

In [4]:

```
1 from sklearn.datasets import fetch_20newsgroups
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 from sklearn.naive_bayes import MultinomialNB, GaussianNB
6 from sklearn.metrics import f1_score
7 from sklearn.pipeline import Pipeline
8 from scipy.sparse import issparse
9 from sklearn.model_selection import GridSearchCV
```

只修改了train_val()函数于Pipeline

版本1, to_array无嵌入, GaussianNB单独

In [2]:

```
1 # def train_val(D, vect, nb, to_array=False):
2 #     X_train, X_test, y_train, y_test=D
3 #     pp=Pipeline([('vect', vect), ('nb', nb)])
4
5 #     if to_array:#高斯的这个差一个to_array
6 #         X_train_vect=vect.fit_transform(X_train)
7 #         X_test_vect=vect.transform(X_test)
8 #         X_train_vect, X_test_vect=X_train_vect.toarray(), X_test_vect.toarray()
9 #         nb.fit(X_train_vect, y_train)
10 #         pred=nb.predict(X_test_vect)
11 #     else:
12 #         pp.fit(X_train, y_train)
13 #         pred=pp.predict(X_test)
14 #     f1_weighted=f1_score(y_test, pred, average='macro')
15 #     # df_X_train=pd.DataFrame(X_train_vect.toarray(), columns=vectorizer.get_feature_names_o
16 #     return f1_weighted
```

版本2, to_array嵌入, 自定义转化器

In [3]:

```
1 # 创建一个自定义转换器, 用于将稀疏矩阵转换为密集矩阵
2 class SparseToDenseTransformer:
3     def fit(self, X, y=None):
4         return self
5     def transform(self, X):
6         if issparse(X): #判断是否为稀疏矩阵
7             return X.toarray()
8         return X
9 def train_val(D, vect, nb, to_array=False):
10     X_train, X_test, y_train, y_test=D
11     if to_array:
12         pp=Pipeline([('vect', vect), ('hidden', SparseToDenseTransformer()), ('nb', nb)])
13     else:
14         pp=Pipeline([('vect', vect), ('nb', nb)])
15     pp.fit(X_train, y_train)
16     pred=pp.predict(X_test)
17     f1_weighted=f1_score(y_test, pred, average='macro')
18     # df_X_train=pd.DataFrame(X_train_ved.toarray(), columns=vectorizer.get_feature_names_out
19     return f1_weighted
```

实验1: 应该使用GaussianNB还是MultinomialNB?

In [4]:

```
1 #由于内存问题, 对比时选用四类数据
2 news = fetch_20newsgroups(subset='all', categories=['alt.atheism', 'talk.religion.misc', 'comp
3 D =train_test_split(news.data, news.target, test_size=0.2, random_state=2023)
4 #控制变量, 使用相同的特征提取器
5 vect=TfidfVectorizer(stop_words='english')
6 #使用两种分类器
7 mnb=MultinomialNB(alpha=0.1)
8 gnb=GaussianNB()
9 #预测且评价
10 print('MultinomialNB的加权F1', train_val(D, vect, mnb))
11 print('GaussianNB的加权F1', train_val(D, vect, gnb, True))
```

MultinomialNB的加权F1 0.9440695360360671

GaussianNB的加权F1 0.9276497366606062

实验1结论: 性能上, MultinomialNB表现更优, 且接收稀疏矩阵传入, 大大减少了时空复杂度, 因此该选用 MultinomialNB。对于后续实验, 选定分类器为MultinomialNB, 此时由于其接收稀疏矩阵的特性, 可以选取所有的数据集, 而不用固定4类

实验2: 使用相同的训练集和测试集, 比较 CountVectorizer和TfidfVectorizer的效果

实验 3: 考察停用词的作用

In [5]:

```
1 news = fetch_20newsgroups(subset='all')
2 D = train_test_split(news.data, news.target, test_size=0.2, random_state=2023)
3 #控制变量，使用相同的分类器
4 mnb=MultinomialNB()
5 #使用两种特征提取器
6 cv_sw=CountVectorizer(stop_words='english')
7 tv_sw=TfidfVectorizer(stop_words='english')
8 cv=CountVectorizer()
9 tv=TfidfVectorizer()
10 #预测且评价
11 print('CountVectorizer的加权F1(无停用词)', train_val(D, cv, mnb))
12 print('TfidfVectorizer的加权F1(无停用词)', train_val(D, tv, mnb))
13 print('CountVectorizer的加权F1(有停用词)', train_val(D, cv_sw, mnb))
14 print('TfidfVectorizer的加权F1(有停用词)', train_val(D, tv_sw, mnb))
```

CountVectorizer的加权F1(无停用词) 0.8354009079661232

TfidfVectorizer的加权F1(无停用词) 0.8341839021793735

CountVectorizer的加权F1(有停用词) 0.8661911054593976

TfidfVectorizer的加权F1(有停用词) 0.8689426147222381

结论：【回答问题二】一方面，在同有停用词或同无停用词的情况下，CountVectorizer与TfidfVectorizer的性能差异不大。**【回答问题三】**另一方面，无论CountVectorizer或TfidfVectorizer都被停用词的引入显著提高了预测效果

实验 4：考察Tf-idf 平滑的作用

In [6]:

```
1 tv=TfidfVectorizer()
2 tv_no_smooth=TfidfVectorizer(smooth_idf=False)
3 print('TfidfVectorizer的加权F1(有平滑)', train_val(D, tv, mnb))
4 print('TfidfVectorizer的加权F1(无平滑)', train_val(D, tv_no_smooth, mnb))
```

TfidfVectorizer的加权F1(有平滑) 0.8341839021793735

TfidfVectorizer的加权F1(无平滑) 0.8339279271757325

结论：有平滑小小的提升

实验五 交叉验证

In [5]:

```
1 news = fetch_20newsgroups(subset='all')
2 X_train,X_test, y_train, y_test =train_test_split(news.data,news.target,test_size=0.2, random_state=42)
3 vect=TfidfVectorizer(stop_words='english')
4 nb=MultinomialNB()
5 parameters = {'nb__alpha': [1,2,3]}
6 pp=Pipeline([('vect',vect), ('nb',nb)])
7 gs = GridSearchCV(pp,
8                   parameters,
9                   scoring = ['accuracy','f1_macro'],
10                  verbose=2,
11                  refit='accuracy',
12                  cv=5,
13                  n_jobs=-1)
14 # 执行多线程并行网格搜索。
15 time_ = gs.fit(X_train, y_train)
16 gs.best_params_, gs.best_score_
17
18 # 输出最佳模型在测试集上的准确性。
19 print(gs.score(X_test, y_test))
20 print(gs.best_params_)
21
```

Fitting 5 folds for each of 3 candidates, totalling 15 fits
0.8843501326259947
{'nb__alpha': 1}

结论:最优参数为1