

This list of exercises simulates an exam for the course *Natuurlijke Taalmodellen en Interfaces*. Answer the questions in the spaces provided. If you run out of space for an answer, continue on the back of the page.

Mobile phones, tablets, computers, e-readers, and other electronic equipments are not allowed. They must be switched off and stored away. Basic calculators (not scientific ones) are allowed, but not required, neither necessary.

Contents

1	Random variables and rules of probabilities	2
2	Categorical distributions	3
3	Markov models	4
4	Hidden Markov models	9
5	Probabilistic context-free grammars	14
6	Deductive systems	18
7	Misc	21

Points

Question:	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	Total
Points:	1	1	3	1	2	1	4	3	3	2	1	2	2	3	1	1	2	4	6	2	45

1 Random variables and rules of probabilities

1. (1 point) Let X be a random variable whose sample space is the English vocabulary Σ and whose mapping to \mathbb{R} is realised by an arbitrary *enumeration* of Σ . Given the partial definitions of X below, mark those that are definitely **invalid**?

$$\bigcirc \quad X(\omega) = \begin{cases} 1 & \text{if } \omega = \{\text{the}\} \\ 2 & \text{if } \omega = \{\text{the}\} \\ 3 & \text{if } \omega = \{\text{cat}\} \\ 4 & \text{if } \omega = \{\text{dog}\} \\ \dots & \end{cases}$$

$$\bigcirc \quad X(\omega) = \begin{cases} 1 & \text{if } \omega = \{\text{the}\} \\ 1 & \text{if } \omega = \{\text{a}\} \\ 2 & \text{if } \omega = \{\text{cat}\} \\ 3 & \text{if } \omega = \{\text{dog}\} \\ \dots & \end{cases}$$

$$\bigcirc \quad X(\omega) = \begin{cases} 1 & \text{if } \omega = \{\text{the}, \text{a}\} \\ 2 & \text{if } \omega = \{\text{cat}\} \\ 3 & \text{if } \omega = \{\text{dog}\} \\ \dots & \end{cases}$$

$$\bigcirc \quad X(\omega) = \begin{cases} 1 & \text{if } \omega = \{\text{the}\} \\ 2 & \text{if } \omega = \{\text{a}\} \\ 3 & \text{if } \omega = \{\text{cat}\} \\ 4 & \text{if } \omega = \{\text{dog}\} \\ \dots & \end{cases}$$

2. (1 point) Number the identities on the right according to the concepts on the left.

1. Chain rule

$$(\text{_____}) \quad P_{A|B}(a|b) = \frac{P_{AB}(a,b)}{P_B(b)}$$

2. Conditional probability

$$(\text{_____}) \quad P_A(a) = \sum_{b \in \mathcal{B}} P_{AB}(a,b)$$

3. Bayes rule

$$(\text{_____}) \quad P_{AB}(a,b) = P_B(b)P_{A|B}(a|b)$$

4. Marginal probability

$$(\text{_____}) \quad P_{B|A}(b|a) = \frac{P_B(b)P_{A|B}(a|b)}{P_A(a)}$$

2 Categorical distributions

3. Let X be a Categorical random variable:

$$X \sim \text{Cat}(\theta_1, \dots, \theta_v)$$

(a) ($\frac{1}{2}$ point) What is the support \mathcal{X} of the random variable?

(b) ($\frac{1}{2}$ point) What is the value of $P_X(x)$?

(c) (1 point) What conditions apply to valid parameters $\langle \theta_1, \dots, \theta_v \rangle$?

(d) (1 point) Given a data set of n i.i.d. observations, what is the maximum likelihood estimate of θ_x ?

Total for Question 3: 3

4. (1 point) Select, out of the list below, vector(s) that constitute(s) **valid** categorical parameters for a categorical random variable that may take on one out of 7 classes.

