

Assignment 1

https://github.com/AdemOdza/CS6320_Assignment1

Group11	Abhirup Mukherjee AXM240026	Adem Odza AXO180008	Buvana Seshathri BXS240020
	Samiha Kuncham S XK240025		

TODO: Fill in your name and details above. If not filled correctly, we will subtract 2pt.

TODO: It is highly suggested to use the template for reports. You are free to use other softwares (e.g. Microsoft Word, Google Doc), but we prefer you follow the structure below.

1 Implementation Details

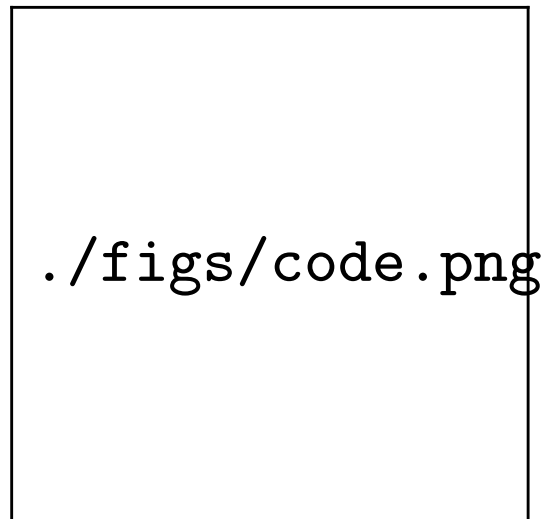


Figure 1: Example of a Code Piece.

```
import numpy
print("this is a piece of code")
```

Listing 1: Example of a Code Piece.

1.1 Unigram and bigram probability computation (15%)

1.1.1 Preprocessing Decisions

We are given two plain-text corpora of reviews where each line is a single review. In order to compute the counts and probabilities for the unigram and bigram, we need to tokenize the words into two flat Python lists with each word of the corpora being a single list element. We chose a list as our data structure since it is easily iterable and sliceable.

These were the preprocessing steps that we performed:

- a. Strip the leading and trailing whitespaces such as spaces, tabs, and newlines.
- b. If the line is empty after removing whitespaces, move to the next line in the file.
- c. Preprocess the current line by converting everything into lowercase.
- d. Remove the punctuations using the `re` library.
- e. If the processed line is not empty, split the line into individual words and add these into the token list.
- f. All words are flattened into a single list from the non-empty processed lines.

The following preprocessing decisions were made, along with the rationale behind each choice:

1. Lowercasing

We converted all the text to lower case to ensure that all the words with the same spelling are treated the same way. For example, “Hotel” and “hotel” are treated the same way (as the same token). This process helped us in reducing the vocabulary size and it also helped the model generalize better.

We used the following line of code to convert the existing text into lower case:

```
line = line.lower()
```

Listing 2: Converting text to lower case.

2. Removed Punctuation

All punctuation is removed so that the focus is on the actual words in the text. This decision was taken since punctuation often doesn’t contribute significant meaning in bag-of-words or similar models and it can introduce unnecessary noise.

The following line of code was introduced to remove punctuation:

```
import re
text = re.sub(r"[^\w\s]", ' ', text).strip()
```

Listing 3: Removing punctuation using regular expressions.

We used the regular expression module from the Python library to remove the punctuation in the text.

- `re.sub(pattern, repl, text)` replaces all occurrences of the regex pattern in text with repl (here, a space “ ”).
- The pattern `r"[^\w\s]"` means:
 - `[...]` : a character class (matches any character inside the brackets)
 - `^` : negation (matches any character not in the set)
 - `\w` : matches any “word” character (letters, digits, or underscore)
 - `\s` : matches any whitespace character (spaces, tabs, newlines)

So, `[^\w\s]` matches any character that is not a word character and not whitespace (i.e., punctuation and symbols). All such characters are replaced with a space. Finally, `.strip()` removes leading and trailing whitespace from the result.

3. Tokenization

Since the data given to us is already tokenized with one review per line and the tokens separated by spaces, we simply split each line into words. This made it straightforward and

easy to process each word individually.

The code used to perform tokenization is as follows:

```
def preprocess_file(input_path):
    tokens = []
    with open(input_path, encoding='utf-8') as fin:
        for line in fin:
            line = line.strip()
            if not line:
                continue
            processed = preprocess_line(line)
            if processed:
                tokens.extend(processed.split())
    return tokens
```

Listing 4: Tokenization function to preprocess reviews.

- **Reading Each Line:** The code reads each line (review) from the input file.
- **Preprocessing:** Each line is lowercased and has punctuation removed via `preprocess_line(line)`.
- **Splitting into Tokens:**
 - `processed.split()` splits the cleaned line into a list of words using spaces as delimiters.
 - `tokens.extend(...)` adds these words to the main tokens list.
- The function returns a flat list of all tokens (words) from the file.

The preprocessing steps resulted in two flat Python lists called `train_tokens` and `val_tokens`.

1.1.2 Unigram Counts and Probabilities

We begin with our `train_tokens` list and append start and stop tokens to it. The start token gives the bigram context of the first word and the stop token makes the bigram grammar a true probability distribution.¹

The unigram is the simplest of language models. The unigram probability of a given word is computed using the formula:

$$P(w_i) = \frac{\text{Count}(w_i)}{N}$$

where:

- $\text{Count}(w_i)$ is the number of times the word w_i appears in the training text.
- N is the total number of words (tokens) in the training text.

