# Assignment 1

| Group 22 | Matthew Evans | Ruochen Meng | Saul Garay |
|---|---|---|---|
| | mre042000 | rxm220120 | SAG180009 |

## 1 Implementation Details

The inline documented code can be found at the GitHub repository given above. The experiments can be run from the entry point in `Main.py`.

### 1.1 Overview

The model builds a vocabulary from the training corpus, and computes n-gram probabilities with add-k smoothing (laplace for k=1). Unknown words are handled by replacing low-frequency words with a special `<UNK>` token. The model computes perplexity on the test corpus using the smoothed probabilities.

### 1.2 Data Structures

The model uses type aliases for clarity:

- `Token` - a string representing a word or symbol.

- `NGram` - a tuple of Tokens representing an n-gram.

- `Dataset` - a list of strings, each string is a line from the corpus.

Further, the follow data structures are used:

- `Vocabulary: set[Token]` - the set of known tokens.

- `Token Counts: Dict[Token, int]` - maps each token to its occurrence count.

- `Probability: float` - represents the probability of an n-gram.

### 1.3 Vocabulary Construction

The vocabulary is constructed by iterating over the training corpus and counting the occurrences of each token.

#### 1.3.1 Unknown word handling

Tokens which occur less than a specified threshold are replaced with the `<UNK>` token to handle unknown words in the test set. Specifically, probability mass of tokens is accumulated until a specified coverage threshold is reached, at which point the remaining tokens are replaced with `<UNK>`. The steps are as follows:

1. Sort tokens in descending order by occurrence count.

2. Iterate over token counts, accumulating mass, $M$, i.e.,

$$\mathrm{M_{cum}} = \sum_{t_i}^{t_n} \mathrm{C}(t_i)$$

If $(\mathrm{M_{cum}}/M_{\mathrm{total}}) < \mathrm{M_{coverage}}$, assign the token's mass to the `<UNK>` token.

### 1.3.2 Implementation

Our implementation, given in 1, computes the vocabulary and applies the unknown token handling as described.

```python
def _get_vocabulary(data: Dataset, coverage: float) -> set[Token]:
    token_counts: dict[str, int] = dict()
    for line in data:
        for token in NGramLanguageModel._parse_tokens(line):
            token_counts[token] = token_counts.get(token, 0) + 1

    token_counts = NGramLanguageModel.apply_unk_by_coverage(token_counts,
                                                             coverage)

    vocabulary = set(token_counts.keys())
    return vocabulary

def apply_unk_by_coverage(token_counts, coverage):
    total_mass = sum(token_counts.values())
    sorted_token_counts = sorted(token_counts.items(),
                                 key=lambda x: x[1],
                                 reverse=True)

    updated_token_counts: dict[str, int] = dict()
    cumulative_mass = 0
    unknown_count = 0
    for token, count in sorted_token_counts:
        if cumulative_mass / total_mass < coverage:
            updated_token_counts[token] = count
            cumulative_mass += count
        else:
            unknown_count += count
    updated_token_counts[UNKNOWN_TOKEN] = unknown_count
    return updated_token_counts
```

Listing 1: Vocabulary construction and unknown token handling.

## 1.4 Probability Computation and Smoothing

The model computes n-gram probabilities using add-k smoothing, supporting $k \geq 0$. The probability of an n-gram is computed as

$$P(w_n|w_{n-1},\ldots,w_1) = \frac{C(w_1,\ldots,w_n) + k}{C(w_1,\ldots,w_{n-1}) + kV}$$

where $C(w_1,\ldots,w_n)$ is the count of the n-gram, $C(w_1,\ldots,w_{n-1})$ is the count of the (n-1)-gram context, $V$ is the vocabulary size, and $k$ is the smoothing parameter.

### 1.4.1 Implementation

Our implementation, given in 2, handles both unigram and higher-order n-grams, with special handling for the unigram case where the denominator is the total token count plus $kV$. The add-k smoothing is captured by the parameter `smoothing`.

```python
def get_probability(self, *tokens: Token) -> float:
    vocabulary_size = len(self.vocabulary)

    # handle unigram case
    if self.n == 1:
        a = (self.ngram_counts.get(tokens, 0) + self.smoothing)
        b = (self.all_tokens_count + (self.smoothing * vocabulary_size))
        return a / b if b > 0 else 0.0

    # handle ngram case
    else:
        preceding = tokens[:-1]

        # Count of the full n-gram tokens
        a = self.ngram_counts.get(tokens, 0) + self.smoothing

        # Count of the (n-1) preceding tokens
        b = self.context_counts.get(preceding, 0)
            + (self.smoothing * vocabulary_size)

        return a / b if b > 0 else 0.0
```

Listing 2: Probability calculation function with add-k smoothing.

## 1.5 Implementation of perplexity

Perplexity is computed as the exponentiation of the negative average log-probability of the test corpus. The formula for perplexity is given by

$$PP(W) = \exp\left(-\frac{1}{N}\sum_{i=1}^{N}\log P(w_i|w_{i-1},\ldots,w_{i-n+1})\right)$$

where $W$ is the test corpus, $N$ is the total number of tokens in the test set, and $P(w_i|w_{i-1},\ldots,w_{i-n+1})$ is the probability of token $w_i$ given its preceding $n-1$ tokens.

### 1.5.1 Implementation

The implementation, given in 3, iterates over each line in the test dataset, computes the log-probability of each n-gram, and accumulates these, taking their mean, to compute the overall perplexity.

```python
import numpy as np
def get_perplexity(self, datum: str) -> float:
    ngrams: list[NGram] = []

    ngrams.extend(NGramLanguageModel._get_ngrams_from_line(datum,
                                                           self.vocabulary,
                                                           self.n))
    nll = - np.log([self.get_probability(*tokens) for tokens in ngrams])
    pp: np.float64 = np.exp(np.mean(nll))

    return pp

def get_mean_perplexity(self, data: Dataset) -> float:
    pps: list[float] = []
    for line in data:
        pp = self.get_perplexity(line)
        pps.append(pp)
    return float(np.mean(pps))
```

Listing 3: Perplexity calculation.

# 2 Evaluation, Analysis and Findings

# 3 Miscellaneous

## 3.1 Python Library Usage

The implementation uses minimal external libraries, primarily relying on Python's standard library. The only external library used is NumPy, which is employed for numerical computations, specifically for calculating logarithms and exponentials in the perplexity calculation. The project was written in Python 3.13.7.

## 3.2 Feedback

- The coding portion was completed within approximately one day, with a few issues resolved via pull requests over the following week or so.

- The assignment was simple enough to complete with minimal experience in NLP, while still fostering practical skills.

- Providing a target perplexity as a benchmark would enhance the clarity of expectations.

- Incorporating automated grading with the option for unlimited submissions would be a valuable addition, allowing students to iteratively refine their implementations based on immediate feedback.

## 3.3 Contributions

- Matthew Evans: Initial project setup, vocabulary construction, unknown token handling, probability computation, smoothing, and perplexity calculation.

- Ruochen Meng: Review, bug fixes in probability calculation, improvements to unknown token handling.

- Saul Garay: Review, bug fixes in vocabulary construction and probability calculation.