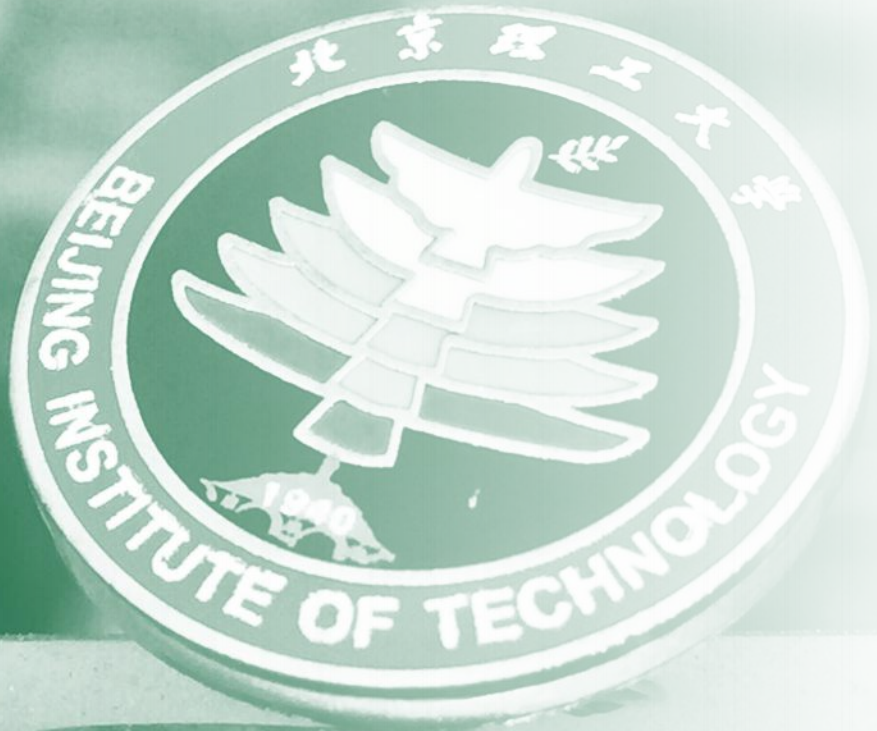


文本聚类

答辩人：钟天声，丛颖，王啸天，杨毅辰，戴铭瑞

德以明理 学以精工



目录 | CONTENTS

- 1 文本聚类简介
- 2 文本预处理和文本表征
- 3 聚类算法介绍
- 4 demo展示和分析
- 5 研究进展

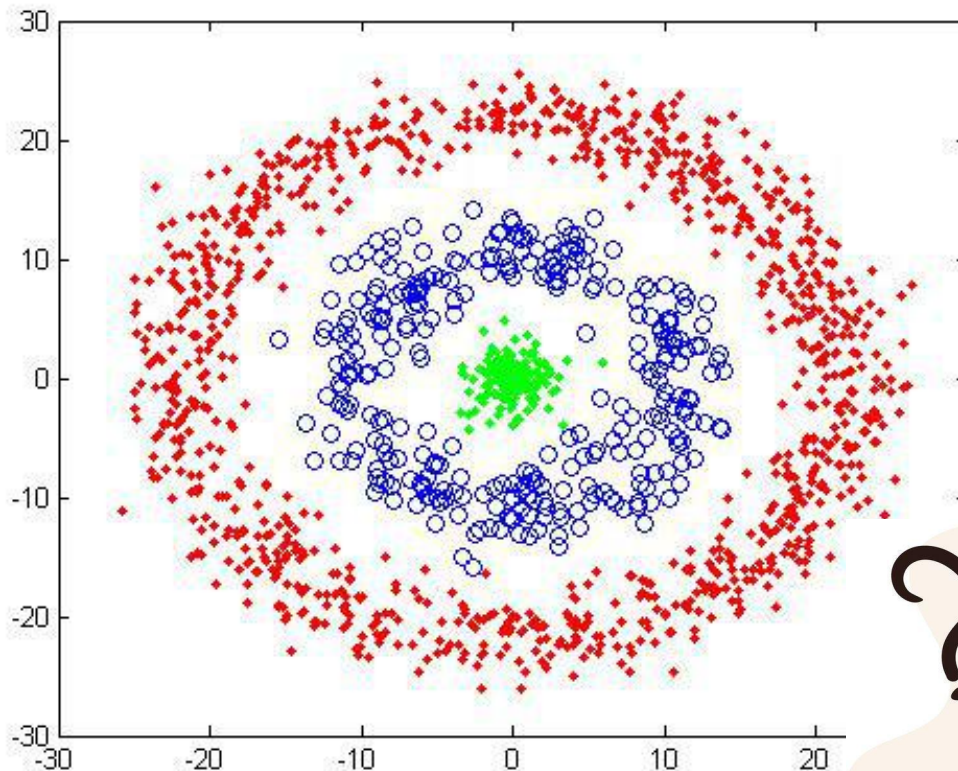


1

文本聚类简介

聚类的概念

聚类是一种流行的无监督学习方法，其任务是将一组数据分为若干个类别，使同一个类别的数据之间的差异更小而不同类别间的数据差异更大。如今聚类已经在各个领域受到广泛应用。





文本聚类

文本聚类即将一系列文本作为聚类的数据，把这些文本分为几类，同一类中的文本相似度较大而不同类别的文本间相似度则较小。

例：同为武侠小说的文章相似度较高，而计算机专业书籍和武侠小说的相似度就比较低。

文本聚类

文本聚类是一种无监督学习，机器不知道也不需要每个数据的真实标签，只是凭借机器对“簇”的理解来对数据进行聚类。


文本分类

文本分类是一种监督式学习，机器针对有标注的数据训练一个分类器，对未标注的文本进行分类。

- Newsblaster是哥伦比亚大学开发的一个多文档文摘系统，其功能是将每天发生的重要新闻文本进行聚类，对其中同主题的文本进行冗余消除、信息融合、文本生成等处理，从而生成一个简洁的摘要文档。

Columbia Newsblaster

New Mexico Forest Fire Mostly Contained




Summary:
Gov. Gary Johnson said the New Mexico blaze started when a resident dumped fireplace ash in a back yard mistakenly thinking the ashes were cold. The fire forced the evacuation of about 1,300 people. A fire was nearly contained Monday, and residents who had been evacuated will be returning home, a National Park Service spokesman said. No one has been injured. The Lincoln County assessor's office pegged property damage from the fire at \$5.2 million in assessed valuation. Last year, 1,649 fires burned 38,890 acres.

Source Articles:

- [Ashes Said to Have Caused NM Wildfire](#) (Lycos 03/26/02)
- [New Mexico Forest Fire Mostly Contained](#) (FOX News 03/26/02)
- [Wired News](#) (Wired 03/26/02)
- [Fires Destroy 36 Homes in New Mexico](#) (Lycos 03/26/02)
- [New Mexico fire nearly contained](#) (CNN 03/26/02)
- [Firefighters tame New Mexico blaze](#) (CNN 03/26/02)
- [Winds Swell Blaze into an Inferno](#) (LA Times 03/26/02)

Columbia Newsblaster

Paper: Pakistan Open to U.S. Troops Crossing Border



Summary:
Taliban and al-Qaida fighters are regrouping in Afghanistan after the recent end of the biggest ground offensive of the war, and are expected to try to mount attacks against U.S. troops there, Vice President Dick Cheney said Sunday. Two U.S. senators visiting soldiers in Afghanistan said on Tuesday that some al Qaeda fighters had fled to Pakistan and raised the possibility of putting U.S. troops on the rugged border to prevent further escapes. Pakistan, once a backer of the ultra-Islamic Taliban movement that ruled most of Afghanistan for six years until it was toppled in late 2001, has become a key ally in the U.S.-led war on terror since the September 11 attacks on the United States. The U.S.-led coalition battling al Qaeda and Taliban forces in Afghanistan is shifting its focus farther south, military officials said Tuesday, responding to unconfirmed reports that Osama bin Laden had been seen in southeastern Afghanistan.

Source Articles:

- [Afghan Authorities Arrest Taliban Commander](#) (Lycos 03/27/02)
- [Senators Want Pakistan to Stop Al Qaeda Fleeing](#) (Lycos 03/27/02)
- [Paper: Pakistan Open to U.S. Troops Crossing Border](#) (Reuters 03/27/02)
- [Sen. Wants Pakistan to Seal Borders](#) (Lycos 03/27/02)

- IBM Watson Explorer和infonetware这两个搜索引擎允许用户输入检索关键词，然后对检索到的文档进行聚类，输出各个类别的简介，这样可以帮助用户缩小检索范围。

Infonetware

Related Topics	Search the Web and Find with RealTerm
<input checked="" type="checkbox"/> RealTerm Technology (1)	<p>What is Infonetware Infonetware is a demonstration of RealTerm technology applied to the Internet search. It submits your query to a traditional Internet search engine and then sorts the results into topics. Now you can find what you want in three mouse clicks. See for yourself!</p> <p>Submitting A Query Type your query into the query box. •capitalization is not important •use quotes to indicate phrases •do not use any regular expressions</p> <p>Using RealTerm Click on a topic to see only relevant documents. For a more powerful research, check some topics into <i>select</i> <input checked="" type="checkbox"/> or <i>exclude</i> <input checked="" type="checkbox"/> state and then click on Drill Down in the tool bar at the top of the page.</p> <p>Contact RealTerm Technology has already been applied to HR, patent search, e-commerce and general corporate search. For further information contact our sales department at sales@infonetware.com.</p>
<input checked="" type="checkbox"/> RealTerm Help (1)	
<input checked="" type="checkbox"/> RealTerm FAQ (1)	
<input checked="" type="checkbox"/> RealTerm White-Paper (1)	
<input checked="" type="checkbox"/> Infogistics Ltd (1)	
<input checked="" type="checkbox"/> Your Feedback (1)	

Copyright © 2000, 2001 Infonetware. All rights reserved.

- 用户的兴趣模式挖掘：将聚类算法用于对用户感兴趣的文档进行聚类，从而发现用户的兴趣模式，有助于定制化推荐。
- 改善文本分类的结果：Fang Y C, Parthasarathy S等人使用PDDP（Principle Direction Divisive Partitioning，一种基于主成分分析的文本聚类方法）改善了文本分类的结果。



1 文本预处理

文本预处理步骤将对文本进行分词和去停用词等处理。



2 文本表征

文本表征指从文本中提取出特征，通过这些特征来表示相应的文本。



3 聚类算法

将上一步得到的用特征表示的文本送入聚类算法进行聚类，得到最终的文本聚类结果。



2 文本预处理与文本表征

文本预处理

分词 (Tokenization)

- 即将 输入的文本 (sequence) 转换为 token 的过程。
- 英文：比较容易实现分词。
- 中文：使用一些分词算法或工具，例如：Jieba分词，百度分词API，阿里云的分词API等。

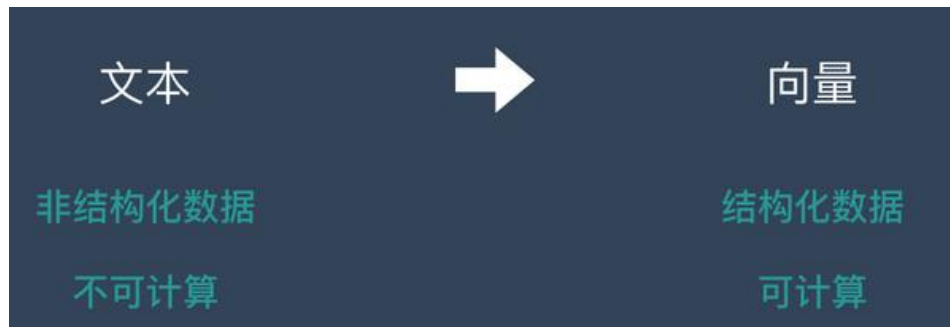
去停用词

- 停用词在这里通常指文本中出现频率很高，但对语义理解帮助不大的词。
- 语气助词、副词、介词、连词等。如中文常见的“的”、“在”、“和”、“接着”之类。英文常见的“of”，“an”
- 停用词表

词性归一化
(tokenization normalization)

- 英文文本处理中特有
- 将相同含义但是不同形式的词转换成同一个词 (token)
- 词干提取 (stemming)
beautiful --> beauty
- 词性还原 (lemmatization)
例如：good, better 和 best 词形还原的结果是 good, good 和 good.

