

贝叶斯分析-多项式朴素贝叶斯文本分类

姓名：熊荣康

学号：SA21229005

摘要

文本分类是NLP应用领域中最常见也最重要的任务类型¹。分类就是人类出于某些需要、按照某些标准，将事物分为若干组的活动。分类可能是我们人类最喜欢做的事情之一。从操作的角度看，分类，就是把事物按照一定的规则分门别类。而文本分类，就是把文本按照一定的规则分门别类。文本分类在搜索引擎、问答系统、会话系统等等重要的信息处理系统中应用非常广泛，可以说无处不在。本文利用朴素贝叶斯分类方法分析了20个网络新闻组语料库，约20000篇新闻的单词出现次数作为特征，用多项式朴素贝叶斯对这些新闻进行分类²。

关键词: 多项式朴素贝叶斯, 文本分类

1. 数据和问题描述

对于文本分类，通常有14种分类算法。8种传统算法：k临近、决策树、多层感知器、朴素贝叶斯（包括伯努利贝叶斯、高斯贝叶斯和多项式贝叶斯）、逻辑回归和支持向量机；4种集成学习算法：随机森林、AdaBoost、lightGBM和xgBoost；2种深度学习算法：前馈神经网络和LSTM³。

多项式朴素贝叶斯通常用于文本分类，其特征都是指待分类文本的单词出现次数或者频次。这里用20个网络新闻组语料库（20 Newsgroups corpus，约20 000篇新闻）的单词出现次数作为特征⁴，使用多项式朴素贝叶斯对这些新闻组进行分类。

2. 模型与数据分析

2.1 贝叶斯分类

朴素贝叶斯模型是一组非常简单快速的分类算法，通常适用于维度非常高的数据集⁵。因为运行速度快，而且可调参数少，因此非常适合为分类问题提供快速粗糙的基本方案。朴素贝叶斯分类器建立在贝叶斯分类方法的基础上，其数学基础是贝叶斯定理（Bayes's theorem）——一个描述统计量条件概率关系的公式⁶。在贝叶斯分类中，我们希望确定一个具有某些特征的样本属于某类标签的概率，通常记为 $P(L|features)$ 。贝叶斯定理告诉我们，可以直接用下面的公式计算这个概率

$$\frac{P(L_1|features)}{P(L_1|features)} = \frac{P(features|L_1)P(L_1)}{P(features|L_2)P(L_2)} \quad (1)$$

现在需要一种模型，帮我们计算每个标签的 $P(features|Li)$ 。这种模型被称为生成模型，因为它可以训练出生成输入数据的假设随机过程（或称为概率分布）。为每种标签设置生成模型是贝叶斯分类器训练过程的主要部分。虽然这个训练步骤通常很难做，但是我们可以通过对模型进行随机分布的假设，来简化训练工作，之所以称为“朴素”或“朴素贝叶斯”，是因为如果对每种标签的生成模型进行非常简单的假设，就能找到每种类型生成模型的近似解，然后就可以使用贝叶斯分类。不同类型的朴素贝叶斯分类器是由对数据的不同假设决定的。

多项式朴素贝叶斯通常用于文本分类，其特征都是指待分类文本的单词出现次数或者频 次。这里用 20 个网络新闻组语料库（20 Newsgroups corpus，约 20 000 篇新闻）的单词出现次数作为特征，用多项式朴素贝叶斯对这些新闻组进行分类。

2.2 高斯朴素贝叶斯

分类器假设每个标签的数据都服从简单的高斯分布，假设有如图(1)所示的数据，代码见附录A.1。

图1:

