Part 2: Text Document Classification

-Bangqi Wang

This project uses two different Naïve Bayes classifiers, Multinomial Model and Bernoulli Model, to classify the text documents.

General Implementation

The train and test files contains the label and the words in each review or topics. Each record is preprocessed and stored as:

```
label word 1:count 1 word 2:count 2 ... word n:count n
```

The basic idea of this project is calculating the conditional probability of each word in each class and the probability of each class. The conditional probability is *likelihood*, P(document | class). The probability of each class is *priors*, P(class). The likelihood of document is represented by a sequence of words in the document, known as *bag of words*, $P(w_i | class)$.

$$P(document \mid class) = P(w_1, \dots, w_n \mid class) = \prod_{i=1}^{n} P(w_i \mid class)$$

prior

spam: 0.33

P(word | spam

the	:	0.0156
to	:	0.0153
and	:	0.0115
of	:	0.0095
you	:	0.0093
а	:	0.0086
with	ı:	0.0080
from	n:	0.0075

P(word | ¬spam)

the :	0.0210
to :	0.0133
of :	0.0119
2002:	0.0110
with:	0.0108
from:	0.0107
and :	0.0105
a :	0.0100

The algorithm will calculate the likelihood of document for each class and find the *maximum likelihood (ML)* as the predicted label for the documents. The main idea is to calculate the maximum likelihood estimate, but different models have different method to calculate the bag of words likelihood. To improve the accuracy and deal with the word that were never seen or seen too few times. The project uses the *Laplacian smoothing*. The smoothing method will introduce below for different models.

Then the project assigns the document to the class with the highest posterior and avoid the underflow by using the logs of probabilities.

$$P(class \mid document) \propto P(class) \prod_{i=1}^{n} P(w_i \mid class)$$
$$L(class \mid document) = \log P(class) + \sum_{i=1}^{n} \log P(w_i \mid class)$$

Multinomial Model

The Multinomial Model calculates the likelihood for each word by calculating the times of occurrences. The algorithm smooths the probabilities by pretending have seen every vocabulary one more time than actually did.

$$P(word \mid class) = \frac{\# \ of \ occurrences \ of \ this \ word \ in \ docs \ from \ this \ class}{total \ \# \ of \ words \ in \ docs \ from \ this \ class + 1}$$

$$P(word \mid class) = \frac{\# \ of \ occurrences \ of \ this \ word \ in \ docs \ from \ this \ class + 1}{total \ \# \ of \ words \ in \ docs \ from \ this \ class + V}$$

$$(V: total \ number \ of \ unique \ words)$$

Bernoulli Model

The Bernoulli model calculates the likelihood for each word by counting whether the word appeared at least once. The algorithm smooths the probabilities by pretending have seen every vocabulary one more times than actually did.

$$P(word \mid class) = \frac{\# \ of \ documents \ this \ word \ appeared \ from \ this \ class}{total \ \# \ of \ documents \ from \ this \ class}$$

$$P(word \mid class) = \frac{\# \ of \ documents \ this \ word \ appeared \ from \ this \ class + 1}{total \ \# \ of \ documents \ from \ this \ class + 2}$$

Part 2.1: Classification Results

The table below are the classification results for part 2.1.

Accuracy & Confusion Matrix

The two tables contain the output for two different datasets. The overall accuracies on the sentiment analysis of movie review task are around 76% and the accuracies on the topical theme classification task are around 93%.

