

# 机器学习Lab4实验报告

---

## Topic-Model——LDA实验

---

PB19020499 桂栋南

### 实验原理

#### 总述

**LDA模型是一种基于文本的生成模型，是基于贝叶斯的话题模型。**

假设的分布：

- 文本由话题的多项分布表示；
- 话题由单词的多项分布表示；
- 文本的话题分布是狄利克雷分布；
- 话题的单词分布是狄利克雷分布。

**LDA是含有因变量的概率图模型：**

其中隐变量为：

- 话题的单词分布；
- 文本的话题分布；
- 文本中各个位置的话题。

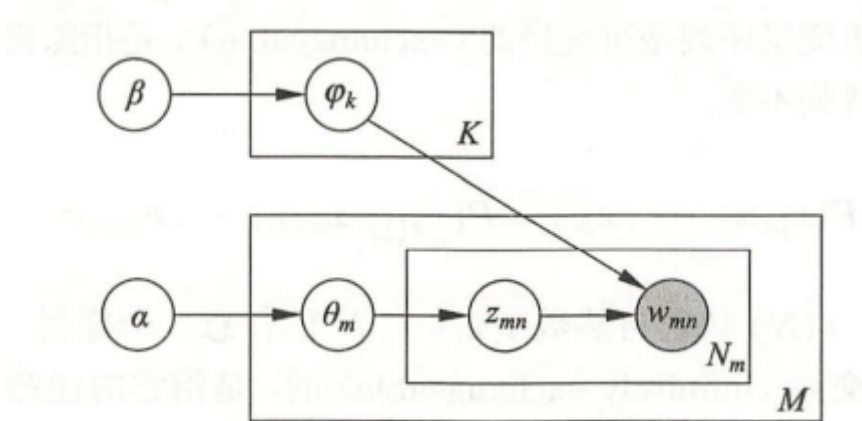
观测变量为：

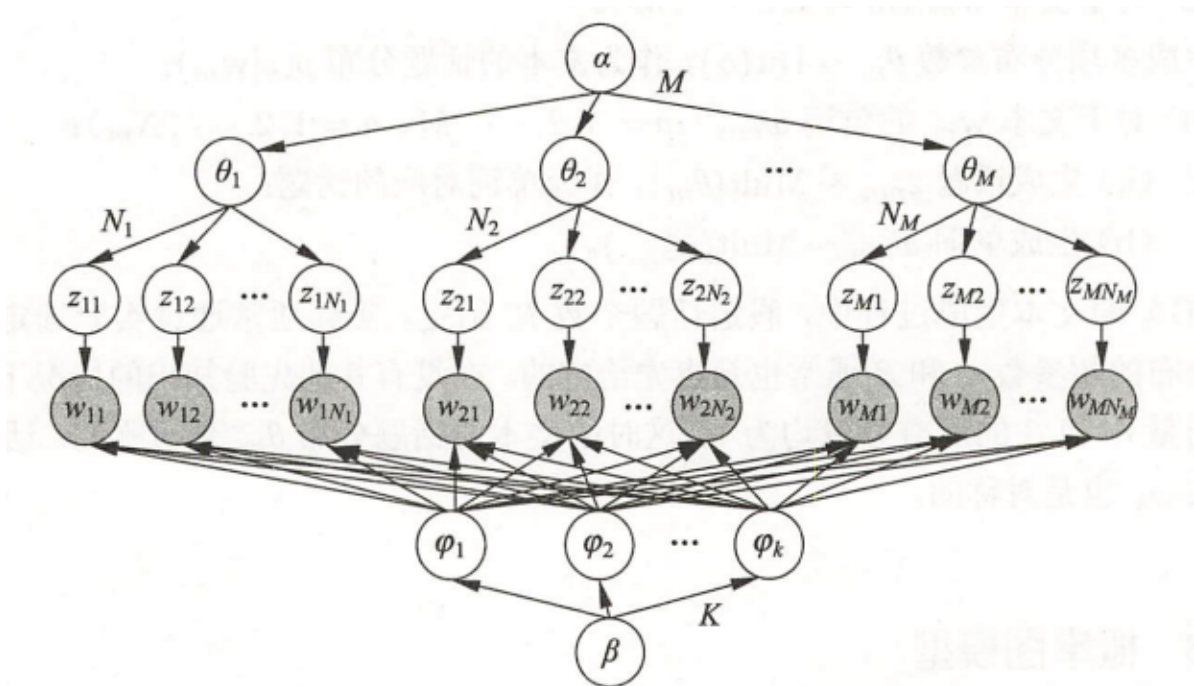
- 文本中各个位置的单词。

#### 模型定义

符号	含义	代码
$V$	单词的总数	V
$W = \{w_1, \dots, w_V\}$	文本集合, $w_v$ 是第 $v$ 个单词	
$M$	文章总数	M
$D = \{\mathbf{w}_1, \dots, \mathbf{w}_M\}$	文本集合, $\mathbf{w}_m$ 是第 $m$ 个文本, $\mathbf{w}_m = (w_{m1}, \dots, w_{mN_m})$	
$K$	话题总数	K
$Z = \{z_1, \dots, z_K\}$	话题集合, $z_k$ 是第 $k$ 个话题	
$\varphi_k = (\varphi_{k1}, \dots, \varphi_{kV})$	$\varphi_{kv}$ 表示话题 $z_k$ 生成单词 $w_v$ 的概率 $p(w z_k)$ , 服从狄利克雷分布	
$\beta = \{\beta_1, \dots, \beta_V\}$	$\varphi_{kv}$ 的狄利克雷分布的超参数	beta
$\theta_m = \{\theta_{m1}, \dots, \theta_{mK}\}$	$\theta_{mk}$ 表示文本 $\mathbf{w}_m$ 生成话题 $z_k$ 的概率 $p(z \mathbf{w}_m)$ , 服从狄利克雷分布	
$\alpha = \{\alpha_1, \dots, \alpha_K\}$	$\theta_{mk}$ 的狄利克雷分布的超参数	alpha
$n_{mk}$	第 $m$ 篇文档中由 $k$ 这个topic产生的单词计数	nmk
$n_{kv}$	第 $k$ 个topic产生单词 $v$ 的计数	nk <sub>v</sub>
$n_m$	第 $m$ 篇文档中产生全部单词的计数	nm
$n_k$	第 $k$ 个topic产生全部单词的计数	nk

