**Overview**

For the second practical we had to write a program which could analyse patterns in which characters occur in a file of text, it was required that the program should count the number of occurrences of all the possible groups of three consecutive characters. It was required that the program ignore punctuation and numerals as well as converting all the letter to lower case. Where a trigram was at the end of a word but consisted of only 2 letters an underscore should be added.

**Design**

**I have put all the program comments in the extension project.**
**Please not the program arguments should be as follows:**
*test.txt nGrams.csv wordGrams.csv dictionary.txt UTF-8*

The design of my project for reading data in from a file bears, close resemblance to my previous practical, due the simplicity and efficiency of the structre in my previous did not feel the need to write a new structure

A new instance of the nGramInfo class is initiated, and the print buffers are set up. A while loop is then entered; the terminating condition is met when the line of the text file is null e.g. end of file. When the while loop is entered the function from the nGram class is called, 'hasNextnGram'. This function returns a Boolean value, which is true when the end of the file is met, otherwise the returned value will be false. Within this function another function is called, findNextnGram. The find findNextnGram creates a new instance of the nGram class, which is a sub class (nested class) of the nGramInfo class. The nGram class contains the following attributes: nGramString, word, predictedWord and occurrences. The instaces of this class hold the information on each n-gram found. The findNextnGram then reads in a line of the text file. The line is split up into words and each word is converted to lower case and all non letter characters are removed. The words are then added to an array. A loop then goes through each word in the array, as long as the item in the array is not just a blank space, and hands the word to numerous other function including: 'findNGram' and 'checkNGramOccurrences'. The find nGram and findWordGram function are quite similar. Both functions take an input string, which is split up into either characters or words. The findNGram function is also recursive for efficiency.  The start character or word is then removed from the front of the string and the input is then handed back into the function. Based on what the user has input (x) each collection of x characters or words are finally handed to another function called 'checkNGramOccurrences'. However, because word-grams in theory could go over multiple lines, I first needed to read every word in the I put file and add it to an array, then with all the words in the array I can look for word-grams. The 'checkNGramOccurrences' function will check an arraylist of object nGram, for object with the same string as the input. If one is found the attribute 'occournces' is incremented for that object. Otherwise a new nGram object is created, where the nGramString attribute is set to the input of the function and then added to the array list. Finally, when all the nGrams and word grams are found the array lists are ordred based on occournecs and the data is output to a CSV file.

The garbage collection system built into the java programming language means the tempory instances of the nGram subclass are automatically deleted if they are not added to the array of nGrams.

I have tried also to implement predictive text functionality, which works by having the user input an n-gram via the console, the program then adds all a dictionary of words each with an assigned score (based on how frequently they are used in common English) to an array. When the user inputs a n-gram the program looks through the array to try and find any words staring with the input n-gram. If there are any matches, the found words are then ordered by score and the top 5 are output to the console.

**Testing**

To test the core functionality of my program I will use the various files supplied as an input and make sure the generated nGrams are firstly the correct size and also that the occurrences of each of the different nGrams found match the actual number of occurrences in the input file.

To carry out erroneous and extreme testing I will do the following tests:
        Try and input an invalid file

Try and input an invalid number when the user input is required.

To test the word prediction functionality, I will output the actual word the program trying to be predicted and the word which was predicted. I can then compare the two words and calculate the accuracy of the predictions over multiple different files.

**Examples**

**Test with errounous input**

| Input | Expected output | Actual output |
|---|---|---|
| nGram value: 0 | Ask user to re-enter a valid number | Please enter an integer >=2 <10:<br>0<br>Please enter an integer >=2 <10: |
| nGram value: 0 | Ask user to re-enter a valid number | Please enter an integer >=2 <10:<br>11<br>Please enter an integer >=2 <10: |
| nGram value: h | Ask user to re-enter a valid number | Please enter an integer >=2 <10:<br>h<br>Please enter an integer >=2 <10: |
| nGram value: <blank> | Ask user to re-enter a valid number | Please enter an integer >=2 <<br><br>\| |
| Corrupt input file | Notify the user the input file is not valid | Please enter an integer >=2 <10:<br>5<br>test1.txt (No such file or directory)<br><br>Process finished with exit code 0 |
| Incorrect number of program arguments | Notify the user the program arguments are incorrect | Please enter an integer >=2 <10:<br>3<br>Program arguments are incorrect...<br><br>Process finished with exit code 0<br>\| |

