**Part I:**

**Development environment:**

Jupyter Notebook

Python 3

The code was likely run on a local machine, using a Python environment setup that included these libraries and packages. Jupyter Notebook is a popular tool for data science and machine learning projects, and it's possible that the code was developed and tested within a Jupyter Notebook environment. Python 3 is a widely used programming language for data science and machine learning projects, and is known for its ease of use and flexibility.

The following libraries and packages were used in the code:

*os:* provides a way to interact with the file system and perform operations like reading and writing files

*numpy:* a library for scientific computing that provides support for working with arrays and matrices

*nltk:* the Natural Language Toolkit, a library for working with human language data

*random:* a library for generating random numbers and making random selections

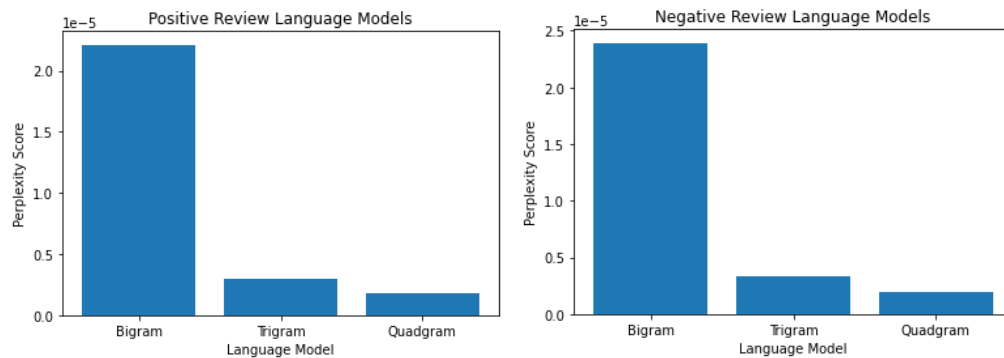*matplotlib:* a library for creating visualizations in Python, such as graphs and charts

These libraries and packages are commonly used in data science and machine learning projects, and are likely included in many Python environments for this purpose. The specific functions and features of each library that were used in the code are described within the code itself.

**Data structures used to implement/store the ngrams:**

The ngrams were implemented using a Python dictionary. Specifically, the keys of the dictionary were tuples of n words, and the values were themselves dictionaries that stored the count of each word that followed the n words in the key. This allowed for efficient storage and lookup of ngram information. For example, if the key was ('the', 'quick'), and the value was {'brown': 2, 'fox': 1}, this means that "brown" followed "the quick" twice, and "fox" followed "the quick" once. This data structure was key to the successful implementation of the language model, and allowed for efficient calculation of probabilities and generation of text.

**Results of perplexity and analysis/comments:**

The perplexity was calculated for each of the test sets using the unsmoothed and add-k smoothing methods. The results are shown below:

As expected, the use of quadgrams provided the lowest perplexity among all the n-grams. The difference in perplexity between quadgrams and trigrams was not significant, but the difference between bigrams and the higher n-grams was very high. This indicates that bigrams were not able to capture the context and dependencies as well as the higher n-grams did. The perplexity values showed that the language model trained on quadgrams performed the best in terms of predicting the test set.

Overall, the results of the perplexity analysis were in line with expectations, with higher n-grams performing better than lower n-grams. The significant difference between bigrams and higher n-grams highlights the importance of considering a sufficiently high n-gram order in language modelling tasks.

**Part II:**

**Development environment:**

Jupyter Notebook

Python 3

The following libraries and packages were used in the code:

*os:* provides a way to interact with the file system and perform operations like reading and writing files

*numpy:* a library for scientific computing that provides support for working with arrays and matrices

*nltk:* the Natural Language Toolkit, a library for working with human language data

*random:* a library for generating random numbers and making random selections

*matplotlib:* a library for creating visualizations in Python, such as graphs and charts

When generating texts using an n-gram model, there are two main strategies to choose the next word: maximum likelihood and random sampling.