文本表示对文本进行**数学建模**，保留语义等需要提取的信息，是后续工作的基础



文本表示模型

■ 离散表示

1. One-hot编码 2. 词袋模型 3. N-gram

■ 分布式表示

经典模型是word2vec，还包括后来的Glove、ELMO、GPT和最近很火的BERT



one-hot编码就是把每个词表示为一个长向量。这个向量的维度是词表大小，向量中只有一个维度的值为1，其余维度为0，这个维度就代表了当前的词。一段文本就应该是多个one-hot组合成的矩阵了。

文本：John likes to watch movies. Mary likes too.

构造词典：

```
{"John": 1, "likes": 2, "to": 3, "watch": 4, "movies": 5, "also":6,  
"football": 7, "games": 8, "Mary": 9, "too": 10}
```

```
["John", "likes", "to", "watch", "movies", "also", "football", "games",  
"Mary", "too " ]
```

One-hot表示：

```
John: [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]  
likes: [0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0]  
too : [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]
```

优点：简洁明了，计算快速。

缺点：无法表述词与词之间的语义关系，同时因为词表太大容易造成维数灾难等。



具体想法是：将每个文档看成一袋子词，忽略每个词出现的顺序，通过统计各个词出现的次数来表述一篇文档，文档的向量表示可以直接将各词的词向量表示加和。

文本：John likes to watch movies. Mary likes too.
==> [1, 2, 1, 1, 1, 0, 0, 0, 1, 1]

优点：这是一个基于统计，能够很快得到文本的表示。

缺点：忽略了文本字词之间的顺序，容易造成信息的丢失or混淆，如“我爱你”和“你爱我”都是同一种表示方法，但他们之间的意思确实不一样的。

2-gram建索引:

"John likes" : 1,

"likes to" : 2,

"to watch" : 3,

"watch movies" :

4,

"Mary likes" : 5,

"likes too" : 6,

"John also" : 7,

"also likes" : 8,

"watch

football" : 9,

"football games":

10,

与词袋模型相比，将连续出现的n个词所构成的词组（N-gram）也作为一个单独的特征放到向量表示中去。这个模型认为，第N个词的出现只与前面N-1个词相关，而与其它任何词都不相关。

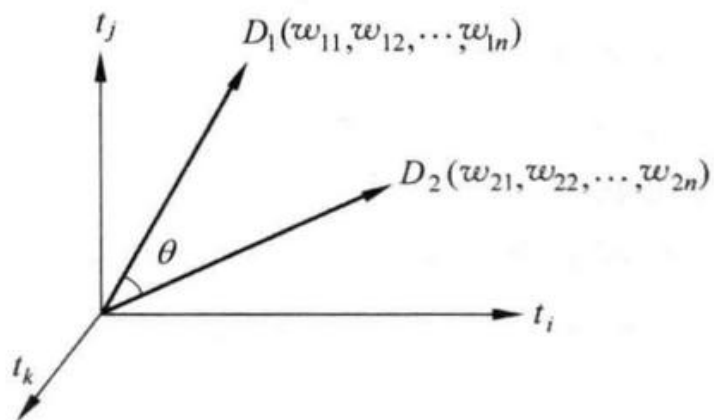
文本: John likes to watch movies. Mary likes too.

[1, 1, 1, 1, 1, 1, 0, 0, 0, 0]

优点: 考虑了词的顺序。

缺点: n越大 数量级越大。

前面的三种表示都是默认每个字词的重要程度是一样的，而在一段文本当中有的词是关键词而有的词就不是很重要，如‘的’、‘吧’、‘啊’等语气词就不是很重要，因此TF-IDF考虑了某个词在文本当中的重要程度，计算公式如下：



$$TF - IDF(t, d) = TF(t, d) * IDF(t)$$

TF：词频，单词t在文档d中出现的频率，反映了单词t对文档的重要性。

IDF：倒文档频度，衡量单词t区别于其他文档的程度。

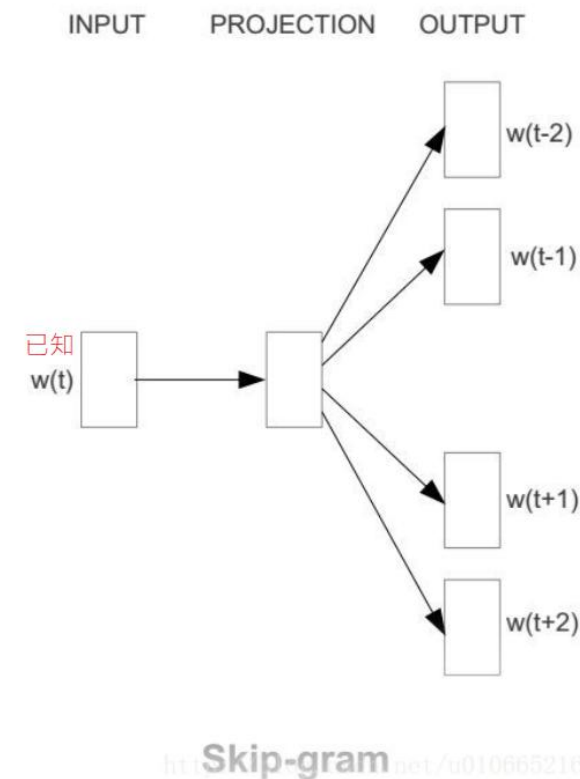
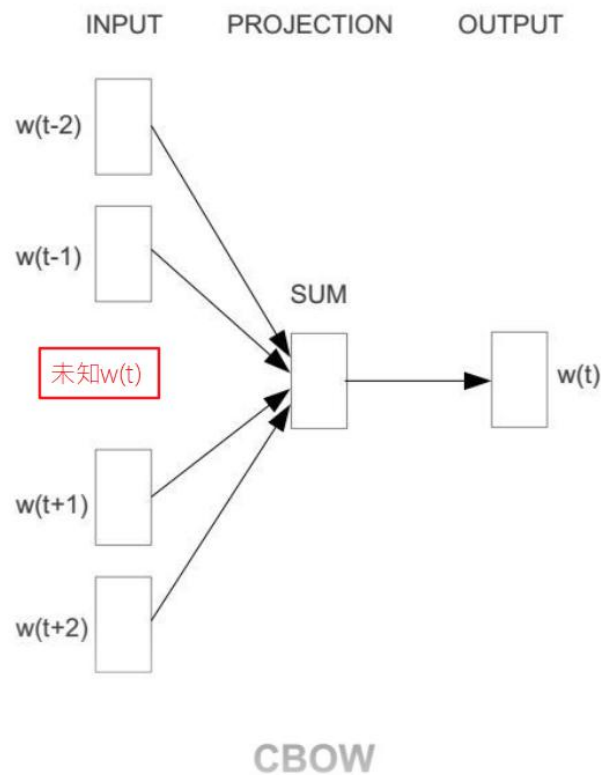
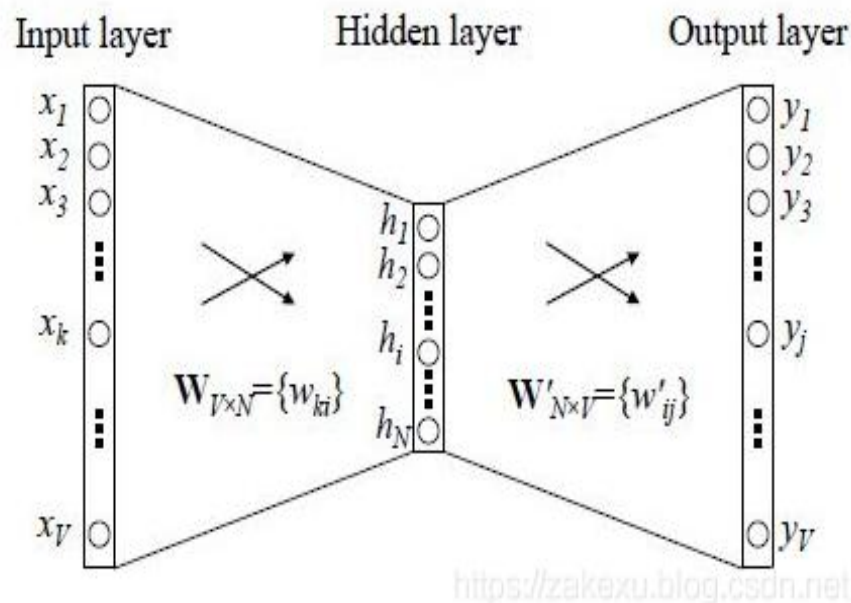
$$IDF(t) = \log \frac{\text{文章总数}}{\text{包含单词的文章总数} + 1}$$

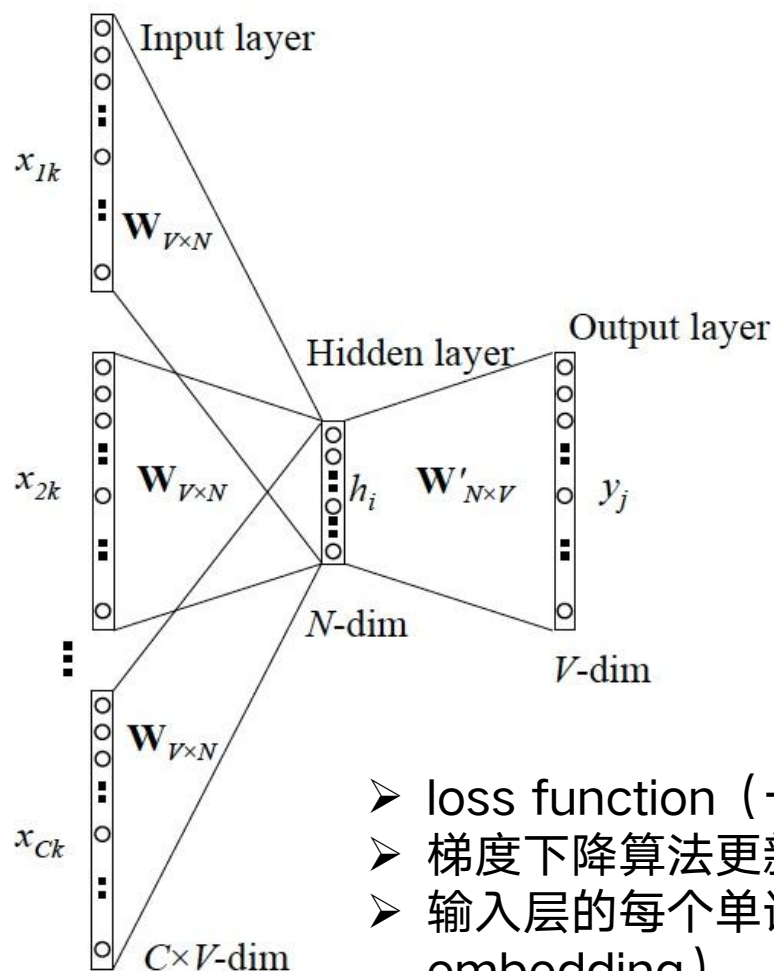
1. CBOW模型

根据中心词 $W(t)$ 周围的词来预测中心词

2. Skip-gram模型

则根据中心词 $W(t)$ 来预测周围词





1. 输入层：上下文单词的onehot. {假设单词向量空间dim为 V ，上下文单词个数为 C }
2. 所有onehot分别乘以共享的输入权重矩阵 W . { $V \times N$ 矩阵， N 为自己设定的数，初始化权重矩阵 W }
3. 所得的向量 相加求平均作为隐层向量, size为 $1 \times N$.
4. 乘以输出权重矩阵 W' { $N \times V$ }
5. 得到向量 { $1 \times V$ } softmax处理得到 V -dim概率分布
6. 概率最大的index所指示的单词为预测出的中间词 (target word) 与true label的onehot做比较, 误差越小越好
6. 根据误差更新权重矩阵

