This assignment was completed in stages.

Stage 1: Building Inverted Index

Here for each document, I traversed through all the words in it. The information regarding the words traversed was stored in a data structure which had the following format:

```
{
{word1 --> (docId, count), (docId,count) ...},
{word2 --> (docId, count), (docId,count) ...}
}
```

There were a few cases that I needed to consider while traversing through words
case1: If the word was not found before.
case2: If the same word was found more than once in the same document.
case3: If the same word was found more than once but in different documents.

I also maintained a dict for storing (docId, noOfWordsInDocuments), which is used in calculating the value of K in the second half of the assignment.

I stored the index in the file in the following format:

```
# word
docId count docId count
```

This made my job of building the index for the second part of the assignment easier.

Stage 2: Calculating BM25 Score

I rebuilt the index again for this part from index.out. After the index was built I calculated the BM25 score using the following formula:

$$\log \left( (N - n_i + 0.5) / (n_i + 0.5) \right) * \{ ( (k_1 + 1) * f_i ) * ( (k_2 + 1) * qf_i ) / (K + f_i) * (k_2 + qf_i) \}$$

I updated the formula from the original one, since we know that R and r = 0
Here,
N = no of documents
$n_i$ = no of documents which contain term i
$k_1$ = 1.2
$f_i$ = no of time term i occurs in the document we are considering
$k_2$ = 100

```
qfi = no of times term 1 occurs in the query
K =  k1 * ( ( 1 - b ) + ( b * float( dl ) / avdl ) )
b = 0.75
dl = length of document, obtainable from the second data structure
mainted as per the first part of the assignment.
avdl = average length of the documents in the corpus.
```