

Previous year (2020) NLP questions for reference

Intro, Regex, Morphology, LM,

Quizzes

1. A book catalogue contains details of book in the following format: (3 marks)

<Name of Book>,<author last name>,<author first name>

Give a regex that matches all lines in which the author's first name is "Martin". For example, "High School English,Wren,Martin" should match.

$^{\wedge}[^,]^*,[^,]^*,[^,]^*\backslash b(\backslash Martin+)\backslash b$

Note: This is one regular expression that matches the given pattern. There could also be other correct regular expressions. Therefore, the evaluation is not restricted to only the expression mentioned as the answer.

2. Write the output of the Porter Stemmer for the following words: (1 mark)

- A. Better
- B. Going
- C. Flying

Better, Go, Fli

3. Consider the sentence: "learn this, learn that, learn all you can but also learn to unlearn."

- A. Give the count of tokens and types for the following sentence.(Specifically mention tokens = x and types = y) (1 mark)

Tokens = 16 (0.5 marks)

Types = 12 (0.5 marks)

- B. Give the number of bigrams. (1 mark)

15

- C. Find the probability of 'that' given 'learn', considering a bigram model with Laplace smoothing. (2 marks)

Ans. $1/7 = 0.1428$

Vocabulary, V = learn, this, that, all, you, can, but, also, to, unlearn

$|V| = 10$

$P(\text{that}|\text{learn}) = (\text{Count}(\text{learn that}) + 1) / (\text{Count}(\text{learn}) + |V|)$

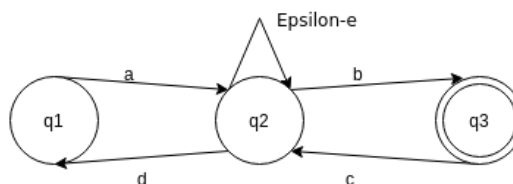
$= (1+1) / (4+10) = 2/14 = 1/7 = 0.1428$

4. For the following epsilon-NFA,

- a. How many states will be there in the equivalent minimal DFA (excluding dead state)? (1 mark)

3

- b. Write regular expression defined by the equivalent minimal DFA (1 mark)



$a(da)*b(c(da)*b)*$

5. Identify the predicate and assign semantic roles in the following sentence. (2 marks)

"John showed Bob his Toyota yesterday."

Predicate: showed (1 mark)

Semantic roles: Agent - John (0.25 marks)

Patient/Receiver - Bob (0.25 marks)

Instrument - Toyota (0.25 marks)

Time - yesterday (0.25 marks)

6. Give two examples of nominalization with different affixes (Note: affixes should not be "-ation", "-ness", "-er", and "-ee"). (2 marks)

Intensity, Movement, and many more

7. Suppose we have a sequence of N independent samples: W = apple apple apple banana banana dates dates eggs eggs eggs frogs grapes grapes. Using Good Turing Smoothing, find the updated probability of

a. drawing an orange. (1 mark)

b. drawing a banana. (1 mark)

W = apple apple apple banana banana dates dates
eggs eggs eggs frogs grapes grapes

N_c = No. of objects with frequency c

count (apple) = 3	}	$N_1 = 1$
count (banana) = 2		$N_2 = 3$
count (dates) = 2		$N_3 = 2$
count (eggs) = 3		
count (frogs) = 1		
count (grapes) = 2		$N = 13$

a) As "orange" is an unseen word,
$$P_{GT}^*(\text{unseen}) = \frac{N_1}{N} = \frac{1}{13}$$

b) Now, "banana" is a seen word,
$$c^* = (c+1) \frac{N_{c+1}}{N_c} \Rightarrow c^* = \frac{(2+1) N_{2+1}}{N_2} = \frac{3 N_3}{N_2}$$
$$= \frac{3 \times 2}{3}$$

$$P_{GT}^*(\text{seen}) = \frac{2}{13}$$

End Sem

8. Construct two-level FSTs (surface \rightarrow intermediate \rightarrow lexical) for the morphological analysis of the following words. [20 marks]