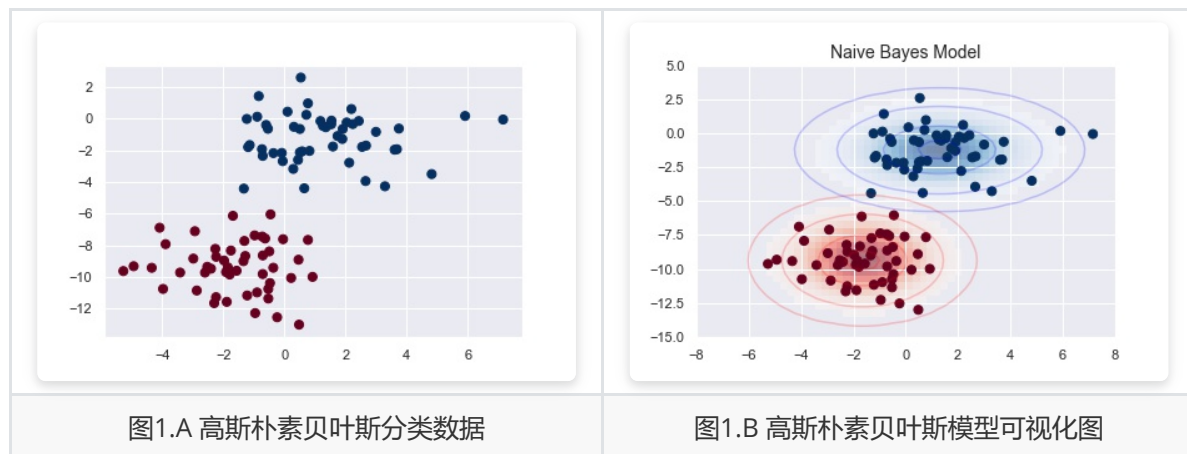
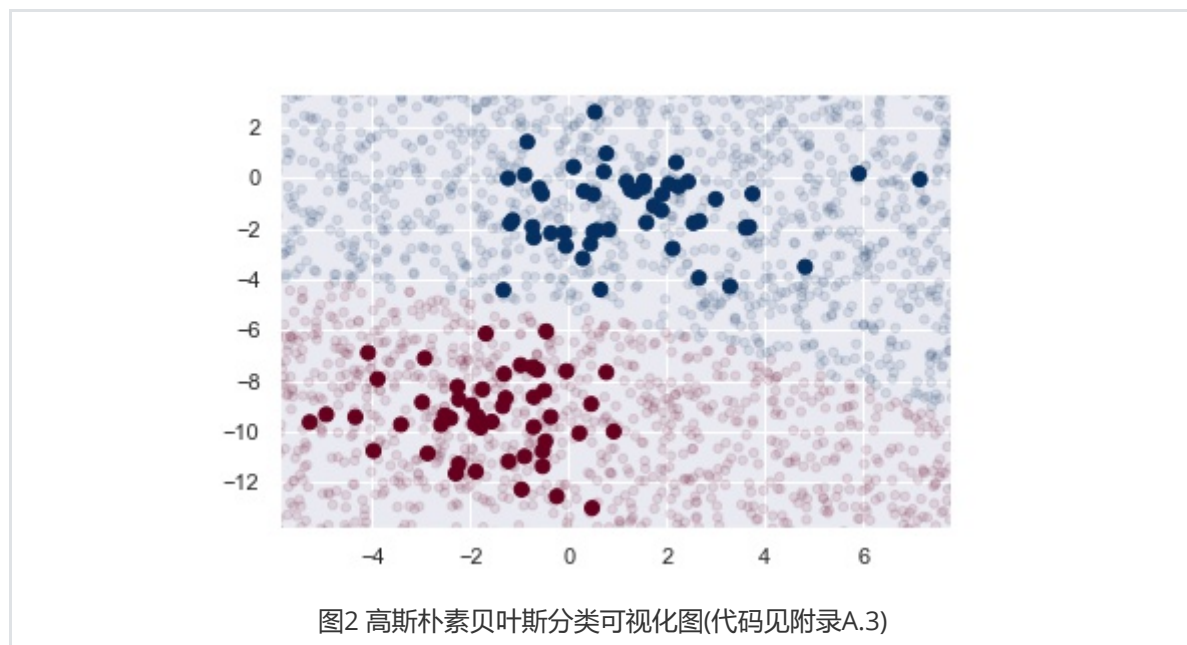


图1.B中的椭圆曲线表示每个标签的高斯生成模型，越靠近椭圆中心的可能性越大。

一种快速创建简易模型的方法是假设数据服从高斯分布，且变量无协方差（线性无关）。只要找出每个标签的所有样本点均值和标准差，再定义一个高斯分布，就可以拟合模型。通过每种类型的生成模型，可以计算出任意数据点的似然估计 $P(features|L_1)$ ，然后根据贝叶斯定理计算出后验概率比值，从而确定每个数据点可能性最大的标签。



