

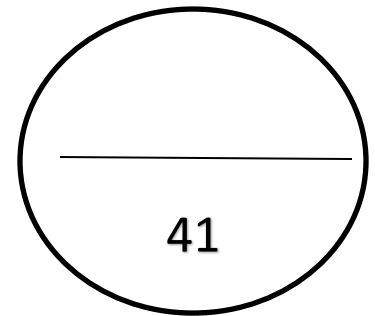
National University of Computer and Emerging Sciences, Lahore Campus



Course:	Natural Language Processing	Course Code:	CS 535
Program:	BS(Computer Science)	Semester:	Spring 2018
Duration:	180 Minutes	Total Marks:	41
Paper Date:	23-May-18	Weight	50%
Section:	ALL	Page(s):	8
Exam:	Final Solution		

Instruction/Notes: Attempt the examination on the question paper and write concise answers. You can use extra sheet for rough work. Do not attach extra sheets used for rough with the question paper. Don't fill the table titled Questions/Marks.

Question	1-4	5-7	8-10	11-14	Total
Marks	/ 9	/ 10	/ 12	/10	/ 41



Q1) You are given the following corpus: [2 + 2 = 4 Marks]

<s> She likes green apples </s>
 <s> Ali likes green apples </s>
 <s> green apples are good for health </s>
 <s> I like red apples </s>

- a) Calculate the probability of following test sentence using trigram language model with linear interpolation. Include <s> and </s> in your counts just like any other token. λ_1 = trigram weight, λ_2 = bigram weight, λ_3 = unigram weight, $\lambda_1 = 0.5$, $\lambda_2 = 0.3$, $\lambda_3 = 0.2$

<s> He likes green apples </s>

Solution:

$$\begin{aligned}
 P(< s > \text{ He likes green apples } < / s >) &= P_1 * P_2 * P_3 * P_4 \\
 &= (7.6 * 10^{-3}) * (1.53 * 10^{-3}) * (2.15 * 10^{-3}) * (2.15 * 10^{-3}) \\
 &= 5.37 * 10^{-11}
 \end{aligned}$$

$$P_1 = \lambda_1 * \text{Count}(< s > \text{ He likes}) / \text{Count}(< s > \text{ He}) + \lambda_2 * \text{Count}(\text{He likes}) / \text{Count}(\text{He}) + \lambda_3 * \text{Count}(\text{likes}) / N$$

- b) Calculate the probability of $P(\text{green} | \text{likes})$ using Kneser Ney smoothing from the corpus given above. $d = \text{discounting factor} = 0.5$

Solution:

$$P_{KN}(w_i | w_{i-1}) = \frac{\max(c(w_{i-1}, w_i) - d, 0)}{c(w_{i-1})} + \frac{1}{c(w_{i-1})} P_{CONTINUATION}(w_i)$$

$$= ((2 - 0.5) / 2) + (0.25) (0.17)$$

- Q2)** Suppose a language model assigns the following conditional n-gram probabilities to a 3-word test set: $1/8, 1/2, 1/6$. What is the perplexity? [2 Marks]

Solution:

$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 \dots w_{i-1})}}$$

- Q3)** $P_{\text{continuation}}(w)$ for a word is defined as follows: [2 Marks]

$$P_{CONTINUATION}(w) = \frac{|\{w_{i-1} : c(w_{i-1}, w) > 0\}|}{\sum_{w'} |\{w'_{i-1} : c(w'_{i-1}, w') > 0\}|}$$

- a) Consider the following incomplete sentence:

"How much wood would a woodchuck chuck would if woodchuck could would chuck"

What is $|\{w_{i-1} : C(w_{i-1} \ w_i) > 0\}|$ for $w_i = \text{"woodchuck"}$?

- i. 0 ii. 1 **iii. 2** iv. 3

- b) Which word is more likely to complete the sentence (follow the last "chuck") based on $P_{\text{continuation}}$?

- i. How ii. wood **iii. would** iv. chuck

- Q4)** Which of the following word pairs, A/B, has A as a hypernym of B? [1 Mark]

- i. Washington/The United States iv. wheel/car
ii. vehicle/car v. None of the above
 iii. Java/programming language

Q5) Consider a trigram HMM tagger with: **[3 Marks]**

- _ The set K of possible tags equal to {D, N, V}
- _ The set V of possible words equal to {the, dog, barks}
- _ The following parameters:

