



CS777

Big Data Analytics

Spring 2021

Professor Kia Teymourian

Assignment 6

Kunfei Chen

U15575304

Introduction

This project aims to apply the MCMC algorithm to achieve the clustering of text documents.

Task 1 – MCMC Algorithm

Given k document clusters, the generative process for n documents is shown in Figure 1.

Figure 1

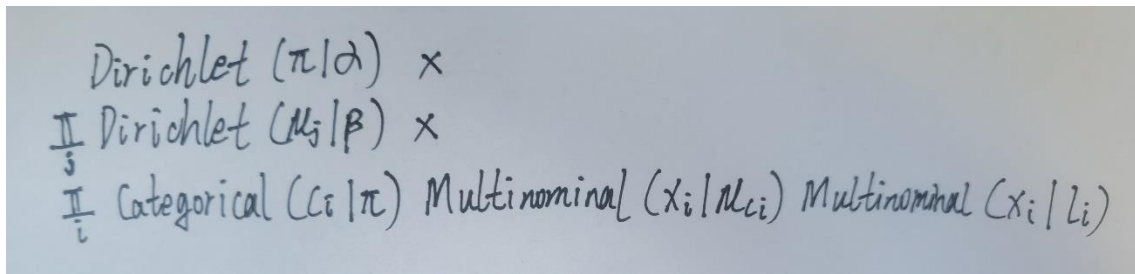
```

 $\pi \sim \text{Dirichlet}(\alpha)$ 
for  $j = 1$  to  $k$  do
     $\mu_j \sim \text{Dirichlet}(\beta)$ 
end for
for  $i = 1$  to  $n$  do
     $c_i \sim \text{Categorical}(\pi)$ 
     $x_i \sim \text{Multinomial}(\mu_{c_i}, l_i)$ 
end for

```

According to this process, the corresponding probability density function is derived as shown in Figure 2.

Figure 2



$$\begin{aligned}
 & \text{Dirichlet}(\pi | \alpha) \times \\
 & \prod_j \text{Dirichlet}(\mu_j | \beta) \times \\
 & \prod_i \text{Categorical}(c_i | \pi) \text{Multinomial}(x_i | \mu_{c_i}) \text{Multinomial}(x_i | l_i)
 \end{aligned}$$

Task 2

In this task, the input documents are transformed into count vectors in which each vector has 20000 entries. All words are converted into lowercase so that the same word with different capitalization are matched and can be counted up as the same word (See Figure 3).

Figure 3

```

sc = SparkContext(appName="A6")
lines = sc.textFile(input_dir)
header_text = lines.map(lambda x: (x[x.index('id')+4:x.index('" url=')], x[x.index('">')+2:][:8].lower()))
regex = re.compile('[^a-zA-Z]')
def get_words(input_val):
    result = []
    words = regex.sub(' ', input_val).split()
    for w in words:
        if len(w) > 2:
            result.append(w)
    return result
header_words = header_text.map(lambda x: (x[0], get_words(x[1])))
top_size = 20000
sorted_words = header_words.flatMap(lambda x: x[1]).map(lambda x: (x,1)).reduceByKey(add).takeOrdered(top_size, key=lambda x: -x[1])
top_words = []
for each in sorted_words:
    top_words.append(each[0])
def countWords(input_val):
    header = input_val[0]
    words = input_val[1]
    numwords = np.zeros(top_size)
    count = 0
    for w in words:
        if w in top_words:
            count += 1
            idx = top_words.index(w)
            numwords[idx] += 1
    return (header, numwords, count)
result = header_words.map(countWords)
result.cache()

```

In addition, for the document “20 newsgroups/comp.graphics/37261”, its number of times that each of the 100 most common dictionary words are identified (See Figure 4). The words list is shown in Figure 5.

Figure 4

```

target = result.filter(lambda x: x[0].split("/)[-1][0:] == "37261").collect()
sort_index = list(np.argsort(target[0][1]))
list.reverse(sort_index)

target_100_words = []
for i in range(177):
    index = sort_index[i]
    target_100_words.append((top_words[index], target[0][1][index]))
print("The 100 most common words appear in document 20 newsgroups/comp.graphics/37261:")
for w in target_100_words:
    print(w)