- ☐ $\langle 0.1, 0.1, 0.1, 0.1, 0.1, 0.2, 0.3 \rangle$
- ☐ $\langle 0.1, 0.1, 0.1, 0.1, 0.1, 0.5 \rangle$
- ☐ $\langle 0.2, 0.2, 0.1, 0.1, 0.1, 0.2, 0.2 \rangle$

3 Markov models

5. Consider the probability of a sentence as given by the following factorisation

$$\begin{aligned} P_S(x_1^n) &= P_N(n)P_{S|N}(x_1^n|n) \\ &= P_N(n) \prod_{i=1}^n P_{X|H}(x_i|x_{<i}) \end{aligned}$$

where S is a random sentence, N a random length, X a random word, and H a random history.

- (a) ($\frac{1}{2}$ point) Select appropriate descriptions for x_1^n

- ☐ an outcome of S
- ☐ a sequence of n random words
- ☐ n outcomes of S

- (b) ($\frac{1}{2}$ point) Select appropriate descriptions for n

- ☐ a random length
- ☐ a random noun
- ☐ the length of the outcome of S

- (c) ($\frac{1}{2}$ point) Select appropriate descriptions for x_i

- ☐ a random word
- ☐ the i th element of the outcome of S
- ☐ the i th random sequence

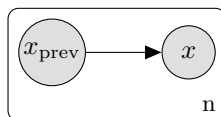
- (d) ($\frac{1}{2}$ point) Select appropriate descriptions for $x_{<i}$

- ☐ a word if $i = 2$
- ☐ a random sequence
- ☐ the i th random history

Total for Question 5: 2

6. (1 point) Let x_1^n be the outcome of a random sentence S , and let $P_{S|N}(x_1^n|n)$ denote its probability value (given length n) under a **unigram** language model. Write down the expression that corresponds to this probability value.

7. Answer questions about the graphical model below, where X is a random variable over exactly v English words.



- (a) ($\frac{1}{2}$ point) Which language model (LM) is this?
 A. unigram LM B. bigram LM C. hidden Markov LM
- (b) ($\frac{1}{2}$ point) How many conditional probability distributions (cpds) are there in the model (ignore the *length* distribution)?
 A. one B. two C. n D. v
- (c) ($\frac{1}{2}$ point) Is $P_{X|X_{\text{prev}}=x_{\text{prev}}}$ a tabular cpd or an inferred distribution?
 A. tabular B. inferred
- (d) ($\frac{1}{2}$ point) Is $P_{S|N=n}$ a tabular cpd or an inferred distribution?
 A. tabular B. inferred
- (e) (1 point) Write down the expression of the probability value $P_S(x_1^n)$ (you may assume appropriate padding exists).

- (f) ($\frac{1}{2}$ point) Assume that the probability value $P_{X|X_{\text{prev}}}(x|x_{\text{prev}})$ can be assessed in constant time. Express the complexity of computing $P_{S|n}(x_1^n|n)$ as a function of sentence length (use *big-O-notation*).

- (g) ($\frac{1}{2}$ point) Suppose we have exactly v words in the vocabulary, and we use a Categorical distribution for each cpd in the model. What is the representation cost of this model (use *big-O-notation*)?

Total for Question 7: 4

8. Consider the following unigram language model, where EOS is a special symbol deterministically added to the end of every sentence, and answer the questions below. In this exercise you are

X	$\text{Cat}(x \boldsymbol{\theta})$
a	θ_a
b	θ_b
c	θ_c
d	θ_d
EOS	θ_{EOS}

expected to pad sentences with a BOS token, which **is not** modelled, and an EOS token, which **is** modelled.

- (a) ($\frac{1}{2}$ point) What is the probability of the sentence a b c a d given its length?

- (b) ($\frac{1}{2}$ point) What is the probability of the sentence a b b d c a a f?

- (c) (1 point) What is the role of smoothing?

- (d) (1 point) Answer true (T) or false (F).

- ___ The sentence a a b c has the same probability as the sentence a b a c.
- ___ The unigram language model is sensitive to word order.
- ___ A smoothed unigram language model has infinite support.
- ___ Without smoothing, and without taking padding into account, the support of the unigram language model above is the set of strings in $\{a, b, c, d\}^*$.

