



学 号 20376310

北京航空航天大学
B E I H A N G U N I V E R S I T Y

深度学习与自然语言处理

大作业 2

文本 LDA 建模与主题分类

院（系）名称	自动化科学与电气工程学院
专业名称	自动化
学生姓名	杨佳木
学 号	20376310

2024 年 5 月

1.摘要

本作业基于 LDA 主题模型和深度学习分类模型，完成给定语料库段落分类任务。首先，基于 LDA 主题模型的构建方法，得到训练文本的主题特征向量的分布。其次，将训练集所得 LDA 模型的状态特征向量输入通过至深度学习分类模型，标签为各个段落对应的小说，获得训练完成的模型。最后将测试集通过 LDA 模型所得的特征向量喂入训练完成的分类器，得到测试集分类结果，计算分类正确率，并通过 K 折交叉验证消除过拟合造成的影响。

本文基于上述的设计流程，通过设计三个实验改变不同主题数量 T ，选择“字”或“词”作为基本单元或选择不同取值长度 K ，并验证短文本和长文本对分类效果的影响。本文最终得到在其他情况不变的时候，主题数 T 越多，训练准确性越高，但是相同 epoch 下，更易出现过拟合现象。token 数越多，训练准确性越高，且能够一定程度上减弱过拟合的现象，深度学习方法相较于传统方法分类有优势。token 数超过 1000，如取 token 数为 3000，以词为基本单元的分类效果好，低于 1000，以字为基本单元的分类效果好，token 为 1000，二者效果差不多。深度学习方法相较于传统方法分类有优势。

2. 问题描述

本文基于 LDA 模型在给定的语料库中抽取 1000 个段落最为数据集（每个段落可以选取 K 个 token, K 可取 20, 100, 500, 1000, 3000），标签为对应段落所属的小说，利用并把每个段落表示为主题分布后进行分类（分类器自由选择）分类。结果使用 10 次交叉验证（900 组数据做训练，剩余 100 组做测试循环十次）。研究问题如下：

- （1）在设定不同的主题个数 T 的情况下，分类性能是否有变化？
- （2）以“词”和以“字”为基本单元下分类结果有什么差异？
- （3）不同的取值的 K 的短文本和长文本，主题模型性能上是否有差异？

3.问题的解决方案

LDA 在主题模型中占有非常重要的地位，常用来文本分类。该方法由 Blei, David M.、Ng, Andrew Y. Jordan 于 2003 年提出，用来推测文档的主题分布。它可以将文档集中每篇文档的主题以概率分布的形式给出，从而通过分析一些文档抽取出它们的主题分布后，便可以根据主题分布进行主题聚类或文本分类。

为了解决所提出的三个问题，解决方案流程如下：

1. 从给定的 16 本金庸小说数据集中，随机、均匀地抽取 1000 个段落，每个段落的标签为对应小说的小说名，每个段落包含 n 个字 ($n \geq 500$)，每个段落作为一个样本；将随机抽取的 k 个段落（即 k 个样本）中的 80% 作为训练样本，剩余 20% 作为测试样本。对样本进行 Jieba 中文分词，将文本转化为一系列词汇，并根据停用词词表过滤无意义词汇和标点符号，保留用用的词汇。构建词袋模型，进行 LDA 建模，设定段落主题数目 T ，得到每个段落的主题特征向量分布。