$q(D *, *) = 1$	$q(N D, V) = 0.3$	$e(\text{dog} N) = 0.4$
$q(N *, D) = 0.5$	$q(\text{STOP} N, V) = 0.6$	$e(\text{barks} N) = 0.6$
$q(V *, D) = 0.5$	$q(\text{STOP} V, N) = 0.4$	$e(\text{dog} V) = 0.1$
$q(V D, N) = 0.7$	$e(\text{the} D) = 1$	$e(\text{barks} V) = 0.9$

with all other parameter values equal to 0. Write down the set of all pairs of sequences $x_1 \dots x_{n+1}, y_1 \dots y_{n+1}$ such that the following properties hold:

- _ $p(x_1 \dots x_{n+1}, y_1 \dots y_{n+1}) > 0$
- _ $x_i \in V$ for all $i \in 1 \dots n$
- _ $y_i \in K$ for all $i \in 1 \dots n$, and $y_{n+1} = \text{STOP}$

Solution:

1. * * The dog barks STOP (D N V)
2. * * The dog barks STOP (D V N)
3. * * The barks dog STOP (D N V)
4. * * The barks dog STOP (D V N)
5. * * The dog dog STOP (D N V)
6. * * The dog dog STOP (D V N)
7. * * The barks barks STOP (D N V)
8. * * The barks barks STOP (D V N)

Q6) Show how following lexicalized grammar rule parameter is decomposed into 2 parameters for learning probabilities from training data. Also show how to use smoothed estimation for the decomposed parameters. **[3 Marks]**

$$q(S(\text{saw}) \rightarrow_2 \text{NP}(\text{man}) \text{VP}(\text{saw}))$$

Solution:

$$\begin{aligned}
 & q(S \rightarrow_2 \text{NP VP} | S, \text{saw}) \\
 = & \lambda_1 \times q_{ML}(S \rightarrow_2 \text{NP VP} | S, \text{saw}) + \lambda_2 \times q_{ML}(S \rightarrow_2 \text{NP VP} | S) \\
 & q(\text{man} | S \rightarrow_2 \text{NP VP}, \text{saw}) \\
 = & \lambda_3 \times q_{ML}(\text{man} | S \rightarrow_2 \text{NP VP}, \text{saw}) + \lambda_4 \times q_{ML}(\text{man} | S \rightarrow_2 \text{NP VP}) \\
 & + \lambda_5 \times q_{ML}(\text{man} | \text{NP})
 \end{aligned}$$

Q7) Write down at least two different parse trees (with different probabilities) for following sentence and PCFG. [4 Marks]

“The boy saw the dog in the park with the telescope”

$S \rightarrow NP VP$ 0.8

$S \rightarrow NP VP PP$ 0.2

$NP \rightarrow DET N$ 0.5

$NP \rightarrow NP PP$ 0.5

$VP \rightarrow V NP$ 1.0

$PP \rightarrow P NP$ 1.0

$N \rightarrow dog$ 0.25

$N \rightarrow boy$ 0.25

$N \rightarrow park$ 0.25

$N \rightarrow telescope$ 0.25

$V \rightarrow saw$ 1.0

$P \rightarrow with$ 0.5

$P \rightarrow in$ 0.5

$DET \rightarrow the$ 1.0

Q8) In the following gloss of different word senses of the words "bank" and "coast" are given. Compute similarity between the words "bank" and "coast" using Lesk algorithm. [4 Marks]

Bank₁: sloping land (especially the slope beside a body of water)

Bank₂: a financial institution that accepts deposits and channels the money into lending activities

Bank₃: a long ridge or pile

Bank₄: an arrangement of similar objects in a row or in tiers

Bank₅: a supply or stock held in reserve for future use (especially in emergencies)