- loss function (一般为交叉熵代价函数)
- 梯度下降算法更新 W 和 W'
- 输入层的每个单词与矩阵 W 相乘得到的向量的就是我们想要的词向量 (word embedding), 矩阵 W 自身也叫做look up table, 免去训练过程直接查表得到单词的词向量了。



3

聚类算法介绍

1. 聚类算法技术
2. 划分算法
3. 层次算法
4. 密度算法
5. 网格算法

什么是聚类分析？

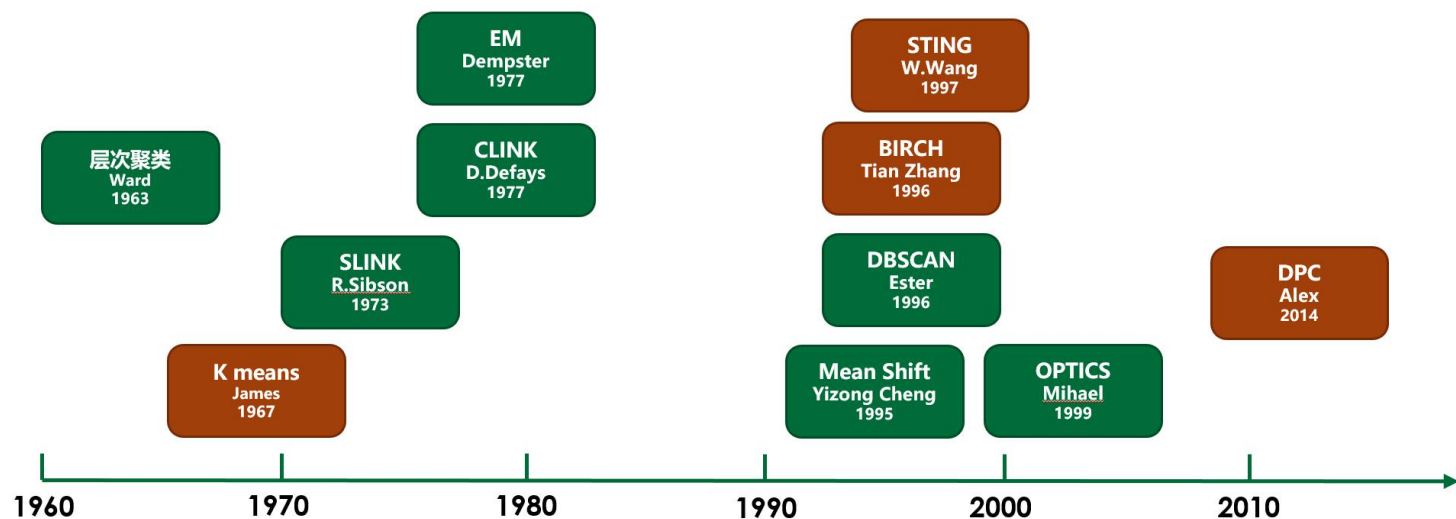
- 聚类分析，是一个把数据对象(或观测)划分成子集的过程。每个子集是一个簇，使得簇中的对象彼此相似，但与其他簇中的对象不相似。由这个过程产生的一系列簇成为一个聚类



- 可以简单概括为：同类相同、异类相异

聚类算法发展

- 聚类算法的历史与有监督学习一样悠久。层次聚类算法出现于1963年，这是非常符合人的直观思维的算法，现在还在使用。



- 聚类算法是最早被用于模式识别及数据挖掘任务的方法之一，并且被用来研究各种应用中的大数据库，因此如今用于大数据的聚类算法受到越来越多的关注



为什么需要聚类算法技术？

常见的机器学习任务大部分都是带有标记的(有监督的)，在真实的世界中，有人工标记的数据仅仅占很少的一部分，对于大数据时代而言甚至可以忽略不计，可以说无监督（不带标记）的学习才是机器学习的终极目标。

聚类算法技术目的？

探究那些无法或者还未标记的数据的内在数据结构或规律，如文本聚类问题。

如何进行聚类？

聚类分析所有的工作都关注在这两者上，划分过程与对象之间相似与否。

聚类过程一般的有：基于划分的、基于层次的、基于密度的以及基于网格的聚类算法

基于划分的聚类算法

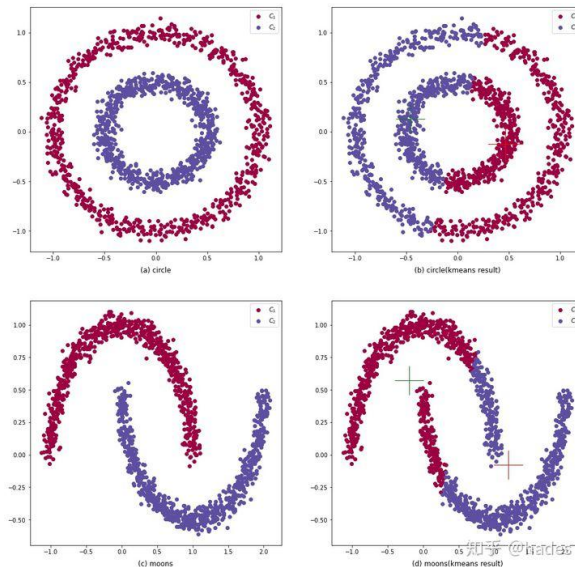
给定一个有 N 个元组或者纪录的数据集，分裂法将构造 K 个分组，每一个分组就代表一个聚类，然后采用一种迭代的重定位技术，尝试通过对对象在划分间移动来改进划分。

常见算法

- K-MEANS算法-k均值聚类算法
- K-MEDOIDS算法
- CLARANS算法

K-Means算法

- 也称为K-平均或者K-均值，一般作为掌握聚类算法的第一个算法。
- 这里的 K 为常数，需事先设定，通俗地说该算法是将没有标注的 M 个样本通过迭代的方式聚集成 K 个簇。

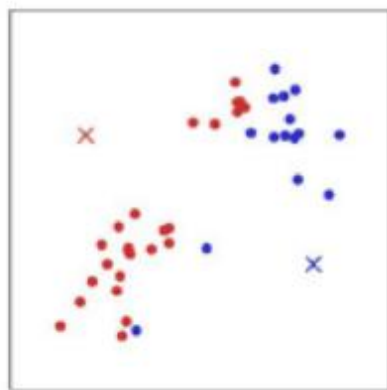




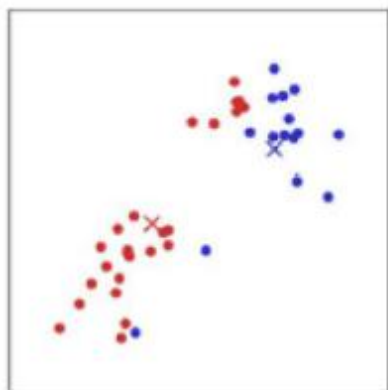
(a)



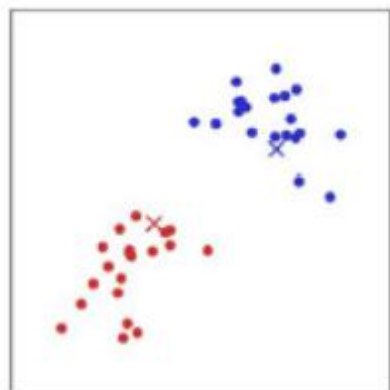
(b)



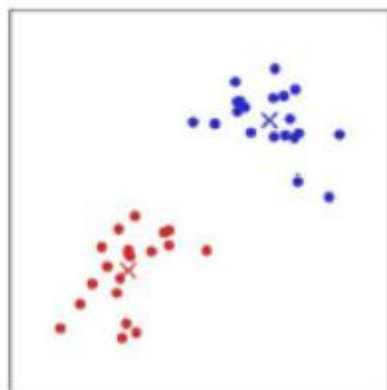
(c)



(d)



(e)



(f)

K-Means算法

1. $K=2$ ，随机两点对应的类别质心
2. 求样本距离第一轮分类
3. 从两类样本点中计算出新聚类中心
4. 重复迭代
5. 标准测度函数收敛



K-Means 算法流程

1. 首先从n个 数据对象任意选择 k 个对象作为初始聚类中心；
2. 对于所剩下其它对象，则根据它们与这些聚类中心的相似度（距离），分别将它们分配给与其最相似的（聚类中心所代表的）聚类；
3. 再计算每个所获新聚类的聚类中心（该聚类中所有对象的均值）
4. 不断重复这一过程直到标准测度函数开始收敛为止。

$$J = \sum_{i=1}^K J_i = \sum_{i=1}^K \sum_{j=1}^N w_{ji} \|X_j - C_i\|^2$$



优点

1. 原理比较简单，实现也是很容易，收敛速度快。
2. 聚类效果较优。
3. 算法的可解释度比较强。
4. 主要需要调参的参数仅仅是簇数 k

缺点

1. K 值的选取不好把握
2. 如果各类别的数据不平衡，则可能聚类效果不佳。
3. 采用迭代方法，得到的结果只是局部最优。
4. 对噪音和异常点比较的敏感。

基于层次的聚类算法

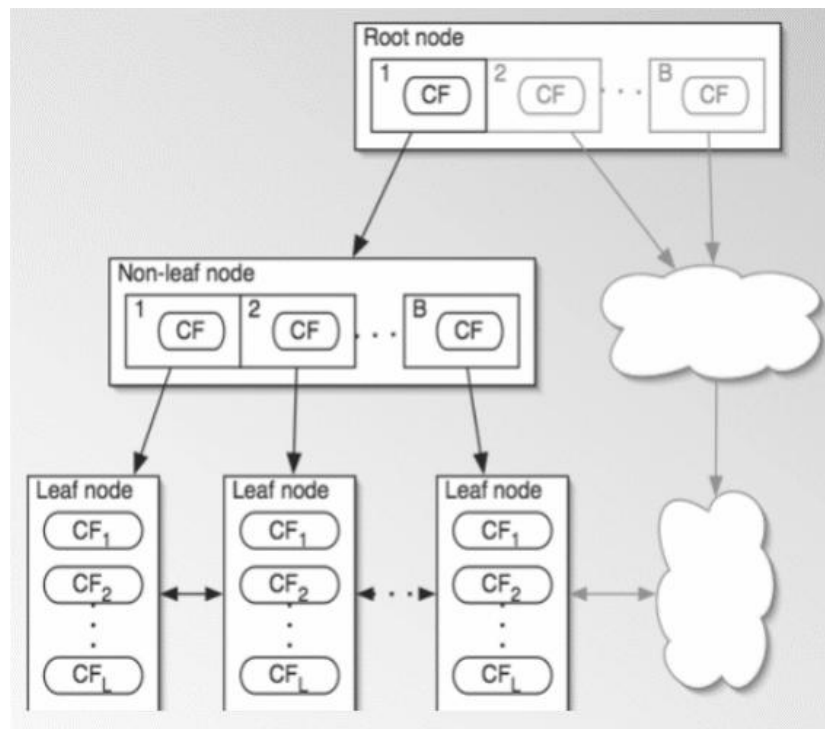
基于层次的聚类方法是指对给定的数据进行层次分解，直到满足某种条件为止。该算法根据层次分解的顺序分为自底向上法和自顶向下法，即凝聚式层次聚类算法和分裂式层次聚类算法。

