

INTRODUCTION TO NLP
ASSIGNMENT 1

REPORT

February 5, 2024

Ashna Dua
2021101072
`ashna.dua@students.iiit.ac.in`

1 Tokenization

The first part of the assignment includes designing a tokenizer using regex. The **Tokenizer** class provides several methods for preprocessing the corpus, replacing certain entities with special tokens, and splitting the text into sentences and words.

The following entity replacements have been implemented:

- URL
- MAILID
- DATE
- TIME
- AGE
- PERCENTAGE
- CURRENCY
- NUM
- MENTION
- HASHTAG

2 N-Grams

This part of the assignment includes designing a class that takes a value of n, and the corpus path to generate n-sized ngram model.

To obtain the probability of the model, the following formula is used:

$$P(ngram|context) = \frac{\text{count}(\text{context} + \text{ngram})}{\text{count}(\text{context})} \quad (1)$$

3 Good Turing Smoothing

In this part of the assignment, smoothing technique, i.e. **Good Turing Smoothing** is implemented with a **3-Gram** model.

3.1 Explanation

Good Turing Smoothing adjusts the distribution of given frequencies to estimate missing probabilities. This involves multiple steps, which are as follows:

3.1.1 Frequency of Frequency

In this step, a structure of frequency of frequency is defined, which is as follows:

$$N_r = \text{count}(w_1, w_2, \dots, w_n \mid \text{freq}(w_1, w_2, \dots, w_n) = r) \quad (2)$$

3.1.2 Calculation of Z

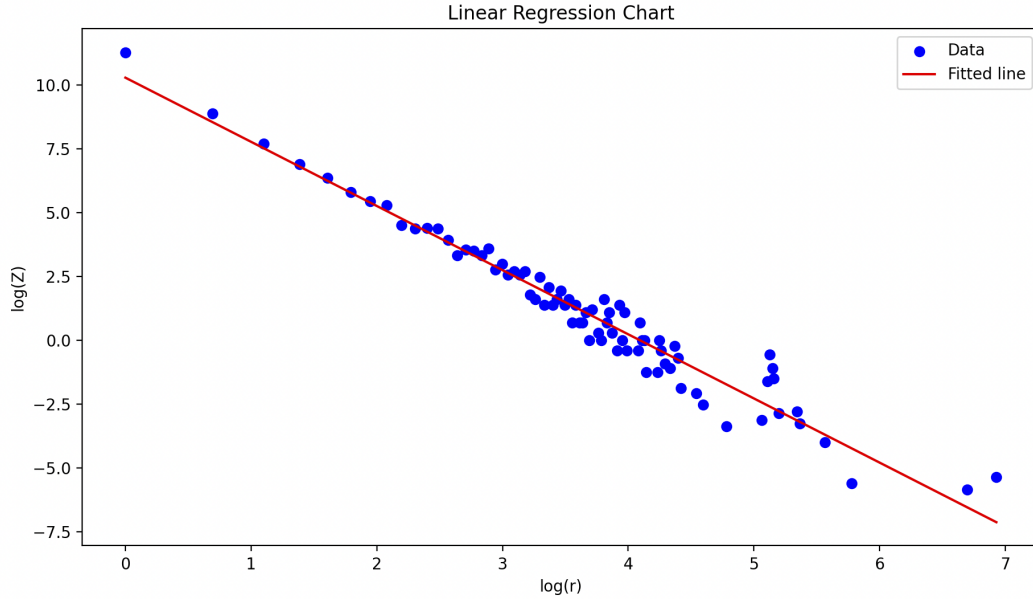
For smoothing the erratic values in N_r for large r , we make a log-log plot of N_r vs r . This causes an issue because for large r , many N_r will be 0. Therefore, we define Z_r .

$$Z_r = \frac{N_r}{0.5 \cdot (t - q)} \quad (3)$$

For where q , r , and t are three consecutive subscripts with non-zero counts N_q , N_r , N_t . For the special case when r is 1, take q to be 0. In the opposite special case, when $r = r_{\text{last}}$ is the index of the last non-zero count, replace the divisor $\frac{1}{2}(t - q)$ with $r_{\text{last}} - q$, so $Z_{r_{\text{last}}} = \frac{N_{r_{\text{last}}}}{r_{\text{last}} - q}$.

$$Z_{r_{\text{last}}} = \frac{N_{r_{\text{last}}}}{r_{\text{last}} - q} \quad (4)$$

Linear Regression is performed over the log-log domain of Z_r and r .



3.1.3 Calculation of S

To calculate the probabilities by assuming that expected occurrences for r , i.e. $\text{count}^*(r)$ is similar to the expected occurrences for $r+1$, i.e. $\text{count}^*(r+1)$.

$$S_r = \frac{(r+1) \cdot S_{n_{r+1}}}{S_{n_r}} \quad (5)$$

Here, S_{n_r} is taken from the regression line, and S_r represents the adjusted count.

3.1.4 Calculation of Probability

To calculate the probability of a token given its context, we use the following formula:

$$P(\text{token}|\text{context}) = \frac{S(\text{context} + \text{token})}{\sum_{t \in V} S(\text{context} + t)} \quad (6)$$

4 Linear Interpolation

In this part of the assignment, smoothing technique, i.e. **Linear Interpolation** is implemented with a **3-Gram** model.

