The first part of the coursework, was quite enjoyable, just to familiarize us with what the basic approaches in cryptography are, I personally found it enjoyable.

The first step I took to solve it was to write a little Utils class which was supposed to be reused for the second part of the coursework, along the other methods in this class I wrote a private class (which later on I realised it was a mistake) to deal with the Frequency Analysis itself. The idea was to have this class hidden in the CryptoUtils class and then manipulate objects of FrequencyAnaliser type, this was achieved at some extend, but my implementation was not a particularly good, I should have put some more thought in it.
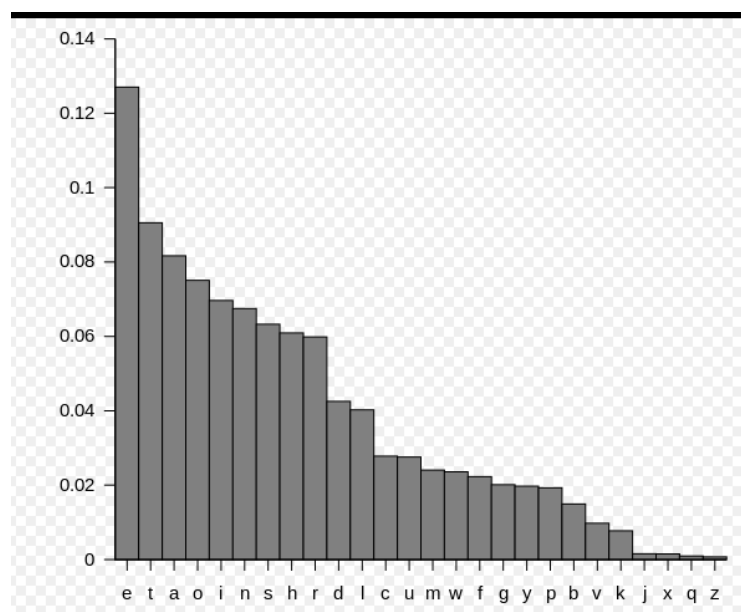
However, all required functionalities are present. I had an issue with the Frequency Analysis, my code was skipping about 30 000 valid characters (A-Zz-A) from the file, later on I realised that I have forgotten to convert the text to lower case. For the implementation of this I decided to use mainly streams, because they make the code more readable and are easy to write, but not all the time of course. The results from frequency analysis were correct with corresponding sources on the web.

Here are my findings :



(Right picture from Wikipedia)

As we can observe, the letter distributions are almost identical.

The second step, was to implement methods which would find bigrams and trigrams in the texts, once I had those, it was made clear that we probably had to deal with a shift cipher. In the cyphered text, we observe a lot of appearances of the string "xli" which in a way implied that it could be corresponding to "the" in the original text, by looking at the bigrams and trigrams for both of the texts, I was able to work out that the cipher used is not a substitution one. Afterwards I worked out the shift values and deciphered the text. Another approach would have been to perform exhaustive search, since the key space is not significant.

An issue that I have encountered was breaking the shift mechanism while I tried to optimize it for the second part of the coursework, which forced me to introduce two functions for character rotation, because it was simply easier than trying to debug it and make it work for both of them. I think it was particularly bad idea to keep using integers for the char values, from the research I have made many people suggest using byte values or hexadecimal values. Nevertheless, I managed to present a working implementation. I lesson I have learned is that a single simple mistake will always result in completely different output from the program.

