

# Einführung in die Computerlinguistik

## Hausaufgabe Language Models, Abgabe 07.05.2019

Laura Kallmeyer

SoSe 2019, Heinrich-Heine-Universität Düsseldorf

**Aufgabe 1** Consider a toy example where the sentences are built over  $\{a, b, c\}$  and our training data consists of the following sentences:

$\langle s \rangle a a b c a b c \langle /s \rangle$

$\langle s \rangle a b a b b a a b b a a \langle /s \rangle$

$\langle s \rangle b b a a b a b a b a \langle /s \rangle$

$\langle s \rangle b b c a b a b \langle /s \rangle$

1. Give the following bigram probabilities and their  $\log_2$  values:

$$P(a|\langle s \rangle) \quad P(a|a) \quad P(\langle /s \rangle|a)$$

2. Using these values, the probability of a sequence  $\langle s \rangle w_1 w_2 \dots w_n$  amounts to

$$P(\langle s \rangle w_1 w_2 \dots w_n) = 2^{\sum_{k=1}^n \log_2 P(w_k | w_{k-1})}$$

while the perplexity value of this sequence is

$$PP(\langle s \rangle w_1 w_2 \dots w_n) = \sqrt[n]{\frac{1}{2^{\sum_{k=1}^n \log_2 P(w_k | w_{k-1})}}} = \sqrt[n]{2^{\sum_{k=1}^n |\log_2 P(w_k | w_{k-1})|}}$$

Compute both, probability and perplexity, of the following sequences:

(1)  $\langle s \rangle a \langle /s \rangle$

(2)  $\langle s \rangle a a a a \langle /s \rangle$

3. Explain why the perplexity value is better for measuring the quality of (1) and (2).

Solution:

$$1. \quad P(a|\langle s \rangle) = \frac{1}{2}, \quad \log_2 \frac{1}{2} = -1 \quad P(a|a) = \frac{1}{4}, \quad \log_2 \frac{1}{4} = -2 \quad P(\langle /s \rangle|a) = \frac{1}{8}, \quad \log_2 \frac{1}{8} = -3$$

$$2. \quad (1): \text{probability } 2^{-1-3} = 2^{-4} = \frac{1}{2^4} = \frac{1}{16}, \text{ perplexity } \sqrt[4]{2^4} = 4$$

$$(2): \text{probability } 2^{-1-2-2-2-3} = 2^{-10} = \frac{1}{2^{10}} = \frac{1}{1024}, \text{ perplexity } \sqrt[10]{2^{10}} = 2^2 = 4$$

3. The probability of a sequence is not a good candidate for measuring its quality since the longer a sequence the lower its probability, i.e., shorter sequences are preferred. This is why (1) has a much higher probability than (2).

The perplexity value does not have this problem since for a sentence of length  $n$ , we have an exponent  $\frac{1}{n}$ , which balances the multiplication of  $n$  probabilities ( $\sqrt[n]{x^n} = x$  for all  $n$ ).

**Aufgabe 2** Consider again the same training data.

1. Modify the training data adopting a treatment of unknown words that changes the first occurrence of every word (here first occurrence of  $a, b, c$ ) into  $\langle \text{UNK} \rangle$ .

Give the following bigram probabilities estimated by this model:

$$P(a|a) \quad P(b|a) \quad P(c|a) \quad P(\langle \text{UNK} \rangle|a) \quad P(a|\langle s \rangle)$$

2. Calculate the probability of the sequence

$$(3) \quad \langle \mathbf{s} \rangle a a d$$

according to this model.

Solution:

1. Modified training data:

$\langle \mathbf{s} \rangle \langle \text{UNK} \rangle a \langle \text{UNK} \rangle \langle \text{UNK} \rangle a b c \langle / \mathbf{s} \rangle$

$\langle \mathbf{s} \rangle a b a b b a a b b a a \langle / \mathbf{s} \rangle$

$\langle \mathbf{s} \rangle b b a a b a b a b a \langle / \mathbf{s} \rangle$

$\langle \mathbf{s} \rangle b b c a b a b \langle / \mathbf{s} \rangle$

$$P(a|a) = \frac{3}{15} = \frac{1}{5} \quad P(b|a) = \frac{9}{15} = \frac{3}{5} \quad P(c|a) = \frac{0}{15} = 0 \quad P(\langle \text{UNK} \rangle | a) = \frac{1}{15} \quad P(a|\langle \mathbf{s} \rangle) = \frac{1}{4}$$

$$2. P(\langle \mathbf{s} \rangle a a d) = \frac{1}{4} \cdot \frac{1}{5} \cdot \frac{1}{15}$$

**Aufgabe 3** Consider again the training data from the first exercise.

This time, we assume that treating unknown words is not necessary. But we do Laplace smoothing in order to cope with unseen  $n$ -grams.

Note that our vocabulary of words that can be the second element of a bigram includes the end of sentence marker  $\langle / \mathbf{s} \rangle$ , i.e.,  $V = \{a, b, c, \langle / \mathbf{s} \rangle\}$ .

1. Calculate the new values we obtain for

$$P(a|\langle \mathbf{s} \rangle) \quad P(a|a) \quad P(c|a) \quad P(a|c) \quad P(\langle / \mathbf{s} \rangle | a)$$

2. Calculate the probabilities of the following sequences using this new model:

$$(4) \quad \langle \mathbf{s} \rangle a c a \langle / \mathbf{s} \rangle$$

$$(5) \quad \langle \mathbf{s} \rangle a a c a c a \langle / \mathbf{s} \rangle$$

Solution:

$$1. P(a|\langle \mathbf{s} \rangle) = \frac{3}{8} \quad P(a|a) = \frac{5}{20} = \frac{1}{4} \quad P(c|a) = \frac{1}{20} \quad P(a|c) = \frac{3}{7} \quad P(\langle / \mathbf{s} \rangle | a) = \frac{3}{20}$$

$$2. P((4)) = \frac{3}{8} \cdot \frac{1}{20} \cdot \frac{3}{7} \cdot \frac{3}{20}$$

$$P((5)) = \frac{3}{8} \cdot \frac{1}{4} \cdot \frac{1}{20} \cdot \frac{3}{7} \cdot \frac{1}{20} \cdot \frac{3}{7} \cdot \frac{3}{20}$$