

# CS 440: NAÏVE BAYES CLASSIFICATION

## PART 1: DIGIT CLASSIFICATION

### 1.1 Single Pixels as Features

Our implementation is quite simple. We created an object that would hold matrices for the binary 0 (white) and 1 (black) for each class. These matrices would be the size of an image. This meant that when we encountered a feature, we could simply find that exact location in a specific matrix and increment the count. For example, if we were looking at pixel (1,1) in a training image that represented the digit class 5 and the pixel was black, we would go to the black matrix for class 5 and increment its count. We followed this procedure for every feature of every training image. After doing this, we compared the testing data against the populated matrices for every single class, using the method described in the Assignment handout. For each testing image, we calculated its posterior probability for each class. We labelled it using the highest posterior probability. Our results are shown below.

We chose to use a smoothing factor of 1. This means that we assume that we have seen every value once more than its actual occurrence value. This made our results more accurate than using no smoothing factor. We chose not to use a larger smoothing factor; we did not want to over-count occurrences by more than 1. Using a smoothing factor of 1 means that if a particular feature never occurred for a given class, we could still operate under the assumption that it did occur once. This allows uncommon features to still be represented. Using a larger smoothing factor would have overrepresented features that do not occur very often. We did, however, build functionality to play around with smoothing factor. We made it possible to vary smoothing factor very easily. In our implementation, smoothing factor can be specified, which allows us to test different smoothing factors very easily. We found the best results with a smoothing factor of 1, as per our hypothesis. We want to represent all features, even if they occur uncommonly, but we do not want to over-represent uncommon features, which larger smoothing factors would do.

### *Results, as referenced above*

Our overall accuracy rate for the digit data set is **77.1%**

Below are the classification rates for each digit. The format is as follows:

*(digit class, classification rate – truncated after 3 decimal points)*

(0, 0.844)  
(1, 0.962)  
(2, 0.776)  
(3, 0.790)  
(4, 0.766)  
(5, 0.673)  
(6, 0.758)  
(7, 0.726)  
(8, 0.601)  
(9, 0.800)

Below is the confusion matrix. Percentages are truncated after 2 decimal points.

	0:	1:	2:	3:	4:	5:	6:	7:	8:	9:
0:	84.44%	0.00%	1.11%	0.00%	1.11%	5.56%	3.33%	0.00%	4.44%	0.00%
1:	0.00%	96.30%	0.93%	0.00%	0.00%	1.85%	0.93%	0.00%	0.00%	0.00%
2:	0.97%	2.91%	77.67%	3.88%	0.97%	0.00%	5.83%	0.97%	4.85%	1.94%
3:	0.00%	2.00%	0.00%	79.00%	0.00%	3.00%	2.00%	6.00%	2.00%	6.00%
4:	0.00%	0.93%	0.00%	0.00%	76.64%	0.00%	2.80%	0.93%	1.87%	16.82%
5:	2.17%	2.17%	1.09%	13.04%	3.26%	67.39%	1.09%	1.09%	2.17%	6.52%
6:	1.10%	6.59%	4.40%	0.00%	4.40%	5.49%	75.82%	0.00%	2.20%	0.00%
7:	0.00%	5.66%	2.83%	0.00%	2.83%	0.00%	0.00%	72.64%	2.83%	13.21%
8:	1.94%	0.97%	2.91%	13.59%	1.94%	5.83%	0.00%	0.97%	60.19%	11.65%
9:	1.00%	1.00%	1.00%	3.00%	9.00%	2.00%	0.00%	2.00%	1.00%	80.00%

Below are the test examples with the highest and lowest posterior probabilities for each digit class. White regions are represented as zeros and colored regions are represented as ones. See below.

Category: 0  
High Frequency: -146.323892153

[illegible]

```
Category:    0
Low Frequency: -552.720115536
```

[illegible]



Category: 1  
High Frequency: -71.6348012014

[illegible]

Category: 1  
Low Frequency: -902.826929454

A 32x100 grid of binary code (0s and 1s) on a black background. The code is arranged to form a stylized 'A' shape. The top 16 rows consist of all 0s. The next 16 rows contain a mix of 0s and 1s, with the 1s forming the vertical bars and the horizontal crossbar of the 'A'. The bottom 16 rows consist of all 0s. The overall effect is a digital, pixelated representation of the letter 'A'.

Category: 2  
High Frequency: -169.157830278

[illegible]



Category: 2  
Low Frequency: -482.154509687

[illegible]



Category: 3  
High Frequency: -145.279996263

[illegible]

```
Category:    3
Low Frequency: -547.318118526
```

[illegible]



Category: 4  
High Frequency: -139.655360285

[illegible]



Category: 4  
Low Frequency: -511.213498046

[illegible]

```
Category:    5
High Frequency: -163.633260013
```

[illegible]



Category: 5  
Low Frequency: -489.904511857

[illegible]



```
Category:    6
High Frequency: -120.651413332
```

[illegible]

Category: 6  
Low Frequency: -619.121045481

[illegible]



Category: 7  
High Frequency: -116.163705357

[illegible]



Category: 7  
Low Frequency: -632.810357967

[illegible]

Category: 8  
High Frequency: -138.566348953

[illegible]



Category: 8  
Low Frequency: -525.03754002

[illegible]



Category: 9  
High Frequency: -125.113036446

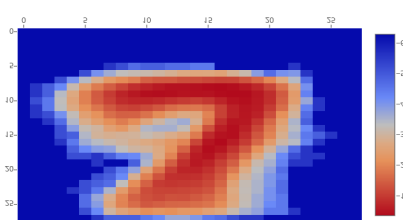
[illegible]

