

DLNLP_2023_Homework3 王彪--ZY2203114

作业三：从上面链接给定的语料库中均匀抽取200个段落（每个段落大于500个词），每个段落的标签就是对应段落所属的小说。利用LDA模型对于文本建模，并把每个段落表示为主题分布后进行分类。验证与分析分类结果，（1）在不同数量的主题个数下分类性能的变化；（2）以"词"和以"字"为基本单元下分类结果有什么差异

1. LDA模型

LDA是一款基于Dirichlet分布的概率主题模型，运用了概率论和词袋模型等知识。从作业内容分析推导这一模型，数据集为金庸的16本小说。本文从小说中随机抽取M个段落，即得到M篇文档，对应第m个文档中有 n_m 个词，假定每篇文档中的词都是从一系列主题中选择出来的，每个主题下有对应的词，文档中一个词 w_{mn} 的产生要经历两个步骤，首先从该文档的主题分布中采样一个主题，然后在这个主题对应的词分布中采样一个词。不断重复直到m个文档都完成上述过程。

随机生成过程的两个步骤都符合Dirichlet分布：

- 首先，抽主题： α 是某文档主题的Dirichlet分布的超参数，采样得到文档m的主题分布 $\theta_m : (p_1, p_2, \dots, p_k)$ ，k是主题个数，采样出某一主题，得到文档m第n个词的主题编号 Z_{mn} ；
- 然后，抽主题下的词： β 是某主题的词语的Dirichlet分布的超参数，采样得到所有主题词语的分布 ϕ ，得到主题1 $(p_{11}, p_{12}, \dots, p_{1v})$ ，.....主题k $(p_{k1}, p_{k2}, \dots, p_{kv})$ ，从第一步主题编号 Z_{mn} 对应的主题分布中抽样得到词语 W_{mn} 。

这一过程中， w_{mn} 是可以观察到的已知变量， α, β 是跟据经验给定的先验参数，其他变量 Z_{mn}, θ, ϕ 都是未知变量，需要根据观察到的变量学习，这里使用Gibbs Sampling算法，其运行方式为，每次抽取概率向量的一个维度，未定其他维度的变量，采样当前维度的值，不断迭代，直到收敛。初始时，随机给文本中的每个词分配主题 $z^{(0)}$ ，然后统计每个主题 z 下词出现的概率，及每个文档 m 下出现主题 z 的数量，以及每个主题下词的总量，据此估计排除当前词的主题分布，然后根据这个分布为该词采样新主题，最终每个文档的主题分布 θ_m ，每个主题的词分布收敛，算法停止。

2. 分类模型

本实验选用支持向量机对文本段落进行多分类，核函数选择线性核。

3. 代码

3.1 数据准备

预处理代码与第一次作业类似，略。

200个段落的选择代码如下：

```

1 def get_segment(file, len, flag):
2     seg = []
3     with open(file, 'r') as f:
4         txt_list = f.readline().split(',')
5         file_path = txt_list[random.randint(0, len(txt_list)-1)] + '.txt'
6         split_words = get_single_corpus(DATA_PATH + file_path, flag)
7         begin = random.randint(0, len(split_words)-len-1)
8         seg.extend(split_words[begin:begin+len])
9     return file_path, seg

```

3.2 LDA模型

3.2.1 初始化

```

1 def initialize():
2     for f_idx, segment in enumerate(segments):
3         temp_word2topic = []
4         for w_idx, word in enumerate(segment):
5             init_topic = random.randint(0, topic_num-1)
6             file2topic[f_idx, init_topic] += 1
7             topic2word[init_topic, word] += 1
8             topic_count[init_topic] += 1
9             temp_word2topic.append(init_topic)
10        word2topic.append(temp_word2topic)

```

3.2.2 困惑度计算

```

1 def compute_perplexity(): # 计算困惑度
2     file_count = np.sum(file2topic, 1)
3     # print(file_count)
4     count = 0
5     perplexity = 0
6     for f_idx, segment in enumerate(segments):
7         for word in segment:
8             perplexity = (perplexity + np.log(((topic2word[:, word] /
9             topic_count) *
10             (file2topic[f_idx, :] /
11             file_count[f_idx])).sum()))
12             count += 1
13
14     return np.exp(perplexity / (-count))

```

3.2.3 gibbs_sample

```

1 def gibbs_sample():
2     global file2topic
3     global topic2word
4     global topic_count
5     new_file2topic = np.zeros([file_num, topic_num])
6     new_topic2word = np.zeros([topic_num, word_num])
7     new_topic_count = np.zeros([topic_num])
8     for f_idx, segment in enumerate(segments):
9         for w_idx, word in enumerate(segment):
10            old_topic = word2topic[f_idx][w_idx]

