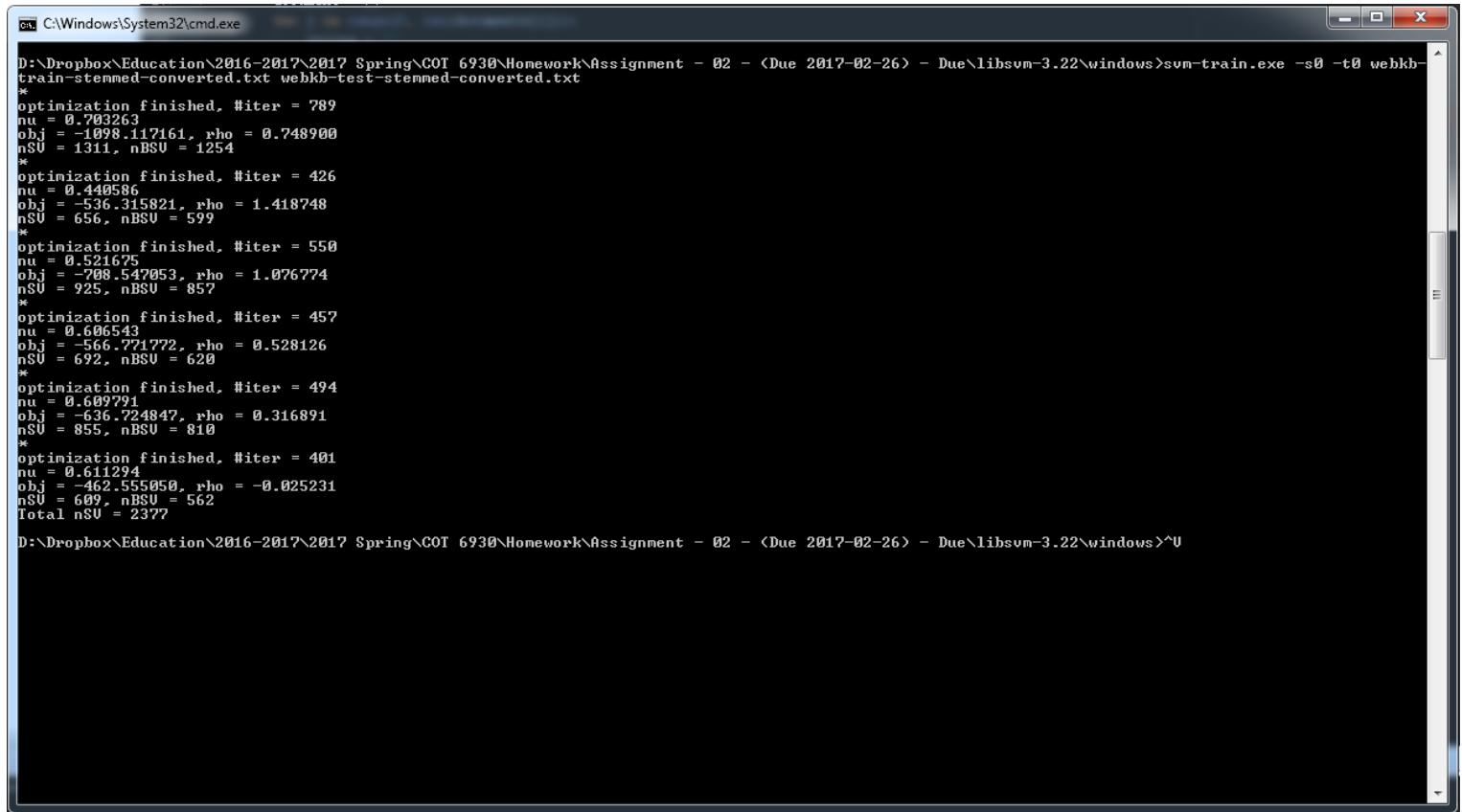


## Maciej Medyk – COT6930 – Natural Language Processing – Homework 02

Question 1 – [10.00pt] – Write a report including the screenshots of generating document-word matrix, input data files for training and testing, command lines you use to train and test, and the output results by libsvm.