Category: 9  
Low Frequency: -649.05931148

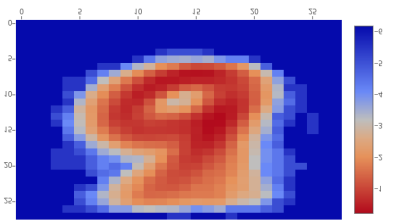
[illegible]



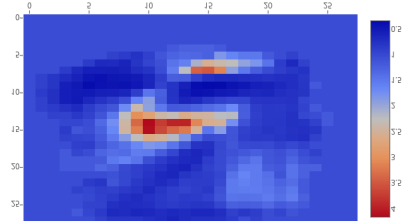
We also looked at odds ratios for the following sets of numbers:  $[(7,9),(4,9),(5,3),(8,3)]$ . The results are below, shown as heat maps of the log likelihoods.



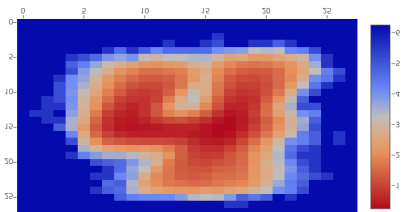
7



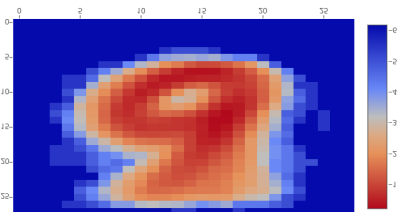
9



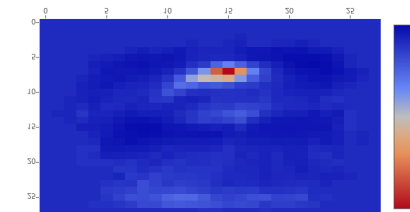
7 over 9



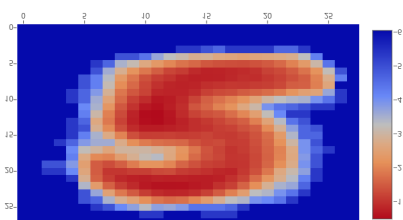
4



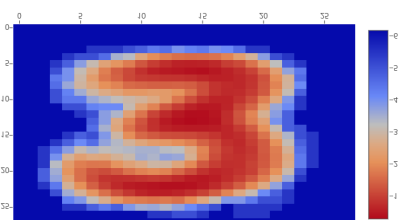
9



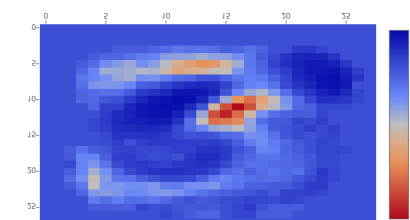
4 over 9



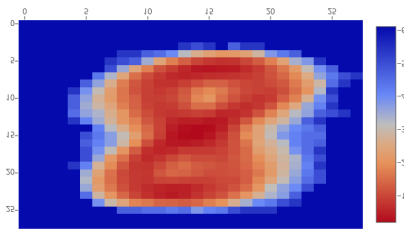
5



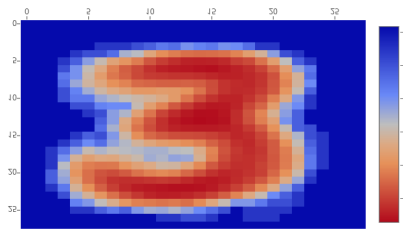
3



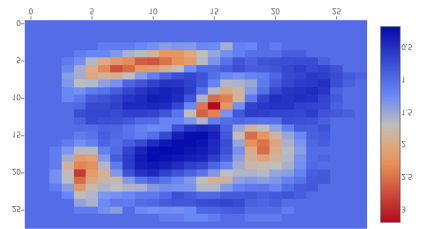
5 over 3



8



3



8 over 3



## Extra Credit for Part 1

We experimented with using ternary features. This is for bonus points. Our implementation was similar, but we had matrices to represent white, gray, and black rather than just white and black. Our results reflected a small improvement.

Our overall accuracy improved to **77.6%**.

Below are the classification rates for each digit. The format is as follows:

*(digit class, classification rate – truncated after 3 decimal points)*

(0, 0.833)  
(1, 0.953)  
(2, 0.766)  
(3, 0.800)  
(4, 0.775)  
(5, 0.684)  
(6, 0.780)  
(7, 0.735)  
(8, 0.621)  
(9, 0.800)

The confusion matrix is shown below.

	0:	1:	2:	3:	4:	5:	6:	7:	8:	9:
0:	83.33%	0.00%	1.11%	0.00%	0.00%	6.67%	4.44%	0.00%	4.44%	0.00%
1:	0.00%	95.37%	0.00%	0.00%	0.00%	1.85%	0.93%	0.00%	1.85%	0.00%
2:	0.97%	2.91%	76.70%	3.88%	0.97%	0.97%	5.83%	1.94%	4.85%	0.97%
3:	0.00%	2.00%	0.00%	80.00%	0.00%	3.00%	2.00%	5.00%	3.00%	5.00%
4:	0.00%	0.00%	0.00%	0.00%	77.57%	0.93%	1.87%	0.93%	1.87%	16.82%
5:	2.17%	1.09%	1.09%	13.04%	3.26%	68.48%	1.09%	1.09%	2.17%	6.52%
6:	0.00%	4.40%	4.40%	0.00%	4.40%	5.49%	78.02%	0.00%	3.30%	0.00%
7:	0.00%	5.66%	2.83%	0.00%	2.83%	0.00%	0.00%	73.58%	2.83%	12.26%
8:	0.97%	1.94%	2.91%	11.65%	1.94%	8.74%	0.00%	0.97%	62.14%	8.74%
9:	1.00%	1.00%	0.00%	2.00%	10.00%	2.00%	0.00%	2.00%	2.00%	80.00%

It appears that ternary features did not make a significant improvement upon classification.