常见算法

- BIRCH算法 → 利用层次方法的平衡迭代规约和聚类
- CURE算法
- CHAMELEON算法

BIRCH算法

- 核心是聚类特征树 Clustering Feature Tree 简称 CF Tree
- 每个节点包括叶子节点都有若干个CF
- 叶子结点由双向链表连接



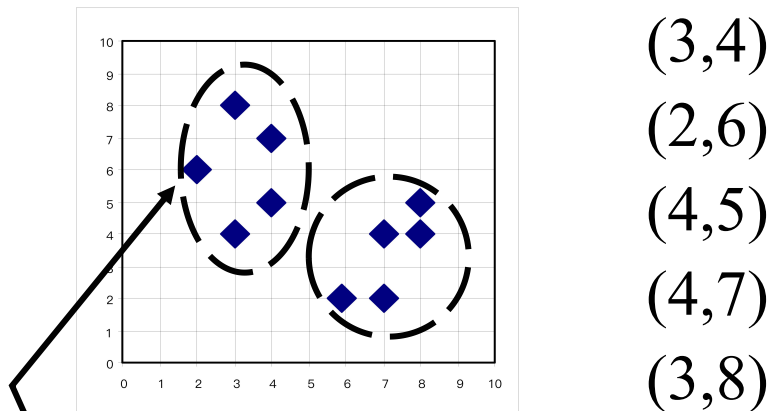


BIRCH算法-CF

■ 聚类特征:

$$CF = (N , LS , SS)$$

■ N: CF中拥有的样本点的数量

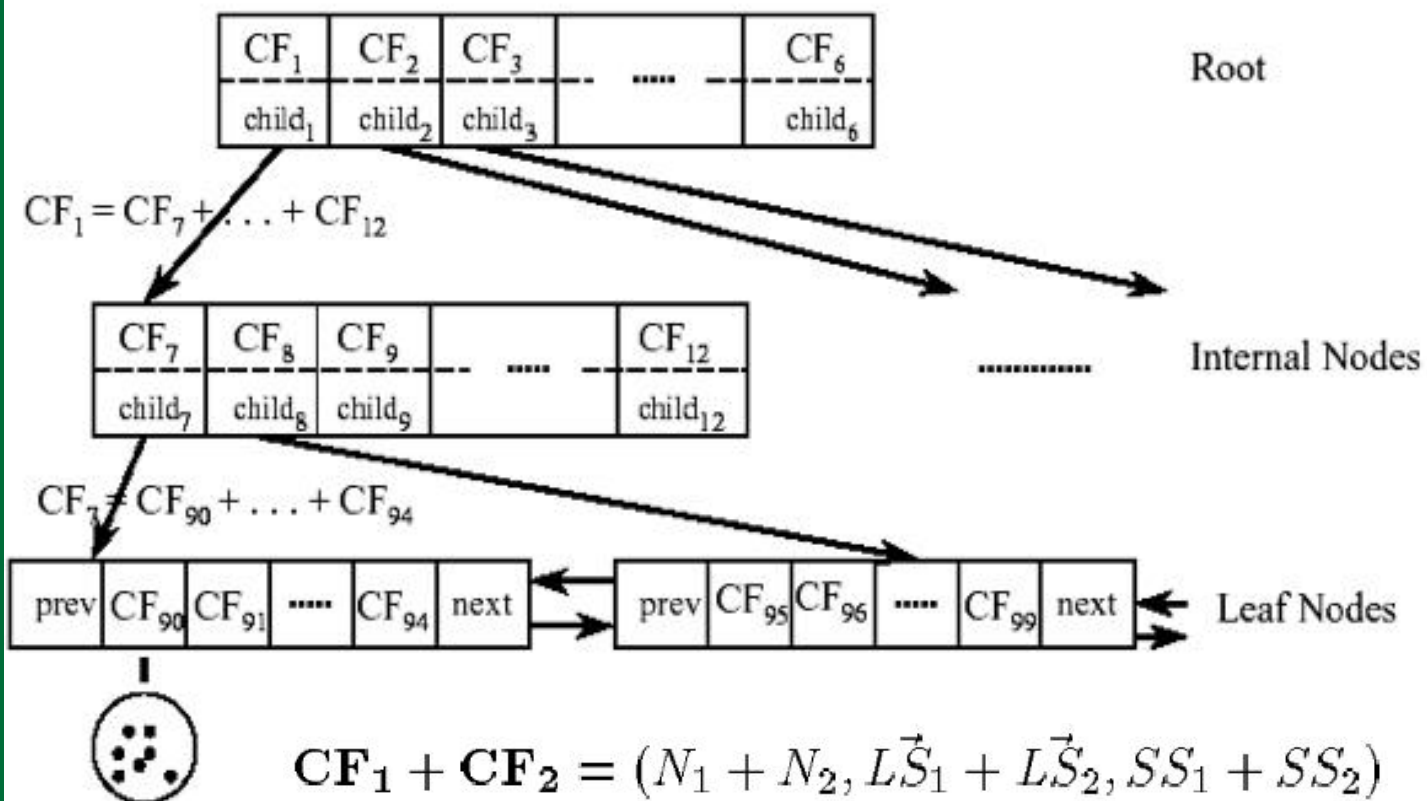
■ LS: CF中拥有的样本点各特征维度的
和向量 $\sum_{i=1}^N X_i$ ■ SS: CF中拥有的样本点各特征维度的
平方和 $\sum_{i=1}^N X_i^2$ 

$$CF = (5, (16,30), (54,190))$$

$$LS = (3+2+4+4+3, 4+6+5+7+8) = (16, 30)$$

$$\begin{aligned} SS &= (9+4+16+16+9+16+36+25+49+64) \\ &= (54+190) \\ &= 244 \end{aligned}$$

B=6,L=5



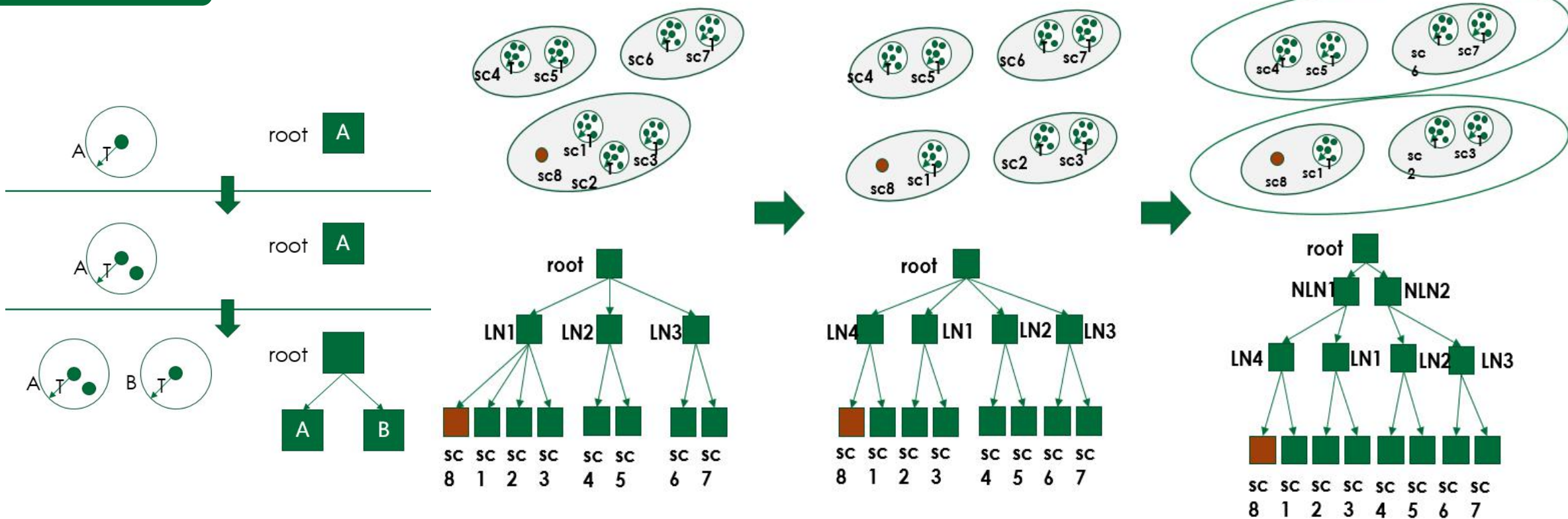
聚类特征树性质

- 满足线性关系
- 叶节点是双向链表
- B: 每个内部节点的最大CF数
- L: 每个叶子节点的最大CF数
- T: 叶节点每个CF的最大样本半径阈值

3.3

基于层次的聚类算法

B=3 L=3





BIRCH算法

BIRCH算法优点：

- 节约内存，所有的样本都在磁盘上，CF Tree仅仅存了**CF节点和对应的指针**。
- **聚类速度快**，只需要一遍扫描训练集就可以建立CF Tree，CF Tree的增删改都很快。
- 可以识别噪音点，还可以对数据集进行初步分类的预处理

BIRCH算法缺点：

- 由于CF Tree对每个节点的**CF个数有限制**，导致聚类的结果可能和真实的类别分布不同。
- 对**高维特征的数据聚类效果不好**。
- 如果数据集的分布簇**不是类似于超球体**，或者说不是凸的，则聚类效果不好。



基于密度的聚类算法

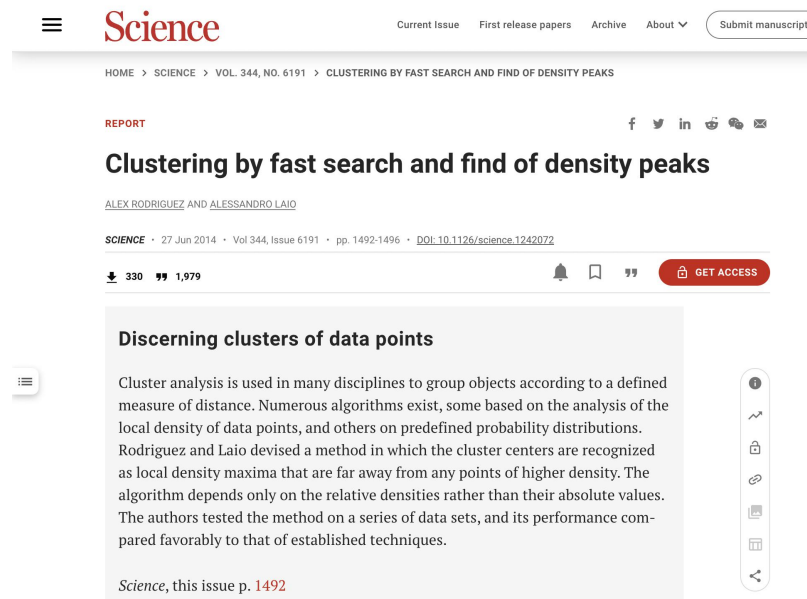
它不是基于各种各样的距离的，而是基于密度的。这样就能克服基于距离的算法只能发现“类圆形”的缺点。只要一个区域中的点的密度大过某个阈值，就把它加到与之相近的聚类中去

常见算法

- DBSCAN 算法
- DPC 算法-基于快速搜索和发现密度峰值的聚类算法
- OPTICS 算法

DPC 算法

- 能够自动地发现簇中心，实现任意形状数据的高效聚类
- 算法分为两步：寻找聚类中心和分配簇标签



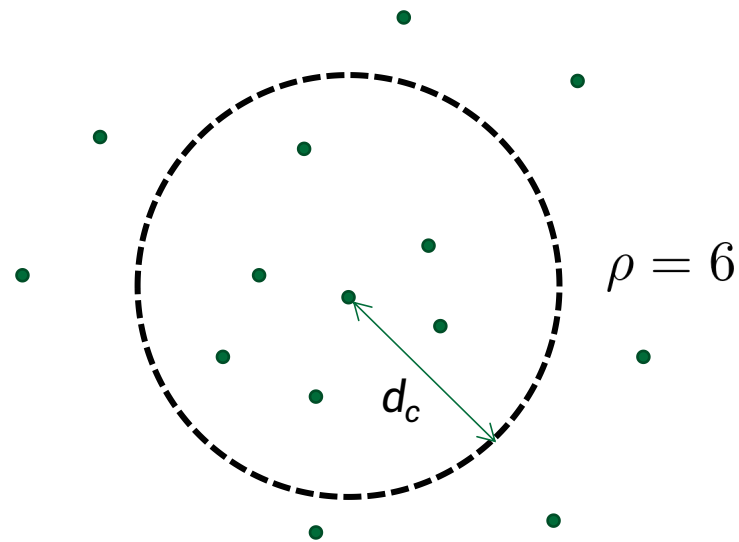
DPC算法-参数

- **密度 ρ** : DPC将每个点的密度定义为, 距离该点 d_c 以内的其它点的个数, 用公式表示为

$$\rho_i = \sum_j \chi(d_{ij} - d_c)$$

- **d_c** : 称为**截断距离**, 是该算法的唯一一个参数, $\chi(x)$ 在 $x < 0$ 时等于1, 在 $x > 0$ 时等于0.