After running svm-train.exe with parameter -t0 for linear kernel and -s0 for C-SVC the results were as follows:

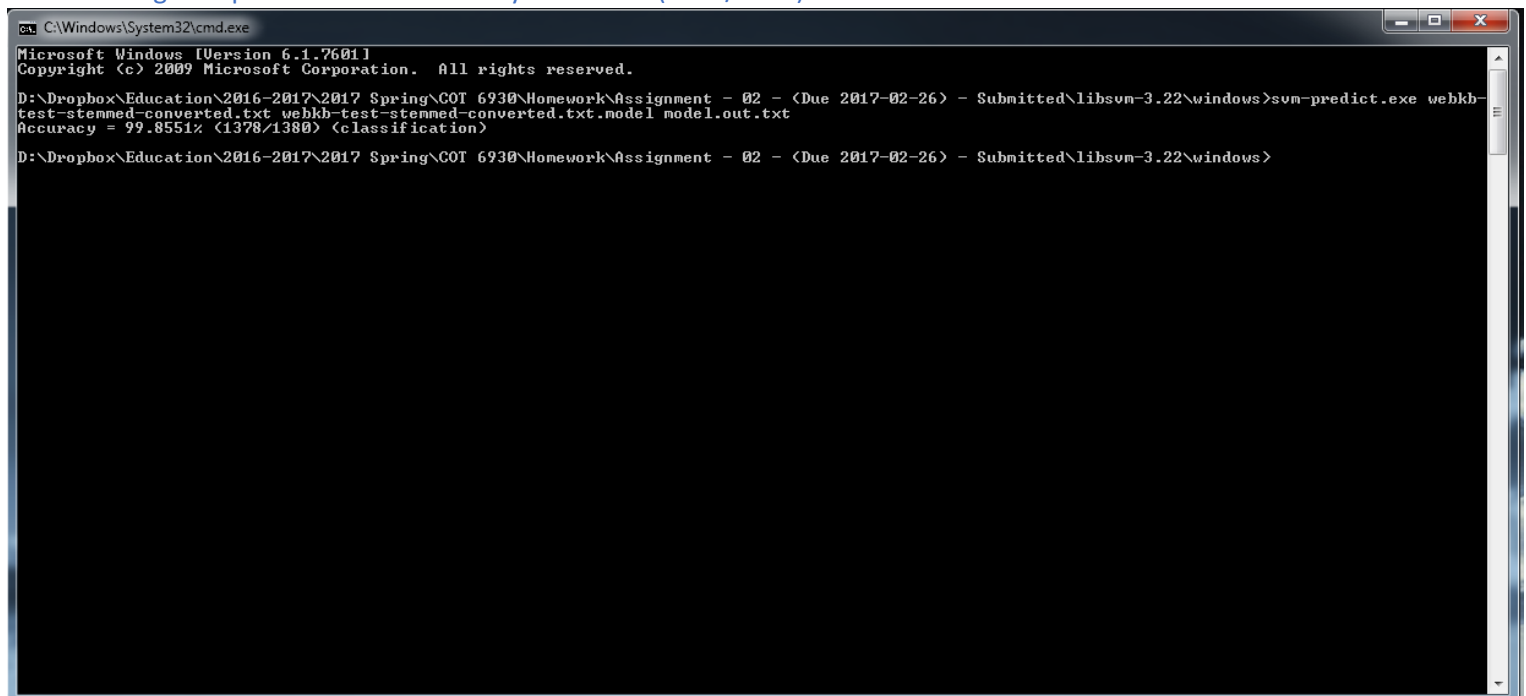


```
C:\Windows\System32\cmd.exe

D:\Dropbox\Education\2016-2017\2017 Spring\COT 6930\Homework\Assignment - 02 - <Due 2017-02-26> - Due\libsvm-3.22\windows>svm-train.exe -s0 -t0 wehkb-train-stemmed-converted.txt wehkb-test-stemmed-converted.txt
*
optimization finished, #iter = 789
nu = 0.703263
obj = -1098.117161, rho = 0.748900
nSV = 1311, nBSV = 1254
*
optimization finished, #iter = 426
nu = 0.440586
obj = -536.315821, rho = 1.418748
nSV = 656, nBSV = 599
*
optimization finished, #iter = 550
nu = 0.521675
obj = -708.547053, rho = 1.076774
nSV = 925, nBSV = 857
