# NLP Assignment - 1

Atharva Pawar - 9427

Atharva Prashant Pawar (9427) — [Batch-D]

## # NLP : Assignment 1 ①

**Q1.** The corpus is as follows:

[BOS] he is my friend [EOS]
[BOS] he is good [EOS]
[BOS] I like good friend [EOS]
[BOS] I like my friend [EOS]

given a sequence of words " [BOS] my friend is good [EOS]"
The vocabulary for the given corpus.

| Word | Frequency |
|------|-----------|
| [BOS] | 4 |
| [EOS] | 4 |
| he | 2 |
| is | 2 |
| my | 2 |
| good | 2 |
| friend | 3 |
| I | 2 |
| like | 2 |

**A)** Calculate the given sequence's probability using bi-gram model (LM1).

$P($ [BOS] my friend is good $)$
$= P($my $|$ [BOS]$) * P($friend $|$my$) * P($is $|$friend$) *$
$\quad P($good $|$ is$) * P($[EOS] $|$good$)$

$= \dfrac{0}{4} \times \dfrac{2}{2} \times \dfrac{0}{3} \times \dfrac{1}{2} \times \dfrac{1}{2} = 0$

**B)** Calculate the probability of the given sequence using tri-gram model (LM2).

$P($[BOS] my friend is good [EOS]$) = P($friend $|$[BOS], my$) *$
$\quad P($is $|$my, friend$) * P($good $|$ friend, is$) * P($[EOS] $|$ is, good$)$

$= \dfrac{0}{0} \times \dfrac{0}{2} \times \dfrac{0}{0} \times \dfrac{1}{1} = 0$

②

③ Calculate the probability of the given sequence using bi-gram model with Laplace smoothing (LM3)

Unique words: 8..(since [BOS] never comes in bigram calculation)

$P([BOS]$ my friend is good $[EOS]) = P(my | [BOS]) * P(friend | my) * P(is | friend) * P(good | is) * P([EOS] | good)$

$= \frac{0+1}{4+8} \times \frac{2+1}{2+8} \times \frac{0+1}{3\times8} \times \frac{1+1}{2+8} \times \frac{1+1}{4+8} = \frac{1}{12} \times \frac{3}{10} \times \frac{1}{11} \times \frac{2}{10} \times \frac{2}{12}$

$= 0.0000767.$

On comparing the result we can see that the bigram model with laplace smoothing (LM3) provides a non-zero probability for the given sequence.

∴ LM3 would be the best model among the three options of the given corpus

* Consider the following training data:

&lt;s&gt; I am Sam &lt;/s&gt;
&lt;s&gt; Sam I am &lt;/s&gt;
&lt;s&gt; Sam I like &lt;/s&gt;
&lt;s&gt; Sam I do like &lt;/s&gt;
&lt;s&gt; do I like Sam &lt;/s&gt;

① The vocabulary of the corpus is as follows:

| Word | Freq. |
|------|-------|
| &lt;s&gt; | 5 |
| &lt;/s&gt; | 5 |
| I | 5 |
| am | 2 |
| Sam | 5 |
| like | 3 |
| do | 2 |

Probability matrix for the bigram model:

|  | \<s\> | \</s\> | I | am | sam | like | do |
|---|---|---|---|---|---|---|---|
| \<s\> | 0 | 0 | 0.2 | 0 | 0.6 | 0 | 0.2 |
| \</s\> | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| I | 0 | 0 | 0 | 0.4 | 0 | 0.4 | 0.2 |
| am | 0 | 0.5 | 0 | 0 | 0.5 | 0 | 0 |
| sam | 0 | 0.4 | 0.6 | 0 | 0 | 0 | 0 |
| like | 0 | 0.66 | 0 | 0 | 0.33 | 0 | 0 |
| do | 0 | 0 | 0.5 | 0 | 0 | 0.5 | 0 |

② Predict the most probable word:

a. \<s\> sam ...
The word with the highest probability after sam ie 'I'
(Probability = 0.6)

b. \<s\> sam I do ...
The words with highest & equal probability after do
are 'I' & 'like' (Probability = 0.5).

c. \<s\> sam I am sam ...
The word with highest probability is 'I'
(Probability = 0.6)

d. \<s\> do I like
The most probable word with highest probability after
like is \</s\> (Probability = 0.66).