可以得到两个标签的后验概率，如果需要评估分类器的不确定性，那么这类贝叶斯方法非常有用。

Table 1: 两类标签的后验概率

$P(L_1 features)$	$P(L_2 features)$
0.	1.
0.	1.
0.	1.
1.	0.
0.	1.
0.	1.
0.08	0.92
0.	1.
0.	1.
1.	0.
1.	0.
1.	0.
0.89	0.11

由于分类的最终效果只能依赖于一开始的模型假设，因此高斯朴素贝叶斯经常得不到非常好的结果。但是，在许多场景中，尤其是特征较多时，这种假设并不妨碍高斯朴素贝叶斯成为一种有用的方法。

2.3 多项式朴素贝叶斯

多项式朴素贝叶斯 (multinomial naive Bayes) 它假设特征是由一个简单多项式分布生成的。多项分布可以描述各种类型样本出现次数的概率，因此多项式朴素贝叶斯非常适合用于描述出现次数或者出现次数比例的特征。这里模型数据的分布不再是高斯分布，而是用多项式分布代替。

- 训练阶段

$$P(C = c) = \frac{\text{属于}c\text{类的文档数}}{\text{训练集文档总数}}$$

$$P(w_i|c) = \frac{\text{词}w_i\text{在属于}c\text{类的所有文档中出现次数}}{\text{属于}c\text{类的所有文档中的词总数}}$$

拉普拉斯平滑（加1平滑）：

$$P(w_i|c) = \frac{\text{词}w_i\text{在属于}c\text{类的所有文档中出现次数} + 1}{\text{属于}c\text{类的所有文档中的词总数}}$$

加1平滑可以认为是采用均匀分布作为先验分布，即每个词项在每个类别中出现一次，然后根据训练数据对得到的结果进行更新。也就是说 $\frac{\text{未登录词的估计值为1}}{\text{词汇表长度}}$

- 预测阶段

$$\begin{aligned} & \arg \max_{c \in C} P(c|w_1, \dots, w_n) \\ &= \arg \max_{c \in C} [P(c)P(w_1, \dots, w_n|c)] \\ &= \arg \max_{c \in C} [P(c)P(w_1|c)P(w_2|c) \cdots P(w_n|c)] = \arg \max_{c \in C} [\log P(c) + \cdots + \log P(w_n|c)] \end{aligned}$$

字符串	预测结果
'Research reported here this week at the annual meeting of the Society of Integrative and Comparative Biology has revealed black-legged ticks infected with the Lyme disease-causing microbe thrive in below-freezing weather and can be active even in winter. The finding suggests the variable winter conditions brought on by climate change could increase ticks' activity, boosting the odds that people will encounter the ticks and come down with Lyme disease.'	'sci.med'
'India, from the earliest days of the pandemic, has reported far fewer COVID-19 deaths than expected given the toll elsewhere—an apparent death “paradox” that some believed was real and others thought would prove illusory. Now, a prominent epidemiologist who contended the country really had been spared the worst of COVID-19 has led a rigorous new analysis of available mortality data and concluded he “got it wrong.” India has “substantially greater” COVID-19 deaths than official reports suggest, says Prabhat Jha of the University of Toronto— close to 3 million, which is more than six times higher than the government has acknowledged and the largest number of any country.'	"talk.politics.guns"
'The Peruvian anchovy is a small fish with a big impact. Only about the size of an index finger, they make up the single largest fish catch in the world—sometimes up to 15% of the global haul. Nearly all the highly nutritious fish are ground up to feed salmon and other farm-raised species that are worth billions of dollars. Now, scientists studying ancient sediments and fossils have shown warming waters once nearly eliminated this valuable resource, raising fears that today's climate change could repeat the disaster.'	'sci.crypt'
“The finding is really concerning,” says Becca Selden, a marine ecologist at Wellesley College who was not involved with the research. The new record of climate change’s impact on fish shows “a complete shift in what that ecosystem looked like,” she says.'sending a payload to the ISS'	'soc.religion.christian'
'discussing islam vs atheism'	'soc.religion.christian`
'determining the screen resolution'	'comp.graphics`