*
optimization finished, #iter = 457
nu = 0.606543
obj = -566.721772, rho = 0.528126
nSV = 692, nBSV = 620
*
optimization finished, #iter = 494
nu = 0.609791
obj = -636.724847, rho = 0.316891
nSV = 855, nBSV = 810
*
optimization finished, #iter = 401
nu = 0.611294
obj = -462.555050, rho = -0.025231
nSV = 609, nBSV = 562
Total nSV = 2377

D:\Dropbox\Education\2016-2017\2017 Spring\COT 6930\Homework\Assignment - 02 - <Due 2017-02-26> - Due\libsvm-3.22\windows>^U
```

After running svm-predict.exe the accuracy is 99.8551 (1378/1380)



```
C:\Windows\System32\cmd.exe

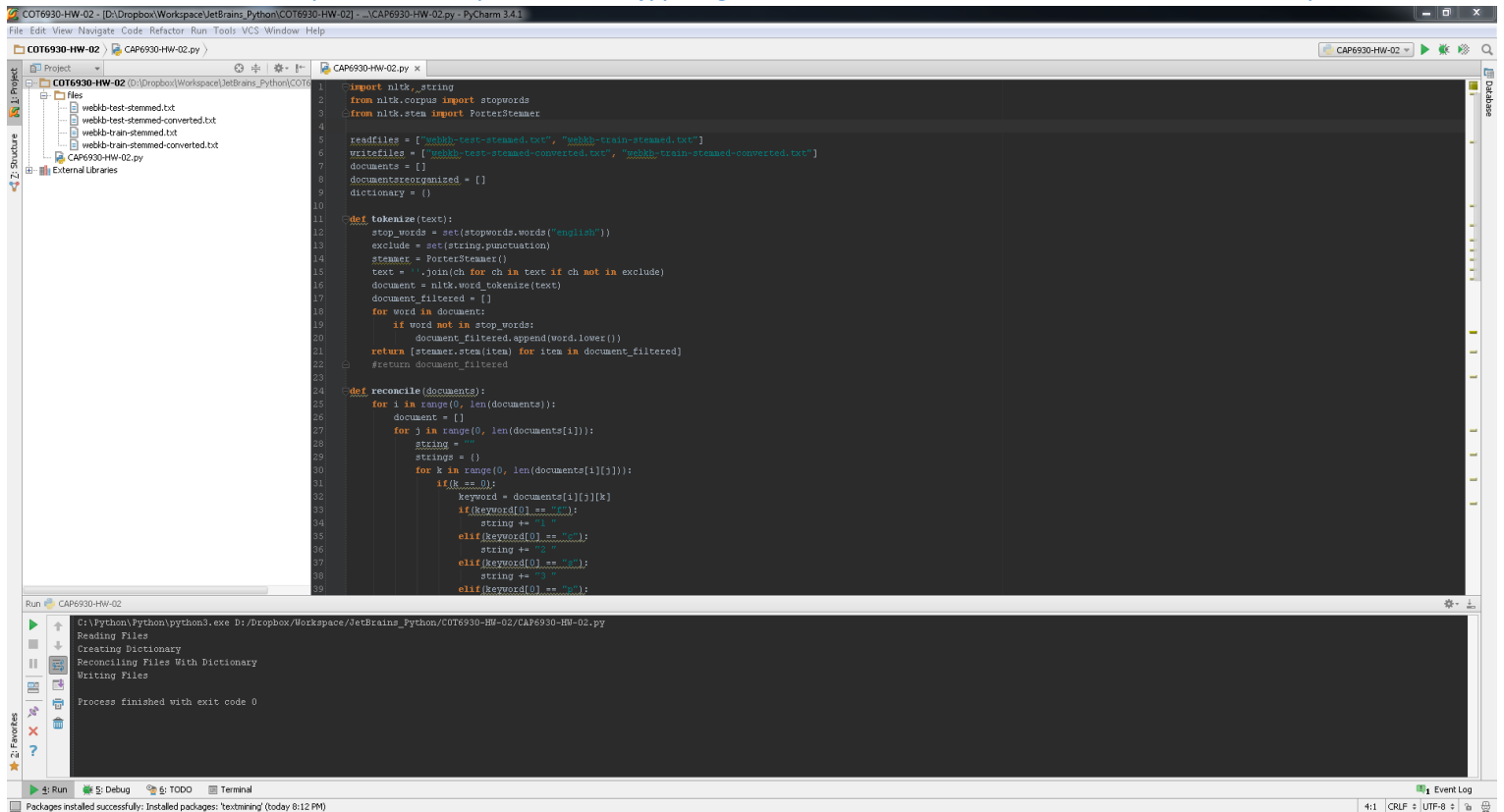
Microsoft Windows [Version 6.1.7601]
Copyright (c) 2009 Microsoft Corporation. All rights reserved.