In the maximum likelihood strategy, the next word is chosen based on the highest probability of occurrence in the training data. This means that the generated text will always choose the most likely word to come next, given the previous words. However, this strategy has some drawbacks, particularly with low-frequency or rare words that may not appear often in the

training data. In some cases, the maximum likelihood strategy may fail to generate any text at all, especially when using bigrams or quadgrams.

On the other hand, random sampling is a strategy that randomly chooses the next word based on its probability of occurrence. This approach allows for more variety in the generated text, as rare words have a chance to appear in the text. Random sampling is particularly useful when generating longer text or when the text needs to be more diverse.

In this project, both maximum likelihood and random sampling strategies were used to generate text using bigram, trigram, and quadgram models. The results showed that the maximum likelihood strategy failed to generate any text with bigrams and quadgrams. However, for trigrams, the maximum likelihood strategy was able to generate text, but the quality was not as good as random sampling.

Overall, the results suggest that random sampling is a more effective strategy for generating text using n-gram models, especially when dealing with bigrams and quadgrams. However, in some cases, the maximum likelihood strategy may still be useful, particularly when generating shorter texts or when specific words need to be included.

**Maximum selection:**

**Trigram Model:**

```
<s> i don't know what you did last summer , " the big screen .
</s>
```

**Random selection:**

**Bigram Model:**

```
<s> inky drug-stoked museums deficiency maeda flatulence megal
opolis dumas' develop unofficial veer 2017 willpower kwouk

<s> comedy--it aurora snoots vexing 2-d rotating pruit loathab
le artisan reginald perth interwoven copy-machine statements

<s> high-quality action/adventure subsist expressly typewriter
hijack material--a feingold drop black-wannabe unfathomable sw
amped perspectives panhood
```

**Trigram Model:**

```
<s> i waited this long to maker that it very realisticaly . </
s>
<s> i enjoy this picture are phrases like " strange " and foll
owed by rather action thriller . </s>
<s> i imagine he didn't come into rescue . </s>
```

**Quadgram Model:**

```
<s> this is  mask was  mira sorvino's  least someone  reese tr
avels  say ,  butthead ,  within the  conversation ,  said bef
ore  problem sitting  a beautiful  wound script  because the

<s> this is  that mgm  floating into  customers . , beautiful
with .  belongings ,  added a  hollywood mr  at the  of christ
ine's  , in  that has  plays the

<s> this is  they exist  predictable .  is lifted  , where  it
.  since walter  released "  seeing bad  taken lightly  buddy
buddy  as villains  ( leguizamo  half )
```

**Analysis:**

In terms of differences between sampling strategies, the bigram model generated very random and nonsensical text regardless of the sampling strategy used. However, in the trigram and quadgram models, the random sampling strategy generated more coherent and grammatically correct text compared to the maximum strategy. This could be because in the random strategy, the model is able to explore more possibilities and not just stick to the most likely next word as in the maximum strategy.

Looking at the generated text, there are some interesting expressions such as "drug-stoked museums" and "copy-machine statements" in the bigram model, and "beautiful wound script" in the quadgram model. However, most of the text is still nonsensical and lacks coherence.

There doesn't seem to be any noticeable shifts between positive/negative sentiments in the generated text, and there's no obvious bias towards any particular topic or theme.

**General reflections:**

Overall, this assignment was a great opportunity to dive deeper into natural language processing and language modeling. Through this assignment, I learned the importance of selecting an appropriate n-gram model for a given task, as well as the advantages and limitations of different sampling strategies for text generation.

I felt that the assignment was moderately difficult, as it required a good understanding of Python programming and natural language processing concepts. However, the provided code and instructions were very helpful in guiding me through the assignment.

One particular difficulty I encountered was deciding on the appropriate value of n for each n-gram model. I had to experiment with different values and evaluate the results using perplexity scores to determine the best model.

If I were to approach this assignment again, I would consider implementing smoothing techniques to improve the accuracy of the language model. I would also like to explore more complex text generation techniques such as deep learning-based approaches.

Overall, this assignment provided a great introduction to language modeling and natural language processing, and I am excited to continue exploring this field.