③ Which of the following sentence is probable:

a. \<s\> Sam I do I like \</s\>
P(\<s\> sam I do I like \</s\> = P (sam | \<s\>) * P(I | sam
    * P(do | I) * P (I | do) * P (like | I) * P (\</s\> | like)
= 0.6 * 0.6 * 0.2 * 0.5 * 0.4 * 0.66
= 0.0095.

b. $\langle s \rangle$ Sam I am $\langle /s \rangle$

$P(\langle s \rangle$ Sam I am $\langle /s \rangle) = P(\text{sam} | \langle s \rangle) * P(I | \text{sam})$

$\quad * P(\text{am} | I) * P(\langle /s \rangle | \text{am})$

$= 0.6 * 0.6 * 0.4 * 0.5$

$= 0.072.$

c. $\langle s \rangle$ I do like sam I am $\langle /s \rangle$

$P(\langle s \rangle$ I do like sam I am $\langle /s \rangle)$

$= P(I | \langle s \rangle) * P(do | I) * P(like | do) * P(sam | like)$

$\quad * P(I | sam) * P(am | I) * P(\langle /s \rangle | am)$

$= 0.2 * 0.2 * 0.5 * 0.33 * 0.6 * 0.4 * 0.5$

$= 0.00079$

$\therefore$ Most Probable sentence is $\langle s \rangle$ Sam I am $\langle /s \rangle$

Q3) Use the same corpus as given in Q2. This time we a bi-gram LM with Laplace Smoothing (Unique words: 6)

1) a. $P(do | I) = \dfrac{1+1}{5+6} = \dfrac{2}{11} = 0.181.$

b. $P(do | sam) = \dfrac{0+1}{5+6} = \dfrac{1}{11} = 0.091$

c. $P(sam | \langle s \rangle) = \dfrac{3+1}{5+6} = \dfrac{4}{11} = 0.364$

d. $P(sam | do) = \dfrac{0+1}{2+6} = \dfrac{1}{8} = 0.125$

e. $P(I | sam) = \dfrac{3+1}{5+6} = \dfrac{4}{11} = 0.36$

f. $P(I | do) = \dfrac{1+1}{2+6} = \dfrac{2}{8} = 0.25$

g. $P(like | I) = \dfrac{2+1}{5+6} = \dfrac{3}{11} = 0.273$

2) &lt;s&gt; do sam I like.

P( &lt;s&gt; do sam I like) = P (do|&lt;s&gt;) * P (sam|do)
        * P(I|sam) * P (like |I)

$$= \frac{1+1}{5+6} * \frac{0+1}{2+6} * \frac{3+1}{5+6} * \frac{2+1}{5+6} = 0.00225$$

&lt;s&gt; Sam do I like

P( &lt;s&gt; Sam do I like.) = P (sam |&lt;s&gt;) * P (do|sam)
        * P(I|do) * P(like |I)

$$= \frac{3+1}{5+6} * \frac{0+1}{2+6} * \frac{1+1}{2+6} * \frac{2+1}{5+6} = 0.00226$$

∴ Most probable sentence is &lt;s&gt; Sam do I like.

Q4) Apply Levensthten minimum edit distance algo. to find the similarity b/w two words 'Honda' & ~~Hyana~~ 'hundai'

| | # | h | u | n | d | a | i |
|---|---|---|---|---|---|---|---|
| # | 0→ | 1→ | 2→ | 3→ | 4→ | 5→ | 6 |
| h | 1↓ | 0→ | 1→ | 2→ | 3→ | 4→ | 5 |
| o | 2↓ | 1↓ | 1→ | 2→ | 3→ | 4→ | 5 |
| n | 3↓ | 2↓ | 2↓ | 1→ | 2→ | 3→ | 4 |
| d | 4↓ | 3↓ | 3↓ | 2↓ | 1→ | 2→ | 3 |
| a | 5 | 4 | 4 | 3 | 2 | 1→ | 2 |

∴ The minimum edit distance is 2

```
h ⊘  o   n   d   a   *
|    |R  |   |   |   |I
h    u   n   d   a   i
```

* Operations :
    R(o, u) ; I(i)