```

Figure 5

```

The 100 most common words appear in document 20 newsgroups/comp.graphics/37261:
('and', 12.0)
('navy', 11.0)
('the', 8.0)
('lipman', 7.0)
('for', 6.0)
('presentations', 6.0)
('seminar', 6.0)
('scientific', 5.0)
('virtual', 5.0)
('visualization', 5.0)
('reality', 5.0)
('will', 4.0)
('robert', 4.0)
('mil', 4.0)
('presentation', 4.0)
('one', 3.0)
('oasis', 3.0)
('bethesda', 3.0)
('should', 3.0)
('maryland', 3.0)
('center', 3.0)
('authors', 2.0)
('length', 2.0)
('related', 2.0)
('work', 2.0)
('information', 2.0)
('david', 2.0)
('warfare', 2.0)
('surface', 2.0)
('call', 2.0)
('carderock', 2.0)

```

Task 3

Then the MCMC algorithm derived in Task 1 is implemented for 200 iterations on the 20 newsgroups data, with the goal of learning the π , μ_j , and c_i values. As shown in Figure 6, after 200 iterations, the 50 most important words with the highest probability in each of the 20 learned mixture components are identified.

Figure 6

```

50 important_words for category 0 : ['the', 'from', 'edu', 'that', 'for', 'subject', 'you', 'and', 'date', 'lines', 'apr', 'gmt', 'about
50 important_words for category 1 : ['the', 'and', 'that', 'for', 'you', 'from', 'this', 'are', 'with', 'have', 'edu', 'subject', 'apr',
50 important_words for category 2 : ['the', 'and', 'that', 'scx', 'chz', 'for', 'kuwait', 'from', 'not', 'this', 'rlk', 'was', 'unw',
50 important_words for category 3 : ['the', 'and', 'that', 'you', 'for', 'not', 'from', 'are', 'have', 'this', 'they', 'was', 'with',
50 important_words for category 4 : ['the', 'for', 'and', 'from', 'edu', 'apr', 'subject', 'lines', 'date', 'gmt', 'with', 'sale', 'you',
50 important_words for category 5 : ['the', 'that', 'and', 'this', 'you', 'was', 'stephanopoulos', 'president', 'for', 'with', 'have', 'n
50 important_words for category 6 : ['the', 'that', 'and', 'you', 'not', 'are', 'this', 'for', 'from', 'have', 'but', 'god', 'with',
50 important_words for category 7 : ['the', 'and', 'pit', 'period', 'tor', 'from', 'edu', 'det', 'bos', 'chi', 'pts', 'for', 'stl',
50 important_words for category 8 : ['max', 'giz', 'bhj', 'the', 'bxn', 'and', 'qax', 'you', 'from', 'edu', 'that', 'com', 'appears',
50 important_words for category 9 : ['edu', 'from', 'subject', 'lines', 'date', 'apr', 'for', 'gmt', 'the', 'and', 'distribution', 'this',
50 important_words for category 10 : ['the', 'and', 'for', 'from', 'edu', 'subject', 'date', 'lines', 'that', 'apr', 'gmt', 'you', 'have
50 important_words for category 11 : ['the', 'and', 'for', 'that', 'you', 'this', 'from', 'are', 'with', 'not', 'have', 'can', 'will',
50 important_words for category 12 : ['the', 'for', 'edu', 'and', 'from', 'you', 'subject', 'are', 'apr', 'gmt', 'lines', 'date', 'com',
50 important_words for category 13 : ['the', 'and', 'edu', 'from', 'for', 'that', 'lines', 'subject', 'apr', 'date', 'you', 'gmt', 'this
50 important_words for category 14 : ['the', 'and', 'that', 'was', 'for', 'you', 'they', 'have', 'from', 'but', 'with', 'this', 'not',
50 important_words for category 15 : ['the', 'and', 'that', 'for', 'from', 'this', 'you', 'apr', 'subject', 'com', 'date', 'lines', 'edu
50 important_words for category 16 : ['air', 'the', 'subject', 'lines', 'from', 'date', 'ahf', 'kjjz', 'khf', 'gmt', 'region', 'apr', 'ne
50 important_words for category 17 : ['the', 'and', 'for', 'that', 'from', 'with', 'you', 'have', 'this', 'date', 'edu', 'lines', 'subje
50 important_words for category 18 : ['from', 'edu', 'date', 'for', 'subject', 'lines', 'gmt', 'and', 'apr', 'ndet', 'loop', 'the', 'uas
50 important_words for category 19 : ['the', 'that', 'you', 'and', 'frank', 'not', 'can', 'are', 'for', 'objective', 'from', 'writes',

```

Task 4

Finally, in order to evaluate the accuracy of the learning algorithm, the number of documents in each of the 20 learned mixture components are printed. In addition, the percentage of the top three real categories for documents assigned to each cluster are also identified (See Figure 7). It can be seen that generally each component has a high prevalence of several related categories with the lowest percentage of 0.16 and the highest percentage of 0.61. In addition, the top categories in the 20 clusters are roughly correspond to the 20 categories present in the data.