Coast₁: the shore of a sea or ocean

Coast₂: a slope down which sleds may coast

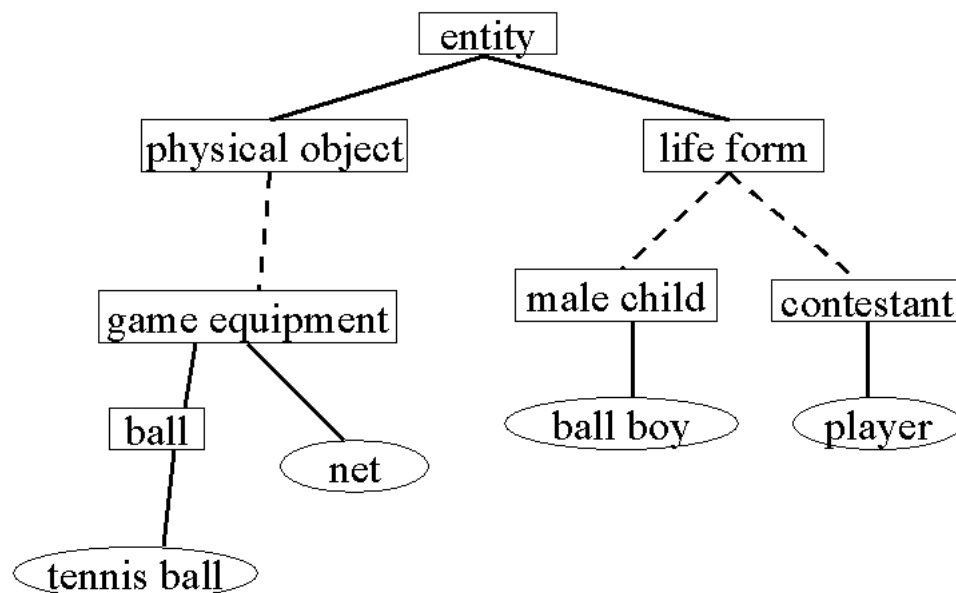
Coast₃: the area within view

Coast₄: the act of moving smoothly along a surface while remaining in contact with it

Solution:

Bank₁ and Coast₂ = 1

Q9) Following is a WordNet hierarchy. The probabilities of words are given in table below: [4 Marks]



Word	Probability
entity	0.395
Physical object	0.167
Life form	0.0231
Game equipment	0.00453
Male child	0.00153
contestant	0.00743
Ball	0.000343
Net	0.00054
Ball boy	0.000113
Player	0.000445
Tennise ball	0.000189

- a) Compute path based similarity between “tennis ball” and “net”

Solution:

1/4

- b) Compute information content based similarity proposed by Lin (Lin Similarity function) between “ball” and “player”

Solution:

$$\log(0.395) / (\log(0.0003) * \log(0.0004))$$

Q10) a) Write down context vectors of words mango and apple using PPMI (Positive Pointwise Mutual Information) of words. [2 Marks]

Counts(w, context)				
	information	data	sweet	Fitness
Banana	0	0	5	3
Apple	3	2	4	6
Mechanical	5	4	1	2
computer	7	6	0	1

Solution:

Probabilities

Apple : $(3/49)=0.06$ 0.04 0.08 0.12

PPMI : 0 0 0.42 0.73

b) Following table gives co-occurrence counts based on syntactic dependencies of words. Write down context vectors of words duty and responsibility using PPMI (Positive Pointwise Mutual Information) of words. (You can assume following table contains all words that can appear as object of a given a word. E.g. total count of words that appear as object of “assert” is 10. Sum of row counts represent total count of the word in collection. E.g. duty appears 22 times in collection. Total words in collection = $N = 100$) [2 Marks]

	Object of assert	Object of assign	Object of avoid	Object of become	Modified by collective	Modified by assumed
duty	3	4	5	3	5	2
responsibility	2	2	7	4	2	7
taxes	0	0	3	0	0	1
danger	0	0	6	0	1	0
control	5	0	0	1	0	0

Q11) Compute value of ROUGE-2 score for following summary. **[2 Marks]**

System Generated Summary: The quake had a preliminary magnitude of 6.9. in an area so isolated there are no roads connecting it to the outside world.

Reference Summary (Human Generated Summary): The quake had a preliminary magnitude of 6.9. An earthquake in the same region in February killed 2300 people and left thousands homeless.

Solution:

7/21

Q12) The first step in query focused multi document summarization is to simplify the sentences. Simplify following sentences using simple rules discussed in class. **[4 Marks]**

- a) Genette's bedroom desk, the biggest disaster area in the house, is a collection of overdue library books, dirty plates, computer components, old mail, cat hair, and empty potato chip bags.
- b) Robbie, a hot-tempered tennis player, charged the umpire and tried to crack the poor man's skull with a racket.
- c) The car began sliding sideways, and then it hit the tree," she said
- d) He died in France, as a matter of fact, and wated to be buried there.

Solution:

- e) Genette's bedroom desk, is a collection of overdue library books, dirty plates, computer components, old mail, cat hair, and empty potato chip bags.
- f) Robbie, charged the umpire and tried to crack the poor man's skull with a racket.
- g) The car began sliding sideways, and then it hit the tree,"
- h) He died in France, and wated to be buried there.

Q13) Word occurrence in sentiment analysis matters more than word frequency. Briefly describe difference between multinomial Naïve Bayes and Boolean Multinomial Naïve Bayes for sentiment analysis. **[2 Marks]**

Solution:

Boolean Multinomial Naïve Bayes clips word counts of all words in all documents at 1.

Q14) Give at least 5 features that can be used to resolve ambiguity in name entity recognition. **[2 Marks]**

Solution:

Identity of word

Neighboring words

Part of speech of word

Part of speech of neighboring words

Uppercase

Shape of word

Presence of hyphen