对于中文文本需要进行分词处理，常见分词工具见Table 3。中文任务分词，一般使用 `jieba` 分词。

Table 3: 常用分词工具

分词工具	地址	简单说明
中科院计算所NLPIR	http://ictclas.nlpir.org/	包括中文分词；英文分词；词性标注；命名实体识别；新词识别；关键词提取；支持用户专业词典与微博分析。NLPIR系统支持多种编码、多种操作系统、多种开发语言与平台
ansj分词器	https://github.com/NLPchina/ansj_seg	开源的Java中文分词工具，基于中科院的ictclas中文分词算法。目前实现了中文分词、中文姓名识别、用户自定义词典、关键字提取、自动摘要、关键字标记等功能
哈工大的LTP	https://github.com/HIT-SCIR/ltp	主页上给过调用接口，每秒请求的次数有限制
清华大学THULAC	https://github.com/thunlp/THULAC	目前已经有Java、Python和C++版本，并且代码开源
斯坦福分词器	https://nlp.stanford.edu/software/segmenter.shtml	作为众多斯坦福自然语言处理中的一个包，目前最新版本3.7.0,Java实现的CRF算法。可以直接使用训练好的模型，也提供训练模型接口
Hanlp分词器	https://github.com/hankcs/HanLP	求解的是最短路径。优点：开源、有人维护、可以解答。原始模型用的训练语料是人民日报的语料，当然如果你有足够的语料也可以自己训练
结巴分词	https://github.com/xywyw/cppjieba	基于前缀词典实现高效的词图扫描，生成句子中汉字所有可能成词情况所构成的有向无环图(DAG)；采用了动态规划查找最大概率路径,找出基于词频的最大切分组合；对于未登录词，采用了基于汉字成词能力的HMM模型，使用了Viterbi算法
KCWS分词器(字嵌入+Bi-LSTM+CRF)	https://github.com/kongweizh/kcws	本质上是序列标注，这个分词器用人民日报的80万语料，据说按照字符正确率评估标准能达到97.5%的准确率
ZPar	https://github.com/francischiang/zpar/releases	新加坡科技设计大学开发的中文分词器，包括分词、词性标注和Parser，支持多语言，据说效果是公开的分词器中最好的，C++语言编写
IKAnalyzer	https://github.com/wkik-analyzer	IK才用了特有的“正向迭代最细粒度切分算法”，支持细粒度和智能分词两种切分模式

3. 结论

由于朴素贝叶斯分类器对数据有严格的假设，因此它的训练效果通常比复杂模型的差。其优点主要体现在以下四个方面

- 训练和预测的速度非常快
- 直接使用概率预测
- 通常很容易解释
- 可调参数（如果有的话）非常少

这些优点使得朴素贝叶斯分类器通常很适合作为分类的初始解。如果分类效果满足要求，那么万事大吉，就获得了一个非常快速且容易解释的分类器。但如果分类效果不够好，那么你可以尝试更复杂的分类模型，与朴素贝叶斯分类器的分类效果进行对比，看看复杂模型的分类效果究竟如何。

朴素贝叶斯分类器非常适合用于以下应用场景。

- 假设分布函数与数据匹配（实际中很少见）
- 各种类型的区分度很高，模型复杂度不重要
- 非常高维度的数据，模型复杂度不重要

在新维度会增加样本数据信息量的假设条件下，高维数据的簇中心点比低维数据的簇中心点更分散。因此，随着数据维度不断增加，像朴素贝叶斯这样的简单分类器的分类效果会和复杂分类器一样，甚至更好——只要你有足够的数据，简单的模型也可以非常强大⁷。

本文的代码已经上传GitHub，请查看<https://github.com/RongkangXiong/Bayes-Gaussian-Classify/settings>⁸

附录

A.1

```
1 %matplotlib inline
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns; sns.set()
5
6 from sklearn.datasets import make_blobs
7 x, y = make_blobs(100, 2, centers=2, random_state=2, cluster_std=1.5)
8 plt.scatter(x[:, 0], x[:, 1], c=y, s=50, cmap='RdBu');
9 plt.savefig('./images/fig1.jpg')
```