```

```

11         p = np.divide(np.multiply(file2topic[f_idx, :], topic2word[:,
word]), topic_count)
12         new_topic = np.random.multinomial(1, p/p.sum()).argmax()
13         word2topic[f_idx][w_idx] = new_topic
14         new_file2topic[f_idx, new_topic] += 1
15         new_topic2word[new_topic, word] += 1
16         new_topic_count[new_topic] += 1
17     file2topic = new_file2topic
18     topic2word = new_topic2word
19     topic_count = new_topic_count

```

3.2.4 画图

```

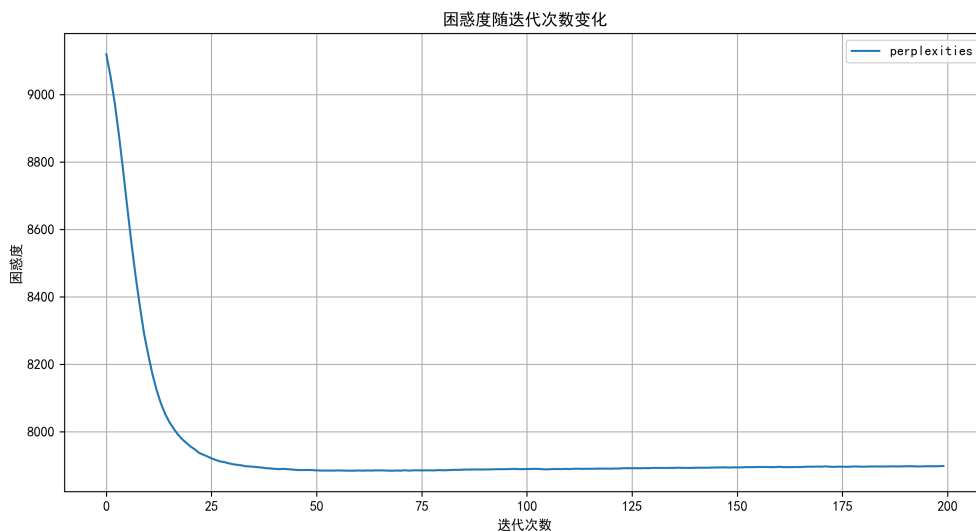
1  def drawfig(x, perplexities, name): # 画图
2      plt.rcParams["figure.figsize"] = (12, 6)
3      plt.rcParams['font.sans-serif'] = ['SimHei']
4      plt.title('困惑度随迭代次数变化')
5      plt.xlabel('迭代次数')
6      plt.ylabel('困惑度')
7      plt.grid()
8      plt.plot(x, perplexities, label='perplexities')
9      plt.legend()
10     plt.savefig(name, dpi=300)

