

Optional tests

Marco Kuhlmann

01

Text classification

(3 points)

- a) Use Maximum Likelihood estimation with add-one smoothing to estimate the class probabilities and word probabilities of a Naive Bayes text classifier from the following document collection. Assume that the vocabulary consists of the set of all words occurring in the documents. Answer with fractions.

document	class
ant	A
ant bear	B
bear camel	B
camel	C

- b) Based on the probabilities just estimated, compute the class-specific scores that the Naive Bayes classifier uses to predict the class for the following document:

ant bear camel

Answer with fractions.

- c) Here are some class frequencies in a document collection:

	class X	class Y	class Z
training data	2,460	2,952	1,968
test data	738	492	615

What is the precision for class X of the most frequent class baseline on the test data? Answer with a fraction.

Sample answers:

a) Estimated probabilities:

$$P(A) = 1/4 \quad P(\text{ant} | A) = 2/4 \quad P(\text{bear} | A) = 1/4 \quad P(\text{camel} | A) = 1/4$$

$$P(B) = 2/4 \quad P(\text{ant} | B) = 2/7 \quad P(\text{bear} | B) = 3/7 \quad P(\text{camel} | B) = 2/7$$

$$P(C) = 1/4 \quad P(\text{ant} | C) = 1/4 \quad P(\text{bear} | C) = 1/4 \quad P(\text{camel} | C) = 2/4$$

b) Class-specific scores:

$$\text{score}(A) = P(A) \cdot P(\text{ant} | A) \cdot P(\text{bear} | A) \cdot P(\text{camel} | A) = \frac{1}{4} \cdot \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{1}{4}$$

$$\text{score}(B) = P(B) \cdot P(\text{ant} | B) \cdot P(\text{bear} | B) \cdot P(\text{camel} | B) = \frac{2}{4} \cdot \frac{2}{7} \cdot \frac{3}{7} \cdot \frac{2}{7}$$

$$\text{score}(C) = P(C) \cdot P(\text{ant} | C) \cdot P(\text{bear} | C) \cdot P(\text{camel} | C) = \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{2}{4}$$

c) $\frac{0}{0}$ (which is mathematically undefined)

02

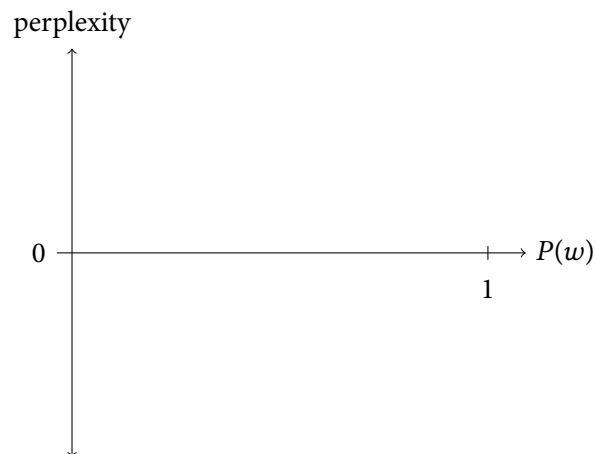
Language modelling

(3 points)

The WikiText language modelling dataset is a collection of 2 million tokens extracted, comprising a vocabulary of 33,000 unique words. We have the following selected counts of unigrams and bigrams:

<i>the</i>	<i>book</i>	<i>first</i>	<i>the book</i>	<i>book the</i>	<i>first book</i>	<i>book first</i>
113,161	611	3,981	200	1	8	0

- a) Estimate the following probabilities using maximum likelihood estimation without smoothing. Answer with fractions.
- $P(\textit{first})$
 - $P(\textit{book} \mid \textit{first})$
- b) Now, use additive smoothing with $k = 0.05$.
- $P(\textit{first})$
 - $P(\textit{first} \mid \textit{book})$
- c) We evaluate a unigram language model on a one-word sentence w . Sketch how the perplexity of the model varies with $P(w)$ by completing the following diagram. What is the minimal value for the perplexity measure?

**Sample answers:**

- $\frac{3981}{2000000}$
 - $\frac{8}{3981}$
- $\frac{3981+0.05}{2000000+0.05 \times 33000}$
 - $\frac{0+0.05}{611+0.05 \times 33000}$
- The graph has the same shape as that for entropy in slide 35 from the slide deck for Unit 2, but the minimal value is 1 instead of 0.