4.1 Explanation

To calculate the $P(\text{token} \text{ given context})$, Linear interpolation uses backoff techniques to give a weighted probability using lower degree models.

$$P(w_3|w_1w_2) = \lambda_1 \hat{P}(w_3) + \lambda_2 \hat{P}(w_3|w_2) + \lambda_3 \hat{P}(w_3|w_1w_2) \quad (7)$$

Also,

$$\lambda_1 + \lambda_2 + \lambda_3 = 1 \quad (8)$$

To find the values of lamdas, we used the following algorithm:

```
set  $\lambda_1 = \lambda_2 = \lambda_3 = 0$ 
foreach trigram  $t_1, t_2, t_3$  with  $f(t_1, t_2, t_3) > 0$ 
    depending on the maximum of the following three values:
        case  $\frac{f(t_1, t_2, t_3)-1}{f(t_1, t_2)-1}$ : increment  $\lambda_3$  by  $f(t_1, t_2, t_3)$ 
        case  $\frac{f(t_2, t_3)-1}{f(t_2)-1}$ : increment  $\lambda_2$  by  $f(t_1, t_2, t_3)$ 
        case  $\frac{f(t_3)-1}{N-1}$ : increment  $\lambda_1$  by  $f(t_1, t_2, t_3)$ 
    end
end
normalize  $\lambda_1, \lambda_2, \lambda_3$ 
```

5 Generation

In this part of the assignment, generation was done for Simple N-Gram models, and Linear Interpolation only, because the probability assigned to unseen n-grams in Good Turing is very, causing the model to generate random values.

The first run is to obtain non-zero probabilities from the n degree model. If no such non-zero probability is received, the model is backed off to a smaller degree model. This is done only in the case of N-Gram model, as Linear Interpolation handles OOV tokens on its own.

6 Results

6.1 Smoothing Results

The following perplexities are received for LM1, LM2, LM3 and LM4 on the train and test corpus:

	Train	Test
LM1	389.19	998.40
LM2	57.43	600.26
LM3	1023.28	2826.79
LM4	234.23	4357.01

6.2 Generation Results

6.2.1 N-Gram Model

In the n -gram model, the probability of a word depends on the previous n words. The generator starts with the highest order n -gram (the context) and falls back to lower order n -grams if it can't find a suitable word to follow the context.

Generally, higher order-ngrams (larger n) can capture more context and produce more fluent text, but also require more data and are prone to overfitting.

The following image shows the result of 5 n -gram models, i.e. 1, 2, 3, 4, and 5. Unigram and Bigram models produce results that are random in nature and not fluent, where trigram, onwards the results make sense and are more fluent in comparison to the lower n -gram models.

```

input sentence:
it was not to be
1
output:
('be',)
('so',) 0.03423566878980892
('in',) 0.03264331210191083
('a',) 0.02945859872611465
2
output:
('to', 'be')
('sure',) 0.038636363636363635
('in',) 0.03636363636363636
('a',) 0.031818181818181815
3
output:
('not', 'to', 'be')
('supposed',) 0.22727272727272727
('a',) 0.09090909090909091
('the',) 0.045454545454545456
4
output:
('was', 'not', 'to', 'be')
('supposed',) 0.5
('.',) 0.125
('had',) 0.125
5
output:
('it', 'was', 'not', 'to', 'be')
('supposed',) 0.8
('.',) 0.2

```

In an Out-of-Data (OOD) scenario, where the context or some words in the context have not been seen in the training data, n-gram models can struggle. This is because n-gram models base their predictions on the frequency of occurrence of n-grams in the training data. If an n-gram has not been seen before, the model will assign it a probability of zero, which can cause issues.

To handle OOD scenarios, n-gram models often use smoothing techniques. These techniques assign a small amount of probability mass to unseen n-grams to ensure that every possible n-gram has a non-zero probability.

For the case of generation, back off technique has been used, where the model falls back to a lower-order n-gram model when it encounters an unseen n-gram.

Also, if no such ngrams are found the EOS token is given a very low probability to mark the end of generation.

6.2.2 Linear Interpolation Model

In the Linear Interpolation model, the probability of a word is calculated as a weighted sum of the probabilities from the unigram, bigram, and trigram models. This approach helps to balance the strengths of different n-gram models.

Linear interpolation can help to mitigate the data sparsity problem that often occurs with higher-order n-gram models. By giving some weight to lower-order n-grams, it ensures that every word in the vocabulary has a non-zero probability, even if it doesn't appear in the context in the training data.

This provides a balance between the context captured by higher-order n-grams and the robustness of lower-order n-grams, leading to more fluent and sensible text generation.

```
input sentence:
it is known
3
output:
('after',) 0.3210541203609964
(',',) 0.08881653081436887
('to',) 0.06408410519791823
```

Figure 1: Output of Linear Interpolation - Pride and Prejudice

```
input sentence:
it is known
3
output:
('by',) 0.2928419131646954
('to',) 0.10615732159242794
('.',) 0.06521614275252895
```

Figure 2: Output of Linear Interpolation - Ulysses

Both the corpus provide a different top prediction with different probabilities. Following are the observations:

- Corpus Influence: The 'Ulysses' model predicts 'by' as the most likely next word, while the 'Pride and Prejudice' model predicts 'after'. This reflects the differences in the writing styles of the two books.
- Probability Distributions: The probabilities assigned to the top predictions are also different in the two models. This is due to differences in the frequency of these word sequences in the respective corpora.