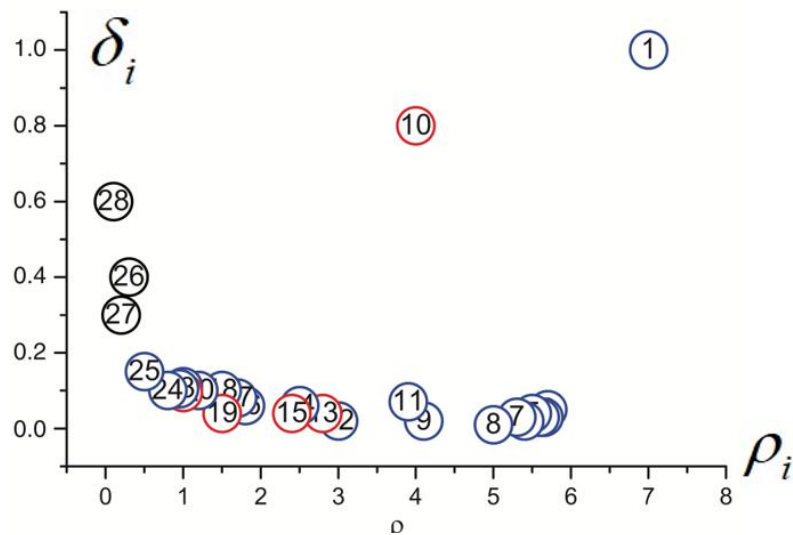
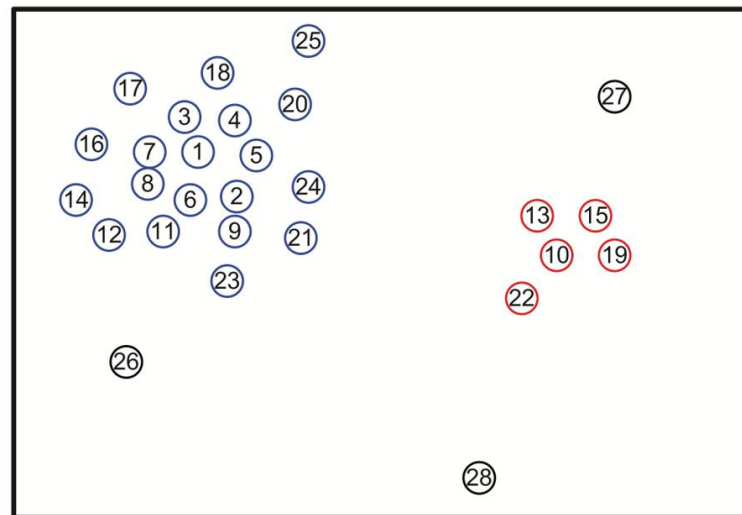
- **δ** : 每个点与最近的密度更高的点的距离, 在所有密度比P点大的点中, 假设离P点最近的那个点是点Q, 则 δ_p 为PQ的距离。计算完密度后, DPC将计算每个点的 δ 值。 用公式表示为:
$$\delta_i = \min_{j: \rho_j > \rho_i} (d_{ij})$$





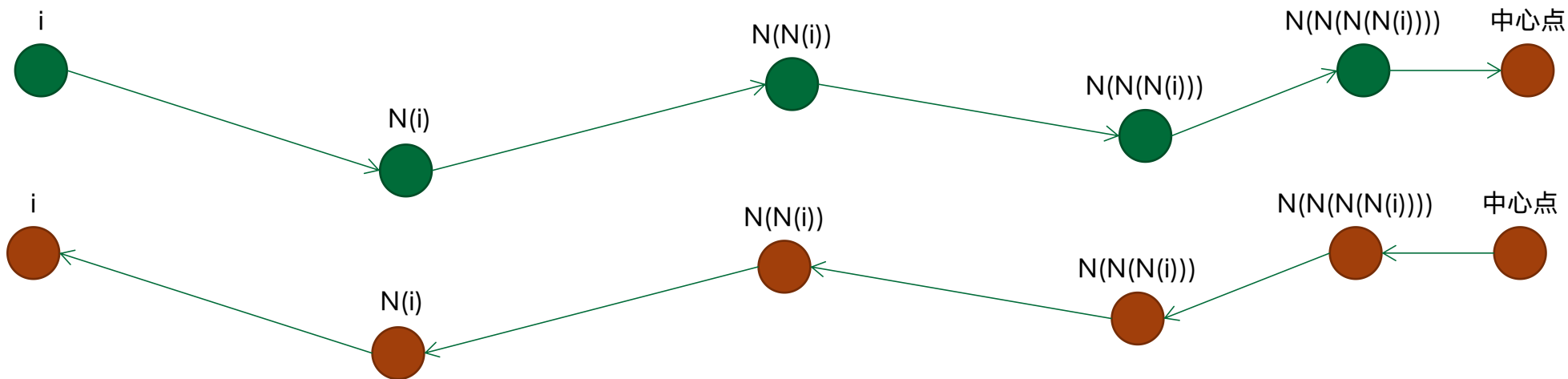
DPC算法-聚类流程

- 1. 根据截断距离算出任意数据点 x_i 的局部密度 ρ_i
- 2. 对于任意数据点 x_i , 计算出 δ_i
- 3. 以 ρ_i 为横轴, 以 δ_i 为纵轴, 画出决策图
- 4. 利用决策图, 将 ρ_i 和 δ_i 都相对较高的点标记为簇中心; 将 ρ_i 相对较低但是 δ_i 相对较高的点标记为噪声点
- 5. 将剩余点进行分配, 分配时, 将每个剩余点分配到它的最近邻且密度比其大的数据点所在簇



DPC算法-聚类流程

- 在选择完中心点后，对数据集中剩余的各点进行簇的分配工作。
- 分配的原则十分简单，在之前计算 δ 那一步，我们已经知道了任意一点 i ，所有密度大于 i 的点中距离 i 最近的点是 N_i 。所以只需将每个点 i 分配到与 N_i 相同的簇中即可，最后找到密度值最大的中心点，这是一个递归的过程，如下图所示。





优点

1. 参数非常少，只有一个。
3. 在服从高斯分布的数据集上非常准确。
3. 实现简单，时间复杂度仅 $O(n^2)$

缺点

1. 在不符合高斯分布的数据集上结果不好，且很难通过调整参数改善。



基于网格的聚类算法

用不同的网格划分方法，将数据空间划分成为有限个单元的网格结构，并对网格数据结构的统计信息进行压缩表达，基于这些统计信息判断高密度网格单元，最后将相连的高密度网格单元识别为簇

常见算法

- STING算法-统计信息网格
- CLIQUE算法
- WAVE-CLUSTER算法

STING算法

Statistical Information Grid-based method

是一种基于网格的多分辨率聚类技术，将空间区域划分为矩形单元。



STING 算法

- **核心思想**：根据属性的相关统计信息进行划分网格，而且网格是分层次的，在一个网格内的数据点即为一个簇。
- 针对不同级别的分辨率，通常存在**多个级别**的矩形单元，这些单元形成了一个层次结构：**高层**的每个单元被划分为**多个低一层**的单元。

STING 算法的两个参数：

- **网格的步长**——确定空间网格划分
- **密度阈值**——网格中对象数量大于等于该阈值表示该网格为稠密网格

STING算法-聚类流程

1. 首先我们先划分一些层次，按层次划分网格；
2. 计算最底层单位网格的统计信息（如均值，最大值和最小值）；
3. 从最底层逐层计算上一层每个父单元格的统计信息，直到最顶层
4. 同时根据密度阈值标记稠密网格。

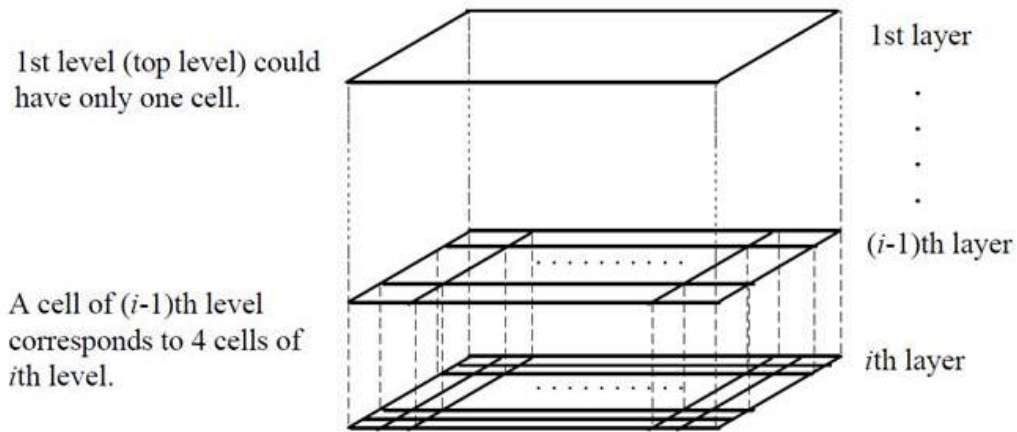
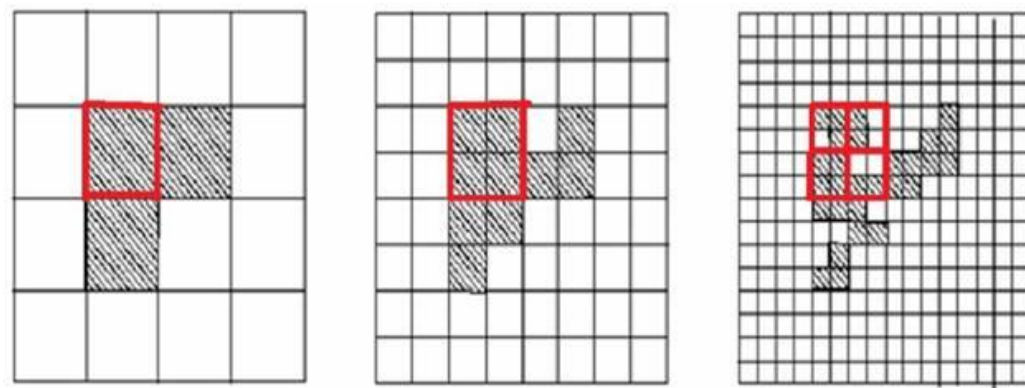