03

Part-of-speech tagging

(3 points)

- a) The evaluation of a part-of-speech tagger produced the confusion matrix shown below. The marked cell gives the number of times the system tagged a word as an adjective (ADJ) whereas the gold standard specified it as a noun (NOUN).

	ADJ	DET	NOUN	VERB
ADJ	1475	0	221	31
DET	5	1835	3	0
NOUN	45	5	3887	167
VERB	28	1	387	2135

Compute the following values. Answer with fractions.

- i. precision on verbs
 - ii. recall on adjectives
- b) Training a Hidden Markov Model (HMM) amounts to estimating two types of probabilities. What is the total number of probability values you need to estimate when training a model with 10 tags and a vocabulary of 29,508 unique words? Answer with a formula that evaluates to a concrete number (example: 2×3). Ignore the beginning-of-sentence and end-of-sentence markers.
- c) One difference between a multi-class perceptron tagger and a tagger based on an HMM is in the feature sets. Which (zero or more) of the following features would you have to choose to provide the multi-class perceptron tagger with the same information that the HMM tagger has access to?
- i. current word
 - ii. word to the left of the current word
 - iii. word to the right of the current word
 - iv. part-of-speech tag of the word to the left of the current word

Sample answers:

- a) i. $\frac{2135}{31+0+167+2135}$ ii. $\frac{1475}{1475+0+221+31}$
 b) $10 \times 10 + 10 \times 29508$
 c) i. and iv.

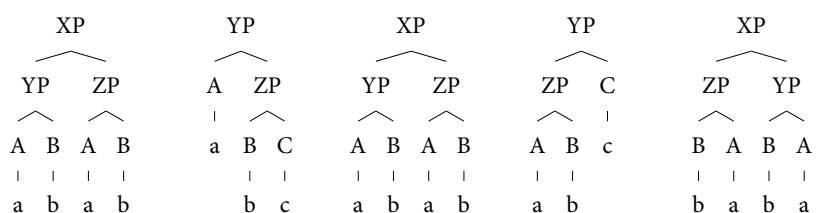
04 Syntactic analysis

(3 points)

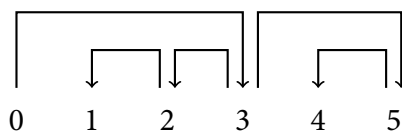
a) You sum up all rule probabilities in a certain probabilistic context-free grammar. Which (zero or more) of the following values can you *not* get as a result?

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

b) Below is a small phrase structure treebank. Read off all rules with left-hand sides XP, YP and ZP and estimate their rule probabilities using maximum likelihood estimation (no smoothing).



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



Sample answers:

a) i. and iii.

b) Rules and estimated probabilities:

$$XP \rightarrow YP ZP \frac{2}{3} \quad XP \rightarrow ZP YP \frac{1}{3}$$

$$YP \rightarrow A B \frac{2}{5} \quad YP \rightarrow A ZP \frac{1}{5} \quad YP \rightarrow ZP C \frac{1}{5} \quad YP \rightarrow B A \frac{1}{5}$$

$$ZP \rightarrow A B \frac{3}{5} \quad ZP \rightarrow B C \frac{1}{5} \quad ZP \rightarrow B A \frac{1}{5}$$

c) Possible answers:

- SH SH SH LA SH LA SH SH LA RA RA
- SH SH SH LA SH SH SH LA RA LA RA

- a) Choose the correct semantic relation: synonym, antonym, hyponym, hypernym?

pigeon	is a/an ... of	animal
big	is a/an ... of	large
parent	is a/an ... of	child
begin	is a/an ... of	start
screwdriver	is a/an ... of	tool

- 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 sentence? Ignore the word *course*, punctuation, and stop words.

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

- c) We read off word vectors from the following co-occurrence matrix (target words correspond to rows, context words correspond to columns):

	<i>caws</i>	<i>dafad</i>
<i>cheese</i>	6	2
<i>sheep</i>	0	4
<i>goat</i>	1	6
<i>bread</i>	5	0

Sort the target words in decreasing degree of semantic similarity (most similar to least similar) to the word *cheese*, assuming that semantic similarity is measured in terms of cosine similarity.

Sample answers:

a) Semantic relations:

pigeon	is a hyponym of	animal
big	is a synonym of	large
parent	is an antonym of	child
begin	is a synonym of	start
screwdriver	is a hyponym of	tool

b) Sense 1 (match with *normal*)

c) *cheese, bread, goat, sheep*