We chose to create a dictionary of the form $\{w_i : \text{freq}(w_i)\}$ where w_i is a word in our vocabulary (key) and $\text{freq}(w_i)$ is its frequency in the training corpus (value). Python dictionaries offer constant time key insertions and lookups ($O(1)$) and are simple to work with. To create this dictionary, we initialized an empty dictionary `unigram_freq` and iterated through the tokens. If a word is not in the dictionary, then we add it to the dictionary with a frequency of 0. Otherwise, we update its frequency.

Next, we initialize another empty dictionary called `unigram_probs`. We iterate through our `unigram_freq` dictionary and keep adding the (key, value) pairs to this new dictionary after dividing the value by N . This gives us our desired dictionary of the form $\{w_i : P(w_i)\}$.

1.1.3 Bigram Counts and Probabilities

The bigram model calculates the probability of a word based on the word that appeared immediately before it. In other words, it computes the conditional probability $P(w_i|w_{i-1})$ of the word w_i , given that the previous word was w_{i-1} . The bigram probability can be computed using the formula:

$$P(w_i|w_{i-1}) = \frac{\text{Count}(w_{i-1}, w_i)}{\text{Count}(w_{i-1})}$$

where:

- $\text{Count}(w_{i-1}, w_i)$ is the number of times the pair of words (the bigram) appears together in that order.
- $\text{Count}(w_{i-1})$ is the unigram count of the preceding word.

Next, we compute the bigram frequencies into a `bigram_freq` dictionary where each key is a bigram and the value is its frequency. To create this dictionary, we initialize the empty `bigram_freq` dictionary and iterate. We use the `zip` function to combine our `train_tokens_with_stop` list and a copy of this list starting from the second word to create a sequence of tuples, where each tuple is a bigram. We then iterate over each tuple and follow the same counting logic as before: if the tuple is not in `bigram_freq`, we assign its value as 0. Otherwise, we increment its count. This gives us a dictionary of the form $\{('w_{i-1}', 'w_i') : \text{bigram frequency}\}$.

Finally, we calculate the bigram probabilities by initializing an empty dictionary `bigram_probs`. We iterate through our `bigram_freq` dictionary and add the (key, value) pairs to this new dictionary after dividing the value by the unigram frequency of the first word in the bigram.

1.1.4 Perplexity

We use the following formula for calculating perplexity:

$$l = \frac{1}{N} \sum_{i=1}^N \log P(w_i|w_{i-1}, \dots, w_{i-n+1})$$
$$PP = 2^{-l}$$

where:

- N is the size of the training set.
- $P(w_i|w_{i-1}, \dots, w_{i-n+1})$ is the probability of the i th N -Gram

This is implemented using two functions: one for unigrams and one for bigrams. Both functions take the `corpus` (Preprocessed corpus), `unigram_counts` (The number of unigrams in the training set), N (Training corpus size), V (Vocab size), and `k` (The smoothing parameter). The bigram perplexity function takes one additional parameter, `bigram_counts` (The number of unigrams in the training set).

The smoothed probabilities are retrieved using the `add_k_smoothing_unigram` function, which passes the value of the `k` parameter and calculates the smoothed probabilities. We get the smoothed probability of every n -gram in the corpus, find their logarithm, and find their sum. We then divide that sum by M to get the l value. We then raise 2 to the power of l before returning.

TODO:

1.2 Smoothing (15%)

1.3 Unknown word handling (15%)

2 Eval, Analysis and Findings

TODO: to add (30%)

```

def calculate_unigram_perplexity(corpus, unigram_counts, N, V, k):

    # M = test corpus size
    M = len(corpus)
    total_log_prob = 0.0

    for token in corpus:
        # set default value of count to 0 if token not found in train set
        c_w = unigram_counts.get(token, 0)

        prob = add_k_smoothing_unigram(c_w, N, V, k)
        total_log_prob += math.log2(prob)

    l = total_log_prob / M
    perplexity = math.exp(-l) # 2.0 ** -l

    return perplexity

```

Listing 5: The Unigram perplexity function.

3 Others

TODO: to add (10%)

```

def calculate_bigram_perplexity(corpus, unigram_counts, bigram_counts, V, k):

    # Create bigrams from the validation corpus
    bigram_corpus = zip(corpus, corpus[1:])

    M = len(corpus)
    total_log_prob = 0.0

    for bigram in bigram_corpus:
        previous_word = bigram[0]

        # set default value of count to 0 if bigram not found in train set
        c_w_v = bigram_counts.get(bigram, 0)

        c_w = unigram_counts.get(previous_word, 0)
        prob = add_k_smoothing_bigram(c_w_v, c_w, V, k)
        total_log_prob += math.log2(prob)

    l = total_log_prob / M
    perplexity = math.exp(-l) # 2.0 ** -l

    return perplexity

```

Listing 6: The Bigram perplexity function.

References

- [1] Chapter 3, Section 3.1.2, Daniel Jurafsky and James H. Martin. *Speech and Language Processing (3rd ed. draft)*. <https://web.stanford.edu/~jurafsky/slp3/>. 2021.