Total for Question 8: 3

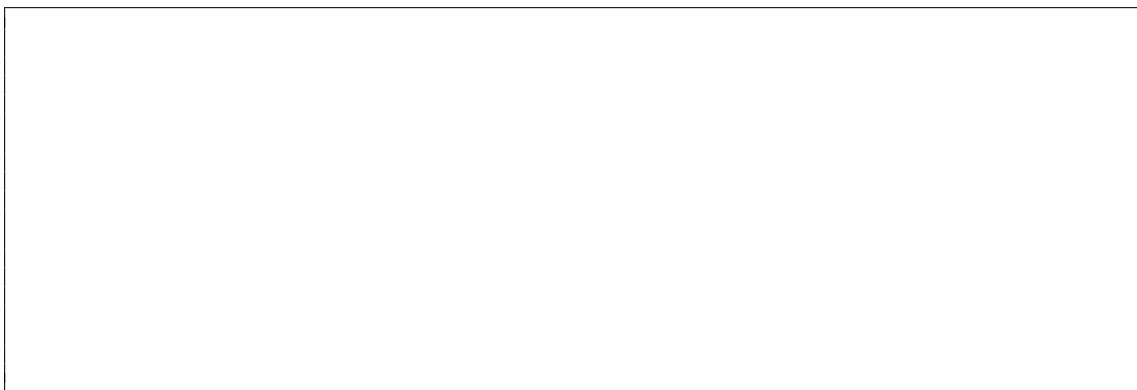
9. Consider the generative story below

$$N \sim P_N$$

$$X_i | X_{i-1} = x_{i-1} \sim \text{Cat}(\theta_1^{(x_{i-1})}, \dots, \theta_v^{(x_{i-1})}) \quad \text{for } i = 1, \dots, n$$

and assume that every distribution in the model gives support to $\{a, b, c, d, \text{UNK}\}$, where UNK is a special token to which we map all unseen words, in addition to a EOS padding symbol that occurs always at the end of strings.

(a) (1 point) Draw the graphical model using plate notation.



(b) Pad the sentence a b c a b and answer the questions below.

i. (1 point) List its bigrams and their counts.

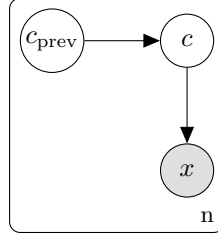


ii. (1 point) What is the probability of the sentence given its length? Express probability values in terms of the parameters shown in the generative story.

Total for Question 9: 3

4 Hidden Markov models

The hidden Markov model (HMM) extends the Markov model with word categories. The graphical model below specifies the conditional independence assumptions of the HMM, where



X is a random word from a vocabulary of v words and C is a random word category (or tag) from a vocabulary of t tags. There are two types of cpds in the HMM. *Transition distributions* used to generate a tag given the tag of the previous word:

$$C|C_{\text{prev}} = c_{\text{prev}} \sim \text{Cat}(\lambda_1^{(c_{\text{prev}})}, \dots, \lambda_t^{(c_{\text{prev}})})$$

And *emission distributions* used to generate a word given its tag:

$$X|C = c \sim \text{Cat}(\theta_1^{(c)}, \dots, \theta_v^{(c)})$$

The joint probability for a sentence x_1^n and tag-sequence c_1^n given length $N = n$ factorises

$$\begin{aligned} P_{X_1^n C_1^n | N}(x_1^n, c_1^n | n) &= P_{X_1^n C_1^n | N}(x_1^n, c_1^n | n) \\ &= \prod_{i=1}^n P_{C|C_{\text{prev}}}(c_i | c_{i-1}) P_{X|C}(x_i | c_i) \end{aligned}$$

in terms of transition and emission probabilities.

Assessing the probability of a sentence, regardless of tag sequence, requires marginalisation

$$\begin{aligned} P_{X_1^n | N}(x_1^n | n) &= \sum_{c_1=1}^t \cdots \sum_{c_n=1}^t P_{X_1^n C_1^n | N}(x_1^n, c_1^n | n) \\ &= \prod_{i=1}^n \sum_{c_{i-1}=1}^t \sum_{c_i=1}^t P_{C|C_{\text{prev}}}(c_i | c_{i-1}) P_{X|C}(x_i | c_i) \end{aligned}$$

10. (2 points) Express the joint probability for a sentence x_1^n and tag-sequence c_1^n given their length as a function of transition and emission parameters.