2. 构建基于深度学习的分类模型，对不同段落特征进行分类训练，并使用验证集验证分类效果。使用 K 折交叉验证避免过拟合问题。

针对问题 1，设计实验 1，通过改变主题数 T 以改变分类效果，得到主题数与分类效果间的关系。

针对问题 2，设计实验 2，通过改变基本单元的类型，得到基本单元的类型与分类效果间的关系。

针对问题 3，设计实验 3，通过改变取值长度 K ，得到取值长度与分类效果间的关系。

4. 基于 LDA 的主题分类模型的建立

基于语料库，使用 LDA 模型进行主题模型建模，具体建模方法如下：

1. 对于每一篇文档 m ，选择一个文档-主题分布

$$\theta_m \sim \text{Dirichlet}(\alpha) \quad (1)$$

对于文档 m 中的每一个词 n ，选择一个主题

$$z_{mn} \sim \text{Multinomial}(\theta_m) \quad (2)$$

并选择一个词，有词的分布先验满足

$$w_{mn} \sim \text{Multinomial}(\beta_{z_{mn}}) \quad (3)$$

其中， α 是文档-主题分布的先验参数， β 是主题-词分布的先验参数。

2. 参数估计

LDA 的目标是估计出 (θ) 和 (β) ，使得整个文档集的概率最大化。

基于这一理论思想，以字为基本单位，假设文章主题数为 10，构建语料库

主题模型如下：

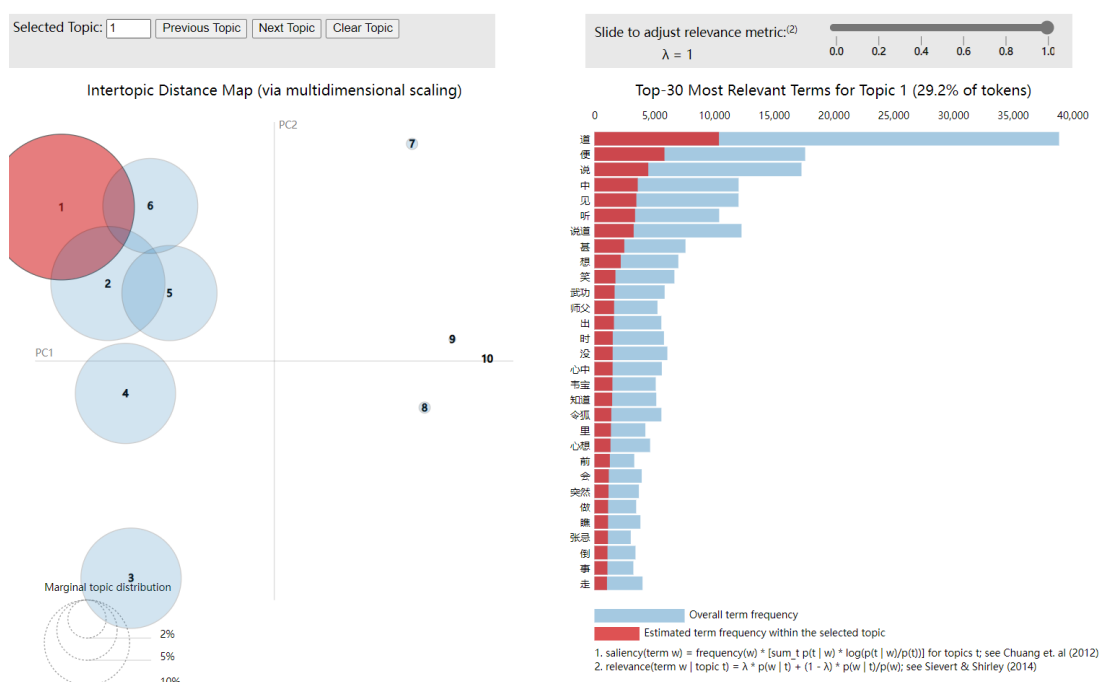


图 1 语料库 LDA 主题模型构建

如图左侧所示为主题数目，右侧所示为每个主题相关词（字）（列举 top-30）。

基于这一主题模型，构建测试集 LDA 主题分布特征矩阵，从每篇小说中抽取的所有段落组合到一起作为训练集，进行训练。首先为每篇文章中的词随机分配一个 topic，然后统计每个段落的 topic 频率。结果如下：

