

Försättsblad till skriftlig tentamen vid Linköpings universitet



Datum för tentamen	2019-06-12
Sal (1)	<u>TER2(3)</u>
Tid	14-18
Utb. kod	TDP030
Modul	TEN1
Utb. kodnamn/benämning Modulnamn/benämning	Språkteknologi Skriftlig tentamen
Institution	IDA
Antal uppgifter som ingår i tentamen	8
Jour/Kursansvarig Ange vem som besöker salen	Marco Kuhlmann
Telefon under skrivtiden	013-284644
Besöker salen ca klockan	endast telefonjour
Kursadministratör/kontaktperson (namn + tfnr + mailaddress)	Veronica Kindeland Gunnarsson, 013-285634, veronica.kindeland.gunnarsson@liu.se
Tillåtna hjälpmedel	inga
Övrigt	
Antal exemplar i påsen	

Exam 2019-06-12

Marco Kuhlmann

This exam consists of two parts:

Part A consists of 5 items, each worth 3 points. These items test your understanding of the basic algorithms that are covered in the course. They require only compact answers, such as a short text, calculation, or diagram.

Part B consists of 3 items, each worth 6 points. These items test your understanding of the more advanced algorithms that are covered in the course. They require detailed and coherent answers with correct terminology.

Note that surplus points in one part do not raise your score in another part.

Grade requirements TDP030:

- Grade 3: at least 12 points in Part A
- Grade 4: at least 12 points in Part A and at least 7 points in Part B
- Grade 5: at least 12 points in Part A and at least 14 points in Part B

Good luck!

The Corpus of Contemporary American English (COCA) is the largest freely-available corpus of English, containing approximately 560 million tokens. In this corpus we have the following counts of unigrams and bigrams:

<i>snow</i>	<i>white</i>	<i>white snow</i>	<i>purple</i>	<i>purple snow</i>
38,186	256,091	122	11,218	0

- a) Estimate the following probabilities using maximum likelihood estimation without smoothing. Answer with fractions.
- $P(\textit{white})$
 - $P(\textit{snow} \mid \textit{white})$
- b) Estimate the following probabilities using maximum likelihood estimation with additive smoothing, $k = 0.01$. Assume that the vocabulary consists of 1,254,100 unique words. Answer with fractions.
- $P(\textit{purple})$
 - $P(\textit{snow} \mid \textit{purple})$
- c) We use maximum likelihood estimation with add- k smoothing to train n -gram models on the COCA corpus, with $n \in \{1, \dots, 5\}$ and $k \in \{0, 0.1, 1\}$. The following table shows the entropy of each trained model on the training data. Which row corresponds to which k -value, and why? Answer with a short text.

	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$
row 1	7.3376	5.9834	6.7332	6.9556	7.0555
row 2	7.3376	3.4269	1.4290	0.5436	0.4171
row 3	7.3376	7.3837	8.4573	8.6577	8.7350

03 Part-of-speech tagging

(3 points)

- a) The evaluation of a part-of-speech tagger produced the following confusion matrix. The marked cell gives the number of times the system classified a document as class NN whereas the gold-standard class for the document was JJ.

	NN	JJ	VB
NN	58	6	1
JJ	5	11	2
VB	0	7	43

Set up fractions for the following values:

- i. accuracy
ii. recall with respect to class JJ
- b) The following matrices specify (parts of) a hidden Markov model. The marked cell specifies the probability for the transition from BOS to AB.

	AB	PN	PP	VB	EOS
BOS	1/11	1/10	1/12	1/11	1/25
AB	1/11	1/11	1/11	1/10	1/14
PN	1/11	1/12	1/12	1/10	1/16
PP	1/13	1/11	1/12	1/14	1/18
VB	1/11	1/10	1/10	1/13	1/15

	she	got	up
AB	1/25	1/25	1/14
PN	1/13	1/25	1/25
PP	1/25	1/25	1/13
VB	1/25	1/14	1/19

Provide a fraction for the probability of the following tagged sentence:

she got up
PN VB PP

Part c) on the next page

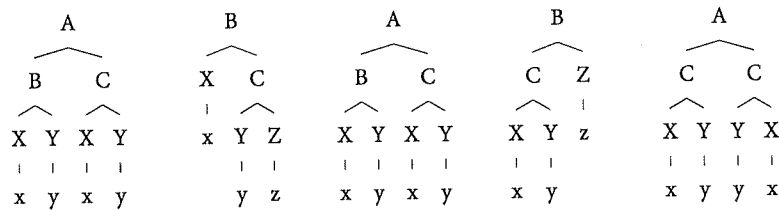
- c) The following table shows a comparison between hidden Markov models and multiclass perceptrons. Complete the three missing cells in this comparison.