Figure 1. Hierarchical Structure



level i level $i+1$ level $i+2$
a cell of $(i-1)$ th level corresponds to 4 cells of (i) th level



STING

■ STING优点:

- 网格结构有利于并行处理和增量更新;
- 效率很高。

■ STING缺点:

- STING的聚类质量取决于网格结构的**最底层的粒度**;
- STING在构建一个父亲单元时没有考虑到子女单元和其他相邻单元之间的联系。所有的簇边界不是水平的，就是竖直的，**没有斜的分界线**。降低了聚类质量。



4 demo展示



数据来源

来源：豆瓣网

内容：书籍简介

数量：85本



聚类方法

- K-means
- Birch
- DPC



调用库

python:

glob、jieba、

sklearn、os等

```

if __name__ == "__main__":
    root = './Ttxtc'
    stopWords = open('Ttxtc1.txt', 'r', encoding='utf-8').read().split('\n')
    docPath = 'doc.txt'
    k = num
    SaveDoc(['', '/Ttxtc'], docPath, stopWords)
    bktitle = [[] for i in range(k)]
    weight = TFIDF(docPath)
    X = PCA(weight, dimension=85) # 将原始权重数据降维
    # plt.scatter(X[:,0],X[:, 1])
    # plt.show()
    # y = kmeans(X, k) # y=聚类后的类标签
    y = birch(X, k)
    # dpc = DPC(density_measure='gauss')
    # y = dpc.fit_predict(X)
    # centers = dpc.get_cluster_center()
    # c_num = len(centers)
    # bktitle = [[] for i in range(c_num)]
    # print("开始聚类:      DPC")
    # print(y)
    plt.scatter(X[:, 0], X[:, 1], c=y)
    # for i in range(len(bookTitle)):
    #     bktitle[y[i] - 1].append(bookTitle[i])
    # for i in range(c_num):
    #     print(bktitle[i], len(bktitle[i]))
    plt.show()
    silhouette_avg, sample_silhouette_values = Silhouette(X, y) # 轮廓系数
    Draw(silhouette_avg, sample_silhouette_values, y, k)

```






第二个fit_trans



语料库

名称

-  建筑.txt
-  科幻.txt
-  历史著作.txt
-  美食.txt
-  武侠小说.txt
-  心理健康论文.txt

修改日期

2021/10/28 22:20
2021/10/28 21:59
2021/10/27 22:35
2021/10/28 22:27
2021/10/28 21:35
2021/10/28 19:21

建筑.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

《认识建筑》两位作者均为执业建筑师，同时又都在欧美知名高校从事建筑理论的研究和教学。他们从自身多年的理论和实践经验出发，坚信只有通过亲身体验才能真正理解建筑。本书以体验为核心，以文字、照片和平面图三者结合的方式引领读者展开一场仿佛亲身参与的“田野调查”。全书以空间、光线、地景和场所等12个和建筑密体验密切相关的主题为框架，精选72座极具代表性的建筑杰作，涵盖全球各种建筑风格。既有古老的埃及金字塔，也有现代化的悉尼歌剧院，既有装饰华美的神圣家族大教堂，也有纯几何形的流水别墅，既有非洲的多贡人村落，也有富有东方风情的伊势神宫。

《造房子》本书是世界建筑最高奖普利兹克奖得主、著名建筑大师王澍的建筑文化随笔集。本书从建筑出发，却不止于建筑，更是一本探讨中国传统文化当代性的著作。传统文化的当代性一直是这些年学界反复思考和讨论的重要课题，王澍以自己的学术素养，以及营造经验，构建出独特的关于东方美学的审美体系，也给出传统文化进入当代的路径，这对于当下有非常重要的学术参考价值。从宋代山水画的意境，到明清园林的审美情趣，作者深入剖析中国传统文化、艺术，更以建筑的角度，从中探寻传统文化、东方哲学的美学价值。王澍的著名建筑作品包括中国美院象山校区、宁波美术馆等，在本书中，从设计开端、建造过程，直至建成后，作者用深入浅出的语言，还原这些作品的诞生历程。从中，我们看到的是作者对于“好的建筑”以及“如何做出重返传统的当代建筑”的深入思考。作者漫谈个人经历、社会与人生，更触及当下人关心的居住空间等话题，大师的成长历程和人文情怀一览无遗。

《建筑的故事》建筑的故事纵向追溯了人类建筑3000年，从原始人寄居洞穴躲避野兽袭击开始，到21世纪英国建筑师修筑草砖房以求与自然和谐共处结尾，而形成了一个完美的闭环，暗含了建筑史与人类史相伴随的命运。金字塔、帕特农神庙、万神庙、巴黎圣母院、水晶宫等16座传世建筑宛如群星闪耀在古典建筑、文艺复兴建筑、哥特式、巴洛克、包豪斯等建筑流派的银河里。每座传世建筑背后不同寻常的故事，再次将读者拉入波谲云诡的大时代，俯瞰建筑流派的演化、时代变迁。

《贝聿铭全集》以时间顺序选择了贝聿铭各个时期担任负责人或建筑设计师的50个建筑项目，图文并茂地介绍了贝氏接受委托的项目，从具体建筑项目中展现贝聿铭的内在建筑思想。书中图片包括建筑实拍图、设计平面图、剖面图、建筑过程图等等；与贝聿铭共事近20年的同事撰写文章，从独到的视角到我们理解贝聿铭和他的建筑作品。全新版本邀请新译者重新翻译，译文力求专业而流畅，不仅订正了外版原稿中的错误，还将此前译本中未译的参考文献、贝聿铭作品名录等，悉数译出，提供给读者更全面翔实的资料信息。

《建筑师》阿斯泰里奥斯波吕普是谁？他有着多种身份——从未建造过建筑的成功建筑师、博学但傲慢的教授、视角独到的美学家以及一位薄情自私的丈夫——但这全是过去的身份。现在，他年过半百，成了自己前半生的影子。然而，一个暴雨夜，闪鸣的雷电将他带上了命运之旅。雷电引起的大火烧毁了他的公寓，也让他抛弃了早已被自己的傲慢毁掉的过去。阿斯泰里奥斯用仅剩的钱买了一张能到达的最远地方的车票。他是在放逐自己，还是在寻找救赎？

《我还未读懂漫山白雪》《我还未读懂漫山白雪》是年轻的写作者、建筑师章程的首部电影随笔集，但文体更广阔，不止于读解电影。他以丰沛而

第 1 行, 第 1 列

100%

Windows (CRLF)

UTF-8



停词表

Txtc1.txt - 记事本

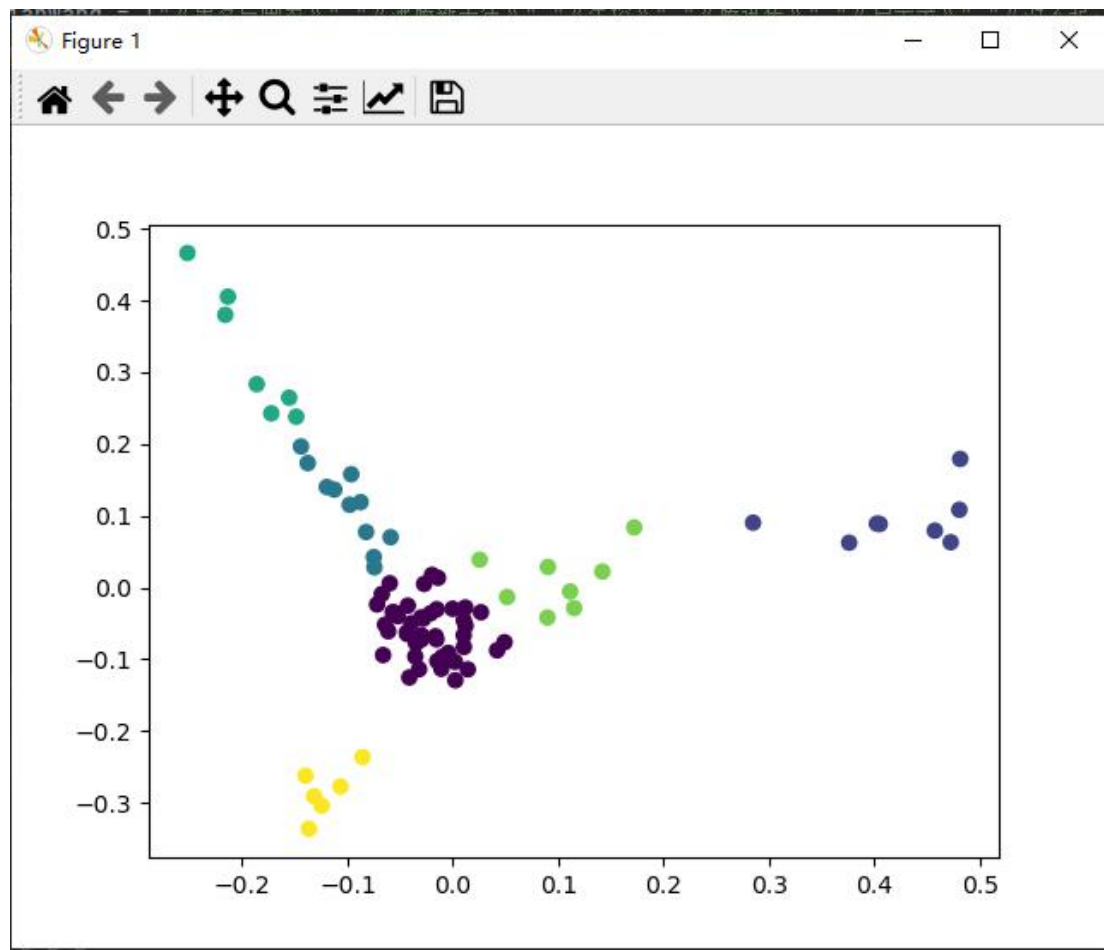
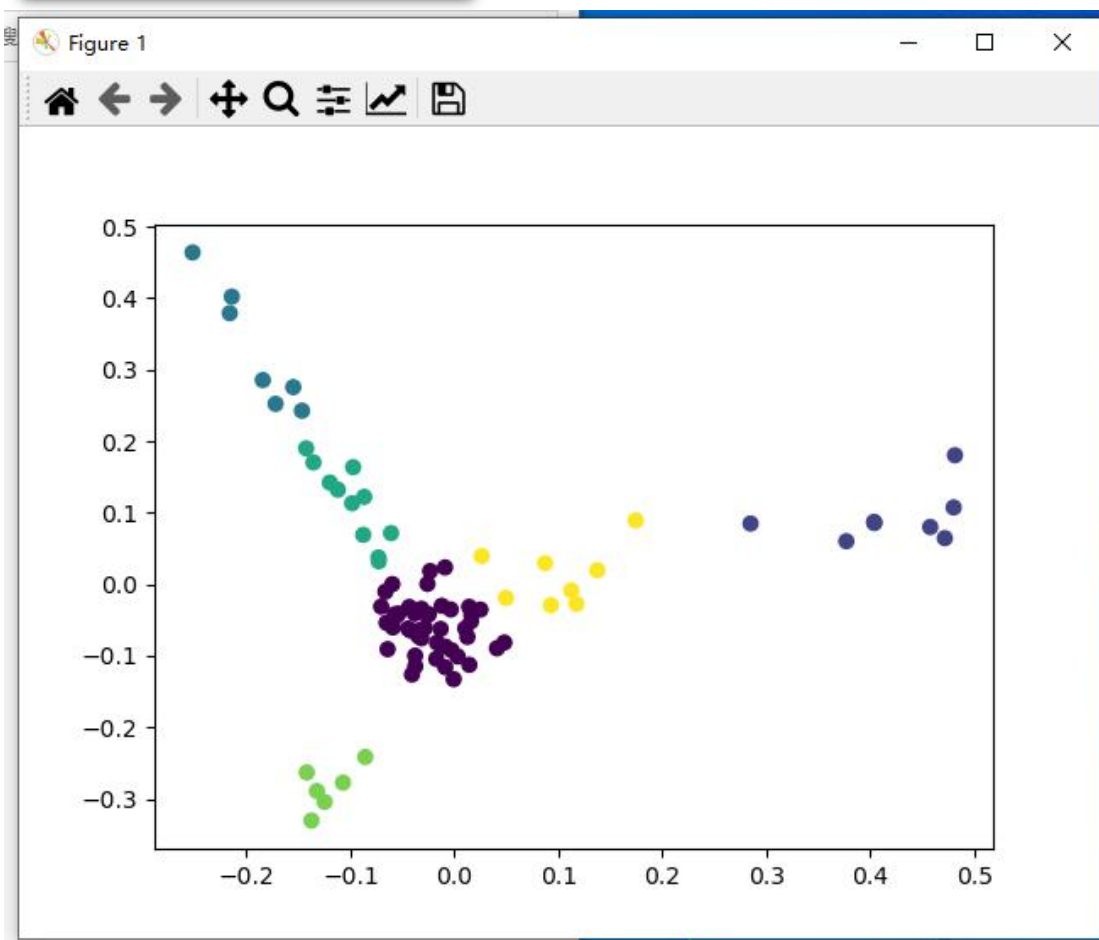
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

