**Q.1]** Backoff vs Interpolation

**(A) Backoff**

- Backoff N-gram modelling is a non-linear method
- We build on N gram model based on (N-1) gram model
- The difference is that in backoff if we have non-zero trigram counts we solely rely on trigram counts & don't interpolate the bigram and unigram counts at all.
- Backoff model in trigram format:

$$P(w_i \mid w_{i-2}\, w_{i-1}) = \begin{cases} \tilde{P}(w_i \mid w_{i-2}\, w_{i-1}) & \text{if } C(w_{i-2}\, w_{i-1}\, w_i) > 0 \\ \alpha(w_{n-2}^{n-1})\, \tilde{P}(w_i \mid w_{i-1}) & \text{if } C(w_{i-2}\, w_{i-1}\, w_i) = 0 \\ & \quad \& \; C(w_{i-1}\, w_i) > 0 \\ \alpha(w_{n-1})\cdot \tilde{P}(w_i) & \text{otherwise} \end{cases}$$

- Doesn't yield valid probability distribution
- Works well for large datasets.


**(B) Interpolation**

- Combines different N-grams by linearly interpolating all 3 models whenever we are computing any trigram.
- Here, we don't train 3 $\lambda$'s as trigram grammar. Instead we make each $\lambda$ a function of the context.
- $\lambda$ terms are used to decide how much to smooth
- $\sum_i \lambda_i = 1$
- Mathematically, $\tilde{P}(w_0 \mid w_{-2}\, w_{-1}) = \lambda_3\, P(w_0 \mid w_{-2}\, w_{-1}) + \lambda_2\, P(w_0 \mid w_{-1})$
  $\qquad\qquad\qquad + \lambda_1 \cdot P(w_0)$
- Can interpolate 'customised' model with general model


**Q.2]** Viterbi algorithm

- It is a variation of the forward algorithm which considers all words simultaneously in order to compute the most likely path.

**Algorithm:**

Input: observations of length T, state-graph of length N

Output: best path

for each state s from 1 to N do

$q[1,s] \leftarrow P(s|s_0) \cdot P(o_1|s)$

backpointers $[1,s] \leftarrow 0$

for each time step t from 2 to T do

for each state s from 1 to N do

$q[t,s] \leftarrow \max\limits_{s'=1}^{N} q[t-1, s'] \cdot P(s|s') \cdot P(o_t|s)$

backpointers $[t,s] \leftarrow \text{argmax}\limits_{s'=1}^{N} q[t-1, s'] \cdot P(s|s')$

$s \leftarrow \text{argmax}\limits_{s'=1}^{N} q[T, s']$

return the backtrace path from the backpointers $[T, s]$

**Example:** Consider a 2 word language: 'fish' & 'sleep'.

Suppose in our training corpus,

'fish' appears 8 times as a noun & 5 times as a verb

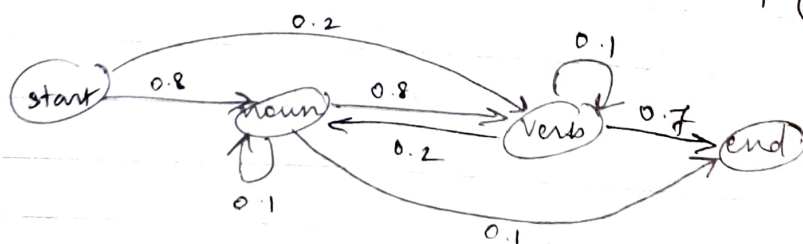'sleep' " 2 " " " & 5 " " " "

⇒ Emission probabilities

○ Noun

— $P(\text{fish} | \text{noun}) = 0.8$

— $P(\text{sleep} | \text{noun}) = 0.2$

○ Verb)

— $P(\text{fish} | \text{verb}) = 0.5$

— $P(\text{fish sleep} | \text{verb}) = 0.5$



0.2

start  0.8

noun  0.8

0.1

Verb  0.7  end

0.2

0.1

0.1

# Token 1: fish

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| start | 1 | 0 | | |
| verb | 0 | $0.2 \times 0.5$ | | |
| noun | 0 | $0.8 \times 0.8$ | | |
| end | 0 | 0 | | |

# Token 2: sleep

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| start | 1 | 0 | 0 | |
| verb | 0 | 0.1 | $0.64 \times 0.8 \times 0.5$ ←max ✓ / $0.1 \times 0.1 \times 0.5$ | |
| noun | 0 | 0.64 | $0.64 \times 0.1 \times 0.2$ ←max ✓ / $0.1 \times 0.2 \times 0.2$ | |
| end | 0 | 0 | . — | |

# Token 3: end

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| start | 1 | 0 | 0 | 0 |
| verb | 0 | 0.1 | 0.256 | - |
| noun | 0 | 0.64 | 0.0128 | - |
| end | 0 | 0 | $0.256 \times 0.7$ ←max ✓ / $0.0128 \times 0.1$ | |

∴ now we can backtrack the most likely path

Q.3] Corpus

&lt;s&gt; I am from DJ &lt;/s&gt;

&lt;s&gt; I am a teacher &lt;/s&gt;

&lt;s&gt; All students are good and intelligent &lt;/s&gt;

&lt;s&gt; Students from DJ Score high marks &lt;/s&gt;

Test data

&lt;s&gt; Students are from DJ &lt;/s&gt;

- Unigram

| &lt;s&gt; | students | are | from | DJ | &lt;/s&gt; |
|---|---|---|---|---|---|
| 4 | 2 | 1 | 2 | 2 | 4 |