Hidden Markov model	Multiclass Perceptron
A	uses weights (arbitrary real numbers)
exhaustive search for the best sequence	B
features: current word, previous word's tag	C

04 Syntactic analysis

(3 points)

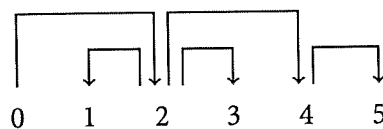
- a) Below is a small phrase structure treebank. Read off all rules whose left-hand sides are either B or C and estimate their rule probabilities using maximum likelihood estimation (no smoothing).



- b) You sum up all rule probabilities in a certain probabilistic context-free grammar. Which (zero or more) of the following values can you *not* get as a result, and why? Answer with a short text.

- i. 0.42 ii. 1 iii. 4.2 iv. 42

- c) State two different sequences of transitions that make the transition-based dependency parser produce the following dependency tree:



a) Draw a partial WordNet-hierarchy for the following synsets:

- | | |
|------------------------------|-----------------------|
| 1 university | 4 hospital |
| 2 institution, establishment | 5 kindergarten |
| 3 educational institution | 6 medical institution |

b) Here are three signatures (glosses and examples) from Wiktionary for different senses of the word *course*:

- 1 A normal or customary sequence. 2 A learning program, as in university. *I need to take a French course.* 3 The direction of movement of a vessel at any given moment. *The ship changed its course 15 degrees towards south.*

Based on these signatures, which of the three senses of the word *course* does the Lesk algorithm predict in the following senses? Ignore the word *course*, stop words, and punctuation.

In the United States, the normal length of a course is one academic term.

c) We read off word vectors for words in Klingon from the following co-occurrence matrix. (Target words correspond to rows, context words correspond to columns.)

	<i>HuSHa'</i>	<i>Ha'DIbaH</i>
<i>qa'vIn</i>	5	1
<i>qurgh</i>	5	5
<i>jonta'</i>	1	0
<i>Dargh</i>	1	4

Sort the four target words by semantic similarity to the word *jonta'*. Start with the most similar target word. Assume that semantic similarity can be quantified in terms of cosine similarity.

Part B

06

Minimum edit distance

(6 points)

- a) Define the concept of the Levenshtein distance between two words. The definition should be understandable even to readers who have not taken this course.
- b) Compute the Levenshtein distance between the two words *game* and *lake* using the Wagner–Fischer algorithm. Your answer should contain both the distance itself and the complete matrix.
- c) Jurafsky and Martin (2018) explain how to augment the Wagner–Fischer algorithm to store backpointers in each cell. Add these backpointers to your matrix. Explain why cells may have several backpointers.

07 Viterbi algorithm

(6 points)

Here is a Hidden Markov model (HMM) specified in terms of costs (negative log probabilities). The marked cell gives the transition cost from BOS to PL.

	PL	PN	PP	VB	EOS
BOS	11	2	3	4	19
PL	17	3	2	5	7
PN	5	4	3	1	8
PP	12	4	6	7	9
VB	3	2	3	3	7

	they	freak	out
PL	17	17	4
PN	3	19	19
PP	19	19	3
VB	19	8	19

When using the Viterbi algorithm to calculate the least expensive (most probable) tag sequence for the sentence 'they freak out' according to this model, one gets the following matrix. Note that the matrix is missing some values.

		they	freak	out
BOS	0			
PL		28	27	A
PN		5	28	35
PP		22	27	B
VB		23	14	36
EOS				C

- Calculate the missing values (marked cells).
- Let m and n denote the number of tags in the HMM and the number of words in the input sentence, respectively. The memory required by the Viterbi algorithm is in $O(mn)$, and the runtime required is in $O(m^2n)$. Explain what these statements mean and how to derive them.
- When one is only interested in the *cost* of the least expensive tag sequence, not in the sequence itself, then the memory required by the Viterbi algorithm is in $O(m)$. Explain this statement. Why does this statement not hold if one wants to reconstruct the actual tag sequence?

Named entity tagging is the task of identifying entities such as persons, organisations, and locations in running text. One way to approach this task is to use the same techniques as in part-of-speech tagging. However, a complicating factor is that named entities can span more than one word. Consider the following sentence:

Alfred Nobel was an inventor from Sweden.

In this example, while the unigram 'Sweden' corresponds to one named entity of type 'location' (LOC), we would also like to identify the bigram 'Alfred Nobel' as a mention of *one* named entity, of type 'person' (PER).

- a) To account for the fact that named entities can span more than one word, we can use the so-called IOB tagging. Explain how this scheme works and show how the example sentence given above would be tagged using IOB tagging.
- b) Named entity taggers often use *gazetteers*. Explain what a gazetteer is and how it can be integrated into a named entity tagger based on the multi-class perceptron.
- c) Jurafsky and Martin (2018) note that evaluation of named entity tagging can be problematic: 'For example, a system that labeled *American* but not *American Airlines* as an organization would cause two errors, a false positive for O and a false negative for I-ORG.' Provide a complete sentence together with gold-standard annotations and system output that illustrates this problem. What accuracy does the system have at the word level? What accuracy does it have at the entity level?