

**CMPE 442 MACHINE LEARNING**  
**Assignment 2**



**Arda Andırın**  
**17921934994**

**Q1)** The size of my dictionary is 150138

**Q2)** The class priors are computed as follows

```
alt.atheism : 0.042425
comp.graphics : 0.051617
comp.os.ms-windows.misc : 0.052236
comp.sys.ibm.pc.hardware : 0.052148
comp.sys.mac.hardware : 0.051087
comp.windows.x : 0.052413
misc.forsale : 0.051706
rec.autos : 0.052501
rec.motorcycles : 0.052855
rec.sport.baseball : 0.052766
rec.sport.hockey : 0.053032
sci.crypt : 0.052590
sci.electronics : 0.052236
sci.med : 0.052501
sci.space : 0.052413
soc.religion.christian : 0.052943
talk.politics.guns : 0.048259
talk.politics.mideast : 0.049850
talk.politics.misc : 0.041100
talk.religion.misc : 0.033322
```

**Q3)** I have used the most frequent 50000 words. I have observed that doing so accuracy rate increased. Also, the probabilities were going near zero so python couldn't make comparisons, this resulted in a low accuracy (around 40%). To solve this I multiplied each new probability of a word with  $10^4$ . This way the accuracy rate got up to 80%.

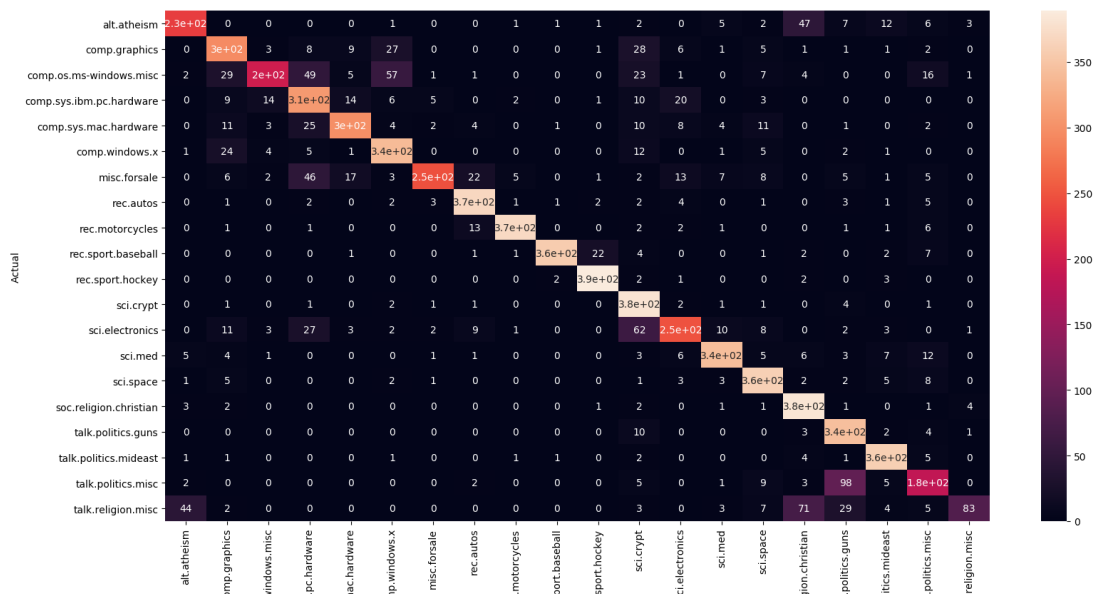
Here is our formula

$$\varphi_{k|y=0} = \frac{\sum_{i=1}^m \sum_{j=1}^{n_i} 1\{X_j^{(i)} = k \wedge Y^{(i)} = 0\} + 1}{\sum_{i=1}^m 1\{Y^{(i)} = 0\} n_i + |V|}$$

Here is the accuracy rate with the confusion matrix.

```
Number of correct predictions = 6087
Number of false predictions = 1445
Total number of tests = 7532

Accuracy rate is = 80.81518852894317 %
```



**Q4)** Yes for example the algorithm confuses the category talk.religion.misc with religion.christian. This is because the subjects are so related that appearance of a word in one of these categories appear also in the other category. The words god, atheism, jesus etc. are all common in religion related newsgroups, it is normal to see false positives. This can also be observed in other categories that are related to each other. For exapmle, some newsgroups in comp categories are also predicted wrong.

**Q5)** When we think of the strong indicators, what comes to the mind first is the most common words. Preprocessor already eliminates common words such as stopwords. What would be a really strong indicator is words that appear exclusively in that newsgroup. I couldn't find a way to find those. So here are the top 10 words for each category.