However, we do believe that weighting gray values differently could be worth further exploration.

We applied our Naïve Bayes Classifier to the face data set. The implementation was very similar to the digit data set, except that there were only two classes and that there were more features per image. This is for bonus points.

Our overall accuracy was 90.6%.

Below is the classification rate for each class. The format is as follows, with 0 being not face and 1 being a face:

*(face class, classification rate – truncated after 3 decimal points)*

(0, 0.883)

(1, 0.931)

The confusion matrix is shown below.

	0:	1:
0:	88.31%	11.69%
1:	6.85%	93.15%

This whole section was for bonus points.



## PART 2: DOCUMENT CLASSIFICATION

### 2.1: For Three-Unit Students

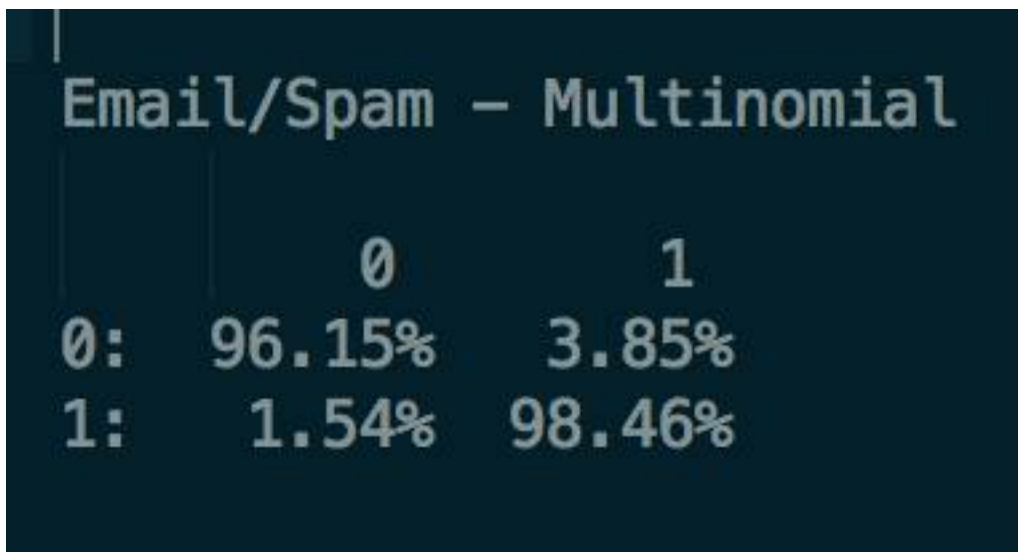
#### *Results*

#### SPAM DATA SET

#### Multinomial

Accuracy 97.30769%

Confusion matrix is below.



A screenshot of a terminal window displaying a confusion matrix for an 'Email/Spam' classification task using a 'Multinomial' model. The matrix shows the following values:

	0	1
0:	96.15%	3.85%
1:	1.54%	98.46%

#### *Top 20 Words for Normal Email*

language	workshop	include
university	email	edu
s	paper	http
linguistic	e	research
de	english	abstract
information	one	address
conference	please	

### *Top 20 Words for Spam Email*

email	program	business
s	send	one
order	free	d
report	money	work
our	list	com
address	receive	nt
mail	name	

### Bernoulli

Accuracy 96.53846%

Confusion matrix is below.

Email/Spam – Bernoulli		
	0	1
0:	93.85%	6.15%
1:	0.77%	99.23%

### *Top 20 Words for Normal Email*

language	please	call
university	e	research
s	follow	www
information	fax	word
linguistic	include	address
http	one	interest
email	english	

### *Top 20 Words for Spam Email*



our  
s  
free  
please  
email  
mail  
one

address  
list  
com  
receive  
http  
us  
send

day  
information  
remove  
here  
over  
want

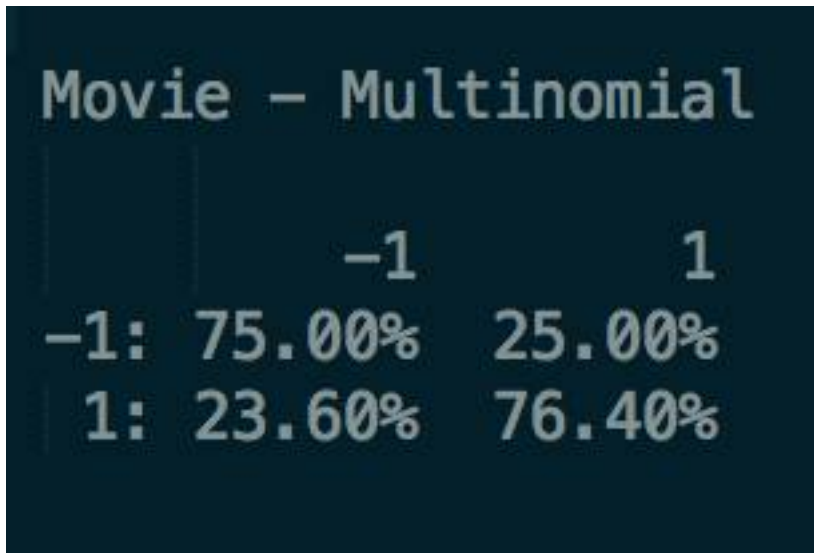


## MOVIE DATA SET

### Multinomial

Accuracy 75.70000%

Confusion matrix is below.



