# Numericals based on N-gram model

Q1. Find the total count of unique bigrams for which likelihood will be estimated.

Alice went to the café

Bob was waiting for Alice

Alice and Bob went to the museum

Total Bigrams = 20

Unique Bigrams = 17

Sol$^n$ – ⟨s⟩ Alice went to the cafe ⟨/s⟩

⟨s⟩ Bob was waiting for Alice ⟨/s⟩

⟨s⟩ Alice and Bob went to the museum⟨/s⟩

| ⟨s⟩ Alice | ⟨s⟩ Bob | ⟨s⟩ Alice |
|---|---|---|
| Alice went | Bob was | Alice and |
| went to | was waiting | and Bob |
| to the | waiting for | Bob went |
| the cafe | for Alice | went to |
| cafe ⟨/s⟩ | Alice ⟨/s⟩ | to the |
| | | the museum |
| | | museum ⟨/s⟩ |
| ⑥ | ⑥ | ⑤ |

Q2.Find the probability of the statement:<S> Michael and Zack played at the playground</s> from the following corpus. Assume a trigram Language model.

<S>the school was open</S>

<S>Michael and Zack went to the school</S>

<S>the playground at the school was huge</S>

<S>Bob and Zack played at the playground</S>

<S>Bob, Michael and Zack were friends</S>

Sol$^n$ – Calculate the trigram's probability

$P(Michael \mid <s>) = \dfrac{1}{5}$

$P(and \mid <s> Michael) = \dfrac{1}{1}$

$P(Zack \mid Michael\ and) = \dfrac{2}{2}$

$P(played \mid and\ zack) = \dfrac{1}{3}$

$P(at \mid Zack\ played) = \dfrac{1}{1}$

$P(the \mid played\ at) = \dfrac{1}{1}$

$P(playground \mid at\ the) = \dfrac{1}{2}$

$P(</s> \mid the\ playground) = \dfrac{1}{2}$

$P(<s>\ Michael\ and\ Zack\ played\ at\ the\ playground\ </s>)$

$= \dfrac{1}{5} \times \dfrac{1}{1} \times \dfrac{2}{2} \times \dfrac{1}{3} \times \dfrac{1}{1} \times \dfrac{1}{1} \times \dfrac{1}{2} \times \dfrac{1}{2}$

$= \dfrac{1}{60}$

Q3.What is the perplexity of the statement in Q2

Sol^n

$P(s) = $ Probability

$n = $ Total token (including $<s>$ and $</s>$)

Perplexity $= P(s)^{-1/n}$

$$= \left(\frac{1}{60}\right)^{-1/9}$$

$$= 1.576$$

Q4.For the given corpus,find the bigram probability estimate by laplace model for
P(do|<s>),P(do|Sam),P(Sam|<s>),P(Sam|do),P(I|Sam),P(I|do),P(like|I)

<S>I am Sam</S>

<S>Sam I am</S>

<S>Sam I like</S>

<S>Sam I do like</S>

<S>do I like Sam</S>

Sol^n

| Tokens | Frequency |
|--------|-----------|
| <s> | 5 |
| </s> | 5 |
| I | 5 |
| am | 2 |
| Sam | 5 |
| like | 3 |
| do | 2 |

Unique tokens = 7

$$P(do|<s>) = \frac{1+1}{5+7} = \frac{2}{12} = \frac{1}{6}$$

$$P(do|Sam) = \frac{0+1}{5+7} = \frac{1}{12}$$

$$P(Sam|<s>) = \frac{3+1}{5+7} = \frac{4}{12} = \frac{1}{3}$$

$$P(Sam|do) = \frac{0+1}{2+7} = \frac{1}{9}$$

$$P(I|Sam) = \frac{3+1}{5+7} = \frac{4}{12} = \frac{1}{3}$$

$$P(I|do) = \frac{1+1}{2+7} = \frac{2}{9}$$

$$P(like|I) = \frac{2+1}{5+7} = \frac{3}{12} = \frac{1}{4}$$

Q5. For a corpus , MLE for bigram "battery life" is 0.27, frequency of "battery" is 800. After applying Laplace smoothening the MLE for "battery life" becomes 0.025. What is the vocabulary size of the corpus?

Sol^n

$$P(w_i \mid w_{i-1}) = \frac{Count(w_{i-1}, w_i)}{Count(w_{i-1})}$$

$$P_{MLE}(Battery\ life) = \frac{f(battery\ life)}{f(battery)}$$

$$0.27 = \frac{f(battery\ life)}{800}$$

$$f(battery\ life) = 800 \times 0.27 = 216$$

With Laplace Smoothening

$$P_{MLE}(life \mid Battery) = \frac{f(battery\ life) + 1}{f(battery) + V}$$

$$0.025 = \frac{216 + 1}{800 + V}$$

$$\therefore \boxed{V = 7880}$$

Q7. Find the edit distance between INTENTION and EXECUTION using dynamic programming.

$$D(i,j) = \min \begin{cases} D(i-1,j)+1 \\ D(i,j-1)+1 \\ D(i-1,j-1) + \begin{cases} 2 & \text{if } S_1(i) \neq S_2(j) \\ 0 & \text{if } S_1(i) = S_2(j) \end{cases} \end{cases}$$

insert −1
delete −1
substitute −1 or 2

Edit distance

| | # | E | X | E | C | U | T | I | O | N | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| N | 9 | 8 | 9 | 10 | 11 | 12 | 11 | 10 | 9 | (8) | 1+1 / 1+1 |
| O | 8 | 7 | 8 | 9 | 10 | 11 | 10 | 9 | 8 | 9 | 0+2 |
| I | 7 | 6 | 7 | 8 | 9 | 10 | 9 | 8 | 9 | 10 | |
| T | 6 | 5 | 6 | 7 | 8 | 9 | 8 | 9 | 10 | 11 | 2+1 / 2+1 |
| N | 5 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 10 | 1+2 |
| E | 4 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 9 | |
| T | 3 | 4 | 5 | 6 | 7 | 8 | 7 | 8 | 9 | 8 | 7+1 |
| N | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 7 | 8 | 7 | 7+1 |
| I | (1) | 2 | 3 | 4 | 5 | 6 | 7 | 6 | 7 | 8 | 6+0 |
| # | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | |

$D(i,j-1)$

$D(i-1,j-1)$     $D(i-1,j)$

Q8 In the sentence-"The only thing we have to fear is fear itself" , the ratio between total number of word token and word types is:

Sol$^n$  Token = 10

Type = 9

Ratio of $\dfrac{Token}{Type} = \dfrac{10}{9}$

Q9 Let the rank of two words W1 and W2 in a corpus be 1600 and 400 respectively. Let m1 and m2 represent the number of meanings of w1 and w2 respectively. The ratio m1:m2 would be?

Sol$^n$  Zipf's law

$$M \propto \sqrt{F} \quad , \quad F \propto \dfrac{1}{rank}$$

$$M_1 = \dfrac{1}{\sqrt{r_1}} = \dfrac{1}{\sqrt{1600}} = \dfrac{1}{40}$$

$$\therefore M_1 : M_2 = 20 : 40$$

$$M_2 = \dfrac{1}{\sqrt{r_2}} = \dfrac{1}{\sqrt{400}} = \dfrac{1}{20}$$

$$\boxed{M_1 : M_2 = 1 : 2}$$

Q10 What is the size of unique words in a document where total number of words=1200 and k=3.71 , β=0.69

Sol$^n$

Heap's Law

$|v|$ = size of vocabulary = $k \cdot N^{\beta}$

$= 3.71(1200)^{0.69}$

$= 494.32$