# HW6 보고서

### 2014105024 컴퓨터학부 김재구

## A. SVM 실행 보고서

1. svm.py를 실행하여 svm\_learn에 대한 input, svm\_classify을 위한 input을 생성한다

svm input을 생성하는 함수는 다음과 같다.

```
def make_svm_input(input_file_name,stopword_list,stemmer,output_file_name,tf_idf_dic,all_word_list):
    input_file = open(input_file_name)
    output_file = open(output_file_name,'w')
    for line in input_file:
        # line list contains => first element : class , last element : newline
        output_str = ''
        line_list = re.split('\t| ', line)
        delete_stopword(line_list, stopword_list)
        UseStemmer(line_list, stemmer)
        word_cnt = len(line_list) - 2
        #extract words from line_list AND make sorted list ascending order
        temp_list = line_list[1:word_cnt+1]
        sorted list = []
        for word in all_word_list:
            if word in temp_list:
                sorted_list.append(word)
        #make output stirng
        if line_list[0] == 'acq':
            output_str += '1'
        else:
            output_str += '-1'
        for word in sorted list:
            index = all_word_list.index(word)+1
            if tf_idf_dic[word] == 0:
                continue
            output_str += ' '+str(index) + ':' + str(tf_idf_dic[word])
        output\_str^-+= '\n'
        output_file.write(output_str)
    return
```

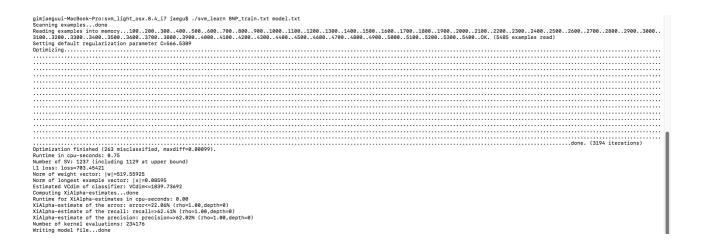
수행결과 다음과 같은 input이 생성된다.