	Sei	ntiment	Aı	nalysis of Mov	rie Review								
	Accur	acy		Confusion Matrix									
	Overall	76.00%			Negative	Positive							
Multinomial	Negative	75.00%		Negative	75.00%	25.00%							
	Positive	77.00%		Positive	23.00%	77.00%							
			<u> </u>										
	Overall	76.00%			Negative	Positive							
Bernoulli	Negative	72.80%		Negative	72.80%	27.20%							
	Positive	79.20%		Positive	20.80%	79.20%							

(Table for dataset 1)

	Binar	y convers	sation topic clas	ssification								
	Accura	acy	Confusion Matrix									
Multinomial	Overall Life Partner Min Wage	90.82% 95.92% 85.71%	Life Partner Min Wage	Life Partner 95.92% 14.29%	Min Wage 4.08% 85.71%							
Bernoulli	Overall Life Partner Min Wage	94.90% 91.84% 97.96%	Life Partner Min Wage	Life Partner 91.84% 2.04%	Min Wage 8.16% 97.96%							

(Table for dataset 2)

Top 10 Words with the Highest Likelihood
The tables below contain the words that are appear a lot in each classes.

	Sentimo	ent Analysis of N	Movi	e Review						
		top 10 words with			d					
	Neg	gative		Positive						
	Movie	0.009317		Film	0.009054					
	Film	0.007300		Movie	0.005952					
	Like	0.005251			0.004337					
	One	0.004610		One	0.003546					
Multinomial		0.003782		Like	0.003166					
Multinoimai	Bad	0.002817		Story	0.003008					
	Story	0.002753		Good	0.002610					
	Much	0.002689		Comedy	0.002659					
	Time	0.002433		Way	0.002564					
	Even	0.002273		Even	0.002438					
	Film	0.03891		Film	0.03931					
	Movie	0.03069		Movie	0.02509					
	One	0.02151		One	0.01533					
	Like	0.01891		Like	0.01380					
Bernoulli		0.01665			0.01366					
Dernoum	Story	0.01137		Story	0.01268					
	Comedy	0.01110		Comedy	0.01143					
	Way	0.01055		Way	0.0101					
	Even	0.01000		Even	0.01073					
	Good	0.00959		Good	0.01031					

	Binary co	nversation topi	c clas	sification	
	· ·	op 10 words with t			d
	Life P	artner		Min	Wage
	Know	0.05446		Know	0.05153
	Yeah	0.04535		Yeah	0.04541
	Uh	0.03034		Like	0.02889
	Like	0.02962		Uh	0.02192
Multinomial	Um	0.02310		Um	0.01960
Multinomiai	Right	0.01933		Right	0.01848
	Just	0.01798		Don	0.01733
	Think	0.01769		Think	0.01707
	Oh	0.01623		Just	0.01649
	Don	0.01590		Oh	0.01482
	Like	0.03780		Um	0.03907
	Know	0.03780		Think	0.03907
	Just	0.03780		Like	0.03907
	Yeah	0.03763		Know	0.03907
Bernoulli	Think	0.03763		Just	0.03907
Dernoum	Don	0.03763		Don	0.03907
	Um	0.03720		Yeah	0.03898
	Right	0.03711		People	0.03898
	Oh	0.03702		Oh	0.03881
	Really	0.03685		Right	0.03872

Top 10 Words with the Highest Odds RationThe tables below contain the words that are more likely to appear in specific class.

	Sentime	ent Analysis of	Mo	vie Review							
		top 10 words with	the	e highest odds ratio)						
	Neş	gative		Positive							
	Flat	15.1695		Disturbing	14.8324						
	Stale	14.1582		Refreshingly	10.8771						
	Dull	13.6526		Haunting	10.8771						
	Tired	12.1356		Grief	10.8771						
Multinomial	Plain	11.1243		Engrossing	10.8771						
Multinomiai	Mediocre	11.1243		Refreshing	9.8882						
	Unfunny	10.1130		Polished	9.8882						
	Poorly 10.1130			Inventive	9.8882						
	Pointless	10.1130		Gripping	9.8882						
	Generic	10.1130		Gem	9.8882						
	Flat	14.7431		Disturbing	14.2439						
	Stale	13.7602	1	Refreshingly	11.1917						
	Dull	12.7774	1	Haunting	11.1917						
	Tired	11.7945	1	Grief	11.1917						
Bernoulli	Plain	10.8116	1	Engrossing	11.1917						
Bernoulli	Mediocre	10.8116	1	Refreshing	10.1742						
	Unfunny	9.8288	1	Polished	10.1742						
	Poorly	9.8288		Inventive	10.1742						
	Pointless	9.8288		Gripping	10.1742						
	Generic	9.8288		Gem	10.1742						