A.2

```
1 from sklearn.datasets import make_blobs
2 x, y = make_blobs(100, 2, centers=2, random_state=2, cluster_std=1.5)
3
4 fig, ax = plt.subplots()
5
6 ax.scatter(x[:, 0], x[:, 1], c=y, s=50, cmap='RdBu')
7 ax.set_title('Naive Bayes Model', size=14)
8
9 xlim = (-8, 8)
10 ylim = (-15, 5)
11
12 xg = np.linspace(xlim[0], xlim[1], 60)
13 yg = np.linspace(ylim[0], ylim[1], 40)
14 xx, yy = np.meshgrid(xg, yg)
```



```

15 xgrid = np.vstack([xx.ravel(), yy.ravel()]).T
16
17 for label, color in enumerate(['red', 'blue']):
18     mask = (y == label)
19     mu, std = X[mask].mean(0), X[mask].std(0)
20     P = np.exp(-0.5 * (Xgrid - mu) ** 2 / std ** 2).prod(1)
21     Pm = np.ma.masked_array(P, P < 0.03)
22     ax.pcolorfast(xg, yg, Pm.reshape(xx.shape), alpha=0.5,
23                  cmap=color.title() + 's')
24     ax.contour(xx, yy, P.reshape(xx.shape),
25               levels=[0.01, 0.1, 0.5, 0.9],
26               colors=color, alpha=0.2)
27
28 ax.set(xlim=xlim, ylim=ylim)
29
30 fig.savefig('./images/gaussian-NB.jpg')

```

A.3

```

1 # 使用Scikit-Learn 的 sklearn.naive_bayes.GaussianNB 评估器
2 from sklearn.naive_bayes import GaussianNB
3 model = GaussianNB()
4 model.fit(X, y)
5
6 # 生成一些新数据来预测标签
7 rng = np.random.RandomState(0)
8 Xnew = [-6, -14] + [14, 18] * rng.rand(2000, 2)
9 ynew = model.predict(Xnew)
10
11 # 将这些新数据画出来, 看看决策边界的位置
12 plt.scatter(X[:, 0], X[:, 1], c=y, s=50, cmap='RdBu')
13 lim = plt.axis()
14 plt.scatter(Xnew[:, 0], Xnew[:, 1], c=ynew, s=20, cmap='RdBu', alpha=0.1)
15 plt.axis(lim);
16
17 # 用predict_proba方法计算样本属于某个标签的概率
18 yprob = model.predict_proba(Xnew)
19 yprob[-20:].round(2)

```

B.1

```

1 from sklearn.datasets import fetch_20newsgroups
2 data = fetch_20newsgroups()
3 data.target_names

```

输出结果:

```

1 ['alt.atheism',
2  'comp.graphics',
3  'comp.os.ms-windows.misc',
4  'comp.sys.ibm.pc.hardware',
5  'comp.sys.mac.hardware',
6  'comp.windows.x',

```



```

7  'misc.forsale',
8  'rec.autos',
9  'rec.motorcycles',
10 'rec.sport.baseball',
11 'rec.sport.hockey',
12 'sci.crypt',
13 'sci.electronics',
14 'sci.med',
15 'sci.space',
16 'soc.religion.christian',
17 'talk.politics.guns',
18 'talk.politics.mideast',
19 'talk.politics.misc',
20 'talk.religion.misc']

```

B.2

```

1  # categories = ['talk.religion.misc', 'soc.religion.christian', 'sci.space',
2  # 'comp.graphics']
3  categories = data.target_names
4  train = fetch_20newsgroups(subset='train', categories=categories)
5  test = fetch_20newsgroups(subset='test', categories=categories)

```

B.3

```

1  From: guykuo@carson.u.washington.edu (Guy Kuo)
2  Subject: SI Clock Poll - Final Call
3  Summary: Final call for SI clock reports
4  Keywords: SI, acceleration, clock, upgrade
5  Article-ID.: shelley.1qvfo9INnc3s
6  Organization: University of Washington
7  Lines: 11
8  NNTP-Posting-Host: carson.u.washington.edu
9
10 A fair number of brave souls who upgraded their SI clock oscillator have
11 shared their experiences for this poll. Please send a brief message
12 detailing
13 your experiences with the procedure. Top speed attained, CPU rated speed,
14 add on cards and adapters, heat sinks, hour of usage per day, floppy disk
15 functionality with 800 and 1.4 m floppies are especially requested.
16
17 I will be summarizing in the next two days, so please add to the network
18 knowledge base if you have done the clock upgrade and haven't answered this
19 poll. Thanks.
20 Guy Kuo <guykuo@u.washington.edu>