D:\Dropbox\Education\2016-2017\2017 Spring\COT 6930\Homework\Assignment - 02 - <Due 2017-02-26> - Submitted\libsvm-3.22\windows>svm-predict.exe wehkb-test-stemmed-converted.txt wehkb-test-stemmed-converted.txt.model model.out.txt
Accuracy = 99.8551% (1378/1380) <classification>

D:\Dropbox\Education\2016-2017\2017 Spring\COT 6930\Homework\Assignment - 02 - <Due 2017-02-26> - Submitted\libsvm-3.22\windows>^U
```

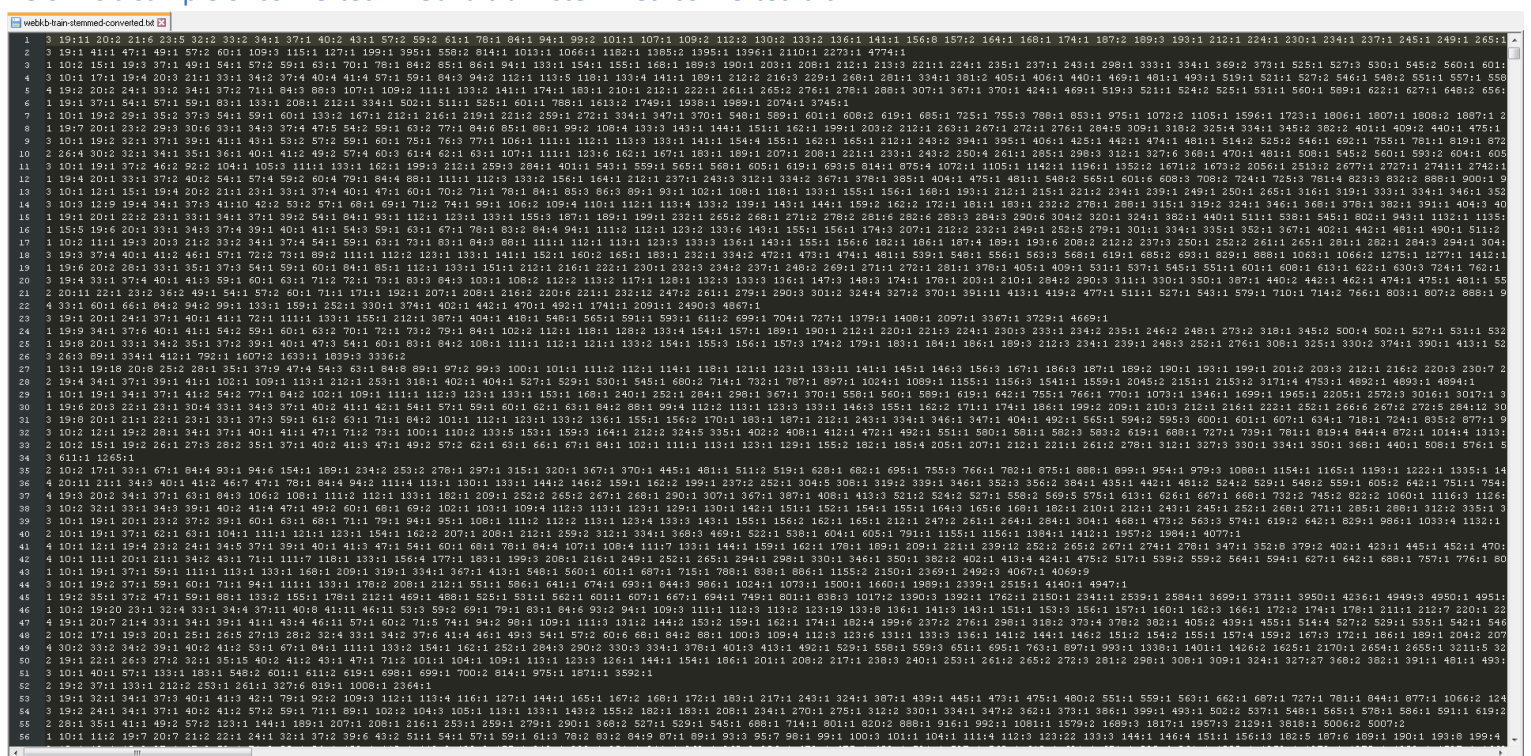
The supplied files required extensive pre-processing that involved python NLTK package and additional custom programming. In this assignment, I wrote a program that loaded each file for preprocessing and then tokenized it. Afterwards, I created dictionary based on the words collected and assigned the word unique numerical value. Then I counted frequency of all words in each sentence and substituted words for dictionary derived values. In this way, I ended up with the desired input files that had a required format by LIBSVM.