去年
nan
安全
食品
一个
中国
自己
可以
我们
他们
美国
没有
问题
这个
生产
进行
万元
现在
很多
不是
如果
这些
进行
这样
福喜

1 历代 政治 得失 作者 专题 演讲 合集 汉 唐 宋 明 清 五代 政府 组织 百官 职权 考试 监察 财经 赋税 兵役 义务 种种 政治 制度 作 提要 勾玄 概观 叙述 因果 演变 指陈 利害得失 高屋建瓴 总括 历史 政治
2 明朝 事儿 讲述 1344 1644 明朝 三百年 间 历史 作品 史料 基础 年代 人物 线 小说 笔法 明朝 十七 帝 王公 权贵 小人物 命运 全景 展示 官场 政治 战争 帝王 心术 著墨 作品 一部 明朝 政治经济 制度
3 欧亚大陆 征服 赶走 杀死 印第安人 澳大利亚人 非洲人 小麦 玉米 牛 猪 世界 作物 牲畜 特定 地区 地区 这部 开创性 著作 中 演化 生物学家 贾 雷德 戴蒙德 揭示 事实上 有助于 历史 环境因素 震撼人心
4 红星 照耀 曾译 西行漫记 1937 初版 畅销 董乐山 译本 中国工农红军 经典 读本 本书 真实 记录 斯诺 1936 西北 革命 根据地 实地 采访 所见所闻 全世界 报道 中国工农红军 红军 领袖 红军 将领 情况
5 第一次世界大战 硝烟 中 迈向 死亡 生命 热烈 生长 威尔士 矿工 少年 失恋 法律系 大学生 穷困潦倒 俄国 兄弟 富有 英俊 英格兰 伯爵 痴情 德国 特工 充满 灰尘 危险 煤矿 闪闪发光 皇室 宫殿 代表 权力
6 德国 战前 欧洲 富裕 强大 经济体 反犹主义 政治 中 处于 边缘 地位 一群 极端分子 恶棍 纳粹党 数年 之中 德国 一党 独裁 国家 极有 教养 民族 引向 道德 物质 文化 废墟 绝境 本书 透过 德国 历史 社会
7 本书 史料 周密 考证 分析 中古 历史 中 门阀 政治 作 探索 中外 学者 习称 魏晋 南北朝 门阀 政治 实际上 东晋 一朝 门阀 政治 皇权 政治 特定 历史 条件 奕态 暂时性 过渡性 形式 门阀 士族 皇权 共治
8 毛泽东 选集 第一卷 包括 毛泽东 同志 革命 时期 中 重要著作 几年 前 地方 出过 几种 版本 毛泽东 选集 著者 审查 体例 颇为 杂乱 文字 错讹 著作 收过去 这部 选集 中国共产党 成立 经历 历史 时期 著作
9 全球 通史 斯塔夫 里 阿诺斯 著 吴象婴 梁赤民 董书慧 王昶译 作者 本书 中 采用 全新 史学观点 方法 世界 看作 不可分割 有机 统一 全球 角度 国家 地区 角度 考察 世界各地 区 人类文明 发展 研究 重点
10 本书 艾柯 论述 中世纪 美学 理论 审美 体验 艺术 实践 作品 中世纪 之美 艾柯 笔下 呈现出 自成一 活力 继承 古希腊 古罗马 传统 教条 思想 环境中 悄然 演变 发展 成熟 批判性 观念 体系 继承 传统 偏
11 建筑 两位 作者 执业 建筑师 欧美 知名 高校 建筑 理论 研究 教学 多年 理论 实践经验 出发 坚信 切身体 理解 建筑 本书 体验 核心 文化 照片 平面图 三者 方式 引领 读者 展开 一场 仿佛 参与 田野 全
12 造 房子 本书 世界 建筑 最高奖 普利兹 克 奖得主 著名 建筑 大师 王澐 建筑 文化 随笔集 本 建筑 出发 建筑 更 是一本 探讨 传统 文化 当代 性 著作 传统 文化 当代 性 学界 反复 思索 讨论 课题 王澐
13 建筑 故事 建筑 故事 纵向 追溯 人类 建筑 3000 原始人 寄居 洞穴 躲避 野兽 袭击 世纪 英国 建筑师 修筑 草 砖房 以求 自然 和谐 共处 结尾 完美 闭环 暗合 建筑史 人类史 相伴 命运 金字塔 帕特农 神
14 贝聿铭 全集 时间 顺序 选择 贝聿铭 时期 担任 负责人 建筑 设计师 建筑 项目 图文并茂 介绍 贝氏 接受 委托 项目 建筑 项目 中 展现 贝聿铭 内在 建筑 思想 书中 图片 包括 建筑 实拍图 设计 平面图 剖面
15 建筑师 阿斯 索 里奥 波吕 普 多种 身份 建造 建筑 成功 建筑师 博学 傲慢 教授 视角 独到 美学家 一位 薄情 自私 丈夫 全是 身份 年过半百 成 前半生 影子 暴雨 夜 闪鸣 雷电 带上 命运 之旅
16 未 读懂 漫山 白雪 未 读懂 漫山 白雪 年轻 写作者 建筑师 章程 首部 电影 随笔集 文体 广阔 解读 电影 丰沛 敏锐 书写 走进 塔 夫斯基 费里尼 阿基考 里斯 马基 伍迪 艾伦 大岛 渚 姜妍 王家卫 侯孝贤 三
17 穿墙 透壁 本书 作者 二十年 古建筑 考察 心得 涵盖 神灵 殿堂 帝王 国度 众生 居所 三个 面向 十六大 类 建筑 探索 五十一 座 经典 个案 时间 秦汉 明清 空间 遍布 中华 大地 无论是 尺度 宏大 宫殿 寺院
18 制造 东京 怀着 梦想 建造 东京 银座 浅草寺 新宿 上野 公园 传统 地标 东京 超越 人类 智慧 怪物 令人 难以置信 明治 时期 东京 走向 衰微 一度 政府 放弃 无数 怀抱 梦想 尝试 失败 推倒重来 商议 制定
19 建筑 好玩 欧洲 篇 本书 活泼 语言 漫画 形式 历史 故事 读者 整体 快速 欧洲 建筑 发展 历史 包括 西方 建筑 老祖宗 古希腊 罗马 角斗场 是否是 帝国 维稳 工具 哥特式 教堂 成 黑暗 化身 本书 适合 欧洲
20 建筑家 安藤忠雄 建筑家 安藤忠雄 四十年 安藤忠雄 从 没 无名 只能 躺 事务所 地板 发呆 打滚 找 空地 发想 建筑 样式 非 学院 出身 建筑师 今日 争相 世界 大学 建筑系 聘请 授课 世界各地 留下 融入 自然
21 安藤忠雄 安藤忠雄 建造 世界 记录 安藤忠雄 年来 代表性 建筑 作品 人生 关键时刻 著文 建筑 中 浓缩 思考 Lens 对话 中 呈现 社会 看法 年轻人 建议 挑战 自由 展开 人生 讲述 住宅 讲透 住宅 原点 光 与影
22 生活 艺术家 手作 私宅 本书 日本 著名 建筑师 中村 好文 拜访 艺术 创作者 宅邸 文集 中村 好文书 中 走访 位 艺术家 住宅 建筑师 眼光 发掘 潜藏在 建筑 中 人性 温度 独特 灵感 书中 宅前 艺术家 那种
23 空间 诗学 本书 初版 1957 现代主义 晚期 建筑 文化 窒息 氛围 中 此书 现象学 象征意义 角度 建筑 展开 独到 思考 想象 作者 空间 填充 物体 容器 人类 意识 居所 建筑学 栖居 诗学 书中 精彩 处 莫过于
24 手绘 紫禁城 遗失 在日本 北京 皇城 建筑 艺术 本书 1901 伊东忠 太 奥山 恒 五郎 学者 北京 紫禁城 相关 建筑 测量 拍照 素描 写生 原件 比例 缩绘后 编著 而成 古代 建筑 结构 装饰 详细 解剖 揭
25 城市 发展史 城市 发展史 起源 演变 前景 著名 城市 理论家 社会 哲学家 刘易斯 芒福德 理论 著作 着重 人文科学 角度 系统地 阐述 城市 起源 发展 展望 远景 城市 发展史 起源 演变 前景 史料 提高 实用性
26 图像 建筑史 一本 理解 古代 建筑 有机 结构 入门 读物 借助 古建筑 典型 实例 照片 图解 阐释 古建筑 结构 体系 三十个 世纪 中 发展 形制 演变 孕育 发祥 史前 时期 发育 成长 汉代 成熟 涅槃 激荡于 唐代
27 心理健康 论文 大学生 人际交往 中 交往 活动中 有时候 两 评价 差距 人会 烦恼 善于 调节 两 评价 全面提高 综合 素质 自我认识 有助于 找到 社会 位置 扮演 社会 角色 人际交往 社会 发展 产物 社会 发展
28 心理健康 论文 大学 大学生 拥有 自由 空间 交往 欲望 友谊 美好 憧憬 结交 朋友 相处 时间 一长 发现 朋友 身上 缺点 毛病 大学生 方式 改造 失败 大学生 疏远 朋友 友谊 停留 表面 大学生 经历 几次 打击
29 心理健康 论文 物质 重要性 强化 社会 不良风气 大学生 影响 大学生 交往 对象 选择 注重 物质条件 带有 功利性 交往 忽视 弱势群体 家庭 困难 内向 学生 忽视 功利性 影响 自卑心理 学生 个人利益 看得 很
30 心理健康 论文 社会交往 中 个体 知识 水平 涵养 影响 交往 效果 应从 点滴 从善如流 勿 善小而 勿以恶小而为之 优化 社交 形象 学生 人际交往 中 社交 恐惧 胆怯 羞怯 自卑 冷漠 孤独 封闭 猜疑 嫉妒
31 心理健康 论文 大学生 人际交往 中 把握 适度 原则 一是 交往 广度 过广 分散 精力 过窄 排他性 二是 交往 深度 交往 对象 浅交 深交 三是 交往 频率 便是 朋友 距离 心理 感 愉悦感 交往 中 大学生 应 学
32 心理健康 论文 大学生 交往 中 建立 人际关系 关键在于 心理 状态 心理 状态 源于 客观 评价 自我 源于 接纳 自我 表现 客观 优势 弱势 自傲 自负 自卑 羞怯 敌视 学会 原谅 包括 失误 过失 注重 自我 修
33 神雕侠侣 金庸 作品集 主人公 杨过 自然而然 走上 非正统 人生道路 入 道流 至情至性 自我 利益 情感 个性 人格 尊严 置于 人生 首位 首要 目标 待人处事 评价 是非 首要 原则 书中 将 杨过 郭靖 杀 父之仇 身
34 射雕 英雄传 金庸 代表作 作于一九五七年 一九五九年 香港 商报 连载 射雕 中 人物 个性 郭靖 诚朴 厚重 黄蓉 机智 狡狴 读者 印象 深刻 这是 传统 小说 戏剧 特征 缺乏 人物 内心世界 复杂性 人物性格 情
35 倚天 屠龙记 金庸 武侠小说 著于 1961 射雕 三部曲 第三部 该书 以元末 群雄 纷起 江湖 动荡 广阔 背景 叙述 武当 弟子 张无忌 江湖 生涯 表现 众 武林 豪杰 质朴 自然 形态各异 精神风貌 展现 人格
36 鹿鼎记 金庸 创作 一部 小说 代表作 小说 讲 扬州 妓院 长大 小孩 韦小宝 武功 姿态 闯江湖 各大 帮会 周旋 皇帝 朝臣 之间 奉旨 远征 云南 俄罗斯 故事 书中 充满 精彩 对白 逆 思考 事件 韦小宝 笑称 不



