

1 准备数据

In [1]:

```
1 # 选取 科技、汽车、娱乐、军事、运动
2 import jieba
3 import pandas as pd
4
5 df_technology = pd.read_csv("./origin_data/technology_news.csv", encoding='utf-8')
6
7 df_technology = df_technology.dropna()
8
9 df_car = pd.read_csv("./origin_data/car_news.csv", encoding='utf-8')
10 df_car = df_car.dropna()
11
12 df_entertainment = pd.read_csv("./origin_data/entertainment_news.csv", encoding='utf-8')
13 df_entertainment = df_entertainment.dropna()
14
15 df_military = pd.read_csv("./origin_data/military_news.csv", encoding='utf-8')
16 df_military = df_military.dropna()
17
18 df_sports = pd.read_csv("./origin_data/sports_news.csv", encoding='utf-8')
19 df_sports = df_sports.dropna()
20
21 df_society = pd.read_csv("./origin_data/society_news.csv", encoding='utf-8')
22 df_society = df_sports.dropna()
```

In [2]:

```
1 # 每类数据取10000条
2 technology = df_technology.content.values.tolist()[0:10000]
3 car = df_car.content.values.tolist()[0:10000]
4 entertainment = df_entertainment.content.values.tolist()[0:10000]
5 military = df_military.content.values.tolist()[0:10000]
6 sports = df_sports.content.values.tolist()[0:10000]
7 society = df_society.content.values.tolist()[0:10000]
```

In [3]:

```
1 ' '.join(jieba.lcut(technology[0]))
```

Building prefix dict from the default dictionary ...

Loading model from cache C:\Users\LIUJIA~1\AppData\Local\Temp\jieba.cache

Loading model cost 0.719 seconds.

Prefix dict has been built successfully.

Out[3]:

'\u3000 \u3000 , 中新网 , 1 月 7 日电 \xa0 恰逢 CES 2017 拉开 大幕 , 却 惊闻 “ AlphaGo 升级版 ” 的 Master 迎来 60 连胜 , 人类 顶尖 围棋 手 在 一周 内 纷纷 败给 这个 谷歌 旗下 DeepMind 团队 打造 的 “ 围棋 大脑 ” , 显然 也 为 聚焦 于 人工智能 的 本届 CES 增添 了 声势 。 而 首次 参展 , 并 致力于 打造 “ 原创 AI 大脑 ” 的 中国 深度 学习 领军 企业 的 商汤 科技 , 在 人工智能 的 浪潮 之 巅 , 及 众多 业界 前辈 和 巨匠 面前 , 将会 交出 一份 怎样 的 答卷 呢 ? ’

2 去停用词

In [4]:

```
1 # 加载停用词转为numpy
2 stopwords=pd.read_csv("origin_data/stopwords.txt", index_col=False, quoting=3, sep="\t", names=['s
3 stopwords=stopwords['stopword'].values #numpy
```

In [5]:

```
1 # 定义一个类, 去停用词
2 sentences=[] #存放所有的数据
3 # sentences_train=[]
4 # sentences_test=[]
5 def preprocess_texts(target_path, catergary, texts):
6     with open(target_path + catergary + '.txt', 'w' ) as f:
7         for text in texts:
8             words=jieba.lcut(text)
9             words=list(filter(lambda x:len(x)>1, words))
10            words=list(filter(lambda x: x not in stopwords, words))
11            if len(words):
12                sentences.append(' '.join(words), catergary))
13                f.write(' '.join(words)+'\n')
14            print(len(sentences))
15 # sentences_train+=sentences[:7000]
16 # sentences_test+=sentences[:3000]
```

In [6]:

```
1 preprocess_texts ('./processed_data/', 'technology', technology )
2 preprocess_texts ('./processed_data/', 'car', car )
3 preprocess_texts ('./processed_data/', 'entertainment', entertainment )
4 preprocess_texts ('./processed_data/', 'military', military )
5 preprocess_texts ('./processed_data/', 'sports', sports )
6 preprocess_texts ('./processed_data/', 'society', society )
```