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	
1	0	0	0.023892	0.037599	0	0.031703	0	0	0	0	0	0	0	0.023814	0	0.096644	0	0	0.011709	0.040585	0	0	0.027454	0.034473	0.011103	0	0	0	
2	0	0	0.003992	0.016183	0.029431	0.014203	0.023217	0	0	0	0	0	0	0.013859	0	0.064715	0	0.014929	0.030656	0.016432	0	0	0	0.038864	0	0	0.015486	0	
3	0	0	0.003964	0.017444	0.043776	0.011732	0.009783	0	0	0	0	0	0	0	0	0.048489	0	0	0	0	0	0	0.018146	0.021396	0.011902	0	0	0	
4	0	0	0	0.017184	0.014546	0.038117	0	0.013885	0	0	0	0	0	0	0	0.0217915	0	0	0.022715	0.01947	0	0	0	0.02025	0.014444	0	0	0.012997	
5	0	0	0.011479	0	0.064621	0.044341	0.048427	0	0	0	0	0	0	0.023523	0	0.05551	0	0	0	0	0.022286	0	0	0	0.02422	0.027269	0	0	0
6	0	0	0.002997	0.016493	0.034558	0	0.074862	0	0	0	0	0	0	0.013929	0	0.067744	0	0.012071	0.017282	0.024569	0	0	0	0.020189	0	0	0	0	
7	0	0	0.002091	0.017126	0.02479	0.020243	0.049389	0	0	0	0	0	0	0.021777	0	0.07883	0	0	0.020772	0.036669	0	0	0.016208	0.015342	0	0.01387	0.018397	0	
8	0	0	0.019477	0.011781	0.020229	0.018905	0.018	0	0	0	0	0	0	0.011672	0	0.012827	0	0.013986	0.02248	0	0	0	0.019986	0.018756	0	0	0	0	
9	0	0	0.021638	0.020703	0.031412	0	0.03154	0	0	0	0	0	0	0	0	0.028781	0	0	0.024918	0.024054	0	0	0.014705	0.027349	0.012227	0	0.0179	0	
10	0.012811	0	0.0218	0.027982	0.020206	0.011987	0.017115	0	0	0	0.011878	0	0	0.013965	0	0.02403	0.02354	0.017485	0.029653	0	0	0	0.020219	0.012135	0.014438	0	0.01305	0	
11	0	0	0.019148	0	0.028906	0.016861	0.028186	0	0	0	0	0	0	0.01709	0	0.024203	0	0	0.017701	0.012169	0	0	0.028887	0.028205	0.012824	0	0	0	
12	0	0	0.020684	0.017982	0.026033	0.014486	0.022297	0	0	0	0.021867	0	0	0.016867	0	0.026379	0.017613	0.022408	0	0	0	0.01628	0.024067	0.017185	0	0.017744	0		
13	0	0	0.016221	0.029311	0.020209	0.017761	0.062522	0	0	0	0	0	0	0.01173	0	0.081647	0	0.012901	0.020302	0	0	0	0.01562	0.028803	0	0	0	0	
14	0	0	0.027791	0.017371	0.027221	0.011024	0.099039	0	0	0	0	0	0	0.013291	0	0.041176	0	0.018601	0.012586	0	0	0	0.020986	0.029676	0.020443	0	0	0	
15	0	0	0.01861	0.014786	0.031796	0.020229	0.060208	0	0	0	0	0	0	0.024102	0	0.057206	0	0	0.033411	0	0	0	0.014744	0.01818	0.020207	0	0	0	
16	0	0	0.018709	0.014144	0.020564	0.014049	0.086024	0	0	0	0	0	0	0.014965	0	0.020962	0	0.010485	0.024451	0	0	0	0.020184	0.027259	0.030801	0	0	0	
17	0.010462	0	0.012108	0.012789	0.013784	0.026112	0.034644	0	0	0	0	0	0	0	0	0.022291	0	0.018774	0.01874	0	0	0	0.020962	0.027136	0	0	0	0	
18	0	0	0.016786	0	0.021293	0	0.07232	0	0	0	0	0	0	0	0	0.04498	0	0	0.018737	0.011288	0	0	0.027226	0	0	0	0.02041	0	
19	0	0	0.018989	0.021432	0.024223	0	0.089852	0	0	0	0	0	0	0.024014	0	0.046999	0	0.012611	0.021048	0	0	0	0.014144	0.017186	0.016284	0	0.018427	0	
20	0	0	0.01521	0.021844	0.020882	0.020243	0.048496	0	0	0	0	0	0	0.024979	0	0.046211	0.016129	0.018187	0.030603	0	0	0	0.028645	0.020536	0	0.010715	0	0	
21	0	0	0.022826	0	0.014671	0.022229	0.029654	0	0	0	0	0	0	0	0	0.033885	0	0	0.015456	0.029964	0	0	0.01877	0.016434	0.016687	0.010083	0	0	