**Test with file: A Midsummer Night's Dream**

| Input | Expected output | Actual output | Expected output | Actual output |
|---|---|---|---|---|
| nGram value:2 | th,2296<br>he,1873<br>er,1302<br>an,1218<br>ou,1050<br>in,973<br>nd,948<br>re,798<br>ha,736<br>en,719 | th,2296<br>he,1873<br>er,1302<br>an,1218<br>ou,1050<br>in,973<br>nd,948<br>re,798<br>ha,736<br>en,719 | i will,47<br>i am,44<br>in the,44<br>and the,26<br>it is,25<br>and i,25<br>i have,24<br>of the,24<br>that i,21<br>i do,19 | i will,47<br>i am,44<br>in the,44<br>and the,26<br>it is,25<br>and i,25<br>i have,24<br>of the,24<br>that i,21<br>i do,19 |
| nGram value:3 | the,1080<br>and,734<br>her,477<br>you,422<br>to_,330<br>hat,319<br>ing,304<br>thi,301<br>his,298<br>of_,267 | the,1080<br>and,734<br>her,477<br>you,422<br>to_,330<br>hat,319<br>ing,304<br>thi,301<br>his,298<br>of_,267 | and i will,8<br>i have a,6<br>pyramus and thisby,6<br>i pray you,6<br>i love thee,6<br>it is not,6<br>and in the,5<br>when thou wakest,5<br>up and down,5<br>in the wood,4 | and i will,8<br>i have a,6<br>pyramus and thisby,6<br>i pray you,6<br>i love thee,6<br>it is not,6<br>and in the,5<br>when thou wakest,5<br>up and down,5<br>in the wood,4 |

**Test with file: Antony & Cleopatra**

| Input | Expected output | Actual output | Expected output | Actual output |
|---|---|---|---|---|
| nGram value:2 | th,3168<br>he,2471<br>an,1949<br>ar,1753<br>er,1686<br>ou,1492<br>ha,1376<br>en,1357<br>at,1286<br>in,1263 | th,3168<br>he,2471<br>an,1949<br>ar,1753<br>er,1686<br>ou,1492<br>ha,1376<br>en,1357<br>at,1286<br>in,1263 | mark antony,256<br>domitius enobarbus,137<br>octavius caesar,122<br>i have,69<br>i am,47<br>of the,42<br>to the,42<br>i will,41<br>my lord,39<br>in the,35 | mark antony,256<br>domitius enobarbus,137<br>octavius caesar,122<br>i have,69<br>i am,47<br>of the,42<br>to the,42<br>i will,41<br>my lord,39<br>in the,35 |
| nGram value:3 | the,1541<br>and,831<br>to_,559<br>you,544<br>hat,514<br>her,508<br>ant,473<br>of_,450<br>our,443<br>nto,418 | the,1541<br>and,831<br>to_,559<br>you,544<br>hat,514<br>her,508<br>ant,473<br>of_,450<br>our,443<br>nto,418 | mark antony i,18<br>mark antony and,17<br>enter mark antony,17<br>lord mark antony,14<br>my lord mark,12<br>enter octavius caesar,12<br>domitius enobarbus aside,11<br>domitius enobarbus i,9<br>octavius caesar i,9<br>cleopatras palace enter,8 | mark antony i,18<br>mark antony and,17<br>enter mark antony,17<br>lord mark antony,14<br>my lord mark,12<br>enter octavius caesar,12<br>domitius enobarbus aside,11<br>domitius enobarbus i,9<br>octavius caesar i,9<br>cleopatras palace enter,8 |

**Test with file: Alls Well That Ends Well**

| Input | Expected output | Actual output | Expected output | Actual output |
|---|---|---|---|---|
| nGram value:2 | he,949<br>th,926<br>an,631<br>in,523<br>nd,446<br>ed,408<br>er,393<br>re,373<br>hi,315 | he,949<br>th,926<br>an,631<br>in,523<br>nd,446<br>ed,408<br>er,393<br>re,373<br>hi,315 | and the,39<br>of the,31<br>in the,25<br>he had,21<br>to the,19<br>the martian,19<br>it was,17<br>the hound,17<br>for a,14<br>into the,14 | and the,39<br>of the,31<br>in the,25<br>he had,21<br>to the,19<br>the martian,19<br>it was,17<br>the hound,17<br>for a,14<br>into the,14 |
| nGram value:3 | the,645<br>and,316<br>ing,230<br>he_,159<br>to_,150<br>of_,147 | the,645<br>and,316<br>ing,230<br>he_,159<br>to_,150<br>of_,147<br>his,121<br>ere,112 | it was a,6<br>he had to,4<br>the man had,4<br>the martian was,4<br>project gutenberg ebook,3 | it was a,6<br>he had to,4<br>the man had,4<br>the martian was,4<br>project gutenberg |

| | his,121<br>ere,112<br>was,100<br>in_,98 | was,100<br>in_,98 | duel on syrtis,3<br>top of the,3<br>a long time,3<br>but he didnt,3<br>it would be,3 | ebook,3<br>duel on<br>syrtis,3<br>top of the,3<br>a long<br>time,3<br>but he<br>didnt,3<br>it would<br>be,3 |