```

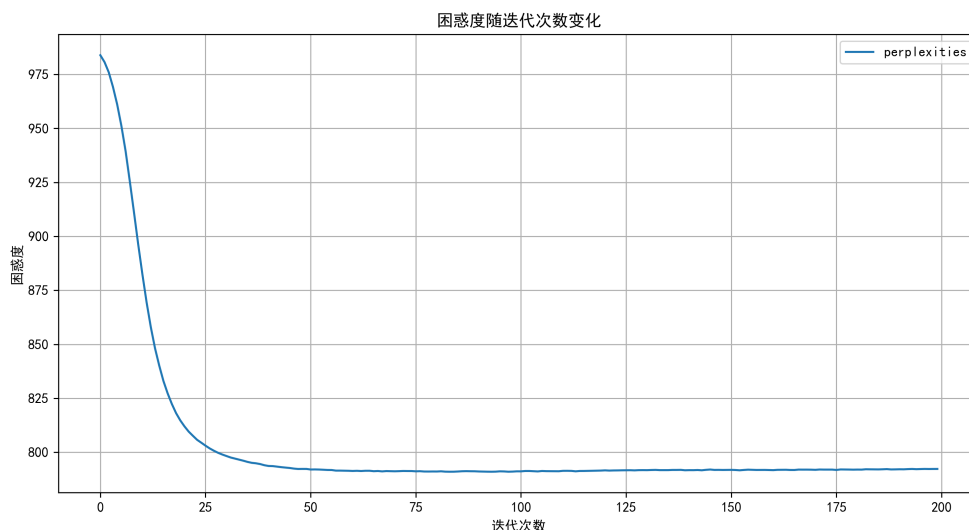
4.实验结果

4.1 16本小说

- 分词：



- 分字符



分词模式下，对每个主题下的词根据概率进行降序排列，部分结果展示如下（分字符不做展示）：

```

1  主题0的高频词为：
2  石破天:0.000719  道:0.000559  孩子:0.000551  不知:0.000496  没:0.000454  知
   道:0.000454  想:0.000395  说:0.000383  武林中:0.000328  这位:0.000320
3  -----
4  主题1的高频词为：
5  心中:0.000723  爹爹:0.000694  见:0.000673  一个:0.000673  两人:0.000589  突
   然:0.000559  快:0.000475  两个:0.000454  伸手:0.000416  站:0.000416
6  -----
7  主题2的高频词为：
8  中:0.001051  众人:0.000820  死:0.000454  住:0.000433  大叫:0.000425  下
   去:0.000366  有人:0.000324  出:0.000307  话:0.000294  周伯通:0.000265
9  -----
10 主题3的高频词为：
11 著:0.000790  麽:0.000774  李文秀:0.000669  走:0.000664  老人:0.000589  汉
   子:0.000580  没:0.000509  说:0.000496  只见:0.000479  师父:0.000467
12 -----
13 主题4的高频词为：
14 道:0.001333  袁承志:0.001030  青青:0.000467  知道:0.000408  问:0.000400  不
   能:0.000349  生死:0.000320  徐天宏:0.000311  话:0.000311  罢:0.000307
15 -----
16 主题5的高频词为：
17 :0.003031  道:0.001829  麽:0.000837  杨过:0.000711  甚:0.000669  郭靖:0.000664
   便:0.000471  夫妇:0.000458  武功:0.000454  说:0.000416

```

可以看到，心中，爹爹，孩子，夫妇，武功，杨过，郭靖，师父等是高频词，效果还不错。

- 分词模型下SVM训练测试结果：

```

1  overall accuracy:0.700000
2  -----
3  accuracy for each class: [0.69230769 1.          1.          0.75
   0.66666667 1.
4  0.          0.58823529 0.73333333 0.77777778 0.75          0.64
5  0.625          0.88888889 0.71428571 0.58823529]
6  -----
7  average accuracy:0.713421
8
9  overall accuracy:0.475000

```

```

10 -----
11 accuracy for each class: [0.5          0.          0.          1.
0.71428571 0.
12 0.          0.25         0.          0.          0.          0.33333333
13 0.4          0.5          1.          0.6          ]
14 -----
15 average accuracy:0.331101

```

- 分字符模型下SVM训练测试结果：

```

1 overall accuracy:0.768750
2 -----
3 accuracy for each class: [0.66666667 0.75          1.          1.
0.62068966 1.
4 0.83333333 0.66666667 1.          1.          0.83333333 0.73333333
5 1.          0.75          0.83333333 0.8          ]
6 -----
7 average accuracy:0.842960
8
9 overall accuracy:0.500000
10 -----
11 accuracy for each class: [1.          0.85714286 0.          0.          0.
0.5
12 1.          0.33333333 1.          0.          1.          0.
13 0.25         0.          ]
14 -----
15 average accuracy:0.424320

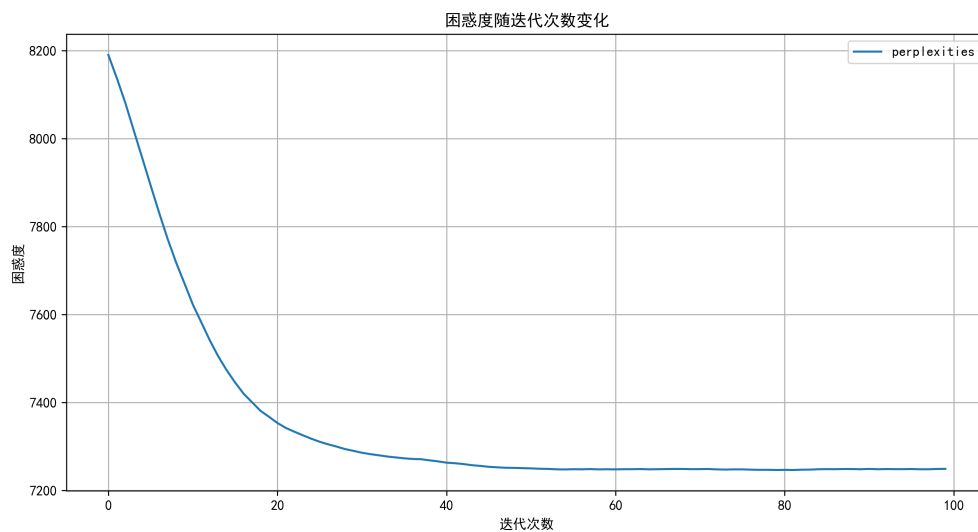
```

