NLP

Assignment - 1

**Q1]** Differentiate between Interpolation and Backoff.

**Ans:** <u>Backoff</u>

- 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 non-zero trigram counts we solely rely on trigram counts and don't interpolate the bigram and unigram counts at all.
- Backoff model in trigram format:

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

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

<u>Interpolation</u>

- This method 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,

$$\tilde{P}(w_0 \mid w_{-2}\, w_{-1}) = \lambda_3 \cdot P(w_0 \mid w_{-2}\, w_{-1})$$
$$+ \lambda_2 \cdot P(w_0 \mid w_{-1})$$
$$+ \lambda_1 \cdot P(w_0)$$

- Can interpolate 'customised' models with 'general' model.

**Q2]** Viterbi algorithm.

Ans: Viterbi algorithm 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

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

$$\text{\& backpointers } [1, s] \leftarrow 0$$

for each time step t from 2 to T do

for each state s from 1 to N do

$$q[t, s] \leftarrow \max_{s'=1}^{N} q[t-1, s'] \cdot p(s|s') \cdot p(o_t/s)$$

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

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

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

**Example:**

Consider a two word language: 'fish' and 'sleep'

Suppose in our training corpus,

'fish' appears 8 times as a noun and 5 times as a verb.

'sleep' appears twice as a noun and 5 times as a verb
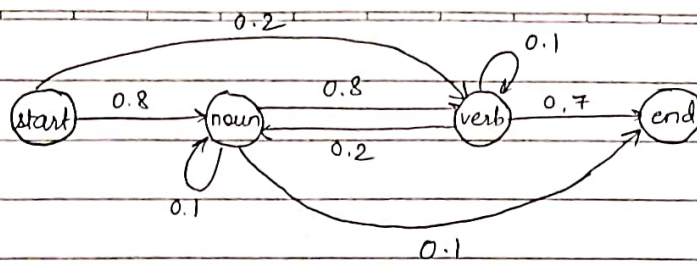
∴ Emission probabilities:

- Noun
  - P(fish | noun): 0.8
  - P(sleep | noun): 0.2

- Verb
  - P(fish | verb): 0.5
  - P(sleep | verb): 0.5

Khushi K. Shah
60004180045
BE - A

**Token 1 : fish**

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| start | 1 | 0 | | |
| verb | 0 | 0.2 × 0.5 | | |
| noun | 0 | 0.8 × 0.8 | | |
| end | 0 | 0 | | |

**Token 2 : sleep**

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| start | 1 | 0 | 0 | |
| verb | 0 | 0.1 | 0.64 × 0.8 × 0.5 ← max ✓ <br> 0.1 × 0.1 × 0.5 | |
| noun | 0 | 0.64 | 0.64 × 0.1 × 0.2 ← max ✓ <br> 0.1 × 0.2 × 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 × 0.7 ← max ✓ <br> 0.0128 × 0.1 |

∴ Now we can backtrack the most likely path

Khushi K. Shah
6000418 0045
BE-A

④

**Q3]** Corpus:

        <s> I am from DJ </s>

        <s> I am a teacher </s>

        <s> All students are good and intelligent </s>

        <s> Students from DJ score high marks </s>

Test data

        <s> students are from DJ </s>

**Ans:** Unigram

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

Bigram

First we find 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 |

Bigram

| | <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 | 2/2 = 1 | 0 |
| DJ | 0 | 0 | 0 | 0 | 0 | 1/2 |
| </s> | 0 | 0 | 0 | 0 | 0 | 0 |

Using MLE to estimate probability for test data.

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

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

↑
hence we need to
apply laplace smoothening.

before applying laplace, we find
   V = count of unique vocabulary in corpus
      = count ({ <s>, </s>, I, am, from, DJ, a, teacher, all,
            students, are, good, and, intelligent, score,
            high, marks })
   = 17

$\therefore P = \left(\dfrac{1+1}{4+17}\right) \times \left(\dfrac{1+1}{2+17}\right) \times \left(\dfrac{0+1}{1+17}\right) \times \left(\dfrac{2+1}{2+17}\right) \times \left(\dfrac{1+1}{2+17}\right)$

$= \dfrac{2}{21} \times \dfrac{2}{19} \times \dfrac{1}{\underset{63}{\cancel{18}}} \times \dfrac{3}{19} \times \dfrac{2}{19} = \boxed{9.257 \times 10^{-6}}$

**Q4]** Corpus:

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

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

 &lt;s&gt; I do not like green eggs and ham &lt;/s&gt;

a) Calculate bigram probability for ① $P(am | Sam)$   ② $P(do | I)$

 ③ $P(am | I)$

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

 ① $P(am | Sam) = \dfrac{P(Sam \, am)}{P(Sam)} = \dfrac{0}{2} = 0$

 ② $P(do | I) = \dfrac{P(I \, do)}{P(I)} = \dfrac{1}{3}$

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

b) Calculate the trigram probability for 'I am Sam'

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

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

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

**soln:** $(\text{&lt;s&gt;}, I), (I, am), (am, Sam), (Sam, \text{&lt;/s&gt;})$

MLE $= P(I | \text{&lt;s&gt;}) \times P(am | I) \times P(Sam | am) \times P(\text{&lt;/s&gt;} | Sam)$

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

Q5] Corpus:

&lt;s&gt; John read Moby Dick &lt;/s&gt;

&lt;s&gt; Mary read a different book &lt;/s&gt;

&lt;s&gt; She read a book by Cher &lt;/s&gt;

a) Calculate MLE for 'John read a book' using bigram

b) Calculate MLE for 'Cher read a book' using bigram

soln: MLE for 'John read a book'

$(\langle s \rangle, John), (John, read), (read, a), (a, book), (book, \langle /s \rangle)$

$MLE = P(John|\langle s \rangle) \times P(read|John) \times P(a|read) \times P(book|a) \times P(\langle /s \rangle|book)$

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

MLE for 'Cher read a book'

$(\langle s \rangle, Cher), (Cher, read), (read, a), (a, book) (book, \langle /s \rangle)$

$MLE = P(Cher|\langle s \rangle) \times P(read|Cher) \times P(a|read) \times P(book|a) \times P(\langle /s \rangle|book)$

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

using add-one smoothing (laplace)

$\frac{count+1}{count+v}$

total no. of unique tokens = v = 11

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

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