Figure 7

```
Total count of cluster 0: 125
Top 3 topics in cluster 0 : [['u' misc.forsale', 0.176], ['u' rec.sport.baseball', 0.128], ['u' talk.religion.miso', 0.104]]
Total count of cluster 1: 3369
Top 3 topics in cluster 1 : [['u' rec.motorcycles', 0.2716], ['u' rec.autos', 0.247], ['u' sci.space', 0.2188]]
Total count of cluster 2: 54
Top 3 topics in cluster 2 : [['u' sci.med', 0.2778], ['u' rec.sport.hockey', 0.2222], ['u' talk.politics.miso', 0.1111]]
Total count of cluster 3: 3036
Top 3 topics in cluster 3 : [['u' talk.politics.mideast', 0.2955], ['u' talk.politics.guns', 0.2915], ['u' talk.politics.miso', 0.2516]]
Total count of cluster 4: 636
Top 3 topics in cluster 4 : [['u' misc.forsale', 0.6148], ['u' sci.med', 0.1164], ['u' sci.electronics', 0.0456]]
Total count of cluster 5: 54
Top 3 topics in cluster 5 : [['u' talk.politics.miso', 0.3704], ['u' sci.med', 0.1852], ['u' talk.politics.guns', 0.1481]]
Total count of cluster 6: 2485
Top 3 topics in cluster 6 : [['u' soc.religion.christian', 0.3646], ['u' alt.atheism', 0.334], ['u' talk.religion.miso', 0.229]]
Total count of cluster 7: 153
Top 3 topics in cluster 7 : [['u' rec.sport.hockey', 0.5621], ['u' alt.atheism', 0.1176], ['u' talk.religion.miso', 0.0784]]
Total count of cluster 8: 99
Top 3 topics in cluster 8 : [['u' talk.politics.miso', 0.2828], ['u' talk.politics.mideast', 0.1818], ['u' talk.religion.miso', 0.1515]]
Total count of cluster 9: 361
Top 3 topics in cluster 9 : [['u' misc.forsale', 0.2188], ['u' comp.windows.x', 0.1108], ['u' comp.graphics', 0.0859]]
Total count of cluster 10: 1780
Top 3 topics in cluster 10 : [['u' comp.graphics', 0.2938], ['u' comp.os.ms-windows.miso', 0.2665], ['u' comp.windows.x', 0.1994]]
Total count of cluster 11: 2329
Top 3 topics in cluster 11 : [['u' sci.crypt', 0.3366], ['u' comp.windows.x', 0.2001], ['u' sci.med', 0.1498]]
Total count of cluster 12: 44
Top 3 topics in cluster 12 : [['u' misc.forsale', 0.1591], ['u' talk.politics.miso', 0.1136], ['u' comp.sys.ibm.pc.hardware', 0.1136]]
Total count of cluster 13: 737
Top 3 topics in cluster 13 : [['u' rec.sport.hockey', 0.46], ['u' rec.sport.baseball', 0.1506], ['u' sci.electronics', 0.0651]]
Total count of cluster 14: 1551
Top 3 topics in cluster 14 : [['u' rec.sport.baseball', 0.5016], ['u' rec.sport.hockey', 0.3217], ['u' sci.med', 0.109]]
Total count of cluster 15: 255
Top 3 topics in cluster 15 : [['u' sci.crypt', 0.2863], ['u' comp.windows.x', 0.1059], ['u' rec.autos', 0.1059]]
Total count of cluster 16: 34
Top 3 topics in cluster 16 : [['u' sci.med', 0.1765], ['u' soc.religion.christian', 0.1765], ['u' talk.politics.guns', 0.0882]]
Total count of cluster 17: 2659
Top 3 topics in cluster 17 : [['u' comp.sys.ibm.pc.hardware', 0.3238], ['u' comp.sys.mac.hardware', 0.3129], ['u' comp.os.ms-windows.miso', 0.1392]]
Total count of cluster 18: 101
Top 3 topics in cluster 18 : [['u' misc.forsale', 0.2277], ['u' comp.windows.x', 0.1089], ['u' soc.religion.christian', 0.0891]]
Total count of cluster 19: 155
Top 3 topics in cluster 19 : [['u' talk.religion.miso', 0.4258], ['u' alt.atheism', 0.3742], ['u' sci.electronics', 0.0581]]
```