## 概率图模型





## LDA生成过程

1. 根据 $\beta$ 生成 $\varphi_k \sim \text{Dir}(\beta)$ , 作为话题 $z_k$ 的单词分布 $p(w|z_k)$
2. 根据 $\alpha$ 生成 $\theta_m \sim \text{Dir}(\alpha)$ , 作为文本 $\mathbf{w}_m$ 的话题分布 $p(z|\mathbf{w}_m)$
3. 根据 $\theta_m$ 生成话题 $z_{mn} \sim \text{Multi}(\theta_m)$ , 作为单词 $w_{mn}$ 对应的话题
4. 根据 $\varphi_k$ 和 $z_{mn}$ 生成单词 $w_{mn} \sim \text{Multi}(\varphi_{z_{mn}})$ , 得到结果

## LDA算法推导

推断目标:

1. 话题序列的集合 $\mathbf{z} = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ 的后验概率分布
2. 参数 $\theta = \{\theta_1, \dots, \theta_M\}$
3. 参数 $\varphi = \{\varphi_1, \dots, \varphi_K\}$

后验概率 $p(\mathbf{z}|\mathbf{w}, \alpha, \beta)$ 的吉布斯抽样

$$p(\mathbf{z}|\mathbf{w}, \alpha, \beta) \propto p(\mathbf{w}, \mathbf{z}|\alpha, \beta) = p(\mathbf{w}|\mathbf{z}, \beta)p(\mathbf{z}|\alpha)$$

其中有:

$$p(\mathbf{w}|\mathbf{z}, \beta) = \prod_{k=1}^K \prod_{v=1}^V \varphi_{kv}^{n_{kv}} = \prod_{k=1}^K \frac{B(n_k + \beta)}{B(\beta)}$$

$$p(\mathbf{z}|\alpha) = \int p(\mathbf{z}|\theta)p(\theta|\alpha)d\theta = \prod_{m=1}^M \frac{B(n_m + \alpha)}{B(\alpha)}$$

所以

$$p(\mathbf{z}|\mathbf{w}, \alpha, \beta) \propto \prod_{k=1}^K \frac{B(n_k + \beta)}{B(\beta)} \cdot \prod_{m=1}^M \frac{B(n_m + \alpha)}{B(\alpha)}$$