	Binary cor	nversation topic	c cla	ssification			
	to	p 10 words with t	he h	ighest odds ratio			
	Life Pa	artner		Min V	Vage		
	Relationship	198.1665		Wage	150.4568		
	Compatibility	119.8154		Minimum	150.3094		
	Communication	116.6344		Welfare	128.7356		
	Marriage	112.1281		Wages	112.7026		
Multinomial	Partner	108.6063		Inflation	90.5393		
Multinomiai	Relationships	96.1350		Waitresses	77.3357		
	Friendship	95.4282		Tax	76.3925		
	Attracted	94.3679		Waitress	68.8476		
	Dating	89.0663		Salary	65.0751		
	Compatible	82.7044		Increase	58.4733		
	Attracted	62.1980		Wage	79.2308		
	Compatibility	60.2543		Wages	70.9990		
	Relationship	55.3951		Inflation	63.7962		
	Relationships	46.9725		Waitresses	56.5934		
Bernoulli	Friendship	46.6485		Minimum	56.5934		
Dernoulli	Attraction	46.6485		Waitress	54.5355		
	Marriage	44.2189		Welfare	52.4775		
	Communication	44.2189		Salary	46.8182		
	Dating	43.0851		Retail	41.1589		
	Qualities	40.8175		Increase	40.1299		

Conclusion

The top 10 words with highest likelihood are almost the same for both classes because the words with the most occurrences are similar to stop words. Those words cannot stand for any class because they have little meaning and they are more likely the general words around the topic. However, the top 10 words with highest odds ratio are high representative. The words with high odds ratio means that the words are more likely to appear in specific class only. The words with high likelihood I discussed above have similar high occurrences in any class. Therefore, the words with high likelihood not necessary have high odds ratio.

Part 2.2: Classification Results

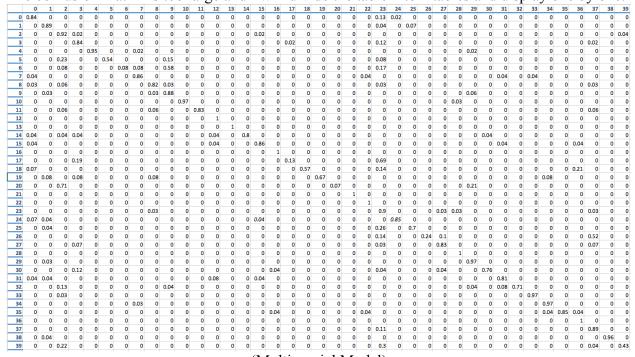
The dataset for this part is pretty large and contains 40 classes. Some classes might have common words but with different frequency. In this case, the multinomial model has more accuracy because the occurrences of words do matter in the topics. E.g. distributed system & database system.

Accuracy

	Full 40 Topic Corpus											
Accuracy												
	Multinomial Bernoulli											
Accuracy	83.88%	50.87%										

Confusion Matrix

The confusion matrix is too large and I will divide the table into 4 tables and display one by one.



(Multinomial Model)

Multinomial Model

Confusion Matrix [0:20][0:20]