#### <BNP\_train.txt>

-1 5:0.006925270059623867 6:0.006101287799200853 8:0.016384554046443005 9:0.016393850816898536 10:0.019789665297985842 12:0.00886691997971888 16:0.016384554046443005 10:0.019789665297985842 12:0.00886691997971888 15:0.0023581534416526914 17:0.006394307471101231 18:0.02174951083698346 -1 8:0.016384554046443005 10:0.019789665297985842 18:0.021749510836983465 25:0.03353862540732 35:0.0074010424230205344 38:0.016767682627646205 4 1 8:0.016384554046443005 9:0.016393850816898536 10:0.019789665297985842 17:0.006394307471101231 18:0.021749510836983465 -1 8:0.016384554046443005 9:0.0163938508168985365 10:0.019789665297985842 17:0.006394307471101231 18:0.021749510836983465 25:0.03353862540732 28: -1 8:0.016384554046443005 9:0.016393859816898536 10:0.019789665297985842 17:0.006394307471101231 12:0.021749510836983465 25:0.033538662540732 28: -1 17:0.006394307471101231 26:0.021749510836983465 25:0.033538662540732 28: -1 17:0.006394307471101231 24:0.009482539841247844 25:0.03353862540732 28:0.009686930293337842 30: -1 8:0.016384554046443005 9:0.016393858016998536 10:0.019789665297985842 17:0.006394307471101231 18:0.021749510836983465 25:0.03353862540732 26: -1 8:0.016384554046443005 9:0.016384554046443005 10:0.019789665297985842 17:0.006394307471101231 18:0.021749510836983465 25:0.03353862540732 26: -1 8:0.016384554046443005 10:0.019789665297985842 17:0.006394307471101231 25:0.03533862540732 28:0.009686930293337842 36: -1 6:0.006101287799200853 8:0.016384554046443005 10:0.019789665297985842 17:0.006394307471101231 25:0.03533862540732 28:0.009686930293337842 32: -1 6:0.006101287799200853 9:0.016393858016998536 17:0.006394307471101231 25:0.012593160355604733 28:0.009686930293337842 32: -1 8:0.01638455404643005 10:0.019789665297985842 17:0.006394307471101231 25:0.03533862540732 28:0.009686930293337842 31: -1 8:0.01638455404643005 10:0.019789665297985842 17:0.006394307471101231 25:0.03533862540732 28:0.009686930293337842 31: -1 8:0.01638455404643005 10:0.019789665297985842 17:0.006394307471101231 25:0.035333862540732 28:0.009686930293337842 31: -1 8:0.016384554046443005 10:0.019789665297985842 18:0.02174951083693465 25:0.035333862540732 28:0.009686930293337842 35: 0.0074010442730025344 31: -1 8:0.016384554046443005 10:0.019789665297985842 18:0.02174951083693465 25:0.035333862540732 28:0.0096866930293337842 35: 0.0074010442730025344 31: -1 8:0.016384554046443005 10:0.019789665297985842 12:0.00886691997971888 17: 0.00639450779200853 7: -0.00835179893452374 8: 0.016384554046443005 10: -0.019789665297985842 12: 0.00886691997971888 17: 0.00639450779200853 7: 0.00835179893452374 8: 0.016384554046443005 9: 0.0163893850816898536 10: 0.019789665297985842 12: 0.00886691997971888 17: 0.00639450779200853 7: 0.0083517989345 -1 8:0.016384554046443005 9:0.016393850816898536 29:0.0019789665297985842 12:0.00886691997971888 17:0.006394307471101231 18:0.021749510836983465 2 1 5:0.006925270059623867 6:0.006101287799200853 8:0.016384554046443005 9:0.016393850816898536 10:0.019789665297985842 12:0.00886691997971888 16:-1 8:0.016384554046443005 18:0.022381698942324506 41:0.0037066875709585263 -1 8:0.016384554046443005 17:0.006394307471101231 18:0.021749510836983465 25:0.03353862540732 26:0.012593168352602473 40:0.019268873882164243 44 1 1:0.0022120041805741065 7:0.00815179893462374 8:0.016384554046443005 9:0.016393850816898536 10:0.019789665297985842 12:0.00886691997971888 25: -1 7:0.00815179893462374 8:0.016384554046443005 10:0.019789665297985842 18:0.021749510836983465 24:0.009482539841247844 25:0.03533862540732 26:0 -1 16:0.003842246855271652 82:0.001189404408621469 96:0.000845177729611017 103:0.008252190416672647 15:0.0004022423421329334 229:0.00068918507 16:0.006101287799200853 7:0.00815179893462374 8:0.016384554046443005 10:0.019789665297985842 14:0.009400377083638752 18:0.021749510836983465 24 1 6:0. 016384554046443005 103:0.00252190416672647 305:0.0033347224597954353 357:0.001931661370331603 420:0.008703935648928405 776:0.00160379451
1:0.0022120041805741065 6:0.006101287799200853 7:0.00815179893462374 8:0.016384554046443005 9:0.016393850816898536 10:0.019789665297985842 17:
-1 1:0.0022120041805741065 7:0.00815179893462374 96:0.0008451727296910197 276:0.0006159909698862033 545:0.0009191817490044648 630:0.000759065239
-1 8:0.016384554046443005 10:0.019789665297985842 18:0.021749510836983465 25:0.035353862540732 38:0.016767682627646205 44:0.013149623110788171 49
-1 7:0.00815179893462374 630:0.0007590652399220948 1374:0.0003548551460239322 3532:8.7077607966050290-05 6785:0.0001050805195912532 7040:9.10823
-1 8:0.016384554046443005 10:0.019789665297985842 18:0.021749510836983465 25:0.03353862540732 28:0.009686930293337842 30:0.0022381698942324506 3 31 32 33 17:0.006394307471101231 24:0.009482539841247844 96:0.0008451727296910197 103:0.000252190416672647 776:0.0006159909698862033 545:0.00091918174 6:0.006101287799200853 7:0.00815179893462374 17:0.006394307471101231 96:0.0008451727296910197 630:0.0007590652399220948 732:0.000552302898173 7:0.00815179893462374 17:0.006394307471101231 30:0.0022381698942324596 96:0.0008451727296910197 226:0.0011315820741369753 276:0.0006159909698 24:0.009482539841247844 84:0.0035845380169542283 103:0.008252190416672647 630:0.0007590652399220948 1374:0.0003548551460239322 3532:8.7077607 6:0.006101287799200853 7:0.00815179893462374 8:0.016384554046443005 10:0.019789665297985842 17:0.006394307471101231 18:0.021749510836983465 2 39 40 41 17:0.00815179893462374 8:0.016384554046443005 10:0.019789665297985842 18:0.021749510836983465 25:0.03353862540732 26:0.012593168352602473 31:0. -1 16:0.003842246855271652 96:0.0008451727296910197 103:0.008252190416672647 135:0.0030232377675683876 545:0.0009191817490044648 732:0.000552302