Code available at link : <https://www.dropbox.com/s/4jyy2k2gv9hewdv/Homework%20-%202002%20-%20Code.zip?dl=0>



```
1 import nltk, string
2 from nltk.corpus import stopwords
3 from nltk.stem import PorterStemmer
4
5 readfiles = ["webkb-test-stemmed.txt", "webkb-train-stemmed.txt"]
6 writefiles = ["webkb-test-stemmed-converted.txt", "webkb-train-stemmed-converted.txt"]
7 documents = []
8 documentsreorganised = []
9 dictionary = {}
10
11
12 def tokenise(text):
13     stop_words = set(stopwords.words('english'))
14     exclude = set(string.punctuation)
15     stemmer = PorterStemmer()
16     text = ''.join(ch for ch in text if ch not in exclude)
17     document = nltk.word_tokenize(text)
18     document_filtered = []
19     for word in document:
20         if word not in stop_words:
21             document_filtered.append(word.lower())
22     return [stemmer.stem(item) for item in document_filtered]
23
24 def reconcile(documents):
25     for i in range(0, len(documents)):
26         document = []
27         for j in range(0, len(documents[i])):
28             string = ""
29             strings = []
30             for k in range(0, len(documents[i][j])):
31                 if k == 0:
32                     keyword = documents[i][j][k]
33                     if(keyword[0] == '@'):
34                         string += " "
35                     elif(keyword[0] == '#'):
36                         string += " "
37                     elif(keyword[0] == '$'):
38                         string += " "
39                     elif(keyword[0] == '%'):
```

Below is a sample of converted "webkb-train-stemmed-converted.txt"



1 139:11 2012 2116 2315 3212 3312 3411 3711 4012 4311 5712 5912 6111 7811 8411 9411 9912 10111 10711 10912 11212 13012 13312 13611 14111 15618 15712 16411 16811 17411 18712 18913 19311 21211 22411 23011 23411 23711 24511 24911 26511 27711 2811 2841 2871 2911 2941 2971 3011 3041 3071 3101 3131 3161 3191 3221 3251 3281 3311 3341 3371 3401 3431 3461 3491 3521 3551 3581 3611 3641 3671 3701 3731 3761 3791 3821 3851 3881 3911 3941 3971 4001 4031 4061 4091 4121 4151 4181 4211 4241 4271 4301 4331 4361 4391 4421 4451 4481 4511 4541 4571 4601 4631 4661 4691 4721 4751 4781 4811 4841 4871 4901 4931 4961 4991 5021 5051 5081 5111 5141 5171 5201 5231 5261 5291 5321 5351 5381 5411 5441 5471 5501 5531 5561 5591 5621 5651 5681 5711 5741 5771 5801 5831 5861 5891 5921 5951 5981 6011 6041 6071 6101 6131 6161 6191 6221 6251 6281 6311 6341 6371 6401 6431 6461 6491 6521 6551 6581 6611 6641 6671 6701 6731 6761 6791 6821 6851 6881 6911 6941 6971 7001 7031 7061 7091 7121 7151 7181 7211 7241 7271 7301 7331 7361 7391 7421 7451 7481 7511 7541 7571 7601 7631 7661 7691 7721 7751 7781 7811 7841 7871 7901 7931 7961 7991 8021 8051 8081 8111 8141 8171 8201 8231 8261 8291 8321 8351 8381 8411 8441 8471 8501 8531 8561 8591 8621 8651 8681 8711 8741 8771 8801 8831 8861 8891 8921 8951 8981 9011 9041 9071 9101 9131 9161 9191 9221 9251 9281 9311 9341 9371 9401 9431 9461 9491 9521 9551 9581 9611 9641 9671 9701 9731 9761 9791 9821 9851 9881 9911 9941 9971 10001 10031 10061 10091 10121 10151 10181 10211 10241 10271 10301 10331 10361 10391 10421 10451 10481 10511 10541 10571 10601 10631 10661 10691 10721 10751 10781 10811 10841 10871 10901 10931 10961 10991 11021 11051 11081 11111 11141 11171 11201 11231 11261 11291 11321 11351 11381 11411 11441 11471 11501 11531 11561 11591 11621 11651 11681 11711 11741 11771 11801 11831 11861 11891 11921 11951 11981 12011 12041 12071 12101 12131 12161 12191 12221 12251 12281 12311 12341 12371 12401 12431 12461 12491 12521 12551 12581 12611 12641 12671 12701 12731 12761 12791 12821 12851 12881 12911 12941 12971 13001 13031 13061 13091 13121 13151 13181 13211 13241 13271 13301 13331 13361 13391 13421 13451 13481 13511 13541 13571 13601 13631 13661 13691 13721 13751 13781 13811 13841 13871 13901 13931 13961 13991 14021 14051 14081 14111 14141 14171 14201 14231 14261 14291 14321 14351 14381 14411 14441 14471 14501 14531 14561 14591 14621 14651 14681 14711 14741 14771 14801 14831 14861 14891 14921 14951 14981 15011 15041 15071 15101 15131 15161 15191 15221 15251 15281 15311 15341 15371 15401 15431 15461 15491 15521 15551 15581 15611 15641 15671 15701 15731 15761 