11. (1 point) The HMM has _____ transition distributions, each of which has _____ parameters, it also contains _____ emission distributions, each of which has _____ parameters, and therefore, the total representation cost of the HMM (in big-O-notation) is _____.

12. (2 points) Write the generative story of the HMM (you may assume appropriate padding exists).

13. (2 points) Consider the following transition and emission distributions.

	$X = 1$	$X = 2$	$X = 3$	\dots	$X = v$		$i = 1$	$i = 2$	$i = 3$
$C = 1$	$\theta_1^{(1)}$	$\theta_2^{(1)}$	$\theta_3^{(1)}$	\dots	$\theta_v^{(1)}$	$C = 1$	$\lambda_1^{(0)}\theta_{x_1}^{(1)}$?	
$C = 2$	$\theta_1^{(2)}$	$\theta_2^{(2)}$	$\theta_3^{(2)}$	\dots	$\theta_v^{(2)}$	$C = 2$	$\lambda_2^{(0)}\theta_{x_1}^{(2)}$		
$C = 3$	$\theta_1^{(3)}$	$\theta_2^{(3)}$	$\theta_3^{(3)}$	\dots	$\theta_v^{(3)}$	$C = 3$	$\lambda_3^{(0)}\theta_{x_1}^{(3)}$		

Transition distributions (left) and emission distributions (right)

We can use an HMM model defined with these cpds to find the best possible way to tag an input sentence $\langle x_1, x_2, x_3 \rangle$. The table below shows 3 cells used to compute the Viterbi recursion $\alpha(i, j)$. What is the value of the Viterbi entry $\alpha(i = 2, j = 1)$?

	$i = 1$	$i = 2$	$i = 3$
$C = 1$	$\lambda_1^{(0)}\theta_{x_1}^{(1)}$?	
$C = 2$	$\lambda_2^{(0)}\theta_{x_1}^{(2)}$		
$C = 3$	$\lambda_3^{(0)}\theta_{x_1}^{(3)}$		

Viterbi table $\alpha(i, j)$: assume $j = 0$ to correspond to the BOS tag.

14. Consider the tagged sequences below where the first sequence occurs n_1 times, the second sequence occurs n_2 times, and the third sequence occurs n_3 times.

BoS — A — EoS
 | |
 x EoS

BoS — A — B — EoS
 | | |
 x y EoS

BoS — A — B — C — A — EoS
 | | | |
 x y t z EoS

- (a) (1 point) Estimate by maximum likelihood the transition distribution given that the previous category is 'A'.

- (b) (1 point) Estimate by maximum likelihood the emission distribution given that the category is 'A'.

(c) (1 point) What is the probability of the second sequence pair, given its length, as a function of maximum likelihood estimates?

[illegible]

Total for Question 14: 3

5 Probabilistic context-free grammars

Let $\mathfrak{G} = \langle \Sigma, \mathcal{V}, S, \mathcal{R} \rangle$ be a context-free grammar (CFG) where

- Σ is the set of terminals
- \mathcal{V} is the set of nonterminals
- $S \in \mathcal{V}$ is the start symbol
- \mathcal{R} is a set of context-free rules

and also assume that the most complex rule has a sequence of a symbols on its right-hand side (RHS).

15. (1 point) What is the general form of a context-free rule in \mathcal{R} ? Make sure to formally specify the set to which left-hand side (LHS) and RHS belong.

16. (1 point) If we know that \mathfrak{G} is in Chomsky normal form (CNF), what can we say about rules in \mathcal{R} ?

17. A probabilistic CFG (PCFG) extends a CFG with a probability distribution over derivations.

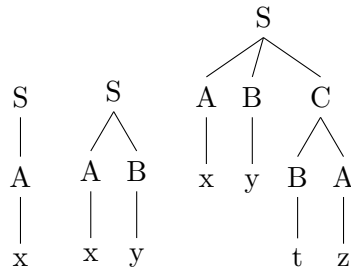
(a) ($\frac{1}{2}$ point) Define a random rule.

(b) ($\frac{1}{2}$ point) Define a random derivation.