于是满条件分布

$$p(z_i|\mathbf{z}_{-i}, \mathbf{w}, \alpha, \beta) = \frac{1}{Z_{z_i}} p(\mathbf{z}|\mathbf{w}, \alpha, \beta) \propto \frac{n_{kv} + \beta_v}{\sum_{v=1}^V (n_{kv} + \beta_v)} \cdot \frac{n_{mk} + \alpha_k}{\sum_{k=1}^K (n_{mk} + \alpha_k)}$$

## 参数估计

$$\theta_{mk} = \frac{n_{mk} + \alpha_k}{\sum_{k=1}^K (n_{mk} + \alpha_k)}$$
$$\varphi_{kv} = \frac{n_{kv} + \beta_v}{\sum_{v=1}^V (n_{kv} + \beta_v)}$$

## LDA-Gibbs算法

输入: 文本的单词序列  $\mathbf{w} = \{\mathbf{w}_1, \dots, \mathbf{w}_m, \dots, \mathbf{w}_M\}$ ,  $\mathbf{w}_m = (w_{m1}, \dots, w_{mn}, \dots, w_{mN_m})$ ;

输出: 文本的话题序列  $\mathbf{z} = \{\mathbf{z}_1, \dots, \mathbf{z}_m, \dots, \mathbf{z}_M\}$ ,  $\mathbf{z}_m = (z_{m1}, \dots, z_{mn}, \dots, z_{mN_m})$  的后验概率分布  $p(\mathbf{z} | \mathbf{w}, \alpha, \beta)$  的样本计数, 模型的参数  $\varphi$  和  $\theta$  的估计值;

参数: 超参数  $\alpha$  和  $\beta$ , 话题个数  $K$ 。

1. 设所有计数矩阵的元素  $n_{mk}, n_{kv}$ , 计数向量的元素  $n_m, n_k$  初值为 0;
2. 对所有文本  $\mathbf{w}_m, m = 1, 2, \dots, M$   
对第  $m$  个文本中的所有单词  $w_{mn}, n = 1, 2, \dots, N_m$ 
  1. 抽样话题  $z_{mn} = z_k \sim \text{Mult}(\frac{1}{K})$ ;
  2. 增加文本-话题计数  $n_{mk} = n_{mk} + 1$ , 增加文本-话题和计数  $n_m = n_m + 1$ , 增加话题-单词计数  $n_{kv} = n_{kv} + 1$ , 增加话题-单词和计数  $n_k = n_k + 1$ ;
3. 循环执行以下操作, 直到进入燃烧期  
对所有文本  $\mathbf{w}_m, m = 1, 2, \dots, M$   
对第  $m$  个文本中的所有单词  $w_{mn}, n = 1, 2, \dots, N_m$ 
  1. 当前的单词  $w_{mn}$  是第  $v$  个单词, 话题指派  $z_{mn}$  是第  $k$  个话题;
  2. 减少计数  $n_{mk} = n_{mk} - 1, n_m = n_m - 1, n_{kv} = n_{kv} - 1, n_k = n_k - 1$ ;
  3. 按照满条件分布进行抽样

$$p(z_i | \mathbf{z}_{-i}, \mathbf{w}, \alpha, \beta) \propto \frac{n_{kv} + \beta_v}{\sum_{v=1}^V (n_{kv} + \beta_v)} \cdot \frac{n_{mk} + \alpha_k}{\sum_{k=1}^K (n_{mk} + \alpha_k)}$$

得到新的第  $k'$  个话题, 分配给  $z_{mn}$ ;

4. 增加计数  $n_{mk'} = n_{mk'} + 1, n_m = n_m + 1, n_{k'v} = n_{k'v} + 1, n_{k'} = n_{k'} + 1$ ;
  5. 得到更新的两个计数矩阵  $N_{K \times V} = [n_{kv}]$  和  $N_{M \times K} = [n_{mk}]$ , 表示后验概率分布  $p(\mathbf{z} | \mathbf{w}, \alpha, \beta)$  的样本计数; 利
4. 用得到的样本计数, 计算模型参数

$$\theta_{mk} = \frac{n_{mk} + \alpha_k}{\sum_{k=1}^K (n_{mk} + \alpha_k)}$$
$$\varphi_{kv} = \frac{n_{kv} + \beta_v}{\sum_{v=1}^V (n_{kv} + \beta_v)}$$