-----  
alt.atheism  
-----

god : 765  
one : 722  
people : 576  
writes : 562  
subject : 542  
line : 538  
would : 501  
organization : 472  
dont : 458  
atheist : 446  
-----

comp.graphics  
-----

line : 774  
image : 758  
subject : 627  
organization : 581  
file : 554  
graphic : 460  
university : 360  
program : 329  
would : 295  
x : 279  
-----

comp.os.ms-windows.misc  
-----

maxaxaxaxaxaxaxaxaxaxaxaxaxax : 3317  
window : 1082  
line : 654  
file : 641  
subject : 618  
organization : 579  
driver : 375  
university : 333  
use : 323  
problem : 318  
-----

comp.sys.ibm.pc.hardware  
-----

drive : 976  
scsi : 704  
line : 637  
subject : 610  
organization : 583  
card : 478  
system : 410  
mb : 388  
one : 379  
disk : 353  
-----

comp.sys.mac.hardware  
-----

line : 652  
subject : 591  
organization : 552  
mac : 546  
apple : 420  
problem : 372  
drive : 357  
one : 350  
university : 308  
nntppostinghost : 295  
-----

comp.windows.x  
-----

x : 5152  
window : 933  
line : 857  
subject : 788  
file : 778  
organization : 599  
program : 553  
use : 508  
widget : 503  
server : 474  
-----

misc.forsale  
-----

line : 628  
subject : 601  
organization : 573  
sale : 560  
new : 329  
university : 326  
offer : 273  
nntppostinghost : 263  
distribution : 252  
email : 246  
-----

rec.autos  
-----

car : 1223  
line : 642  
subject : 625  
organization : 589  
writes : 484  
article : 453  
would : 430  
one : 362  
like : 327  
dont : 322  
-----

rec.motorcycles  
-----

bike : 688  
line : 638  
organization : 612  
subject : 611  
writes : 503  
article : 473  
dod : 455  
one : 394  
like : 334  
nntppostinghost : 303  
-----

|   |  |  |  |
|---|--|--|--|
| rec.sport.baseball<br>-----<br>line : 648<br>subject : 618<br>organization : 601<br>year : 592<br>game : 558<br>writes : 468<br>team : 432<br>article : 391<br>player : 364<br>run : 343<br>-----<br>rec.sport.hockey<br>-----<br>team : 954<br>game : 917<br>line : 707<br>subject : 635<br>organization : 611<br>hockey : 594<br>player : 522<br>play : 491<br>year : 438<br>would : 429<br>-----<br>sci.crypt<br>-----<br>key : 1449<br>encryption : 840<br>chip : 839<br>would : 707<br>clipper : 705<br>line : 693<br>system : 668<br>subject : 665<br>one : 642<br>organization : 621 | sci.electronics<br>-----<br>line : 782<br>subject : 672<br>organization : 581<br>one : 481<br>use : 371<br>would : 367<br>writes : 301<br>university : 291<br>like : 275<br>nntppostinghost : 268<br>-----<br>sci.med<br>-----<br>subject : 649<br>line : 619<br>organization : 610<br>one : 567<br>article : 462<br>writes : 439<br>would : 419<br>people : 322<br>msg : 312<br>dont : 311<br>-----<br>sci.space<br>-----<br>space : 1200<br>line : 646<br>subject : 635<br>organization : 631<br>would : 553<br>writes : 452<br>one : 413<br>nasa : 407<br>article : 402<br>launch : 377 | soc.religion.christian<br>-----<br>god : 1477<br>one : 817<br>would : 775<br>christian : 755<br>subject : 671<br>people : 656<br>line : 636<br>jesus : 618<br>organization : 559<br>say : 520<br>-----<br>talk.politics.guns<br>-----<br>gun : 1231<br>would : 819<br>people : 657<br>line : 608<br>subject : 584<br>organization : 555<br>one : 531<br>writes : 507<br>article : 492<br>right : 487<br>-----<br>talk.politics.mideast<br>-----<br>armenian : 1268<br>people : 996<br>one : 883<br>israel : 879<br>israeli : 791<br>turkish : 712<br>would : 711<br>subject : 661<br>jew : 630<br>line : 626 | talk.politics.misc<br>-----<br>would : 692<br>people : 684<br>writes : 611<br>q : 575<br>article : 573<br>line : 525<br>one : 515<br>dont : 508<br>organization : 506<br>subject : 500<br>-----<br>talk.religion.misc<br>-----<br>god : 516<br>one : 463<br>subject : 418<br>line : 418<br>people : 410<br>organization : 403<br>christian : 401<br>jesus : 390<br>would : 387<br>writes : 357 |
|---|--|--|--|

**Q6)** There are some ways to increase the accuracy of the classifier. One of them is with using weights. One of the most popular weight implementation in text classification is “term frequency- inverse term frequency.” It assumes that, how important a word is inversely proportional to how often it occurs across all documents.

Another way to increase the accuracy is removing common words. These words do not have any meaning for us because it appears in all the categories. For example one, would, writes, etc.