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
0	0.84	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0.89	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0.92	0.02	0	0	0	0	0	0	0	0	0	0	0	0.02	0	0	0	0
3	0	0	0	0.84	0	0	0	0	0	0	0	0	0	0	0	0	0	0.02	0	0
4	0	0	0	0	0.95	0	0	0.02	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0.23	0	0	0.54	0	0	0	0.15	0	0	0	0	0	0	0	0	0	0
6	0	0	0.08	0	0	0	0.08	0.08	0	0.58	0	0	0	0	0	0	0	0	0	0
7	0.04	0	0	0	0	0	0	0.86	0	0	0	0	0	0	0	0	0	0	0	0
8	0.03	0	0.06	0	0	0	0	0	0.82	0.03	0	0	0	0	0	0	0	0	0	0
9	0	0.03	0	0	0	0	0	0	0.03	0.88	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0.97	0	0	0	0	0	0	0	0	0
11	0	0	0.06	0	0	0	0	0	0.06	0	0	0.83	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
14	0.04	0	0.04	0.04	0	0	0	0	0	0	0	0	0.04	0	0.8	0	0	0	0	0
15	0.04	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0.86	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
17	0	0	0	0.19	0	0	0	0	0	0	0	0	0	0	0	0	0	0.13	0	0
18	0.07	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.57	0
19	0	0.08	0	0.08	0	0	0	0	0.08	0	0	0	0	0	0	0	0	0	0	0.67

(Multinomial Model)

Confusion Matrix [0:20][20:40]

	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39
0	0	0	0	0.13	0.02	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0.04	0	0.07	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.04
3	0	0	0	0.12	0	0	0	0	0	0	0	0	0	0	0	0	0	0.02	0	0
4	0	0	0	0	0	0	0	0	0	0.02	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0.08	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0.17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0.04	0	0	0	0	0	0	0	0	0.04	0	0.04	0	0	0	0	0	0
8	0	0	0	0.03	0	0	0	0	0	0	0	0	0	0	0	0	0	0.03	0	0
9	0	0	0	0	0	0	0	0	0	0.06	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0.03	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.06	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0	0	0.04	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0.69	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0.14	0	0	0	0	0	0	0	0	0	0	0	0	0.21	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.08	0	0	0	0	0

(Multinomial Model)

Confusion Matrix [20:40][0:20]

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
20	0	0	0.71	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0.03	0	0	0	0	0	0	0	0	0	0	0
24	0.07	0.04	0	0	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0	0
25	0	0.04	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
27	0	0	0	0.07	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
29	0	0.03	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
30	0	0	0	0.12	0	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0
31	0.04	0.04	0	0	0	0	0	0	0	0	0	0	0.08	0	0	0.04	0	0	0	0
32	0	0	0.13	0	0	0	0	0	0	0.04	0	0	0	0	0	0	0	0	0	0
33	0	0	0.03	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
34	0	0	0	0	0	0	0	0.03	0	0	0	0	0	0	0	0	0	0	0	0
35	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0
36	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
37	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
38	0	0.04	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
39	0	0	0.22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

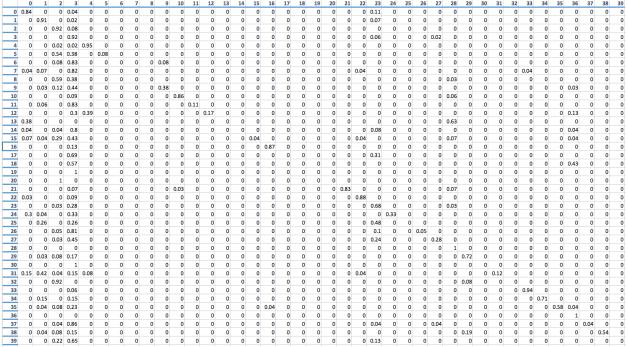
(Multinomial Model)

Confusion Matrix [20:40][20:40]