## 实验实现

### 总体框架

```
1 class LDA():
2
3     def __init__(self, alpha=6, beta=0.01, iter=50, topics_num=8):
4         self.alpha = alpha
5         self.beta = beta
6         # 最大迭代次数
```

```

7         self.iter = iter
8         # 所有doc的列表和停用词列表
9         documents, stopwords = self.read_file()
10        print("Reading file has finished!")
11        # docs是二维列表, docs[i][j]是第i篇文章第j个单词的对应的编号
12        # word2id是字典, 记录关键字是单词, 值为对应编号的转换关系
13        # id2word是字典, 记录关键字是编号, 值为对应单词的转换关系
14        self.docs, self.word2id, self.id2word = self.preprocessing(
15            documents, stopwords)
16        print("Preprocessing has finished!")
17        # M是文章数量
18        self.M = len(self.docs)
19        # V是所有单词的数量
20        self.V = len(self.word2id)
21        # Z是二维列表, z[i][j]记录第i篇文章第j个单词的topic
22        self.Z = []
23        # K是topic数量
24        self.K = topics_num
25        # 初始化
26        # nmk是array(M*K), nmk[m][k]表明第m篇文档中由k这个topic产生的单词计数
27        self.nmk = np.zeros([self.M, self.K]) + self.alpha
28        # nkV是array(K*V), nkV[k][v]表明第k个topic产生单词v的计数
29        self.nkv = np.zeros([self.K, self.V]) + self.beta
30        # nk是array(K), nk[k]表明第k个topic产生全部单词的计数
31        self.nk = np.zeros([self.K]) + self.V * self.beta
32
33        # 读取文件
34        def read_file(self):
35            pass
36
37        # 预处理(分词, 去停用词, 为每个word赋予一个编号, 文档使用word编号的列表表示)
38        def preprocessing(self, documents, stopwords):
39            pass
40
41        # 进行初始化的多项分布的分配
42        def multinomial(self):
43            pass
44
45        # gibbs采样
46        def gibbsSampling(self):
47            pass
48
49        def perplexity(self):
50            pass
51
52        def run(self):
53            pass
54

```

- `init` 初始化读文件并进行预处理, 定义模型参数。
- `run` 函数调用初始化的 `multinomial` 函数对话题采样, 并调用进行循环的 `gibbsSampling` 函数进行话题采样, 并调用 `perplexity` 函数存储生成匹配情况。

## 预处理函数

```
1  # 预处理(分词, 去停用词, 为每个word赋予一个编号, 文档使用word编号的列表表示)
2  def preprocessing(self, documents, stopwords):
3      word2id = {}
4      id2word = {}
5      docs = []
6      currentDocument = []
7      currentWordId = 0
8
9      flag_list = ['n', 'nz', 'vn']
10     for document in documents:
11         # 分词
12         segList = psg.cut(document)
13         for seg_word in segList:
14             word = seg_word.word
15             # 单词长度大于1, 并且不包含数字和单词, 并且不是停止词, 并且词性符合要求
16             if len(word) > 1 and not re.search('[0-9]', word) \
17                 and word not in stopwords and not re.search('[a-z]',
18 word) and not re.search('[A-Z]', word) \
19                 and seg_word.flag in flag_list:
20                 # word出现过, docs直接增加其编号
21                 if word in word2id:
22                     currentDocument.append(word2id[word])
23                 else:
24                     # 增加一个新的word-ID, 加入docs
25                     currentDocument.append(currentWordId)
26                     # 记录word->id的对应关系
27                     word2id[word] = currentWordId
28                     # 记录id->word的对应关系
29                     id2word[currentWordId] = word
30                     # 准备下一个word-ID
31                     currentWordId += 1
32             docs.append(currentDocument)
33             currentDocument = []
34     return docs, word2id, id2word
```

- 该函数实现预处理功能。
  - 实现去除停用词和单词数字;
  - 记录单词word和训练的id之间的匹配关系。

## 初始多项分布的话题分配

```

1 def multinomial(self):
2     for d, doc in enumerate(self.docs):
3         # d是文档的序号, doc是文档内容相应单词的编号
4         zCurrentDoc = []
5         for w in doc:
6             z = np.random.multinomial(1, [1/self.K]*self.K).argmax()
7             zCurrentDoc.append(z)
8             self.nmk[d, z] += 1
9             self.nkv[z, w] += 1
10            self.nk[z] += 1
11            self.nm[d] += 1
12        self.Z.append(zCurrentDoc)