# Bigram occurrence count

| | <s> | students | are | from | DJ | </s> |
|---|---|---|---|---|---|---|
| <s> | 0 | 1 | 0 | 0 | 0 | 0 |
| students | 0 | 0 | 1 | 1 | 0 | 0 |
| are | 0 | 0 | 0 | 0 | 0 | 0 |
| from | 0 | 0 | 0 | 0 | 2 | 0 |
| DJ | 0 | 0 | 0 | 0 | 0 | 1 |
| </s> | 0 | 0 | 0 | 0 | 0 | 0 |

| | <s> | students | are | from | DJ | </s> |
|---|---|---|---|---|---|---|
| <s> | 0 | 1/4 | 0 | 0 | 0 | 0 |
| students | 0 | 0 | 1/2 | 1/2 | 0 | 0 |
| are | 0 | 0 | 0 | 0 | 0 | 0 |
| from | 0 | 0 | 0 | 0 | 2/2¹ | 0 |
| DJ | 0 | 0 | 0 | 0 | 0 | 1/2 |
| </s> | 0 | 0 | 0 | 0 | 0 | 0 |

Using MLE to estimate probability of test data

$$P = P(\text{students} \mid \text{<s>}) \cdot P(\text{are} \mid \text{students}) \cdot P(\text{from} \mid \text{are}) \cdot$$
$$P(\text{DJ} \mid \text{from}) \cdot P(\text{</s>} \mid \text{DJ})$$

$$= \frac{1}{4} \times \frac{1}{2} \times 0 \times 1 \times \frac{1}{2}$$

hence we need to apply
Laplace smoothing

To apply laplace smoothing

V = count of unique vocabulary

= count({ <s>, </s>, I, am, from, DJ, a, teacher, all, students, are, good, and, intelligent, score, high, marks })

= 17

$$\therefore P = \left(\frac{1+1}{4+17}\right) \cdot \left(\frac{1+1}{2+17}\right) \cdot \left(\frac{0+1}{1+17}\right) \cdot \left(\frac{2+1}{2+17}\right) \cdot \left(\frac{1+1}{2+17}\right) = 9.257 \times 10^{?}$$

**Q4]** Corpus

> $\langle s \rangle$ I am Sam $\langle /s \rangle$
>
> $\langle s \rangle$ Sam I am $\langle /s \rangle$
>
> $\langle s \rangle$ I do not like green eggs & ham $\langle /s \rangle$

(a) Bigram probability (i) P(am | Sam) (ii) P(do | I)

(iii) P(am | I)

- $P(w_n | w_{n-1}) = \dfrac{C(w_{n-1}, w_n)}{C(w_{n-1})}$

(i) $P(am | Sam) = \dfrac{P(Sam\ am)}{P(am)} = \dfrac{0}{2} = 0$

(ii) $P(do | I) = \dfrac{P(I\ do)}{P(I)} = \dfrac{1}{3}$

(iii) $P(am | I) = \dfrac{P(I\ am)}{P(I)} = \dfrac{2}{3}$

(b) Trigram probability 'I am Sam'

- $P(w_n | w_{n-2}\ w_{n-1}) = \dfrac{C(w_{n-2}\ w_{n-1}\ w_n)}{C(w_{n-2}\ w_{n-1})}$

$P(Sam | I\ am) = \dfrac{C(I\ am\ Sam)}{C(I\ am)} = \dfrac{1}{2}$

(c) MLE for 'I am Sam' using bigram

- $(\langle s \rangle, I), (I, am), (am, Sam), (Sam, \langle /s \rangle)$

MLE = $P(I | \langle s \rangle) \cdot P(am | I) \cdot P(Sam | am) \cdot P(\langle /s \rangle | Sam)$

$= \dfrac{2}{3} \times \dfrac{2}{3} \times \dfrac{1}{2} \times \dfrac{1}{2} = \dfrac{1}{9}$

Q.5] Corpus

    \<s\> John read Moby Dick \</s\>

    \<s\> Mary read a different book \</s\>

    \<s\> she read a book by Cher \</s\>

(a) MLE for 'John read a book'.

(\<s\>, John), (John, read), (read, a), (a, book), (book,

$MLE = P(John | \<s\>) \cdot P(read | John) \cdot P(a | read) \cdot P(book)$

$P(\</s\> | book)$

$= \frac{1}{3} \times \frac{1}{1} \times \frac{2}{3} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{18} = 0.056$

(b) MLE for 'Cher read a book'

(\<s\>, cher), (cher, read), (read, a), (a, book), (book, \</s\>)

$MLE = P(Cher | \<s\>) \cdot P(read | cher) \cdot P(a | read) \cdot P(book | a)$

$P(\</s\> | book)$

$= \frac{0}{3} \times \frac{0}{1} \times \frac{2}{3} \times \frac{1}{2} \times \frac{1}{2}$

Using add-one smoothing (Laplace)
Total unique tokens = 11

$\left( \frac{Count + 1}{Count + V} \right)$

$= \frac{0+1}{3+11} \cdot \frac{0+1}{1+11} \cdot \frac{2+1}{3+11} \cdot \frac{1+1}{2+11} \cdot \frac{1+1}{2+11}$

$= \frac{1}{14} \cdot \frac{1}{12} \cdot \frac{3}{14} \cdot \frac{2}{13} \cdot \frac{2}{13} = 3.019 \times 10^{-5}$