- (c) (1 point) Write down the probability distribution of a derivation r_1^m given its length as a function of the factor $P_{\text{RHS}|\text{LHS}}$.

Total for Question 17: 2

18. Consider the treebank below where the first tree occurs n_1 times, the second tree occurs n_2 times, and the third tree occurs n_3 times.



- (a) (1 point) Use this treebank to derive the minimal set of context-free rules that could reconstruct it.

- (b) (1 point) Consider we extend this grammar to a PCFG, write down the maximum likelihood estimates for all pre-terminal rules.

- (c) (1 point) Write down a derivation (as an ordered sequence of rule applications) for the second tree.

- (d) (1 point) What is the probability of the second tree under a PCFG estimated via maximum likelihood using the given treebank.

Total for Question 18: 4

6 Deductive systems

19. In the HMM model we often need to represent the space of all possible analyses of a sentence x_1^n , for example, we need that space of options in order to characterise the marginal probability $P_{X_1^n}(x_1^n)$ as well as in order to find the best tag sequence. Below, we have a deductive system that compactly represents the weighted set of all possible analyses.

INPUT	tagset $\{1, \dots, t\}$ and sentence x_1^n
ITEM	$[c, i]$ where $c \in \{1, \dots, t\} \cup \{\text{BoS}, \text{EOS}\}$ and $i \in \{0, n+1\}$
GOAL	$[\text{EOS}, n+1]$
AXIOMS	$[\text{BoS}, 0]$
TAG	$\frac{[c, i]}{[c', i+1]} \quad i < n \text{ and } c' \in \{1, \dots, t\}$
CONCLUDE	$\frac{[c, n]}{[\text{EOS}, n+1]}$

In this program an item $[c, i]$ refers to word x_i being tagged with tag c . We augment the tag set $\{1, \dots, t\}$ with two special symbols $\{\text{BoS}, \text{EOS}\}$ that helps us track the beginning and the end of the tag sequence.

- (a) (1 point) How many items can we prove for an input x_1^n (use big-O-notation)?

- (b) (1 point) How many inferences are valid for an input x_1^n (use big-O-notation)?

- (c) (2 points) For $t = 2$ draw a graph where each state corresponds to an instantiation of a valid item and where each inference corresponds to an edge.



- (d) (2 points) Each path of execution of the deductive system corresponds to one complete analysis, we can call it a *derivation* since it stands for a way to *derive* or *prove* the goal item. We can easily extend the system to assign a weight to each inference. Let us assume a parameterisation of our HMM generative story

$$C_i | C_{i-1} = c_{i-1} \sim \text{Cat}(\lambda_1^{(c_{i-1})}, \dots, \lambda_t^{(c_{i-1})})$$

$$X_i | C_i = c_i \sim \text{Cat}(\theta_1^{(c_i)}, \dots, \theta_v^{(c_i)})$$

in terms of transition and emission distributions (ignoring length). An AXIOM is a trivial inference, thus we give it a dummy weight (not to interfere with the total)

$$\text{AXIOMS} \quad [\text{BoS}, 0] : \bar{1}$$

CONCLUDE is an inference which serves the purpose of ending the tag sequence and the word sequence with special EOS tokens positioned at $i = n + 1$. This requires a transition to EOS and an emission of EOS, thus the weight is

$$\text{CONCLUDE} \quad \frac{[c, n]}{[\text{EOS}, n + 1] : \lambda_{\text{EOS}}^{(c)} \times \theta_{\text{EOS}}^{(\text{EOS})}}$$

What is the weight of the TAG rule?

Total for Question 19: 6

7 Misc

20. (2 points) Compare n -gram LMs, hidden Markov models, and probabilistic CFGs in terms of what they can or cannot capture, cost of representation, and algorithmic complexity of assessing important quantities such as marginals and argmax.

Assessment

Question	Points	Score
1	1	
2	1	
3	3	
4	1	
5	2	
6	1	
7	4	
8	3	
9	3	
10	2	
11	1	
12	2	
13	2	
14	3	
15	1	
16	1	
17	2	
18	4	
19	6	
20	2	
Total:	45	