**Test with file: Duel on Syrtis**

| Input | Expected output | Actual output | Expected output | Actual output |
|---|---|---|---|---|
| nGram value:2 | th,7522<br>he,6802<br>in,4052<br>er,3859<br>an,3758<br>ha,3442<br>ou,3307<br>re,3296<br>on,2950<br>at,2930 | th,7522<br>he,6802<br>in,4052<br>er,3859<br>an,3758<br>ha,3442<br>ou,3307<br>re,3296<br>on,2950<br>at,2930 | of the,434<br>in the,250<br>it was,170<br>it is,156<br>i have,146<br>the moor,139<br>sir henry,139<br>upon the,133<br>to the,131 | of the,434<br>in the,250<br>it was,170<br>it is,156<br>i have,146<br>the moor,139<br>sir henry,139<br>upon the,133<br>to the,131 |
| nGram value:3 | the,4466<br>and,1919<br>of_,1599<br>ing,1408<br>to_,1407<br>hat,1381<br>tha,1250<br>you,1068<br>his,1052<br>her,1023 | the,4466<br>and,1919<br>of_,1599<br>ing,1408<br>to_,1407<br>hat,1381<br>tha,1250<br>you,1068<br>his,1052<br>her,1023 | upon the<br>moor,45<br>that it was,27<br>of the moor,26<br>out of the,23<br>it was a,22<br>on the moor,22<br>sir henry<br>baskerville,22<br>i could not,20<br>in front of,19<br>it is a,19 | upon the<br>moor,45<br>that it was,27<br>of the moor,26<br>out of the,23<br>it was a,22<br>on the moor,22<br>sir henry<br>baskerville,22<br>i could not,20<br>in front of,19<br>it is a,19 |

**Test with file: Hound of the Baskervilles**

| Input | Expected output | Actual output | Expected output | Actual output |
|---|---|---|---|---|
| nGram value:2 | th,6728<br>he,6363<br>in,4612<br>er,3952<br>re,3157<br>ha,3055<br>ed,3024<br>ou,2987<br>an,2903<br>at,2719 | th,6728<br>he,6363<br>in,4612<br>er,3952<br>re,3157<br>ha,3055<br>ed,3024<br>ou,2987<br>an,2903<br>at,2719 | of the,302<br>in the,228<br>it was,180<br>on the,160<br>to the,140<br>it is,124<br>mrs<br>inglethorp,115<br>i was,109<br>at the,99<br>had been,97 | of the,302<br>in the,228<br>it was,180<br>on the,160<br>to the,140<br>it is,124<br>mrs<br>inglethorp,115<br>i was,109<br>at the,99<br>had been,97 |
| nGram value:3 | the,3777<br>ing,1963<br>and,1448 | the,3777<br>ing,1963<br>and,1448<br>hat,1346 | that it was,32<br>i did not,27<br>out of the,27 | that it was,32<br>i did not,27<br>out of the,27<br>one of the,24 |

| | hat,1346<br>to_,1339<br>her,1324<br>of_,1269<br>you,1126<br>tha,1083 | `to_,1339`<br>`her,1324`<br>`of_,1269`<br>`you,1126`<br>`tha,1083` | one of the,24<br>there was a,23<br>i do not,23<br>i dont know,22<br>what do you,19<br>it was a,18 | `there was a,23`<br>`i do not,23`<br>`i dont know,22`<br>`what do you,19`<br>`it was a,18` |
| --- | --- | --- | --- | --- |

**Test word prediction**

| Input | Expected output | Actual output |
| --- | --- | --- |
| Input nGram: t | the<br>to<br>to<br>that<br>this | `Enter quit to end program or otherwise enter an n-gram:`<br>`t`<br>`The predicted words based on the input 't'(in order of relevance) are:`<br>`the`<br>`to`<br>`to`<br>`that`<br>`this` |
| Input nGram: th | the<br>that<br>this<br>they<br>that | `Enter quit to end program or otherwise enter an n-gram:`<br>`th`<br>`The predicted words based on the input 'th'(in order of relevance) are:`<br>`the`<br>`that`<br>`this`<br>`they`<br>`that` |
| Input nGram: he | he<br>her<br>here<br>her<br>help | `Enter quit to end program or otherwise enter an n-gram:`<br>`he`<br>`The predicted words based on the input 'he'(in order of relevance) are:`<br>`he`<br>`her`<br>`here`<br>`her`<br>`help` |
| Input nGram: 0 | Program should wait for valid input | `Enter quit to end program or otherwise enter an n-gram:`<br>`0`<br>`Enter quit to end program or otherwise enter an n-gram:` |
| Input nGram: gobbledygook | No predictons found | `Enter quit to end program or otherwise enter an n-gram:`<br>`gobbledygook`<br>`No predictions found` |
| Input nGram: quit | Program exits | `Enter quit to end program or otherwise enter an n-gram:`<br>`quit`<br><br>`Process finished with exit code 0` |

**Evaluation**

After testing my program, I pleased to say it works as expected and is able to process large data files quickly and efficiently. I have used multiple classes and universal methods which take in various parameters making the program overall more modular and adaptable.

The output files that are generated contain the correct required data as well as other information that is generated based on the input files.


**Conclusion**

In conclusion I'm satisfied with end program, it is able to carry out the required tasks with ease and is able to handle erroneous data without crashing. The added extensions massively increase the overall functionality and makes the program far more flexible. The word prediction extension did however prove troublesome due to the ambiguity of the extension requirement.