16分类下SVM在测试集的表现非常差。

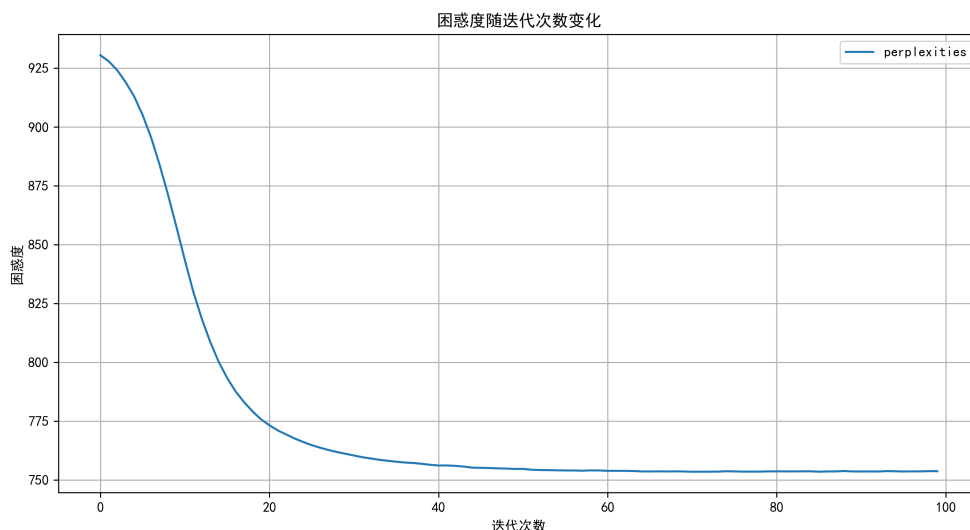
4.2 5本小说

降低分类难度，选择5本小说。

- 分词



- 分字符



分词模式下，对每个主题下的词根据概率进行降序排列，部分结果展示如下（分字符不做展示）：

- 1 主题0的高频词为：
- 2 萧峰:0.000606 一声:0.000450 伸手:0.000433 却是:0.000408 想:0.000374 起
- 3 来:0.000362 怎地:0.000349 金花婆婆:0.000332 杀:0.000328 有人:0.000324
- 4 主题1的高频词为：
- 5 道:0.003566 说:0.000782 难道:0.000471 罢:0.000450 没:0.000433 二人:0.000416
- 6 不肯:0.000353 定:0.000324 两位:0.000320 一句:0.000320
- 7 主题2的高频词为：
- 8 道:0.003036 便:0.001846 说:0.001472 说道:0.000669 弟子:0.000555 这
- 9 位:0.000391 和尚:0.000391 包不同:0.000391 想:0.000379 星宿:0.000379
- 10 主题3的高频词为：
- 11 郭靖:0.001985 黄蓉:0.001535 两人:0.000808 欧阳锋:0.000799 洪七公:0.000669 一
- 12 个:0.000547 吃:0.000534 功夫:0.000526 黄蓉道:0.000463 转身:0.000379
- 13 主题4的高频词为：
- 14 中:0.000871 慕容复:0.000702 丐帮:0.000589 段誉:0.000475 帮主:0.000459 王语
- 15 嫣:0.000429 心下:0.000404 段:0.000395 一个:0.000379 杨:0.000374
- 16 主题5的高频词为：
- 17 师父:0.001653 武功:0.000694 黄药师:0.000589 中:0.000505 心中:0.000505 梅超
- 风:0.000429 弟子:0.000404 点:0.000370 少林寺:0.000341 欧阳克:0.000332

可以看出，许多我们耳熟能详的词作为高频词被选了出来，效果非常好。

- 分词模型下SVM训练测试结果：

```

1 overall accuracy:0.793750
2 -----
3 accuracy for each class: [0.81818182 0.83333333 0.7755102  0.70588235
4 0.80645161]
5 -----
6 average accuracy:0.787872
7
8 overall accuracy:0.625000
9 -----
10 accuracy for each class: [0.85714286 1.          0.36842105 1.          0.6
11 ]
12 -----
13 average accuracy:0.765113

```

测试集的准确率达到0.765.

- 分字符模型下SVM训练测试结果:

```

1 overall accuracy:0.925000
2 -----
3 accuracy for each class: [0.91428571 0.96666667 0.90322581 0.96428571
4 0.88888889]
5 -----
6 average accuracy:0.927471
7
8 overall accuracy:0.925000
9 -----
10 accuracy for each class: [0.9          1.          0.91666667 1.          0.8
11 ]
12 -----
13 average accuracy:0.923333

```

测试集的准确率达到0.923，效果比分词模型更佳.

5. 结论

LDA模型能够较好地解决一词多义和多词一意的的问题，实验说明了LDA的有效性，并通过SVM进行了验证，针对金庸小说，有效的分类通常是一些有密切关系的人物，或者有关联的动作等。在分类类别非常多时，效果比较差，而在5分类时，效果很好，且分字符模式的效果优于分词模式。