22	0	0	0.012347	0.016784	0.023143	0	0.018909	0	0	0	0.01221	0	0	0.018201	0	0.010826	0	0.010109	0	0.024613	0	0	0	0.020193	0	0	0	0	0
23	0	0	0.016534	0.020884	0.011032	0	0.010386	0	0	0	0	0	0	0.024932	0	0.068777	0	0.018939	0.021353	0.020097	0	0	0.028931	0.029413	0	0	0	0	
24	0	0	0.024051	0.012226	0.026293	0.070522	0	0	0	0	0	0	0	0.057267	0	0	0	0.019961	0.026779	0	0	0	0.016484	0.025767	0	0	0.010149	0	
25	0	0.016621	0	0.048489	0.015987	0.051744	0	0	0	0	0.021137	0	0	0.074211	0	0.01245	0	0.014686	0	0	0	0.019893	0.028969	0.012281	0	0.020207	0	0	
26	0	0	0.014928	0.034684	0.010407	0.098222	0	0	0	0	0	0	0	0.008993	0	0.018238	0.013932	0.017964	0	0	0	0.025802	0.022448	0	0	0	0	0	
27	0	0	0.015539	0.019547	0	0.072572	0	0	0	0	0	0	0	0.010584	0	0.067723	0	0.035148	0.015417	0	0	0.021519	0.01684	0	0	0	0	0	
28	0	0.020322	0.011821	0	0.011357	0.061642	0	0	0	0	0	0	0	0.020284	0	0.019554	0.011114	0.010162	0	0	0	0	0.015247	0	0	0	0	0	
29	0	0.022161	0.027473	0.019024	0.042111	0.038487	0	0	0	0	0	0	0	0.020767	0	0.048229	0	0.01187	0.018722	0.024726	0	0	0.044331	0.019799	0	0	0	0	
30	0	0.018789	0.027876	0.01737	0.021147	0.034543	0	0	0	0	0	0	0	0.029636	0	0.03389	0	0.014783	0.019797	0	0	0.024886	0.015743	0	0	0.012317	0	0	
31	0	0.011104	0.022372	0.024652	0	0.057796	0	0	0	0	0	0	0	0.050538	0	0	0	0.017486	0.010204	0	0	0	0.017486	0.010204	0	0	0	0	
32	0	0.026075	0.02548	0.021065	0.026339	0.030342	0	0	0	0	0	0	0	0.026694	0	0.043907	0	0.022279	0.029068	0	0	0.011871	0.020091	0.010787	0	0.014138	0	0	
33	0	0.012341	0	0.020779	0.021444	0.020322	0	0	0	0	0	0	0	0	0	0.047793	0	0	0.020261	0	0	0.020879	0.018391	0.013754	0.012504	0	0	0	
34	0	0.012887	0.019625	0.027923	0	0.058362	0	0	0	0	0	0	0	0.044963	0	0.044923	0.024371	0.020139	0	0	0	0.020255	0	0	0	0	0	0	
35	0	0.014125	0.018348	0.04023	0	0.028997	0	0	0.021897	0	0.011121	0	0.012906	0	0.020279	0.011584	0.02848	0	0.015988	0.014986	0	0	0.015988	0.014986	0	0.017873	0	0	
36	0.011313	0	0.012579	0.016253	0.020284	0	0.05941	0	0	0	0.018786	0	0.025973	0	0.020028	0.015177	0.029829	0	0.024889	0.018957	0	0	0	0.024889	0.018957	0	0	0	
37	0	0.003098	0.014037	0.029973	0	0.054164	0	0	0	0	0.016878	0	0.020287	0	0.012984	0.038186	0.02223	0	0	0	0.021581	0	0	0	0	0	0	0	
38	0	0.0106	0.01328	0.028223	0	0.049588	0	0	0	0	0	0	0	0.027622	0	0.043421	0	0.012646	0.022975	0	0	0.020561	0	0.020179	0	0	0	0	
39	0	0.013307	0.019184	0.012263	0	0.024701	0	0	0	0	0	0	0	0.025542	0	0.018951	0	0.015029	0	0	0	0.026547	0.028205	0	0	0.012847	0	0	
40	0.012187	0	0.011842	0.027482	0.029062	0.018886	0.027482	0	0	0	0	0	0	0.011999	0	0.016388	0	0	0.022185	0	0	0	0.029618	0.013111	0	0	0	0	
41	0	0.013731	0.017739	0.021208	0.013387	0.045618	0	0	0	0	0	0	0	0.023342	0	0.016208	0	0	0.012703	0	0	0	0.014027	0.019333	0.033779	0	0	0	
42	0	0.02098	0.023589	0.020253	0	0.060775	0	0	0	0	0	0	0	0.025185	0	0.054721	0	0.019997	0.019487	0	0	0	0.014827	0.019333	0	0	0	0	
43	0	0.01028	0.013911	0.028302	0.01115	0.048489	0	0	0	0	0	0	0	0.025005	0	0.04485	0.012186	0.011811	0	0	0	0.013737	0.020093	0	0.011432	0	0	0	
44	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	