15791 15821 15851 15881 15911 15941 15971 16001 16031 16061 16091 16121 16151 16181 16211 16241 16271 16301 16331 16361 16391 16421 16451 16481 16511 16541 16571 16601 16631 16661 16691 16721 16751 16781 16811 16841 16871 16901 16931 16961 16991 17021 17051 17081 17111 17141 17171 17201 17231 17261 17291 17321 17351 17381 17411 17441 17471 17501 17531 17561 17591 17621 17651 17681 17711 17741 17771 17801 17831 17861 17891 17921 17951 17981 18011 18041 18071 18101 18131 18161 18191 18221 18251 18281 18311 18341 18371 18401 18431 18461 18491 18521 18551 18581 18611 18641 18671 18701 18731 18761 18791 18821 18851 18881 18911 18941 18971 19001 19031 19061 19091 19121 19151 19181 19211 19241 19271 19301 19331 19361 19391 19421 19451 19481 19511 19541 19571 19601 19631 19661 19691 19721 19751 19781 19811 19841 19871 19901 19931 19961 19991 20021 20051 20081 20111 20141 20171 20201 20231 20261 20291 20321 20351 20381 20411 20441 20471 20501 20531 20561 20591 20621 20651 20681 20711 20741 20771 20801 20831 20861 20891 20921 20951 20981 21011 21041 21071 21101 21131 21161 21191 21221 21251 21281 21311 21341 21371 21401 21431 21461 21491 21521 21551 21581 21611 21641 21671 21701 21731 21761 21791 21821 21851 21881 21911 21941 21971 22001 22031 22061 22091 22121 22151 22181 22211 22241 22271 22301 22331 22361 22391 22421 22451 22481 22511 22541 22571 22601 22631 22661 22691 22721 22751 22781 22811 22841 22871 22901 22931 22961 22991 23021 23051 23081 23111 23141 23171 23201 23231 23261 23291 23321 23351 23381 23411 23441 23471 23501 23531 23561 23591 23621 23651 23681 23711 23741 23771 23801 23831 23861 23891 23921 23951 23981 24011 24041 24071 24101 24131 24161 24191 24221 24251 24281 24311 24341 24371 24401 24431 24461 24491 24521 24551 24581 24611 24641 24671 24701 24731 24761 24791 24821 24851 24881 24911 24941 24971 25001 25031 25061 25091 25121 25151 25181 25211 25241 25271 25301 25331 25361 25391 25421 25451 25481 25511 25541 25571 25601 25631 25661 25691 25721 25751 25781 25811 25841 25871 25901 25931 25961 25991 26021 26051 26081 26111 26141 26171 26201 26231 26261 26291 26321 26351 26381 26411 26441 26471 26501 26531 26561 26591 26621 26651 26681 26711 26741 26771 26801 26831 26861 26891 26921 26951 26981 27011 27041 27071 27101 27131 27161 27191 27221 27251 27281 27311 27341 27371 27401 27431 27461 27491 27521 27551 27581 27611 27641 27671 27701 27731 27761 27791 27821 27851 27881 27911 27941 27971 28001 28031 28061 28091 28121 28151 28181 28211 28241 28271 28301 28331 28361 28391 28421 28451 28481 28511 28541 28571 28601 28631 28661 28691 28721 28751 28781 28811 28841 28871 28901 28931 28961 28991 29021 29051 29081 29111 29141 29171 29201 29231 29261 29291 29321 29351 29381 29411 29441 29471 29501 29531 29561 29591 29621 29651 29681 29711 29741 29771 29801 29831 29861 29891 29921 29951 29981 30011 30041 30071 30101 30131 30161 30191 30221 30251 30281 30311 30341 30371 30401 30431 30461 30491 30521 30551 30581 30611 30641 30671 30701 30731 30761 30791 30821 30851 30881 30911 30941 30971 31001 31031 31061 31091 31121 31151 31181 31211 31241 31271 31301 31331 31361 31391 31421 31451 31481 31511 31541 31571 31601 31631 31661 31691 31721 31751 31781 31811 31841 31871 31901 31931 31961 31991 32021 32051 32081 32111 32141 32171 32201 32231 32261 32291 32321 32351 32381 32411 32441 32471 32501 32531 32561 32591 32621 32651 32681 32711 32741 32771 32801 32831 32861 32891 32921 32951 32981 33011 33041 33071 33101 33131 33161 33191 33221 33251 33281 33311 33341 33371 33401 33431 33461 33491 33521 33551 33581 33611 33641 33671 33701 33731 33761 33791 33821 33851 33881 33911 33941 33971 34001 34031 34061 34091 34121 34151 34181 34211 34241 34271 34301 34331 34361 34391 34421 34451 34481 34511 34541 34571 34601 34631 34661 34691 34721 34751 34781 34811 34841 34871 34901 34931 34961 34991 35021 35051 35081 35111 35141 35171 35201 35231 35261 35291 35321 35351 35381 35411 35441 35471 35501 35531 35561 35591 35621 35651 35681 35711 35741 35771 35801 35831 35861 35891 35921 35951 35981 36011 36041 36071 36101 36131 36161 36191 36221 36251 36281 36311 36341 36371 36401 36431 36461 36491 36521 36551 36581 36611 36641 36671 36701 36731 36761 36791 36821 36851 36881 36911 