### *Top 20 Words for Negative Reviews*

movie	much	comedy
film	time	never
like	even	nothing
one	characters	makes
--	good	plot
bad	little	make
story	would	

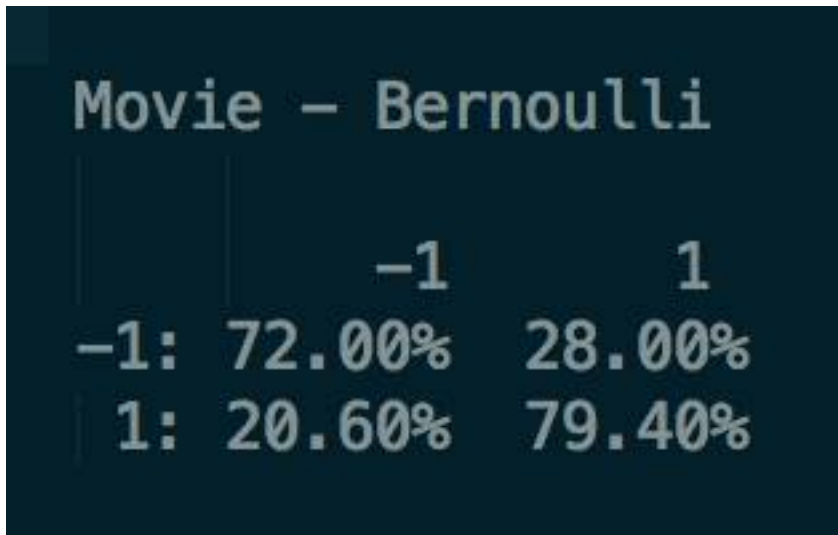
### *Top 20 Words for Positive Reviews*

film	comedy	funny
movie	way	make
--	even	life
one	time	us
like	best	makes
story	much	characters
good	performances	

## Bernoulli

Accuracy 75.70000%

Confusion matrix is below.



### *Top 20 Words for Negative Reviews*

movie	bad	comedy
film	time	nothing
like	even	makes
one	characters	plot
story	little	never
much	good	make
--	would	

### *Top 20 Words for Positive Reviews*

film	way	funny
movie	even	makes
one	good	life
like	best	make
--	time	characters
story	much	work
comedy	performances	



## Part 2.2: For Four-Unit Students (Extra Credit for Us)

Though we are three-unit students, we attempted the four-unit problem. This is for bonus points.

### Results

#### 8 CATEGORY DATA SET

#### Multinomial

Accuracy 92.01521%

Confusion matrix is below.

8 Category – Multinomial									
	0	1	2	3	4	5	6	7	
0:	97.06%	0.00%	0.00%	0.00%	2.94%	0.00%	0.00%	0.00%	
1:	0.00%	96.97%	0.00%	0.00%	0.00%	0.00%	0.00%	3.03%	
2:	0.00%	0.00%	100.00%	0.00%	0.00%	0.00%	0.00%	0.00%	
3:	0.00%	3.57%	3.57%	82.14%	0.00%	3.57%	0.00%	7.14%	
4:	2.13%	2.13%	2.13%	0.00%	93.62%	0.00%	0.00%	0.00%	
5:	0.00%	20.00%	0.00%	0.00%	0.00%	80.00%	0.00%	0.00%	
6:	0.00%	0.00%	2.17%	0.00%	0.00%	2.17%	95.65%	0.00%	
7:	0.00%	13.79%	0.00%	3.45%	0.00%	6.90%	0.00%	75.86%	

#### Top 20 Words for sci.space

space	subject	time
nt	like	data
would	us	first
one	system	orbit
launch	also	edu
nasa	writes	mission
earth	could	

#### Top 20 Words for comp.sys.ibm.pc.hardware

drive	scsi	nt
-------	------	----

ide  
one  
card  
drives  
controller  
system

disk  
subject  
use  
would  
edu  
hard

bus  
get  
m  
data  
also

*Top 20 Words for rec.sport.baseball*

nt  
would  
year  
edu  
writes  
one  
game

good  
team  
subject  
last  
article  
think  
players

like  
baseball  
games  
better  
well  
time

*Top 20 Words for comp.windows.x*

x  
window  
use  
nt  
subject  
file  
server

also  
available  
get  
edu  
motif  
version  
system

program  
sun  
c  
one  
m  
windows

*Top 20 Words for talk.politics.misc*

nt  
would  
people  
q  
one  
mr  
think

writes  
president  
article  
government  
stephanopoulos  
know  
us

edu  
like  
subject  
going  
right  
get

*Top 20 Words for misc.forsale*

new  
edu  
dos  
sale  
appears  
art  
subject

wolverine  
shipping  
cover  
price  
one  
list  
comics

drive  
nt  
hulk  
good  
vs  
system

*Top 20 Words for rec.sport.hockey*

nt  
game  
team  
hockey  
would  
play  
subject

period  
season  
nhl  
games  
one  
first  
year

think  
players  
get  
la  
edu  
like

*Top 20 Words for comp.graphics*

image  
jpeg  
edu  
nt  
file  
images  
data

also  
graphics  
software  
available  
use  
one  
program

files  
format  
get  
version  
system  
ftp

**Bernoulli**

Accuracy 84.41065%

Confusion matrix is below.



