

Final Exam NATURAL LANGUAGE PROCESSING

8th of September 2015

Remember to fill in your name on all pages
GOOD LUCK

PROBLEM 1 (1 point). Compute the Good-Turing smoothing of the words *perch*, *trout* and *bass* in a corpus that contains only the following words (words are given with their counts in the table below):

<i>carp</i>	<i>perch</i>	<i>whitefish</i>	<i>trout</i>	<i>salmon</i>	<i>eel</i>
10	3	2	1	1	1

PROBLEM 2 (2 points). Suppose we have the following short movie reviews, each labeled with a genre, either **comedy** or **action**.

1. *fun, couple, love, love* (**comedy**)
2. *fast, furious, shoot* (**action**)
3. *couple, fly, fast, fun, fun* (**comedy**)
4. *furious, shoot, shoot, fun* (**action**)
5. *fly, fast, shoot, love* (**action**)

Using a simple Naïve Bayes approach with Laplace smoothing, how would we classify the following review: *fun, couple, shoot, action*?

PROBLEM 3 (3 points). In a set of 806,791 documents, we get the following data on a few terms and a few documents:

term	document frequency	Doc 1	Doc 2	Doc 3
car	18,165	27	4	24
auto	6,723	3	33	0
insurance	19,241	0	39	29
best	25,235	14	9	17

Compute the tf-idf value for these terms and documents. What is the cosine similarity between query “best car best insurance” and Doc 1, 2 and 3, respectively (use tf-idf weighting - nnc.ltc variation). Which of these three documents would be ranked first by a search engine using the nnc.ltc scheme?

PROBLEM 4 (1 point). What is the mean average precision (MAP) for the following sequence of retrieved documents, where R denotes a relevant document and N denotes an irrelevant document? (Assume there are 20 relevant documents in the collection)

N R R N R R N R N R N R N

PROBLEM 5 (3 points). Simulate the run of the Earley Parser for grammar G on the sentence *The large can can hold the water*, where

$G = (\{S, NP, VP, D, J, N, V\}, \{the, large, can, hold, water\}, S, R)$,

and the set R of rules is given by

$R = \{S \rightarrow NP VP, NP \rightarrow D J N \mid D N \mid J N, VP \rightarrow V VP \mid V NP\} \cup \{D \rightarrow the, J \rightarrow large, N \rightarrow can|water, V \rightarrow can|hold\}$

Does the grammar G accept this sentence?

DEFINITION 1. Let w be a word that appears c times in the corpus, N_c = the count of things we've seen c times, and N = the total number of tokens in the corpus. Then,

$$c^*(w) = \begin{cases} \frac{(c+1)N_{c+1}}{N_c}, & \text{if } c > 0 \\ N_1, & \text{if } c = 0 \end{cases}$$

$$P_{GT}^*(w) = \frac{c^*(w)}{N}$$

DEFINITION 2. Let S be a sentence, C a class of possible labels and V the vocabulary. Then,

$$c_{NB}(S) = \operatorname{argmax}_{c \in C} P(c) \prod_{w \in S} P(w | c),$$

where the prior probability of the class is $P(c) = \text{doccount}(c)/N_{doc}$, and the probability of a word w given the class c , computed using Laplace (add-1) smoothing with unknown words, is:

$$P(w | c) = \frac{\text{count}(w, c) + 1}{\sum_{v \in V} \text{count}(v, c) + |V| + 1}$$

DEFINITION 3. SMART Notation: denotes the combination in use in an engine, with the notation ddd.qqq, using the acronyms from the following table:

Term frequency	Document frequency	Normalization
n (natural): $tf_{t,d}$	n (no): 1	n (none): 1
l (logarithm): $1 + \log_{10}(tf_{t,d})$	t (idf): $\log_{10}(\frac{N}{df_t})$	c (cosine): $\frac{1}{\sqrt{\sum_i x_i^2}}$

DEFINITION 4. The Mean Average Precision (MAP) is the average of the precision value obtained for the top k documents, each time a relevant document is retrieved.

Precision = TP/(TP+FP)

TP = true positives (the number of gold answers that were correctly guessed)

FP = false positives (the number of guessed answers that were incorrect)

FN = false negative (the number of gold answers that were not correctly guessed)

DEFINITION 5. Start by building a sequence of state sets called *Earley sets*

Input: $x_1 x_2 \dots x_n$

Sequence of state sets: S_0, S_1, \dots, S_n

A set S_i contains various “tokens”: $[A \rightarrow \alpha \bullet \beta, j]$

- *Scan*: $[A \rightarrow \dots \bullet a \dots, j] \in S_i$, $a = x_{i+1}$,
add $[A \rightarrow \dots a \bullet \dots, j]$ to S_{i+1}

- *Predict*: $[A \rightarrow \dots \bullet B \dots, j] \in S_i$
add $[B \rightarrow \bullet \alpha, i]$ to S_i for all rules $B \rightarrow \alpha$

- *Complete*: $[A \rightarrow \dots \bullet, j] \in S_i$
add $[B \rightarrow \dots A \bullet \dots, k]$ to S_i for all tokens $[B \rightarrow \dots \bullet A \dots, k]$ in S_j

Accept $x_1 x_2 \dots x_n$ if in S_n there exist a rule of type $[S \rightarrow \alpha \bullet, 0]$