```

C.1

```

1  from sklearn.feature_extraction.text import TfidfVectorizer
2  from sklearn.naive_bayes import MultinomialNB
3  from sklearn.pipeline import make_pipeline
4  model = make_pipeline(TfidfVectorizer(), MultinomialNB())
5

```

```

6 # 通过这个管道，就可以将模型应用到训练数据上，预测出每个测试数据的标签
7 model.fit(train.data, train.target)
8 labels = model.predict(test.data)
9
10 # 这样就得到每个测试数据的预测标签，可以进一步评估评估器的性能了
11 from sklearn.metrics import confusion_matrix
12 mat = confusion_matrix(test.target, labels)
13 sns.heatmap(mat.T, square=True, annot=True, fmt='d', cbar=False,
14 xticklabels=train.target_names, yticklabels=train.target_names)
15 plt.xlabel('true label')
16 plt.ylabel('predicted label');

```

C.2

```

1 def predict_category(s, train=train, model=model):
2     pred = model.predict([s])
3     return train.target_names[pred[0]]

```

D.文本特征

分类特征

一种常见的非数值数据类型是分类数据，例如，浏览房屋数据的时候，除了看到“房价”（price）和“面积”（rooms）之类的数值特征，还会有“地点”（neighborhood）信息，数据可能像这样：

```

1 In[1]: data = [
2     {'price': 850000, 'rooms': 4, 'neighborhood': 'Queen Anne'},
3     {'price': 700000, 'rooms': 3, 'neighborhood': 'Fremont'},
4     {'price': 650000, 'rooms': 3, 'neighborhood': 'Wallingford'},
5     {'price': 600000, 'rooms': 2, 'neighborhood': 'Fremont'}
6 ]

```

可以把分类特征用映射关系编码成整数：

```

1 In[2]: {'Queen Anne': 1, 'Fremont': 2, 'Wallingford': 3}

```

但是，在 Scikit-Learn 中这么做并不是一个好办法：这个程序包的所有模块都有一个基本假设，那就是数值特征可以反映代数量（algebraic quantities）。因此，这样映射编码可能会让人觉得存在 `Queen Anne < Fremont < Wallingford`，甚至还有 `Wallingford - Queen Anne = Fremont`，这显然是没有意义的。

面对这种情况，常用的解决方法是独热编码。它可以有效增加额外的列，让 0 和 1 出现在对应的列分别表示每个分类值有或无。当你的数据是像上面那样的字典列表时，用 ScikitLearn 的 DictVectorizer 类就可以实现：

```

1 In[3]: from sklearn.feature_extraction import DictVectorizer
2 vec = DictVectorizer(sparse=False, dtype=int)
3 vec.fit_transform(data)
4 out[3]: array([[ 0,  1,  0, 850000,  4],
5 [  1,  0,  0, 700000,  3],
6 [  0,  0,  1, 650000,  3],
7 [  1,  0,  0, 600000,  2]], dtype=int64)

```

`neighborhood` 字段转换成三列来表示三个地点标签，每一行中用 1 所在的列对应一个地点。当这些分类特征编码之后，你就可以和之前一样拟合 Scikit-Learn 模型了。如果要看每一列的含义，可以用下面的代码查看特征名称：

```
1 In[4]: vec.get_feature_names()
2 Out[4]: ['neighborhood=Fremont',
3 'neighborhood=Queen Anne',
4 'neighborhood=wallingford',
5 'price',
6 'rooms']
```

这种方法也有一个显著的缺陷：如果你的分类特征有许多枚举值，那么数据集的维度就会急剧增加。然而，由于被编码的数据中有许多 0，因此用稀疏矩阵表示会非常高效：

```
1 In[5]: vec = DictVectorizer(sparse=True, dtype=int)
2 vec.fit_transform(data)
3 Out[5]: <4x5 sparse matrix of type '<class 'numpy.int64'>'
4 with 12 stored elements in Compressed Sparse Row format>
```