图 2 段落主题频率

将段落名作为标签，构建训练集。同理，选取 10%的数据作为验证集。基于这一样本设计分类算法。

采用基于自注意力机制优化的多层感知机模型进行训练，该模型属于深度学习网络，网络结构如下：

```
class SelfAttnrtn(nn.Module):
    def __init__(self, embed_dim):
        super(SelfAttnrtn, self).__init__()
        self.query = nn.Linear(embed_dim, embed_dim)
        self.key = nn.Linear(embed_dim, embed_dim)
        self.value = nn.Linear(embed_dim, embed_dim)

    def forward(self, x):
```

```

        q = self.query(x)
        k = self.key(x)
        v = self.value(x)
        attn_weight = torch.matmul(q, k.transpose(1, 2))
        attn_weight = nn.functional.softmax(attn_weight, dim=-1)
        attned_values = torch.matmul(attn_weight, v)
        return attned_values

class Net(nn.Module):
    def __init__(self, in_dim, n_hidden_1, n_hidden_2, out_dim):
        super().__init__()
        self.attention = SelfAttnrntion(in_dim)
        self.layer1 = nn.Sequential(
            nn.Linear(in_dim, n_hidden_1), nn.ReLU(True), nn.Dropout(0.3))
        self.layer2 = nn.Sequential(
            nn.Linear(n_hidden_1, n_hidden_2), nn.ReLU(True), nn.Dropout(0.3))
        self.layer4 = nn.Sequential(nn.Linear(n_hidden_2, out_dim))

    def forward(self, x):
        x = self.attention(x)
        x = self.layer1(x)
        x = self.layer2(x)
        x = self.layer4(x)
        return x

```

基于这一网络结构，使用交叉熵损失函数进行训练，得到训练结果。并和基于 SVM 的训练结果进行对比。

```

100%|██████████| 90/90 [00:00<00:00, 402.55it/s]
100%|██████████| 10/10 [00:00<00:00, 1666.46it/s]
train acc = 39.222%, loss = 1.831393297513326
epoch = 29, valid acc = 29.00%, loss = 2.188531219959259

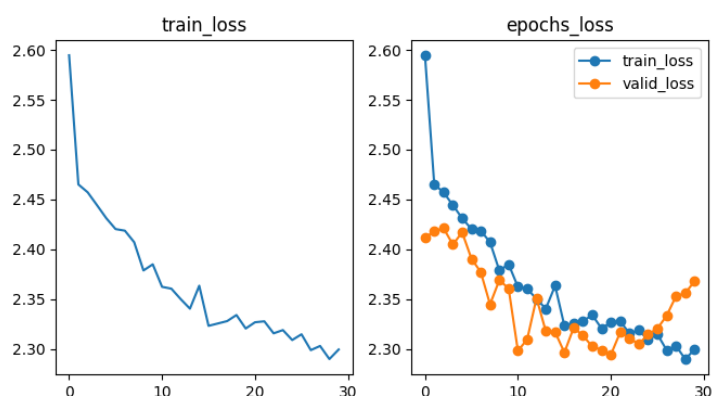
```