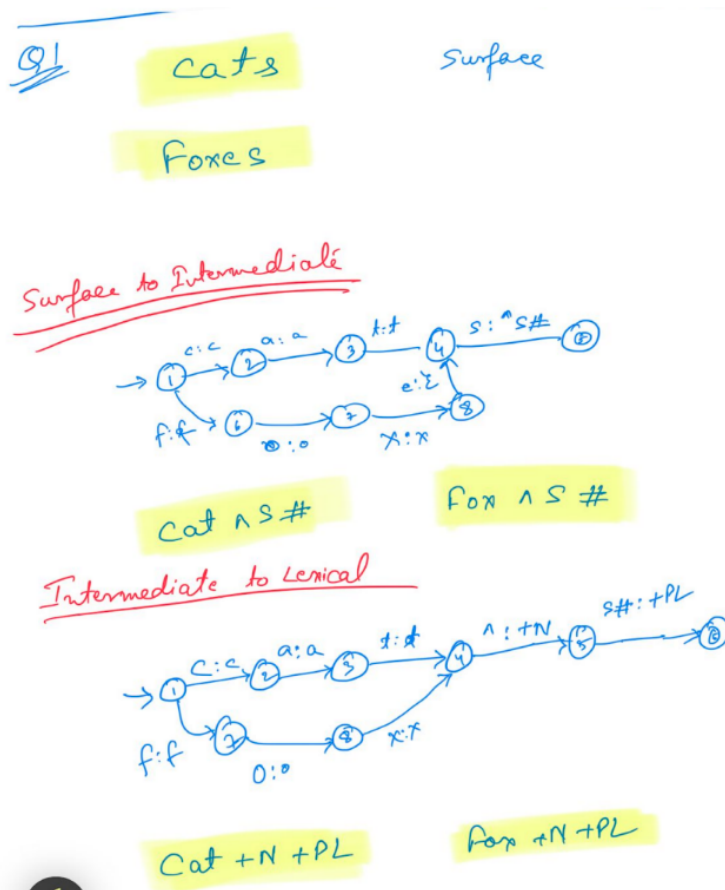
Inputs: cats,

foxes

Outputs: cat + N + PL,

fox + N + PL

[Note: 5 marks each for surface \rightarrow intermediate for the two words and 5 marks each for intermediate \rightarrow lexical for the two words]



9. For a trigram language model, find the next probable word (blank) in the sequence considering Laplace smoothing. Show the computation. [4 marks]

Corpus: ["He is apple of my eye.", "He hit my eye."]

Vocabulary: [He, is, apple, of, my, eye, hit, ate, .]

Sequence: He ate my _____

$$|\text{Vocabulary}| = |V| = 9$$

For a trigram model, below are the probabilities of the next probable word with Laplace smoothing,

$$P(\text{He} \mid \text{ate my}) = \frac{\text{count}(\text{ate my He}) + 1}{\text{count}(\text{ate my}) + 9} = \frac{0 + 1}{0 + 9} = \frac{1}{9}$$

$P(\text{is} \mid \text{ate my}) = \text{count}(\text{ate my is})+1 / \text{count}(\text{ate my})+9 = 0+1 / 0+9 = 1/9$

$P(\text{apple} \mid \text{ate my}) = \text{count}(\text{ate my apple})+1 / \text{count}(\text{ate my})+9 = 0+1 / 0+9 = 1/9$

$P(\text{of} \mid \text{ate my}) = \text{count}(\text{ate my of})+1 / \text{count}(\text{ate my})+9 = 0+1 / 0+9 = 1/9$

$P(\text{my} \mid \text{ate my}) = \text{count}(\text{ate my my})+1 / \text{count}(\text{ate my})+9 = 0+1 / 0+9 = 1/9$

$P(\text{eye} \mid \text{ate my}) = \text{count}(\text{ate my eye})+1 / \text{count}(\text{ate my})+9 = 0+1 / 0+9 = 1/9$

$P(\text{hit} \mid \text{ate my}) = \text{count}(\text{ate my hit})+1 / \text{count}(\text{ate my})+9 = 0+1 / 0+9 = 1/9$

$P(\text{ate} \mid \text{ate my}) = \text{count}(\text{ate my ate})+1 / \text{count}(\text{ate my})+9 = 0+1 / 0+9 = 1/9$

All the probabilities are the same because the word “ate” is not present in the corpus.

Hence, the next probable word can be any word present in the vocabulary because the probability of any word is the same.

One of the heuristics is sorting the vocabulary and then picking the first word from the vocabulary.

Therefore, Next predicted word is “apple”.

Note: The final solution must identify a word with any heuristic.

-----END-----

Indraprastha Institute of Information Technology Delhi (IIIT-Delhi)

CSE556: NLP

Quiz 2