

Name: Arroj

Quiz 1(NLP)

Roll Number: 21L-5619

Q1. Write down the Vocabulary, number of tokens and number of types in the following text:

23.5+1.5
3.2

Solution: [2+2+1 marks]

Vocabulary:

{

بلی کو سمجھانے آئے، چونچے، کٹی، ہزار
نے، اک، بات، نہ، مانی، روئے
زادہ زار، سنو، گپ، شپ، ناویش
ندی، ڈوب، چلی، شیر، اور، بکری،
مل، کڑ، بیٹھے، گھوڑا، گھاس، کھائے، تینوں
اپنی، ضد، کئے، پورے، کون، کسے، سمجھانے

Tokens:

42

Types:

38

بلی کو سمجھانے آئے

چونچے کٹی ہزار

بلی نے اک بات نہ مانی

روئے زار و زار

سنو گپ شپ سنو گپ شپ

ناویش ندی ڈوب چلی

شیر اور بکری مل کر بیٹھے

گھوڑا گھاس نہ کھائے

تینوں اپنی ضد کے پورے

کون کسے سمجھانے

consider
as 1
token
and
type

Q2. Assume you have the following training corpus: [10+5 Marks]

<s> I am from Vellore </s>

<s> I am a teacher </s>

<s> Students are good and are from various cities </s>

<s> Students from Vellore do engineering </s>

i) Find the Bigram probability of the given test sentence, including <s> & </s> as a token. ii) Compute the perplexity of this bigram.

Test data:

<s> Students are from Vellore cities </s>

Note: Use Laplace (Add-1) Smoothing if needed. Write down vocabulary as well if you use it.

Vocabulary: <s>, I, am, from, Vellore, a, teacher, Students,
are, good, and, ~~from~~, various, cities,
do, engineering, </s>

$|V| = 16$

(a)

Solution:

<S>Students are from Vellore cities </S>

We use Add-1 smoothing as.

$$P(\text{cities}|\text{Vellore}) = 0$$

$$\text{Formula: } P(w_i|w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + |V|}$$

$$\Rightarrow P(\text{students}|\text{<S>}) = \frac{2+1}{4+16} = \frac{3}{20}$$

$$\Rightarrow P(\text{are}|\text{students}) = \frac{1+1}{2+16} = \frac{2}{18}$$

$$\Rightarrow P(\text{from}|\text{are}) = \frac{1+1}{2+16} = \frac{2}{18}$$

$$\Rightarrow P(\text{Vellore}|\text{from}) = \frac{2+1}{3+16} = \frac{3}{19}$$

$$\Rightarrow P(\text{cities}|\text{Vellore}) = \frac{0+1}{2+16} = \frac{1}{18}$$

$$\Rightarrow P(\text{</S>}|\text{cities}) = \frac{1+1}{1+16} = \frac{2}{17}$$

$$P(\text{<S>Students are from Vellore cities </S>}) =$$

$$\left(\frac{3}{20}\right) \left(\frac{2}{18}\right) \left(\frac{2}{18}\right) \left(\frac{3}{19}\right) \left(\frac{1}{18}\right) \left(\frac{2}{17}\right)$$

$$= 1.91 \times 10^{-6}$$

$$(b) \text{ perplexity} = P(\text{sentence})^{-1/N}$$

$$= (1.91 \times 10^{-6})^{-1/N}$$

$$= 8.97$$

$$N = 6 + 1$$

<S></S>

Name: Aveej

Quiz 1(NLP)

Roll Number: 212-5619

Q3: Use byte pair encoding on training corpus to generate suitable vocabulary and apply it on our testing corpus and describe how well it segments our testing data. [3+2 Marks]

Training Corpus:

bat cat mat rat hat pat sat fat bat cat

Testing Set:

bat cat mat lat dat

Solution:

Vocabulary :- b, a, t, c, m, r, h, p,

2 bat-

2 cat-

1 mat-

1 rat-

1 hat-

1 pat-

1 sat-

1 fat-

s, f, at, at-, bat-, cat-

By merging
a and t

By merging
at and -

By merging
b and at-

By merging
c and at-

Testing set: let's assume we have:

bat-, cat-, mat-, lat-, dat-

It segment our testing data as:

bat as ~~cat as~~
~~bat-~~, ~~cat-~~, m ~~at-~~ > mat-

we don't have lat and dat- in
our vocabulary (even not by segmenting)

we don't have l and d in our
vocabulary only at-