36941 36971 37001 37031 37061 37091 37121 37151 37181 37211 37241 37271 37301 37331 37361 37391 37421 37451 37481 37511 37541 37571 37601 37631 37661 37691 37721 37751 37781 37811 37841 37871 37901 37931 37961 37991 38021 38051 38081 38111 38141 38171 38201 38231 38261 38291 38321 38351 38381 38411 38441 38471 38501 38531 38561 38591 38621 38651 38681 38711 38741 38771 38801 38831 38861 38891 38921 38951 38981 39011 39041 39071 39101 39131 39161 39191 39221 39251 39281 39311 39341 39371 39401 39431 39461 39491 39521 39551 39581 39611 39641 39671 39701 39731 39761 39791 39821 39851 39881 39911 39941 39971 40001 40031 40061 40091 40121 40151 40181 40211 40241 40271 40301 40331 40361 40391 40421 40451 40481 40511 40541 40571 40601 40631 40661 40691 40721 40751 40781 40811 40841 40871 40901 40931 40961 40991 41021 41051 41081 41111 41141 41171 41201 41231 41261 41291 41321 41351 41381 41411 41441 41471 41501 41531 41561 41591 41621 41651 41681 41711 41741 41771 41801 41831 41861 41891 41921 41951 41981 42011 42041 42071 42101 42131 42161 42191 42221 42251 42281 42311 42341 42371 42401 42431 42461 42491 42521 42551 42581 42611 42641 42671 42701 42731 42761 42791 42821 42851 42881 42911 42941 42971 43001 43031 43061 43091 43121 43151 43181 43211 43241 43271 43301 43331 43361 43391 43421 43451 43481 43511 43541 43571 43601 43631 43661 43691 43721 43751 43781 43811 43841 43871 43901 43931 43961 43991 44021 44051 44081 44111 44141 44171 44201 44231 44261 44291 44321 44351 44381 44411 44441 44471 44501 44531 44561 44591 44621 44651 44681 44711 44741 44771 44801 44831 44861 44891 44921 44951 44981 45011 45041 45071 45101 45131 45161 45191 45221 45251 45281 45311 45341 45371 45401 45431 45461 45491 45521 45551 45581 45611 45641 45671 45701 45731 45761 45791 45821 45851 45881 45911 45941 45971 46001 46031 46061 46091 46121 46151 46181 46211 46241 46271 46301 46331 46361 46391 46421 46451 46481 46511 46541 46571 46601 46631 46661 46691 46721 46751 46781 46811 46841 46871 46901 46931 46961 46991 47021 47051 47081 47111 47141 47171 47201 47231 47261 47291 47321 47351 47381 47411 47441 47471 47501 47531 47561 47591 47621 47651 47681 47711 47741 47771 47801 47831 47861 47891 47921 47951 47981 48011 48041 48071 48101 48131 48161 48191 48221 48251 48281 48311 48341 48371 48401 48431 48461 48491 48521 48551 48581 48611 48641 48671 48701 48731 48761 48791 48821 48851 48881 48911 48941 48971 49001 49031 49061 49091 49121 49151 49181 49211 49241 49271 49301 49331 49361 49391 49421 49451 49481 49511 49541 49571 49601 49631 49661 49691 49721 49751 49781 49811 49841 49871 49901 49931 49961 49991 50021 50051 50081 50111 50141 50171 50201 50231 50261 50291 50321 50351 50381 50411 50441 50471 50501 50531 50561 50591 50621 50651 50681 50711 50741 50771 50801 50831 50861 50891 50921 50951 50981 51011 51041 51071 51101 51131 51161 51191 51221 51251 51281 51311 51341 51371 51401 51431 51461 51491 51521 51551 51581 51611 51641 51671 51701 51731 51761 51791 51821 51851 51881 51911 51941 51971 52001 52031 52061 52091 52121 52151 52181 52211 52241 52271 52301 52331 52361 52391 52421 52451 52481 52511 52541 52571 52601 52631 52661 52691 52721 52751 52781 52811 52841 52871 52901 52931 52961 52991 53021 53051 53081 53111 53141 53171 53201 53231 53261 53291 53321 53351 53381 53411 53441 53471 53501 53531 53561 53591 53621 53651 53681 53711 53741 53771 53801 53831 53861 53891 53921 53951 53981 54011 54041 54071 54101 54131 54161 54191 54221 54251 54281 54311 54341 54371 54401 54431 54461 54491 54521 54551 54581 54611 54641 54671 54701 54731 54761 54791 54821 54851 54881 54911 54941 54971 55001 55031 55061 55091 55121 55151 55181 55211 55241 55271 55301 55331 55361 55391 55421 55451 55481 55511 55541 55571 55601 55631 55661 55691 55721 55751 55781 55811 55841 55871 55901 55931 55961 55991 56021 56051 56081 56111 56141 56171 56201 56231 56261 56291 56321 56351 56381 56411 56441 56471 56501 56531 56561 56591 56621 56651 56681 56711 56741 56771 56801 56831 56861 56891 56921 56951 56981 57011 57041 57071 57101 57131 57161 57191 57221 57251 57281 57311 57341 57371 57401 57431 57461 57491 57521 57551 57581 57611 57641 57671 57701 57731 57761 57791 57821 57851 57881 57911 57941 57971 58001 58031 58061 58091 58121 58151 58181 58211 58241 58271 58301 58331 58361 58391 58421 58451 58481 58511 58541 58571 58601 58631 58661 58691 58721 58751 58781 58811 58841 58871 58901 58931 58961 58991 59021 59051 59081 59111 59141 59171 59201 59231 59261 59291 59321 59351 59381 59411 59441 59471 59501 59531 59561 59591 59621 59651 59681 59711 59741 59771 59801 59831 59861 59891 59921 59951 59981 60011 60041 60071 60101 60131 60161 60191 60221 60251 60281 60311 60341 60371 60401 60431 60461 60491