8 Category – Bernoulli								
	0	1	2	3	4	5	6	7
0:	82.35%	0.00%	2.94%	0.00%	14.71%	0.00%	0.00%	0.00%
1:	0.00%	100.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
2:	0.00%	0.00%	100.00%	0.00%	0.00%	0.00%	0.00%	0.00%
3:	0.00%	25.00%	3.57%	64.29%	0.00%	0.00%	0.00%	7.14%
4:	0.00%	0.00%	2.13%	0.00%	97.87%	0.00%	0.00%	0.00%
5:	0.00%	40.00%	0.00%	0.00%	10.00%	50.00%	0.00%	0.00%
6:	0.00%	2.17%	2.17%	0.00%	0.00%	0.00%	95.65%	0.00%
7:	3.45%	48.28%	0.00%	0.00%	6.90%	0.00%	0.00%	41.38%

*Top 20 Words for sci.space*

subject	like	new
would	could	much
nt	also	way
space	get	m
writes	think	edu
article	time	see
one	us	

*Top 20 Words for comp.sys.ibm.pc.hardware*

subject	know	system
nt	article	edu
one	card	drive
would	also	work
writes	like	problem
use	m	could
get	two	

*Top 20 Words for rec.sport.baseball*

subject	would	baseball
nt	one	good
writes	last	think
article	year	get
edu	like	time

m  
know

game  
team

first

*Top 20 Words for comp.windows.x*

subject  
x  
nt  
use  
writes  
get  
window

article  
using  
one  
like  
would  
also  
problem

know  
code  
m  
set  
email  
help

*Top 20 Words for talk.politics.misc*

subject  
nt  
writes  
article  
people  
would  
one

like  
edu  
us  
even  
think  
m  
get

government  
could  
know  
make  
time  
much

*Top 20 Words for misc.forsale*

subject  
sale  
edu  
new  
shipping  
please  
email

price  
nt  
one  
get  
condition  
like  
list

want  
used  
good  
use  
sell  
etc

*Top 20 Words for rec.sport.hockey*

subject  
nt  
team

game  
hockey  
writes

would  
one  
article

like  
play  
first  
think

go  
get  
nhl  
year

games  
time  
last

*Top 20 Words for comp.graphics*

subject  
nt  
one  
would  
writes  
also  
like

article  
graphics  
edu  
use  
get  
know  
computer

need  
could  
think  
program  
m  
two



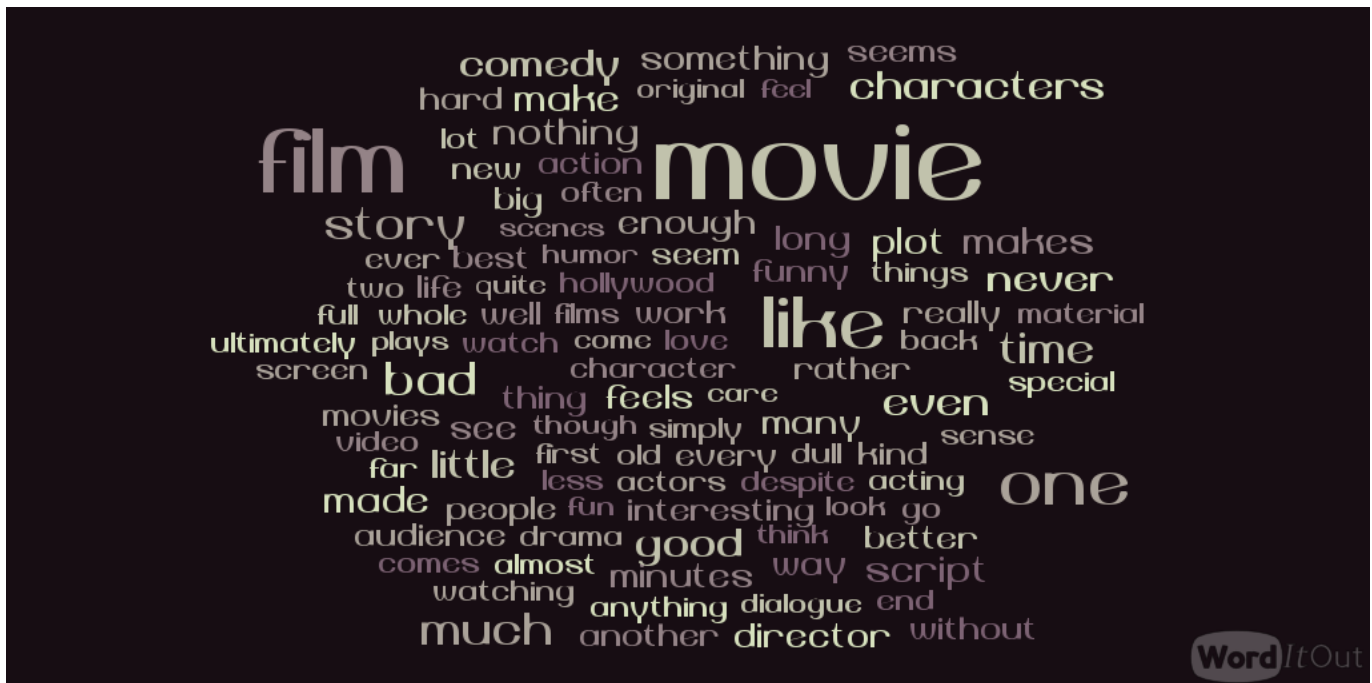
We did word clouds for every single category for ever dataset. This is for bonus points. They are shown below.

