

corpus:

dogs bark often
cats meow sometimes
dogs don't meow
cats don't bark

Perplexity Example

unigram:

dogs	bark	often	cats	meow	sometimes	don't
$\frac{2}{12}$	$\frac{2}{12}$	$\frac{1}{12}$	$\frac{2}{12}$	$\frac{2}{12}$	$\frac{1}{12}$	$\frac{2}{12}$

test sentences:

dogs don't bark
cats meow

PP(dogs don't bark cats meow)

$$= \sqrt[N]{\prod_i^N \frac{1}{P(w_i)}}$$

$$= 5 \sqrt{\frac{1}{P(\text{dogs}) \cdot P(\text{don't}) \cdot P(\text{bark}) \cdot P(\text{cats}) \cdot P(\text{meow})}}$$

$$= 5 \sqrt{\frac{1}{\frac{2}{12} \cdot \frac{2}{12} \cdot \frac{2}{12} \cdot \frac{2}{12} \cdot \frac{2}{12}}}$$

$$= \left(\frac{1}{\frac{2}{12} \cdot \frac{2}{12} \cdot \dots \cdot \frac{2}{12}} \right)^{\frac{1}{5}}$$

$$= \left(\frac{2}{12} \cdot \frac{2}{12} \cdot \dots \cdot \frac{2}{12} \right)^{-\frac{1}{5}}$$

$$= 6$$

Unigram PP calculation using \log_{10} probs

dogs bark often cats meow sometimes don't

$$\log\left(\frac{1}{6}\right) \quad \log\left(\frac{1}{6}\right) \quad \log\left(\frac{1}{12}\right) \quad \log\left(\frac{1}{6}\right) \quad \log\left(\frac{1}{6}\right) \quad \log\left(\frac{1}{12}\right) \quad \log\left(\frac{1}{6}\right)$$
$$\approx -0.78 \quad -0.78 \quad -1.08 \quad -0.78 \quad -0.78 \quad -1.08 \quad -0.78$$

PP (dogs don't bark cats meow)

$$= \sqrt[N]{\prod_i \frac{1}{P(w_i)}}$$

$$= \sqrt[N]{\frac{1}{\sum_{10} \log(P(w_i))}}$$

$$= \sqrt[5]{\frac{1}{10^{-0.78 - 0.78 - 0.78 - 0.78 - 0.78}}}$$

$$= \left(\frac{1}{10^{-3.9}} \right)^{\frac{1}{5}}$$

$$= \left(10^{-3.9} \right)^{-\frac{1}{5}}$$

$$= 10^{+ \frac{3.9}{5}}$$

$$= 6$$

$$PP = 10^{-\frac{\sum \log(P(w_i))}{N}}$$

$\langle s \rangle$ dogs bark often $\langle /s \rangle$
 $\langle s \rangle$ cats meow sometimes $\langle /s \rangle$
 $\langle s \rangle$ dogs don't meow $\langle /s \rangle$
 $\langle s \rangle$ cats don't bark $\langle /s \rangle$

bigram:

unigram counts

	dogs	bark	often	cats	meow	sometimes	don't	$\langle /s \rangle$	
$\langle s \rangle$	2/4			2/4					4
dogs		1/2					1/2		2
bark			1/2					1/2	2
often								1	1
cats					1/2		1/2		2
meow						1/2		1/2	2
smtms								1	1
don't		1/2			1/2				2

test sentences: $\langle s \rangle$ dogs don't bark $\langle /s \rangle$
 $\langle s \rangle$ cats meow $\langle /s \rangle$

= # n-grams
 $N = 4 + 3 = 7$

$$P(\langle s \rangle \text{ dogs don't bark } \langle /s \rangle \langle s \rangle \text{ cats meow } \langle /s \rangle)$$

$$= \sqrt[7]{\frac{1}{P(\text{dogs} | \langle s \rangle) \cdot P(\text{don't} | \text{dogs}) \cdot P(\text{bark} | \text{don't}) \cdot P(\langle /s \rangle | \text{bark}) \cdot P(\text{cats} | \langle s \rangle) \cdot P(\text{meow} | \text{cats}) \cdot P(\langle /s \rangle | \text{meow})}}$$

$$= \sqrt[7]{\frac{1}{\frac{2}{4} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{2}{4} \cdot \frac{1}{2} \cdot \frac{1}{2}}}$$

$$= \left(\frac{1}{.0078} \right)^{\frac{1}{7}}$$

$$= (.0078)^{-\frac{1}{7}}$$

$$= 2$$

(as expected, less than perplexity of the unigram model)

	dogs	bark	often	cats	meow	sometimes	don't	</s>
<s>	-0.3			-0.3				
dogs		-0.3					-0.3	
bark			-0.3					-0.3
often								0
cats					-0.3		-0.3	
meow						-0.3		-0.3
smtms								0
don't		-0.3			-0.3			

using base 10:

$10^x = a$ ← $P(w_2|w_1)$ in the probability matrix,
 replace each $P(w_2|w_1)$ with its log: $\log_{10}(P(w_2|w_1))$
 $\log_{10}(a) = x$

PP(<s> dogs don't bark </s> <s> cats meow </s>)

$$= \sqrt[N]{\prod_i^N \frac{1}{P(w_2|w_1)}}$$

$$= \sqrt[N]{\frac{1}{10^{\sum \log(P(w_2|w_1))}}}$$

$$= \left(\frac{1}{10^{\sum \log(P(w_2|w_1))}} \right)^{\frac{1}{N}}$$

$$= \left(10^{\sum \log(P(w_2|w_1))} \right)^{-\frac{1}{N}}$$

$$= 10^{-\frac{\sum \log(P(w_2|w_1))}{N}}$$

$$= - \frac{(\log(P(\text{dogs}|\text{<s>})) + \log(P(\text{don't}|\text{dogs})) + \log(P(\text{bark}|\text{don't})) + \log(P(\text{</s>}|\text{bark})) + \log(P(\text{cats}|\text{<s>})) + \log(P(\text{meow}|\text{cats})) + \log(P(\text{</s>}|\text{meow}))}{N}$$

$$= 10^{-\frac{(-0.3 - 0.3 - 0.3 - 0.3 - 0.3 - 0.3 - 0.3)}{7}} = 10^{-\left(-\frac{2.1}{7}\right)} = 2$$