#### <BNP\_test.txt>

### 2. BNP\_train.txt로 svm\_learn 작업을 수행한다.



svm train 작업 결과로 생성된 model.txt는 다음과 같다.

```
model.txt ~
SVM-light Version V6.02
0 # kernel type
3 # kernel parameter -d
1 # kernel parameter -g
1 # kernel parameter -s
1 # kernel parameter -r
empty# kernel parameter -u
8172 # highest feature index
5485 # number of training documents
1238 # number of support vectors plus 1
1.1544428 # threshold b, each following line is a SV (starting with alpha*y)
-566.53888571579159361135680228472 8:0.016384553 10:0.019789666 18:0.021749511
25:0.033538625 38:0.016767683 44:0.013149623 50:0.0055565289 54:0.01810213
96:0.00084517273 124:0.00035734026 126:0.0059264963 127:0.0013702284 133:0.0026865716
147:0.0022198537 200:0.0052741929 201:0.014689588 299:0.0022997088 305:0.0033347225
313:0.0029094089 358:0.0010465623 430:0.0012855664 458:0.0020801506 459:0.0016814971
479:0.0034885367 529:0.00045332962 542:0.00629802 608:0.0036667094 688:0.0010800606
730:0.0013370522 735:0.0015095737 765:0.00076288299 770:0.0021074559 805:0.0017250171
842:0.0023054637 904:0.0010132007 943:0.0016537941 952:0.003111905 996:0.0025563736
1083:0.00030468579 1164:0.00024531482 1214:0.0017381095 1436:7.5814693e-05
1461:0.00084322499 1469:0.00013129524 1602:0.00028902653 1608:0.0013359392
1862:0.0020361268 1920:0.00016296659 2125:0.00019040664 2322:0.00025224363
2529:0.00013224162 2530:0.00013224162 3287:0.00022425821 3449:5.3055381e-05
3457:0.00023670334 5767:8.3212457e-05 5974:8.8161076e-05 #
566.53888571579159361135680228472 5:0.00692527 6:0.0061012879 8:0.016384553 10:0.019789666
14:0.0094003771 18:0.021749511 25:0.033538625 26:0.012593169 35:0.0074010426
38:0.016767683 40:0.019268874 54:0.01810213 77:0.0016064595 80:0.0052598831
96: 0.00084517273 \ 131: 0.0044333953 \ 132: 0.0012673686 \ 197: 0.0035134612 \ 201: 0.014689588
212:0.00053734303 267:0.005499742 291:0.0027949237 327:0.0027215697 348:0.0031609824
357:0.0017931662 363:0.0021962137 514:0.0021314137 550:0.0075062169 556:0.0032876073
```

### 3. BNP\_test.txt와 model.txt로 svm\_classify 작업을 수행

gimjaeguui-MacBook-Pro:svm\_light\_osx.8.4\_i7 jaegu\$ ./svm\_classify BNP\_test.txt model.txt output.txt
Reading model...OK. (1237 support vectors read)
Classifying test examples.100 ...000...300 ...600 ...600 ...600 ...600 ...900 ...1000 ...1200 ...1200 ...1300 ...1400 ...1500 ...1600 ...1700 ...1800 ...1900 ...2000 ...2000 ...2100 ...done
Runtime (without IO) in cpu-seconds: 0.00
Accuracy on test set: 95.48% (2090 correct, 99 incorrect, 2189 total)
Precision/recall on test set: 97.61%/67.93%

작업 수행결과 95.48%의 정확도가 나옴을 확인할 수 있다.

## B. Adaboost 구현코드 및 실행결과

#### 1. model1, model2 생성

input 파일을 생성하기 전에 svm\_classify의 결과로 생성된 output.txt 바탕으로 에러를 계산 해야 하므로 output.txt를 내용으로 리스트를 생성한다.

```
def make_output_list(input_file_name):
    input_file = open(input_file_name)
    output_list = []

    for output_val in input_file:
        if float(output_val) > 0:
            output_list.append(1)
        else:
            output_list.append(-1)

    return output_list
```

그리고 가중치를 곱한후 새로운 input 파일생성하고 model2를 생성한다. weight를 구하는 함수는 다음과 같다.