9954
19898
29841
39768
49692
59616

In [300]:

```
1 # sentences_train=[]
2 # sentences_test=[]
3 # for i in range(5):
4 #     sentences_train+=sentences[i*10000:i*10000+7000]
5 #     sentences_test+=sentences[7000+i*10000:(i+1)*10000]
```

In [7]:

```
1 len(sentences)
```

Out[7]:

59616

In [8]:

```
1 for i in sentences:
2     if len(i[0])==0: print( "false")
```

In [9]:

```
1 import random
2 random.shuffle(sentences)
```

3 划分数据集

In [10]:

```
1 from sklearn.model_selection import train_test_split
2 x_train_, y_train_ = zip(*sentences)
3 # x_test, y_test = zip(*sentences_test)
4 # print(len(x_train))
5 x_train, x_val, y_train, y_val= train_test_split(x_train_, y_train_, random_state=2020)
```

4 训练和预测数据

In [11]:

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 from gensim.models import Word2Vec
4 vec = CountVectorizer(
5     analyzer='word',      # 以词为单位
6     max_features=4000,    # 选取4000个主要词构建预料
7 )
8 tfv = TfidfVectorizer(lowercase=True, decode_error='ignore')
9 # w2v = Word2Vec(x_train, size=300, window=5, iter=15, workers=12, seed=2020)
10
11 tfv.fit(x_train)
12 vec.fit(x_train)
```

Out[11]:

```
CountVectorizer(analyzer='word', binary=False, decode_error='strict',
               dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
               lowercase=True, max_df=1.0, max_features=4000, min_df=1,
               ngram_range=(1, 1), preprocessor=None, stop_words=None,
               strip_accents=None, token_pattern='(?u)\\b\\w\\w+\\b',
               tokenizer=None, vocabulary=None)
```

In [12]:

```
1 # one-hot 预测分数
2 from sklearn.naive_bayes import MultinomialNB
3 classifier = MultinomialNB()
4 classifier.fit(vec.transform(x_train), y_train)
5 classifier.score(vec.transform(x_val), y_val)
```

Out[12]:

0.6565351583467526

In [26]:

```
1 # tf-idf预测分数
2 from sklearn.naive_bayes import MultinomialNB
3 classifier = MultinomialNB()
4 classifier.fit(tfv.transform(x_train), y_train)
5 # classifier.score(tfv.transform(x_train), y_train)
6 classifier.score(tfv.transform(x_val), y_val)
```

Out[26]:

0.8236491322240115

In [14]:

```
1 # 构建word2vec模型
2 from gensim.models import Word2Vec
3 w2v = Word2Vec(x_train_, size=60, iter=10, seed=2020, min_count=5 )
```

In [15]:

```
1 import numpy as np
2 x_train_w2v=np.zeros((len(x_train),60))
3 x_train_w2v_ave=np.zeros((len(x_train),60))
4 i=0
5 for words in list(x_train):
6     n=0
7     for word in words:
8         if word in w2v:
9             n+=1
10            vect=w2v[word]
11            x_train_w2v[i,:]+=vect
12            x_train_w2v_ave[i,:]=x_train_w2v[i,:]/n
13            i+=1
```

D:\ProgramData\anaconda\envs\tf2.0\lib\site-packages\ipykernel_launcher.py:8: DeprecationWarning: Call to deprecated `__contains__` (Method will be removed in 4.0.0, use self.wv.__contains__() instead).

D:\ProgramData\anaconda\envs\tf2.0\lib\site-packages\ipykernel_launcher.py:10: DeprecationWarning: Call to deprecated `__getitem__` (Method will be removed in 4.0.0, use self.wv.__getitem__() instead).

Remove the CWD from sys.path while we load stuff.

In [16]:

```
1 x_val_w2v=np.zeros((len(x_val),60))
2 x_val_w2v_ave=np.zeros((len(x_val),60))
3 i=0
4 for words in list(x_val):
5     if len(words)==0: print( i,"false")
6     n=0
7     for word in words:
8         if word in w2v:
9             n+=1
10            vect=w2v[word]
11            x_train_w2v[i,:]+=vect
12 #         if n==0: print("false")
13 x_val_w2v_ave[i,:]=x_train_w2v[i,:]/n
14 i+=1
```

D:\ProgramData\anaconda\envs\tf2.0\lib\site-packages\ipykernel_launcher.py:8: DeprecationWarning: Call to deprecated `__contains__` (Method will be removed in 4.0.0, use self.wv.__contains__() instead).

D:\ProgramData\anaconda\envs\tf2.0\lib\site-packages\ipykernel_launcher.py:10: DeprecationWarning: Call to deprecated `__getitem__` (Method will be removed in 4.0.0, use self.wv.__getitem__() instead).

Remove the CWD from sys.path while we load stuff.

In [17]:

```
1 # x_train_df=pd.DataFrame(x_train_w2v_ave)
2 # y_train_df=pd.DataFrame(y_train)
3 # x_val_df=pd.DataFrame(x_val_w2v_ave)
4 # y_val_df=pd.DataFrame(y_val)
```

In [18]:

```
1 # from sklearn.preprocessing import LabelEncoder
2 # labelEncoder = LabelEncoder()
3 # y_train_label = labelEncoder.fit_transform(y_train_df)
4 # y_train[0:10]
5 # y_train_label=pd.DataFrame(y_train_label)
6 # y_train_label
```

In [19]:

```
1 from sklearn.linear_model import LogisticRegression
2 classifier = LogisticRegression()
3 classifier.fit(x_train_w2v_ave, y_train)
4 classifier.score(x_val_w2v_ave, y_val)
```

Out[19]:

0.42706655931293613

In [20]:

```
1 from sklearn.naive_bayes import GaussianNB
2 classifier = GaussianNB()
3 classifier.fit(x_train_w2v_ave, y_train)
4 classifier.score(x_val_w2v_ave, y_val)
```

Out[20]:

0.3882850241545894

In [21]:

```
1 #
2 from sklearn.model_selection import StratifiedKFold
3 from sklearn.metrics import accuracy_score, precision_score
4 import numpy as np
5
6 def stratifiedkfold_cv(x, y, clf_class, shuffle=True, n_folds=5, **kwargs):
7     stratifiedk_fold = StratifiedKFold(n_splits=n_folds, shuffle=shuffle)
8     y_pred = y[:]
9     for train_index, test_index in stratifiedk_fold.split(x, y):
10         X_train, X_test = x[train_index], x[test_index]
11         y_train = y[train_index]
12         clf = clf_class(**kwargs)
13         clf.fit(X_train, y_train)
14         y_pred[test_index] = clf.predict(X_test)
15     return y_pred
```

In [35]:

```
1 x_train_, y_train_ = zip(*sentences)
2 # x_train_=x_train_[:-1]
3 # y_train_=y_train_[:-1]
4 # len(x_train_)
```

In [36]:

```
1 # one-hot交叉验证
2 NB = MultinomialNB
3 print(precision_score(y_train_, stratifiedkfold_cv(vec.transform(x_train_), np.array(y_train_), N
```

0.6830925308820235

In [38]:

```
1 # tfidf交叉验证
2 NB = MultinomialNB
3 print(precision_score(y_train_, stratifiedkfold_cv(tfv.transform(x_train_), np.array(y_train_), N
```

0.7269982223098465

In [47]:

```
1 # w2v交叉验证
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.model_selection import cross_val_score
4 LR = LogisticRegression()
5
6 score_ndarray = cross_val_score(LR, x_train_w2v_ave, np.array(y_train), cv=5)
7 print(score_ndarray.mean())
8 # print(precision_score(y_train, stratifiedkfold_cv( x_train_w2v_ave, np.array(y_train), LR),
```

0.6345501404018736