	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39
20	0.07	0	0	0	0	0	0	0	0	0.21	0	0	0	0	0	0	0	0	0	0
21	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
22	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0.9	0	0	0	0.03	0.03	0	0	0	0	0	0	0	0	0.03	0	0
24	0	0	0	0	0.85	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0.26	0	0.7	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0	0.14	0	0	0.24	0.1	0	0	0	0	0	0	0	0	0	0.52	0	0
27	0	0	0	0.03	0	0	0	0.83	0	0	0	0	0	0	0	0	0	0.07	0	0
28	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
29	0	0	0	0	0	0	0	0	0	0.97	0	0	0	0	0	0	0	0	0	0
30	0	0	0	0.04	0	0	0	0.04	0	0	0.76	0	0	0	0	0	0	0	0	0
31	0	0	0	0	0	0	0	0	0	0	0	0.81	0	0	0	0	0	0	0	0
32	0	0	0	0	0	0	0	0	0	0.04	0	0.08	0.71	0	0	0	0	0	0	0
33	0	0	0	0	0	0	0	0	0	0	0	0	0	0.97	0	0	0	0	0	0
34	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.97	0	0	0	0	0
35	0	0	0.04	0	0	0	0	0	0	0	0	0	0	0	0.04	0.85	0.04	0	0	0
36	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
37	0	0	0	0.11	0	0	0	0	0	0	0	0	0	0	0	0	0	0.89	0	0
38	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.96	0
39	0	0	0	0.3	0	0	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0.43

(Multinomial Model)

Bernoulli Model



(Bernoulli Model)

Confusion Matrix [0:20][0:20]

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
0	0.84	0	0	0.04	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0.91	0	0.02	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0.92	0.08	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0.92	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0.02	0.02	0.95	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0.54	0.38	0	0.08	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0.08	0.83	0	0	0	0	0	0.08	0	0	0	0	0	0	0	0	0	0
7	0.04	0.07	0	0.82	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0.59	0.38	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0.03	0.12	0.44	0	0	0	0	0	0.38	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0.09	0	0	0	0	0	0	0.86	0	0	0	0	0	0	0	0	0
11	0	0.06	0	0.83	0	0	0	0	0	0	0	0.11	0	0	0	0	0	0	0	0
12	0	0	0	0.3	0.39	0	0	0	0	0	0	0	0.17	0	0	0	0	0	0	0
13	0.38	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
14	0.04	0	0.04	0.8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
15	0.07	0.04	0.29	0.43	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0	0
16	0	0	0	0.13	0	0	0	0	0	0	0	0	0	0	0	0	0.87	0	0	0
17	0	0	0	0.69	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0.57	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

(Bernoulli Model)

Confusion Matrix [0:20][20:40]

	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39
0	0	0	0	0.11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0.07	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0.06	0	0	0	0.02	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0.04	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0.03	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.03	0	0	0
10	0	0	0	0	0	0	0	0	0.06	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.13	0	0	0
13	0	0	0	0	0	0	0	0	0.63	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0.08	0	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0
15	0	0	0.04	0	0	0	0	0	0.07	0	0	0	0	0	0	0	0.04	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0.31	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.43	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

(Bernoulli Model)

Confusion Matrix [20:40][0:20]

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
20	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0.07	0	0	0	0	0	0	0.03	0	0	0	0	0	0	0	0	0
22	0.03	0	0	0.09	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0.03	0.28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
24	0.3	0.04	0	0.33	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0.26	0	0.26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0.05	0.81	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
27	0	0	0.03	0.45	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
29	0	0.03	0.08	0.17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
30	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
31	0.15	0.42	0.04	0.15	0.08	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
32	0	0	0.92	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
33	0	0	0	0.06	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
34	0	0.15	0	0.15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
35	0	0.04	0.08	0.23	0	0	0	0	0	0	0	0	0	0	0	0	0.04	0	0	0
36	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
37	0	0	0.04	0.86	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
38	0	0.04	0.08	0.15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
39	0	0	0.22	0.65	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

(Bernoulli Model)

Confusion Matrix [20:40][20:40]