[illegible]

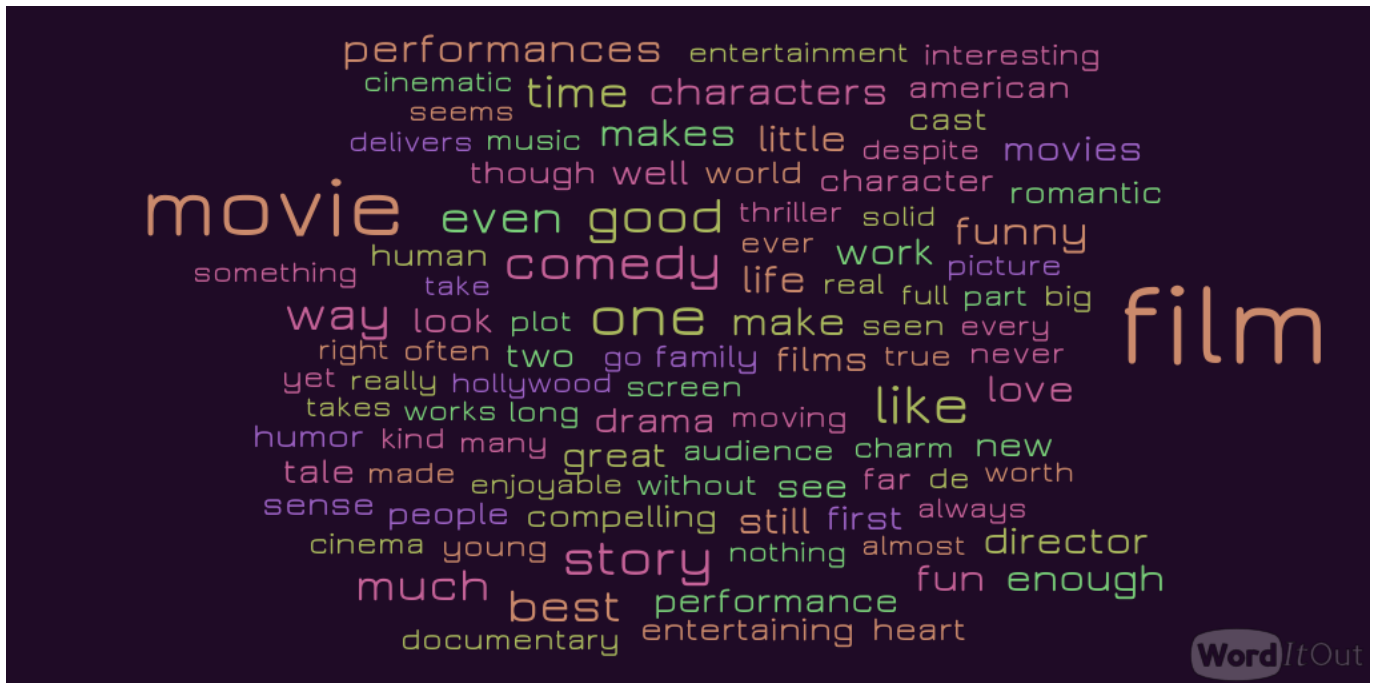
*Spam Email*



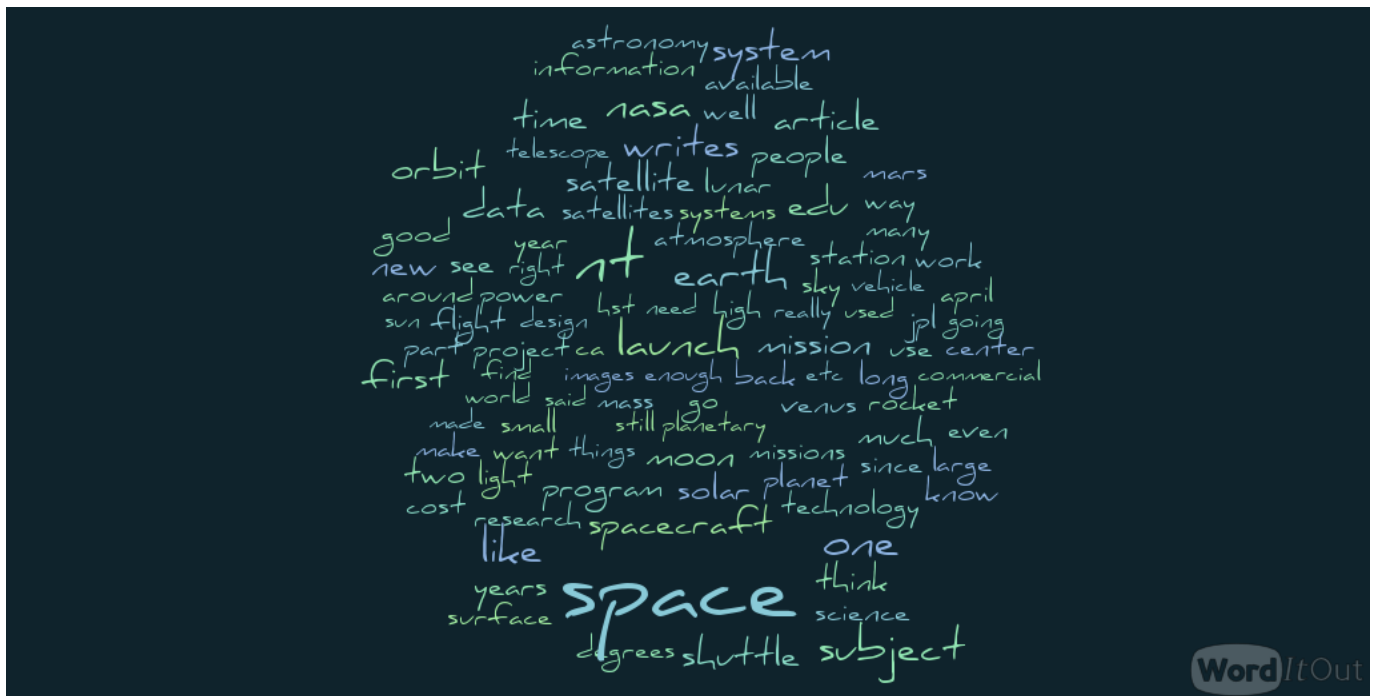
## Negative Movie Review



### Positive Movie Reiview



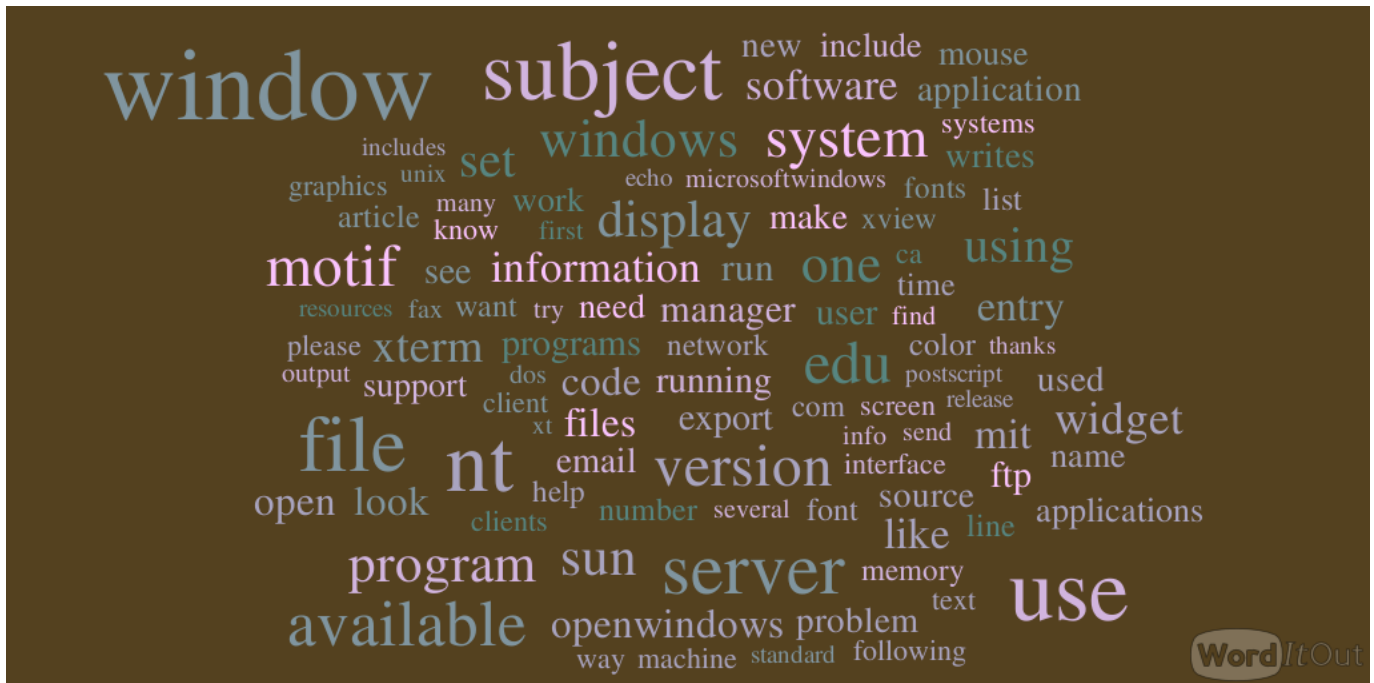
*sci.space*



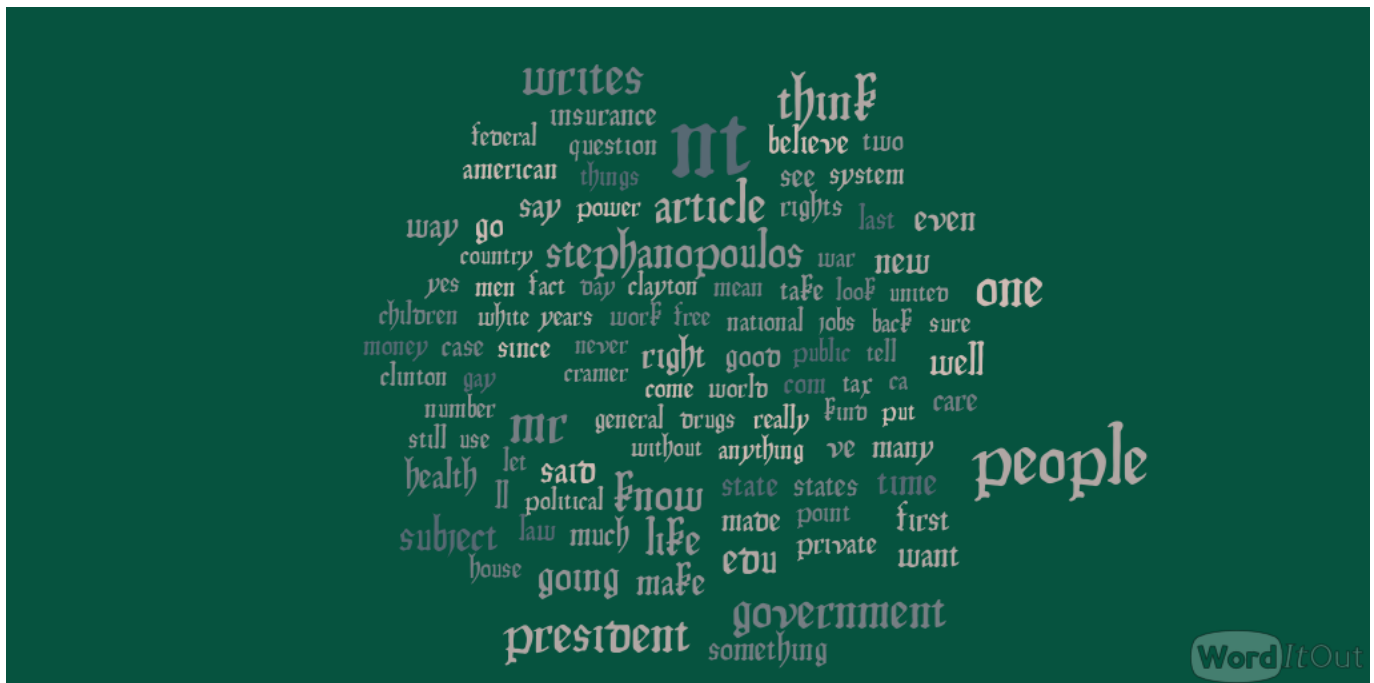