```

- 该函数实现初始化的多项分布的分配。

## Gibbs采样

```

1 # gibbs采样
2 def gibbsSampling(self):
3     # 为每个文档中的每个单词重新采样topic
4     for d, doc in enumerate(self.docs):
5         # d是文档的序号, doc是文档内容相应单词的编号
6         for index, w in enumerate(doc):
7             # index是单词w在该篇文章doc中的序号位置
8             z = self.Z[d][index]
9             # 将当前文档当前单词原topic相关计数减1
10            self.nmk[d, z] -= 1
11            self.nkv[z, w] -= 1
12            self.nk[z] -= 1
13            self.nm[d] -= 1
14            # 重新计算当前文档当前单词属于每个topic的概率
15            pz = np.multiply((self.nkv[:, w]+self.beta),
16                             (self.nmk[d, :] + self.alpha)) / np.sum(self.nmk[d, :] + self.alpha)
17            # 按照计算出的分布进行采样
18            z = np.random.multinomial(1, pz / pz.sum()).argmax()
19            self.Z[d][index] = z
20            # 将当前文档当前单词新采样的topic相关计数加上1
21            self.nmk[d, z] += 1
22            self.nkv[z, w] += 1
23            self.nk[z] += 1
24            self.nm[d] += 1

```

- 该函数实现吉布斯采样过程。
  - 每次采样是对于每篇文档的每个单词的主题的重新分配，分配基于推导得到的参数的多项分布。

## run函数

```

1 def run(self):
2     self.multinomial()
3     # 记录匹配度
4     self.perplexity_score = []
5     for i in range(0, self.iter):

```

```

6         self.gibbsSampling()
7         self.perplexity_score.append(self.perplexity())
8         # 打印iter信息
9         print(time.strftime('%X'), "Iteration: ", i, " completed",
10              " Perplexity: ", self.perplexity_score[-1])
11         if i < 10:
12             continue
13         # 如果匹配度大于前十次的平均值, 停止迭代
14         if self.perplexity_score[-1] >
np.average(self.perplexity_score[-10:]):
15             break

```

- 该函数调用 `multinomial` 和 `gibbsSampling` 进行LDA算法的主要功能
  - `gibbsSampling` 的停止条件是perplexity大于前十次的平均值, 表明该值不在下降, 迭代停止。

## get\_topic 函数

```

1 def get_topic(self, maxTopicWordsNum=10):
2     topicwords = []
3     for z in range(0, self.K):
4         # 对于每个主题, 对在该主题下出现单词的数量排序
5         ids = self.nkv[z, :].argsort()
6         topicword = []
7         # 根据id->word的转换关系转为id
8         for j in ids:
9             topicword.insert(0, self.id2word[j])
10            # 返回需求的每个主题的相应数量的单词
11            topicwords.append(topicword[0: maxTopicWordsNum])
12    return topicwords

```

- 该函数返回需要的每个主题的前maxTopicWordsNum个主题单词。

## 实验结果

### 主题情况

```

array([[ '研究', '科学家', '人类', '时间', '人员', '流感', '技术', '消息', '科技', '利用', 科技
        '太空', '航天飞机', '流感病毒', '海豚', '科学'],
       [ '学生', '经济', '大学', '学校', '工作', '政府', '计划', '教育', '专业', '记者', '情况', 教育
        '银行', '孩子', '能力', '国家'],
       [ '网站', '手机', '互联网', '公司', '网络', '信息', '服务', '网游', '用户', '业务', 彩票
        '视频', '内容', '彩票', '记者', '全国'],
       [ '专家', '网友', '老师', '压力', '分析', '走势', '突破', '黄金', '成本', '股票', '新浪', 股票
        '大盘', '建议', '调整', '机会'],
       [ '比赛', '主队', '赔率', '球队', '主场', '火箭', '客场', '数据', '球员', '公司', '联赛', 体育
        '时间', '情况', '篮板', '奇才'],
       [ '游戏', '比赛', '电子竞技', '玩家', '海选', '总决赛', '冠军', '赛事', '战队', '赛区', 游戏
        '世界', '项目', '奖金', '时间', '星际争霸'],
       [ '项目', '市场', '发展', '投资', '生活', '建筑', '企业', '活动', '地产', '国际', '新浪', 房产
        '文化', '空间', '产品', '设计'],
       [ '电影', '主持人', '票房', '影片', '观众', '演员', '作品', '朋友', '故事', '娱乐', 娱乐
        '现场', '新浪', '感觉', '电影节', '角色']], dtype='<U4')

```

- 红字是根据 `data.xlsx` 得到的主题对应情况。



## 可视化

简单根据gibbs采样50次的运行结果作出下图：