实验 1 不同主题数与分类结果间的关系

设置主题分类的数目分别为 $T=10$ 、20、50、100、200，固定 token 数为 1000，字为单位，分类结果如下

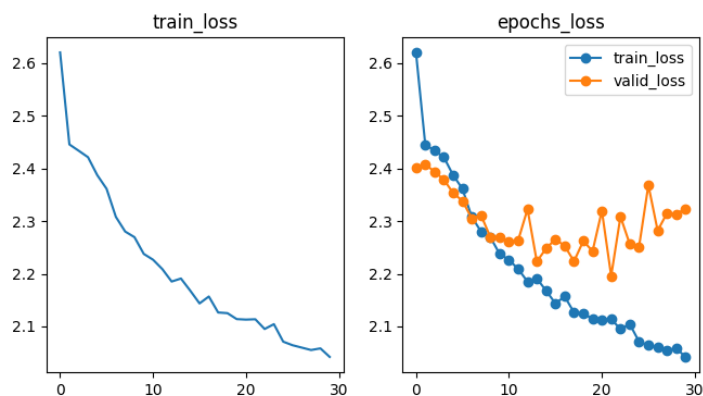
1. 分类主题数 $T=10$

训练集准确率 21.333%，测试集准确率 22.000%，损失函数如下：



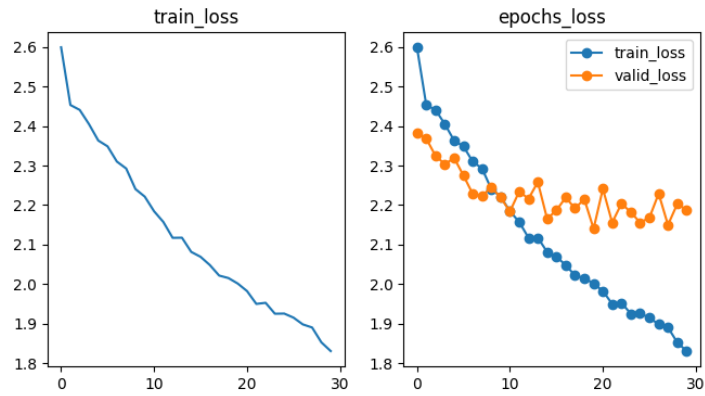
2. 分类主题数 $T=20$

训练集准确率 30.566%，测试集准确率 23.000%，损失函数如下：



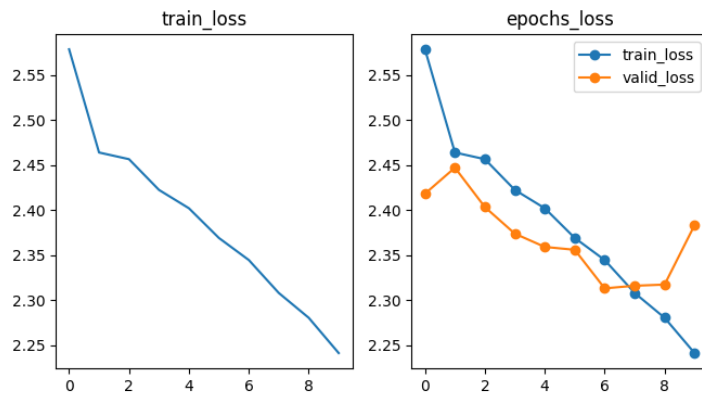
3. 分类主题数 $T=50$

训练集准确率 39.226%，测试集准确率 29.000%，损失函数如下：



4. 分类主题数 $T=100$

训练集准确率 45.245%，测试集准确率 39.000%，损失函数如下：

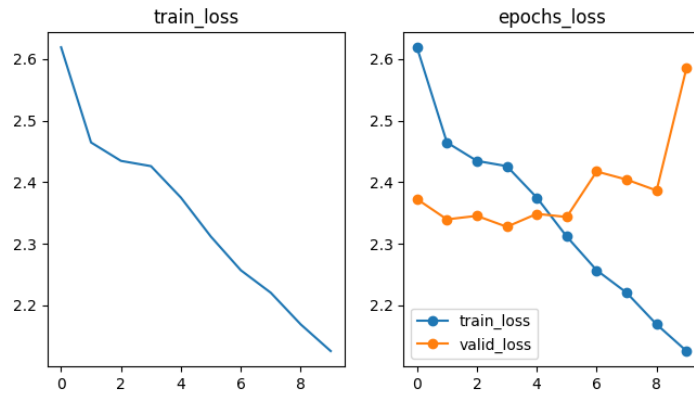


(epochs 过高，会出现过拟合)

5. 分类主题数 $T=200$

训练集准确率 65.937%，测试集准确率 32.000%，损失函数如下：

(epoch 为 10 的情况下，出现过拟合)



主题数 \ 准确率	Self Attention-MLP 训练集准确率/%	Self Attention-MLP 验证集准确率/%	SVM 验证集准确率/%	SVM 验证集准确率/%
10	21.333	22.000	12.564	11.423
20	30.566	23.000	18.342	15.342
50	39.226	29.000	23.231	19.258
100	45.245	39.000	27.543	25.379
200	65.397	32.000	35.437	32.456