在拟合和评估模型时，Scikit-Learn 的许多（并非所有）评估器都支持稀疏矩阵输入。

`sklearn.preprocessing.OneHotEncoder` 和 `sklearn.feature_extraction.FeatureHasher` 是 Scikit-Learn 另外两个为分类特征编码的工具。

文本特征

一种常见的特征工程需求是将文本转换成一组数值，绝大多数社交媒体数据的自动化采集，都是依靠将文本编码成数字的技术手段。数据采集最简单的编码方法之一就是单词统计：给你几个文本，让你统计每个词出现的次数，然后放到表格中。例如下面三个短语：

```
1 In[6]: sample = ['problem of evil',
2 'evil queen',
3 'horizon problem']
```

面对单词统计的数据向量化问题时，可以创建一个列来表示单词“problem”、单词“evil”和单词“horizon”等。虽然手动做也可以，但是用 Scikit-Learn 的 `CountVectorizer` 更是可以轻松实现：

```
1 In[7]: from sklearn.feature_extraction.text import CountVectorizer
2 vec = CountVectorizer()
3 x = vec.fit_transform(sample)
4 x
5 Out[7]: <3x5 sparse matrix of type '<class 'numpy.int64'>'
6 with 7 stored elements in Compressed Sparse Row format>
```

结果是一个稀疏矩阵，里面记录了每个短语中每个单词的出现次数。如果用带列标签的 `DataFrame` 来表示这个稀疏矩阵就更方便了：

```
1 In[8]: import pandas as pd
2 pd.DataFrame(x.toarray(), columns=vec.get_feature_names())
3 Out[8]: evil horizon of problem queen
4 0 1 0 1 1 0
5 1 1 0 0 0 1
6 2 0 1 0 1 0
```

不过这种统计方法也有一些问题：原始的单词统计会让一些常用词聚集太高的权重，在分类算法中这样并不合理。解决这个问题方法就是通过 **TF-IDF** (term frequency-inversedocument frequency, 词频逆文档频率, IDF的大小与一个词的常见程度成反比)，通过单词在文档中出现的频率来衡量其权重。计算这些特征的语法和之前的示例类似：


某一特定词语的IDF，可以由总文件数目除以包含词语的文件数目，再将得到的商取对数得到

$$IDF = \log \frac{1 + n}{1 + docs(w, D)} + 1$$

n 是文档总数, w 是词, $docs(w, D)$ 是词 w 出现的文件数

```
1 In[9]: from sklearn.feature_extraction.text import TfidfVectorizer
2 vec = TfidfVectorizer()
3 x = vec.fit_transform(sample)
4 pd.DataFrame(x.toarray(), columns=vec.get_feature_names())
5 Out[9]: evil horizon of problem queen
6 0 0.517856 0.000000 0.680919 0.517856 0.000000
7 1 0.605349 0.000000 0.000000 0.000000 0.795961
8 2 0.000000 0.795961 0.000000 0.605349 0.000000
```

参考文献

1. Text Classification in Natural Language Processing, 2020. . Analytics Vidhya. URL <https://www.analyticsvidhya.com/blog/2020/12/understanding-text-classification-in-nlp-with-movie-review-example-example/> (accessed 1.10.22). 
2. Jake VanderPlas, 2018. Python数据科学手册. 中国工信出版集团. 
3. <https://www.heywhale.com/mw/project/5cbbe1668c90d7002c810f79> 
4. sklearn.datasets.fetch_20newsgroups [WWW Document], n.d. . scikit-learn. URL https://scikit-learn/stable/modules/generated/sklearn.datasets.fetch_20newsgroups.html (accessed 1.10.22). 
5. Naive Bayes Classifier - an overview | ScienceDirect Topics [WWW Document], n.d. URL <https://www.sciencedirect.com/topics/engineering/naive-bayes-classifier> (accessed 1.10.22). 
6. Hoff, P.D., 2009. A First Course in Bayesian Statistical Methods, Springer Texts in Statistics. Springer New York, New York, NY. <https://doi.org/10.1007/978-0-387-92407-6> 
7. Chauhan G., 2018. All about Naive Bayes [WWW Document]. Medium. URL <https://towardsdatascience.com/all-about-naive-bayes-8e13cef044cf> (accessed 1.10.22). 
8. <https://github.com/RongkangXiong/Bayes-Gaussian-Classify/settings> 