-1 5:0.006928274358863014 6:0.006103934641004152 8:0.016391661943006126 9:0.016400962746557862 10:0.019798250389380954 12:0.008870766599521327v1 -1 8:0.016391661943006126 10:0.019798250389380954 12:0.008870766599521327 15:0.002359176448485573 17:0.006397081429792359 18:0.02175894614250895 -1 8:0.016391661943006126 10:0.019799250389380954 18:0.021758946142508995 25:0.03355317503006766 35:0.007404253120354925 38:0.0167749567318665503 1 8:0.01639243605003192 9:0.016401737292820226 10:0.019799185374788662 16:0.0038440952182735924 17:0.006397383536174041 18:0.02175997372306436 2 1 8:0.01639243005003192 9:0.016400962746557862 10:0.0197981537/786052 10:0.08440952182735924 17:0.000397383536174041 18:0.02175994312300436 2 - 18:0.016391661943006126 9:0.016400962746557862 10:0.0197998250389380954 17:0.006397081429792359 18:0.021758946142508995 25:0.0335317503006766 - 1 17:0.006397081429792359 26:0.012598631481949694 38:0.016774956731865503 65:0.0026074378983894915 87:0.003361028420266549 89:0.004043540670926 1 8:0.016392430605003192 10:0.019799185374788662 17:0.006397383536174041 24:0.009487101540812498 25:0.033554759600846314 28:0.00969159931759735 3 - 1 8:0.016391661943006126 9:0.016400962746557862 10:0.019798250389380954 17:0.006397081429792359 18:0.021758946142508995 25:0.03355317503006766 - 1 8:0.016391661943006126 10:0.019798250389380954 16:0.00384391368919302 17:0.006397081429792359 18:0.021758946142508995 20:0.00241241313008968 -1 6:0.006103934641004152 8:0.016391661943006126 10:0.019798250389380954 17:0.006397081429792359 25:0.03355317503006766 28:0.009691132647478456 -1 6:0.006103934641004152 9:0.016400962746557862 17:0.006397081429792359 26:0.012598631481949694 30:0.0022391408501738694 44:0.01315532763961761 -1 8:0.016391661943006126 10:0.019798250389380954 18:0.021758946142508995 25:0.03355317503006766 28:0.009691132647478456 35:0.007404253120354925  $16: 0.006103934641004152 \ 7: 0.008155335322826092 \ 8: 0.016391661943006126 \ 9: 0.016400962746557862 \ 10: 0.019798250389380954 \ 12: 0.008870766599521327 \ 17-1 \ 8: 0.016391661943006126 \ 9: 0.016400962746557862 \ 29: 0.002125236619345801 \ 197: 0.0035149854270207395 \ 633: 0.0031477602615541827 \ 644: 0.00173329987065 \ 10: 0.016400962746557862 \ 10: 0.01640096274657862 \ 10: 0.01640096274657862 \ 10: 0.01640096274657862 \ 10: 0.01640096274657862 \ 10: 0.01640096274657$ 1 8:0.01639243605003192 9:0.016401737292820226 10:0.019799185374788662 12:0.008871185527313563 17:0.006397383536174041 18:0.02175997372306436 24  $\frac{-1}{6:0.006103934641004152} \ 8:0.021758946142508995 \ 25:0.03355317503006766 \ 26:0.0127598631481949694 \ 29:0.002125236619345801 \ 30:0.002239140059173896 \ -1 \ 6:0.006103934641004152 \ 8:0.016391661943006126 \ 10:0.019798250389380954 \ 18:0.021758946142508995 \ 25:0.03355317503006766 \ 28:0.009691132647478456$ 1 8:0.016391661943006126 9:0.016400962746557862 10:0.019798250389380954 14:0.009404455126152173 16:0.003843913686919302 25:0.03355317503006766 2 1 6:0.006103934641004152 8:0.016391661943006126 17:0.006397081429792359 18:0.021758946142508995 25:0.03355317503006766 26:0.012598631481949694 3 -1 8:0.016391661943006126 9:0.016400962746557862 10:0.019798250389380954 12:0.008870766599521327 16:0.003843913686919302 18:0.021758946142508995 -1 8:0.016391661943006126 10:0.019798250389380954 18:0.02175894614258095 25:0.03355317593006766 28:0.009691132647478456 35:0.008870766599521327 16