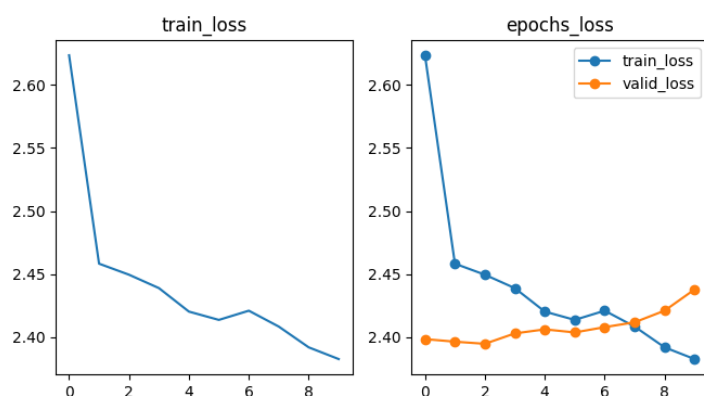
实验结论：主题数 T 越多，训练准确性越高，但是相同 epoch 下，更易出现过拟合现象，深度学习方法相较于传统方法分类有优势。

实验 2 不同 token 与分类结果间的关系

选取 token 数为 20, 100, 500, 1000, 3000, 以字为单位, K=100, 进行分类, 结果如下:

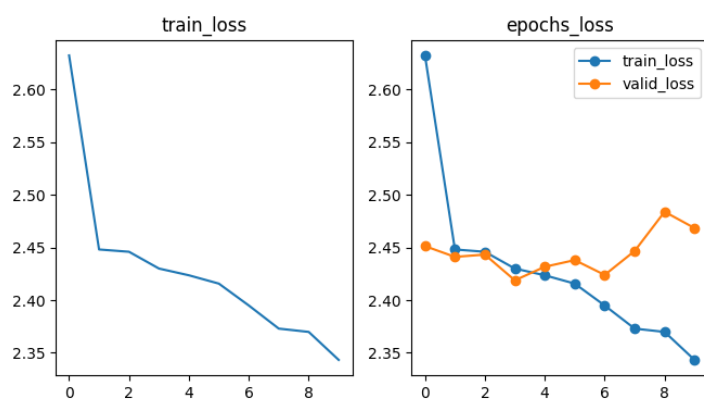
1. token 数 T=20

训练集准确率 17.333%, 测试集准确率 13.000%, 损失函数如下:



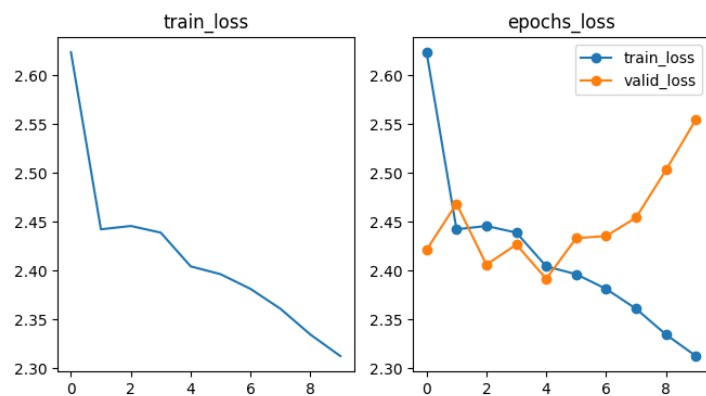
2. token 数 T=50

训练集准确率 19.667%, 测试集准确率 17.000%, 损失函数如下:



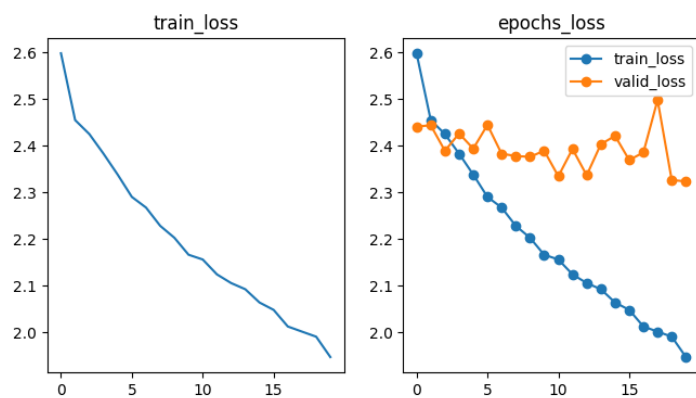
3. token 数 T=100

训练集准确率 25.556%, 测试集准确率 19.000%, 损失函数如下:



