how* you arrive at the cost you enter in each cell of the lattice.

[5 marks]

- (c) Using the language and channel models from **Appendix B**, diagram the Viterbi search for the most likely translation into English for the French phrase "les bonnes dorment". All numbers are *costs*, that is, negative log probabilities, so you can a) sum them, rather than multiplying and b) you are looking to *minimise* the total cost.

Don't forget both the initial (from sentence start) and final (to sentence end) transitions.

Show your work. That is, make clear *how* you arrive at the cost you enter in each cell of the lattice.

[8 marks]

3. Best-first chart parsing

Using the probabilistic CF-PSG from **Appendix C** simulate a best-first top-down chart parser parsing the sentence “time flies”.

Use the technique from Jurafsky & Martin, also used in lectures, of drawing actual edges in your chart for inactive edges and partially complete active edges, but just writing the dotted rule for empty active edges below the relevant vertex.

Show both the chart and the agenda, crossing items off the agenda as you draw them into the chart, and *numbering* both the agenda entry and the chart entry as you do so to show the order in which things were done.

Start with the following two edges already in the chart:

$_0\text{time}_1$
 $_1\text{flies}_2$

with both edges having 0 cost.

Start with the following entry in the agenda:

$_0\text{Top} \longrightarrow \cdot S_0$

with 0 cost and the best possible figure of merit.

Be sure to show the cost of every edge in both chart and agenda, and the figure of merit for edges in the agenda.

Mark each edge in the agenda other than the first (top-down initiator) one as either TD (top-down) or F (fundamental), to show which rule they were 'built' by.

Explain how you are calculating costs and the figure of merit you are using, and why. [15 marks]

Appendix A: Sample corpus data

These extracts from the Brown corpus are for use in answering Question 1 in Part B. The paragraphs below correspond to single lines in the originals.

Extract 1

Hemphill said that the Hughes Steel Erection Co. contracted to do the work at an impossibly low cost with a bid that was far less than the " legitimate " bids of competing contractors .

The Hughes concern then took " shortcuts " on the project but got paid anyway , Hemphill said .

The Controller's charge of rigging was the latest development in an investigation which also brought these disclosures Tuesday :

The city has sued for the full amount of the \$172,400 performance bond covering the contract .

The Philadelphia Transportation Co. is investigating the part its organization played in reviewing the project .

Extract 2

Such a little thing to start with -- the car registration .

" Ida , where is the car license " ? She asked .

" I can't find it in the glove compartment " .

" Vera must have it " , I answered readily enough , recalling her last visit .

" Vera " , she was frowning .

" Why should Vera have it " ?

Had she forgotten she had signed the car away , that whatever they mutually owned had been divided among the children ?

I was silent .

I didn't want to stir things up .

Appendix B: Language and channel models

These models are for use in answering Question 2 in Part B.

		the	them	good	maids	sleep	\$
Language model	.	3.1	15.8	10.3	20.2	14.9	∞
	the	13.7	19.0	10.8	16.6	16.2	13.2
	them	5.8	15.4	11.4	∞	13.9	2.3
	good	7.7	∞	10.3	∞	12.1	4.0
	maids	6.7	∞	∞	∞	4.5	3.2
	sleep	6.8	12.6	11.6	∞	10.3	2.0

		the	them	good	maids	sleep
Channel model	les	2.3	1.1	∞	∞	∞
	le	2.1	∞	∞	∞	∞
	la	1.7	∞	∞	∞	∞
	eux	∞	1.4	∞	∞	∞
	elles	∞	2.6	∞	∞	∞
	bonnes	∞	∞	4.0	5.8	∞
	bonne	∞	∞	1.9	∞	∞
	honnête	∞	∞	6.6	∞	∞
	servantes	∞	∞	∞	3.3	∞
	dormir	∞	∞	∞	∞	1.8
	dormez	∞	∞	∞	∞	6.9
	dorment	∞	∞	∞	∞	3.0

*Note: read these tables as follows: The language model has start states down the left and destination states across the top, so for example the cost of a transition from **the** to **good** is 10.8; the channel model has English across the top and French down the left, so for example the cost of seeing the word 'bonnes' in state **maids** is 5.8. In the language model, full stop ('.') is used for sentence start and dollar-sign ('\$') for sentence end.*

The language model costs were taken from the BNC: the unexpectedly low cost for the bigram "the the" is due almost entirely to errors in the texts. Whether those errors are in the originals, or were introduced by the transcription process, is not clear.

Appendix C: Probabilistic CF-PSG

These rules are for use in answering Question 3 in Part B.

Note: the numbers given are costs

Rule	Cost
$S \rightarrow NP\ VP$	2.3
$S \rightarrow VP$.3
$VP \rightarrow V0$	1.6
$VP \rightarrow V1\ NP$.6
$V1 \rightarrow time$	15
$V0 \rightarrow flies$	11
$NP \rightarrow time$	16.5
$NP \rightarrow flies$	12.5