

Chapter 3: N-gram Language Models

Exercises

$$(Q3.1) \quad P(w_n | w_{n-2} w_{n-1}) = \frac{C(w_{n-2} w_{n-1} w_n)}{C(w_{n-2} w_{n-1})}$$

<S><S> I am Sam </S>

<S><S> Sam I am </S>

<S><S> I do not like green eggs and ham </S>

$$P(I | <S><S>) = 2/3$$

$$P(am | <S>I) = 1/2$$

$$\begin{aligned}(Q3.2) \quad P(i \text{ want chinese food}) &= P(i | <S>) P(\text{want} | i) \\ &\quad P(\text{chinese} | \text{want}) P(\text{food} | \text{chinese}) \\ &\quad P(</S> | \text{food}) \\ &= 0.25 \times 0.33 \times 0.0065 \times 0.52 \times 0.68 \\ &= 0.0001896\end{aligned}$$

$$\begin{aligned}P(i \text{ want chinese food}) &= P(i | <S>) P(\text{want} | i) P(\text{chinese} | \text{want}) \\ &\quad P(\text{food} | \text{chinese}) P(</S> | \text{food}) \\ &= 0.19 \times 0.21 \times 0.0029 \times 0.052 \times 0.4 \\ &= 0.00002406\end{aligned}$$

(Q3.3) The unsmoothed probability is higher because the bigrams used in the sentences are very common and has probabilities. However, in the smoothed case, their probabilities are distributed among not-so-common bigrams which are not used in our test statement.

(Q3.4)

	<s>	I	am	Sam	do	not	like	green	eggs	and	</s>
<s>	1	4	1	2	1	1	1	1	1	1	1
I	1	1	4	1	2	1	1	1	1	1	1
am	1	1	1	3	1	1	1	1	1	1	2
Sam	1	2	1	1	1	1	1	1	1	1	4
do	"	1	"	"	"	2	"	"	"	"	1
not	"	"	"	"	"	1	2	"	"	"	"
like	"	"	"	"	"	"	1	2	"	"	"
green	"	"	"	"	"	"	"	1	2	"	"
eggs	"	"	"	"	"	"	"	"	1	2	"
and	"	"	"	2	"	"	"	"	"	1	"
</s>	"	"	"	1	"	"	"	"	"	"	"

$$P(\text{Sam} | \text{am}) = \frac{C(\text{am Sam})}{C(\text{am})} = \frac{3}{14} = .214$$

(Q3.5)

	<s>	a	b
<s>	0	2	2
a	0	1	1
b	0	1	1

- $P(a a) = P(a | <s>) P(a | a) = 0.5 \times 0.5 = 0.25$
 $P(a b) = P(a | <s>) P(b | a) = 0.5 \times 0.5 = 0.25$
 $P(b b) = P(b | <s>) P(b | b) = 0.5 \times 0.5 = 0.25$
 $P(b a) = P(b | <s>) P(a | b) = 0.5 \times 0.5 = 0.25$

$$P(s \in \{a, b\}^L) = 1$$

$$\begin{aligned}
 (\text{Q3.6}) \quad P(\omega_3 | \omega_1 \omega_2) &= \frac{C(\omega_1 \omega_2 \omega_3) + 1}{\sum_{\omega} (C(\omega_1 \omega_2 \omega) + 1)} \\
 &= \frac{C(\omega_1 \omega_2 \omega_3) + 1}{C(\omega_1 \omega_2) + 9}
 \end{aligned}$$

$$\begin{aligned}
 (\text{Q3.7}) \quad P(\text{sam} | \text{am}) &= d_1 P(\text{sam}) + d_2 (P(\text{sam} | \text{am})) \\
 &= \frac{1}{2} \times \frac{42}{25} + \frac{1}{2} \times \frac{2}{3} \\
 &= \frac{2}{25} + \frac{1}{3} = 0.41
 \end{aligned}$$

$$\begin{aligned}
 (\text{Q3.12}) \quad PP(\omega) &= \sqrt[N]{\prod_{i=1}^N \frac{1}{P(\omega_i)}} \\
 &= \sqrt[10]{\frac{(100)^{10}}{(91)^9}} \\
 &= 1.0726
 \end{aligned}$$

$$P(0) = \frac{91}{100}$$

$$P(1) = P(2) \dots P(9) = \frac{1}{100}$$