-1 8:0.01639243605003192 18:0.02175997372306436 25:0.033554759600846314 26:0.012599226460615034 30:0.002239246595074532 41:0.0037084707213973153 -1 8:0.016391661943006126 17:0.006397081429792359 18:0.021758946142508995 25:0.03355317503006766 26:0.012598631481949694 40:0.01927723304543802 1:0.00221296378538655 7:0.008155335322826092 8:0.016391661943006126 9:0.016400962746557862 10:0.019798250389380954 12:0.008870766599521327 25: -1 7:0.008155335322826092 8:0.016391661943006126 10:0.019798250389380954 18:0.021758946142508995 24:0.00948653527354661 25:0.03355317503006766 -1 16:0.003843913686919302 82:0.0011899203923634882 96:0.0008455393799107074 103:0.008255770356397258 156:0.000402416841661763 229:0.00068948405 1 6:0.006103934641004152 7:0.008155335322826092 8:0.016391661943006126 10:0.019798250389380954 14:0.009404455126152173 18:0.021758946142508995 2 1 8:0.016391661943006126 103:0.008255770356397258 305:0.003336169109775581 357:0.001793944042825924 420:0.008707711563372538 776:0.0016094927370 1:0.00221296378538655 6:0.006103934641004152 7:0.008155335322826092 8:0.016391661943006126 9:0.016400962746557862 10:0.019798250389380954 17:0.016400962746557862 7:0.008155335322826092 96:0.0008455393799107074 276:0.0006162581971836541 545:0.0009195805055880198 630:0.0007593945352 -1 8:0.016391661943006126 10:0.019798250389380954 18:0.021758946142508995 25:0.03355317503006766 38:0.016774956731865503 44:0.013155327639617614 -1 7:0.0081553355322826092 630:0.0007593945352568799 1374:0.0003550090881859025 3532:8.711538372463764e-05 6785:0.00010512610529727464 7040:9.112 -1 8:0.016391661943006126 10:0.019798250389380954 18:0.021758946142508995 25:0.03355317593006766 28:0.009691132647478456 30:0.009486653527354661 96:0.0008455393799107074 103:0.008255770356397258 276:0.0006162581971836541 545:0.00091958050 -1 6:0.006103934641004152 7:0.008155335322826092 17:0.006397081429792359 96:0.0008455393799107074 630:0.0007593945352568799 732:0.00055254249650 -1 7:0.008155335322826092 17:0.006397081429792359 30:0.0082391408501738694 96:0.0008455393799107074 226:0.0011320729735726784 276:0.000616258197 -1 24:0.009486653527354661 84:0.003586093050150668 103:0.008255770356397258 630:0.0007593945352568799 1374:0.0003550090881859025 3532:8.71153837  $-1 \;\; 6:0.006103934641004152 \;\; 7:0.008155335322826092 \;\; 8:0.016391661943006126 \;\; 10:0.019798250389380954 \;\; 17:0.006397081429792359 \;\; 18:0.0217589461425089951 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.0217589461425089951 \;\; 18:0.0217589461425089951 \;\; 18:0.021758946142508961 \;\; 18:0.021758946142508961 \;\; 18:0.021758961 \;\; 18:0.021758961 \;\; 18:0.021758961 \;\; 18:0.0217589619$ -1 7:0.008155335322826092 24:0.00948653357354661 26:0.01259863148194694 96:0.0008455393799107074 135:0.0094245492992255993 180:0.00845393799107074 135:0.00832557034661 26:0.01259863148194694 96:0.0008455393799107074 135:0.00944549992255993 180:0.00845393799107074 103:0.0084545393799107074 103:0.008455393799107074 103:0. 17:0.008155335322826092 8:0.016391661943006126 10:0.019798250389380954 18:0.021758946142508995 25:0.03355317503006766 26:0.012598631481949694 3 -1 16:0.003843913686919302 96:0.0008455393799107074 103:0.008255770356397258 135:0.0030245492992265993 545:0.0009195805055880198 732:0.000552542

# 2. model1, model2 의 알파값을 바탕으로 새로운 결과값을 계산하여 예측값과 비교한다.

```
gimjaeguui-MacBook-Pro:svm_light_osx.8.4_i7 jaegu$ ./svm_classify new_test.txt model.txt new_output.txt
Reading model...OK. (1237 support vectors read)
Classifying test examples..100..200..300..400..500..600..700..800..900..1000..1100..1200..1300..1400..1500..1600..1700..1800..1900..2000..2100..done
Runtime (without IO) in cpu-seconds: 0.01
Accuracy on test set: 95.48% (2000 correct, 99 incorrect, 2189 total)
Precision/recall on test set: 97.61%/87.93%
```

작업 수행결과 95.48%의 정확도가 나옴을 확인할 수 있다.