	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39
20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0.83	0	0	0	0	0	0	0.07	0	0	0	0	0	0	0	0	0	0	0
22	0	0	0.88	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0.68	0	0	0	0	0.03	0	0	0	0	0	0	0	0	0	0	0
24	0	0	0	0	0.33	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0.48	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0	0.1	0	0	0.05	0	0	0	0	0	0	0	0	0	0	0	0	0
27	0	0	0	0.24	0	0	0	0.28	0	0	0	0	0	0	0	0	0	0	0	0
28	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
29	0	0	0	0	0	0	0	0	0	0.72	0	0	0	0	0	0	0	0	0	0
30	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
31	0	0	0.04	0	0	0	0	0	0	0	0	0.12	0	0	0	0	0	0	0	0
32	0	0	0	0	0	0	0	0	0	0.08	0	0	0	0	0	0	0	0	0	0
33	0	0	0	0	0	0	0	0	0	0	0	0	0	0.94	0	0	0	0	0	0
34	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.71	0	0	0	0	0
35	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0.58	0.04	0	0	0
36	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
37	0	0	0	0.04	0	0	0	0.04	0	0	0	0	0	0	0	0	0	0.04	0	0
38	0	0	0	0	0	0	0	0	0	0.19	0	0	0	0	0	0	0	0	0.54	0
39	0	0	0	0.13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

(Bernoulli Model)

Confused Topic in Multinomial Model The table below shows the most confused topic for each topic.

Multinomial

The overall accuracy for multinomial model is around 83.88%.

Real Topic	Accuracy	Confused Topic	Percentage
0	84.44%	23	13.33%
1	89.29%	25	7.14%
2	91.84%	39	4.08%
3	83.67%	23	12.24%
4	95.45%	7	2.27%
5	53.85%	2	23.08%
6	58.33%	23	16.67%
7	85.71%	0	3.57%
8	82.35%	2	5.88%
9	88.24%	29	5.88%
10	97.14%	28	2.86%
11	83.33%	2	5.56%
12	100.00%	0	0.00%
13	100.00%	0	0.00%
14	80.00%	0	4.00%
15	85.71%	0	3.57%
16	100.00%	0	0.00%
17	68.75%	3	18.75%
18	57.14%	36	21.43%
19	66.67%	1	8.33%
20	71.43%	29	21.43%
21	100.00%	0	0.00%
22	100.00%	0	0.00%
23	90.00%	8	2.50%
24	85.19%	0	7.41%
25	69.57%	23	26.09%
26	52.38%	26	23.81%
27	82.76%	3	6.90%
28	100.00%	0	0.00%
29	97.22%	1	2.78%
30	76.00%	3	12.00%
31	80.77%	12	7.69%
32	70.83%	2	12.50%
33	96.88%	2	3.13%
34	97.06%	7	2.94%
35	84.62%	16	3.85%
36	100.00%	0	0.00%
37	89.29%	23	10.71%
38	96.15%	1	3.85%
39	43.48%	23	30.43%

Bernoulli The overall accuracy for Bernoulli model is 50.87%

Real Topic	Accuracy	Confused Topic	Percentage
0	84.44%	23	11.11%
1	91.07%	23	7.14%
2	91.84%	3	8.16%
3	91.84%	23	6.12%
4	95.45%	2	2.27%
5	53.85%	3	38.46%
6	83.33%	2	8.33%
7	82.14%	1	7.14%
8	58.82%	3	38.24%
9	44.12%	9	38.24%
10	85.71%	3	8.57%
11	83.33%	11	11.11%
12	39.13%	3	30.43%
13	62.50%	0	37.50%
14	80.00%	23	8.00%
15	42.86%	2	28.57%
16	86.96%	3	13.04%
17	68.75%	23	31.25%
18	57.14%	36	42.86%
19	100.00%	0	0.00%
20	100.00%	0	0.00%
21	83.33%	3	6.67%
22	87.50%	3	9.38%
23	67.50%	3	27.50%
24	33.33%	3	33.33%
25	47.83%	1	26.09%
26	80.95%	23	9.52%
27	44.83%	27	27.59%
28	100.00%	0	0.00%
29	72.22%	3	16.67%
30	100.00%	0	0.00%
31	42.31%	0	15.38%
32	91.67%	29	8.33%
33	93.75%	3	6.25%
34	70.59%	1	14.71%
35	57.69%	3	23.08%
36	100.00%	0	0.00%
37	85.71%	2	3.57%
38	53.85%	29	19.23%
39	65.22%	2	21.74%