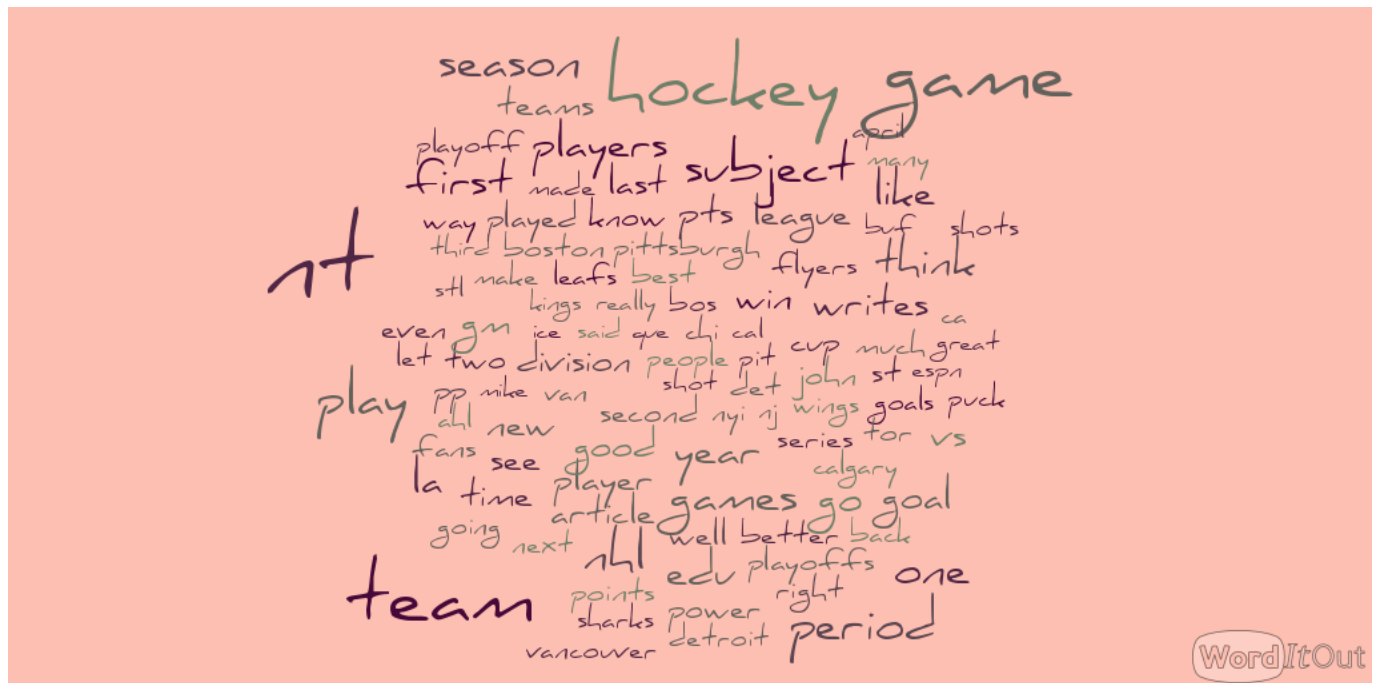


*comp.windows.x*



*talk.politics.misc*







*comp.graphics*



## INDIVIDUAL CONTRIBUTIONS

### Abhishek Nigam

I setup the infrastructure for part 1 and developed the training model for part 2. I generated the Top 20 Words list. We worked on the rest together.

### Jakub Klapacz

I setup the infrastructure for part 2 and developed the training model for part 1. I generated the confusion matrices. We worked on the rest together.