[web-test-datasets.comverted.html](#)

[illegible]

Figure 1. The effect of the number of nodes on the number of nodes in the network. The number of nodes in the network is plotted against the number of nodes in the network. The number of nodes in the network is plotted against the number of nodes in the network.

```
import nltk, string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer

readfiles = ["webkb-test-stemmed.txt", "webkb-train-stemmed.txt"]
writefiles = ["webkb-test-stemmed-converted.txt", "webkb-train-stemmed-converted.txt"]
documents = []
documentsreorganized = []
dictionary = {}
```

```

def tokenize(text):
    stop_words = set(stopwords.words("english"))
    exclude = set(string.punctuation)
    stemmer = PorterStemmer()
    text = ''.join(ch for ch in text if ch not in exclude)
    document = nltk.word_tokenize(text)
    document_filtered = []
    for word in document:
        if word not in stop_words:
            document_filtered.append(word.lower())
    return [stemmer.stem(item) for item in document_filtered]

def reconcile(documents):
    for i in range(0, len(documents)):
        document = []
        for j in range(0, len(documents[i])):
            string = ""
            strings = {}
            for k in range(0, len(documents[i][j])):
                if(k == 0):
                    keyword = documents[i][j][k]
                    if(keyword[0] == "f"):
                        string += "1 "
                    elif(keyword[0] == "c"):
                        string += "2 "
                    elif(keyword[0] == "s"):
                        string += "3 "
                    elif(keyword[0] == "p"):
                        string += "4 "
                else:
                    if dictionary[documents[i][j][k]] in strings.keys():
                        count = strings[dictionary[documents[i][j][k]]]
                        strings[dictionary[documents[i][j][k]]] = count + 1
                    else:
                        strings[dictionary[documents[i][j][k]]] = 1
            maximum = 0
            for key in strings:
                if key > maximum:
                    maximum = key;
            for k in range(0, maximum):
                if k in strings:
                    string += str(k) + ":" + str(strings[k]) + " "
            string = string[:-1]
            if len(string) > 3:
                document.append(string)
            documentsreorganized.append(document)

def createdictionary():
    tracker = 10
    for i in range(0, len(documents)):
        for j in range(0, len(documents[i])):
            for k in range(0, len(documents[i][j])):
                if documents[i][j][k] in dictionary.keys():
                    continue
                else:
                    dictionary[documents[i][j][k]] = tracker
                    tracker += 1

def loadfile(readfiles):
    for i in range(0, len(readfiles)):
        with open("files\\" + str(readfiles[i]), "r") as file:
            document = file.read().split("\n")
    context = []
    for i in range(0, len(document)):
        sentence = tokenize(document[i])
        context.append(sentence)
    documents.append(context)

def readfile(path):
    with open(path, "r") as file:
        document = file.read().split("\n")
    return document

def writefile():
    for i in range(0, len(writefiles)):
        with open("files\\" + str(writefiles[i]), "w") as file:
            for j in range(0, len(documentsreorganized[i])):
                file.write(documentsreorganized[i][j] + "\n")

print("Reading Files")
loadfile(readfiles)

```

```
print("Creating Dictionary")
createdictionary()
print("Reconciling Files With Dictionary")
reconcile(documents)
print("Writing Files")
writefile()
```