Here we can see the relations of letter distribution in the encrypted and plain text for exercise 1:

```
  Printing all digrams fou
Counter is: 321081
th - 12012 - 3.74 %
 he - 10980 - 3.42 %                  Printing all digrams found for file Resources
 in - 7670 - 2.39 %          Counter is: 1137
 er - 7181 - 2.24 %          xl - 32 - 2.81 %
 an - 6511 - 2.03 %           vi - 30 - 2.63 %
 ha - 5772 - 1.80 %           li - 29 - 2.54 %
 ou - 5701 - 1.78 %           iv - 28 - 2.45 %
 re - 5627 - 1.75 %           mr - 28 - 2.45 %
 nd - 4914 - 1.53 %           ih - 22 - 1.92 %
 at - 4896 - 1.52 %           er - 21 - 1.83 %
 on - 4493 - 1.40 %           sy - 18 - 1.57 %
 en - 4268 - 1.33 %           ir - 18 - 1.57 %
 ed - 4264 - 1.33 %           iw - 18 - 1.57 %
 hi - 4157 - 1.29 %           mw - 18 - 1.57 %
 is - 3875 - 1.21 %           mx - 18 - 1.57 %
 to - 3771 - 1.17 %           ev - 16 - 1.39 %
 it - 3753 - 1.17 %           le - 16 - 1.39 %
 as - 3627 - 1.13 %           qi - 15 - 1.30 %
 ve - 3472 - 1.08 %           lm - 15 - 1.30 %
 es - 3454 - 1.08 %           rh - 14 - 1.21 %
 ng - 3421 - 1.07 %           sr - 14 - 1.21 %
 or - 3298 - 1.03 %           ew - 14 - 1.21 %
 ar - 3207 - 1.00 %           wi - 14 - 1.21 %
 me - 3147 - 0.98 %           xs - 14 - 1.21 %
 te - 3146 - 0.98 %           rk - 13 - 1.12 %
 st - 3128 - 0.97 %           wx - 12 - 1.03 %
 se - 2943 - 0.92 %           xi - 12 - 1.03 %
 le - 2780 - 0.87 %           ai - 11 - 0.95 %
 of - 2766 - 0.86 %           ri - 11 - 0.95 %
 ne - 2514 - 0.78 %           rx - 11 - 0.94 %
                              ex - 11 - 0.94 %
```

Plain text                                    Ciphered text

The pictures above, give us hints what could be the possible shift value for the characters. Then the trigrams, give us even further proof:

Ciphered text                    Plaintext

```
 Printing all trigrams for fil
xli - 18 - 2.20 %
 mrk - 10 - 1.22 %          the - 6644 - 2.99 %
 erh - 9 - 1.10 %           and - 3184 - 1.43 %
 lmw - 9 - 1.10 %           ing - 2694 - 1.21 %
 ivi - 8 - 0.97 %           hat - 2134 - 0.96 %
 ali - 6 - 0.73 %           her - 1944 - 0.87 %
 lex - 6 - 0.73 %           tha - 1843 - 0.83 %
 xle - 6 - 0.73 %           you - 1720 - 0.77 %
 irx - 5 - 0.61 %           ere - 1534 - 0.69 %
 aew - 5 - 0.60 %           his - 1502 - 0.68 %
 iww - 5 - 0.60 %           was - 1295 - 0.58 %
 qmr - 5 - 0.60 %           ver - 1222 - 0.55 %
 evi - 5 - 0.60 %           for - 1180 - 0.53 %
 liv - 5 - 0.60 %           ave - 1123 - 0.51 %
 ipc - 4 - 0.48 %           thi - 1100 - 0.49 %
 viw - 4 - 0.48 %           ter - 1081 - 0.49 %
 jsv - 4 - 0.48 %           ent - 1072 - 0.48 %
 jvs - 4 - 0.48 %           ith - 1021 - 0.46 %
 riv - 4 - 0.48 %           all - 1019 - 0.46 %
 aiv - 4 - 0.48 %           wit - 999 - 0.45 %
 qiw - 4 - 0.48 %           hav - 972 - 0.44 %
 amx - 4 - 0.48 %           our - 959 - 0.43 %
 lir - 4 - 0.48 %           ion - 948 - 0.43 %
 swx - 4 - 0.48 %           whi - 929 - 0.42 %
 syv - 4 - 0.48 %           oul - 888 - 0.40 %
 ziv - 4 - 0.47 %           not - 867 - 0.39 %
 mrh - 4 - 0.47 %           one - 837 - 0.38 %
 msr - 4 - 0.47 %           uld - 817 - 0.37 %
 xih - 4 - 0.47 %           ght - 813 - 0.37 %
 vih - 3 - 0.35 %           ear - 809 - 0.36 %
```

Another thing that I observed was while working with the bigrams and trigrams, I had certain values, but colleagues of mine had totally different ones, so I used a couple of online tools, and there was no consistent result amongst them which was quite surprising.

To run the code just execute, the Runner class from the package containing exercise one, no user interaction is needed, a new file "analysis" followed by the current date and time will appear, where any crypto and frequency analysis results can be seen. A file with the deciphered text will appear the file name is "ceaserDecrypt" followed by the current time and date. All these are written under the resources folder, within the project. Comment are available in the code to make it easier to understand.