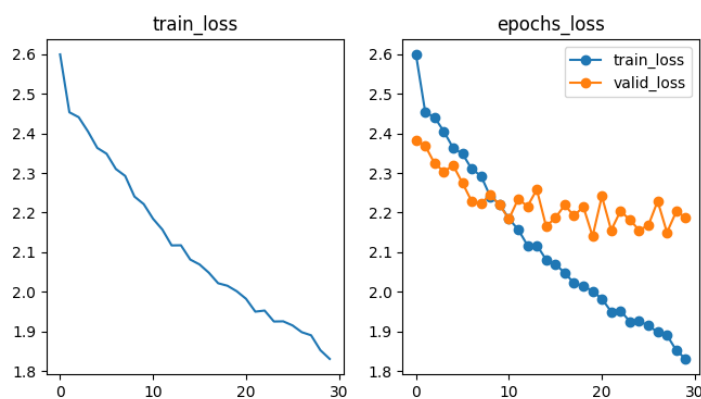
4. token 数 T=500

训练集准确率 34.123%，测试集准确率 31.000%，损失函数如下：



6. token 数 1000

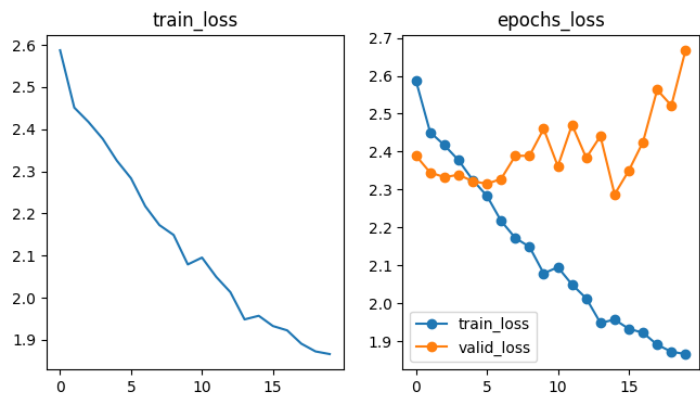
训练集准确率 41.582%，测试集准确率 36.000%，损失函数如下：



7. token 数 3000

训练集准确率 47.942%，测试集准确率 46.000%，损失函数如下：

8.



token 数 \ 准确率	Self Attention-MLP 训练集准确率 /%	Self Attention-MLP 验证集准确率 /%	SVM 验证集准 确率/%	SVM 验证集准 确率/%
20	17.333	13.000	13.663	12.426
50	19.667	17.000	15.462	15.331
100	25.556	19.000	21.233	19.258
500	34.123	31.000	27.543	25.379
1000	41.582	36.000	36.433	32.856
3000	47.942	46.000	40.231	37.538

实验总结：token 数越多，训练准确性越高，且能够一定程度上减弱过拟合的现象，深度学习方法相较于传统方法分类有优势。

实验 3 字与词与分类结果间的关系

选取 token 数为 1000，以字或词为单位，K=100，进行分类，结果如下：

准确率 字/词	Self Attention-MLP 训练集准确率 /%	Self Attention-MLP 验证集准确率 /%	SVM 验证集准 确率/%	SVM 验证集准 确率/%
字	41.582	36.000	36.433	32.856
词	41.268	37.000	35.462	31.331

选取 token 数为 3000，以字或词为单位，K=100，进行分类，结果如下：

准确率 字/词	Self Attention-MLP 训练集准确率 /%	Self Attention-MLP 验证集准确率 /%	SVM 验证集准 确率/%	SVM 验证集准 确率/%
字	47.942	46.000	40.231	37.538
词	48.387	42.000	37.462	32.331

选取 token 数为 100，以字或词为单位，K=100，进行分类，结果如下：

准确率 字/词	Self Attention-MLP 训练集准确率 /%	Self Attention-MLP 验证集准确率 /%	SVM 验证集准 确率/%	SVM 验证集准 确率/%
字	25.556	19.000	21.233	19.258
词	22.345	17.000	17.462	15.331

实验结论：token 数超过 1000，如取 token 数为 3000，以词为基本单元的分类效果好，低于 1000，以字为基本单元的分类效果好，token 为 1000，二者效果差不多。

总结

1. 主题数 T 越多，训练准确性越高，但是相同 epoch 下，更易出现过拟合现象，深度学习方法相较于传统方法分类有优势。
2. token 数越多，训练准确性越高，且能够一定程度上减弱过拟合的现象，深度学习方法相较于传统方法分类有优势。
3. token 数超过 1000，如取 token 数为 3000，以词为基本单元的分类效果好，低于 1000，以字为基本单元的分类效果好，token 为 1000，二者效果差不多。