K-means





K-means

```

开始聚类:           Kmeans
[0 0 0 0 0 0 0 0 0 5 1 1 1 1 5 5 5 0 1 1 1 1 5 5 1 0 5 4 4 4 4 4 0 0 0 0 0
 0 0 0 0 0 0 3 2 0 0 2 3 3 3 0 3 2 3 2 3 2 2 3 0 3 3 2 5 3 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0]

['《中国历代政治得失》', '《明朝那些事儿》', '《红星照耀中国》', '《枪炮、病菌与钢铁》', '《巨人的陨落》', '《第三帝国的到来》', '《东晋门阀政治》', '《毛泽东选集》', '《全球通史》', '《制造东京》', '《城市发展史》', '《神雕侠侣》', '《射雕英雄传》', '《认识建筑》', '《造房子》', '《建筑的故事》', '《贝聿铭全集》', '《建筑也可以很好玩：欧洲篇》', '《建筑家安藤忠雄》', '《安藤忠雄》', '《手绘紫禁城：遗失在日本的北京皇城建筑艺术》'] 8
['《三体》', '《齐马蓝》', '《银河帝国：基地七部曲》', '《三体II》', '《三体III》', '《沙丘》', '《惨败》'] 7
['《索拉里斯星》', '《星之继承者》', '《未来学大会》', '《永恒的终结》', '《小镇奇谈》', '《仿生人会梦见电子羊吗？》', '《你一生的故事》', '《其主之声》', '《球状闪电》', '《机器人大师》', '《环界1》'] 11
['《心理健康论文一》', '《心理健康论文二》', '《心理健康论文三》', '《心理健康论文四》', '《心理健康论文五》', '《心理健康论文六》'] 6
['《中世纪之美》', '《建筑师》', '《我还未读懂漫山白雪》', '《穿墙透壁》', '《生活艺术家的手作私宅》', '《空间的诗学》', '《图像中国建筑史》', '《雪崩》'] 8

计算轮廓系数:
0.5733305955530371
It's over!!

```

```

开始聚类:           Kmeans
[0 0 0 0 0 0 0 0 0 4 1 1 1 1 4 4 4 0 1 1 1 1 4 4 1 0 4 5 5 5 5 5 0 0 0 0 0
 0 0 0 0 0 2 3 0 0 3 2 2 2 0 2 3 2 3 2 3 3 2 0 2 2 3 4 2 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0]

['《中国历代政治得失》', '《明朝那些事儿》', '《红星照耀中国》', '《枪炮、病菌与钢铁》', '《巨人的陨落》', '《第三帝国的到来》', '《东晋门阀政治》', '《毛泽东选集》', '《全球通史》', '《制造东京》', '《城市发展史》', '《神雕侠侣》', '《认识建筑》', '《造房子》', '《建筑的故事》', '《贝聿铭全集》', '《建筑也可以很好玩：欧洲篇》', '《建筑家安藤忠雄》', '《安藤忠雄》', '《手绘紫禁城：遗失在日本的北京皇城建筑艺术》'] 8
['《索拉里斯星》', '《星之继承者》', '《未来学大会》', '《永恒的终结》', '《小镇奇谈》', '《仿生人会梦见电子羊吗？》', '《你一生的故事》', '《其主之声》', '《球状闪电》', '《机器人大师》', '《环界1》'] 11
['《三体》', '《齐马蓝》', '《银河帝国：基地七部曲》', '《三体II》', '《三体III》', '《沙丘》', '《惨败》'] 7
['《中世纪之美》', '《建筑师》', '《我还未读懂漫山白雪》', '《穿墙透壁》', '《生活艺术家的手作私宅》', '《空间的诗学》', '《图像中国建筑史》', '《雪崩》'] 8
['《心理健康论文一》', '《心理健康论文二》', '《心理健康论文三》', '《心理健康论文四》', '《心理健康论文五》', '《心理健康论文六》'] 6

计算轮廓系数:
0.5664446892166937
It's over!!

```



K-means

《心理健康论文一》', '《心理健康论文二》', '《心理健康论文三》', '《心理健康论文四》', '《心理健康论文五》', '《心理健康论文六》

《中世纪之美》', '《建筑师》', '《我还未读懂漫山白雪》', '《穿墙透壁》', '《生活艺术家的手作私宅》', '《空间的诗学》', '《图像中国建筑史》', '《雪崩》

《三体》', '《齐马蓝》', '《银河帝国：基地七部曲》', '《三体II》', '《三体III》', '《沙丘》', '《惨败》

《认识建筑》', '《造房子》', '《建筑的故事》', '《贝聿铭全集》', '《建筑也可以很好玩：欧洲篇》', '《建筑家安藤忠雄》', '《安藤忠雄》', '《手绘紫禁城：遗失在日本的北京皇城建筑艺术》

《索拉里斯星》', '《星之继承者》', '《未来学大会》', '《永恒的终结》', '《小镇奇谈》', '《仿生人会梦见电子羊吗？》', '《你一生的故事》', '《其主之声》', '《球状闪电》', '《机器人大师》', '《环界1》

《中国历代政治得失》', '《明朝那些事儿》', '《红星照耀中国》', '《枪炮、病菌与钢铁》', '《巨人的陨落》', '《第三帝国的到来》', '《东晋门阀政治》', '《毛泽东选集》', '《全球通史》', '《制造东京》', '《城市发展史》', '《神雕侠侣》', '《射雕英雄传》', '《倚天屠龙记》', '《鹿鼎记》', '《小李飞刀：多情剑客无情剑》', '《侠客行》', '《逝去的武林：一代形意拳大师口述历史》', '《武林：一代形意拳大师口述历史》', '《蜀山剑侠传》', '《笑傲江湖》', '《献给阿尔吉侬的花束》', '《克莱因壶》', '《平面国》', '《醉步男》', '《鱼翅与花椒》', '《中国人超会吃》', '《雅舍谈吃》', '《行走的柠檬》', '《四口吃遍江户》', '《随园食单》', '《料理图鉴》', '《美食与文明》', '《五味》', '《于谦：人间烟火》', '《肚子饿万岁》', '《厨艺的常识》', '《鸭川食堂》', '《这本书好吃吗》', '《人间滋味》', '《至味在人间》', '《流动的餐桌》', '《川菜》', '《金牌主厨的面点课》', '《世间味道》



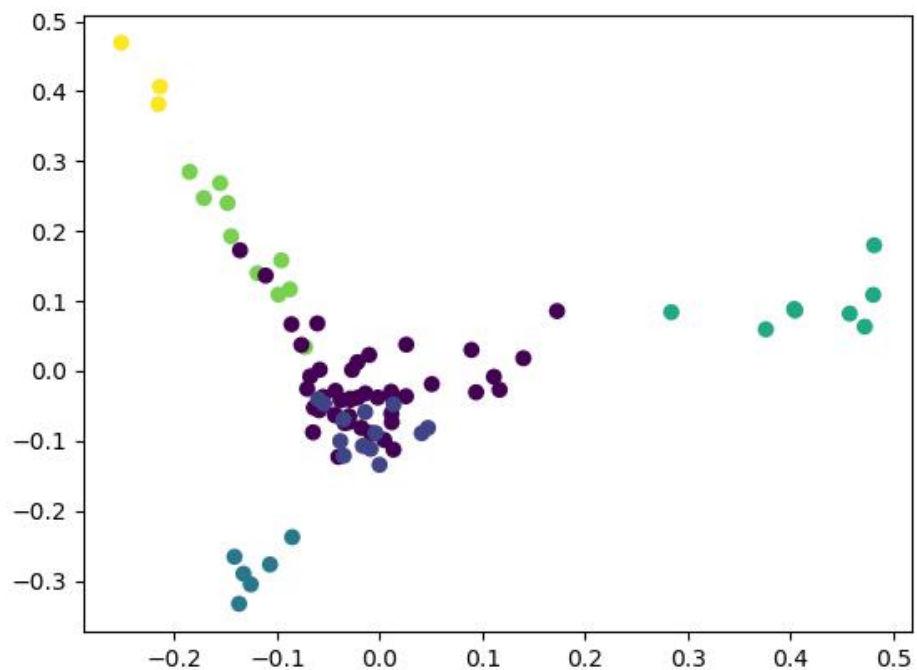
Birch

```

开始聚类:          birch
输出聚类结果:
[0 0 0 0 0 0 0 0 0 0 3 3 3 3 0 0 0 0 3 3 3 0 0 3 0 0 2 2 2 2 2 2 0 0 0 0 0
 0 0 0 0 0 4 5 0 0 4 4 4 4 0 0 4 0 5 4 5 4 4 0 0 0 4 0 0 1 1 0 0 0 1 1 0 1
 1 0 1 1 1 1 1 0 1 1 1]
['《中国历代政治得失》', '《明朝那些事儿》', '《红星照耀中国》', '《枪炮、病菌与钢铁》', '《巨人的陨落》', '《第三帝国的到来》',
'《鱼翅与花椒》', '《中国人超会吃》', '《随园食单》', '《料理图鉴》', '《五味》', '《于谦：人间烟火》', '《厨艺的常识》',
'《心理健康论文一》', '《心理健康论文二》', '《心理健康论文三》', '《心理健康论文四》', '《心理健康论文五》', '《心理健康论文六》',
'《认识建筑》', '《造房子》', '《建筑的故事》', '《贝聿铭全集》', '《建筑也可以很好玩：欧洲篇》', '《建筑家安藤忠雄》', '《建筑家安藤忠雄》',
'《索拉里斯星》', '《齐马蓝》', '《星之继承者》', '《未来学大会》', '《永恒的终结》', '《银河帝国：基地七部曲》', '《你一生的故事》',
'《三体》', '《三体II》', '《三体III》'] 3
计算轮廓系数:
0.014277223302942345
It's over!!

```

Figure 1





Birch

《心理健康论文一》', '《心理健康论文二》', '《心理健康论文三》', '《心理健康论文四》', '《心理健康论文五》', '《心理健康论文六》

《认识建筑》', '《造房子》', '《建筑的故事》', '《贝聿铭全集》', '《建筑也可以很好玩：欧洲篇》', '《建筑家安藤忠雄》', '《安藤忠雄》', '《手绘紫禁城：遗失在日本的北京皇城建筑艺术》

《索拉里斯星》', '《齐马蓝》', '《星之继承者》', '《未来学大会》', '《永恒的终结》', '《银河帝国：基地七部曲》', '《你一生的故事》', '《沙丘》', '《其主之声》', '《惨败》

《三体》', '《三体 II》', '《三体 III》

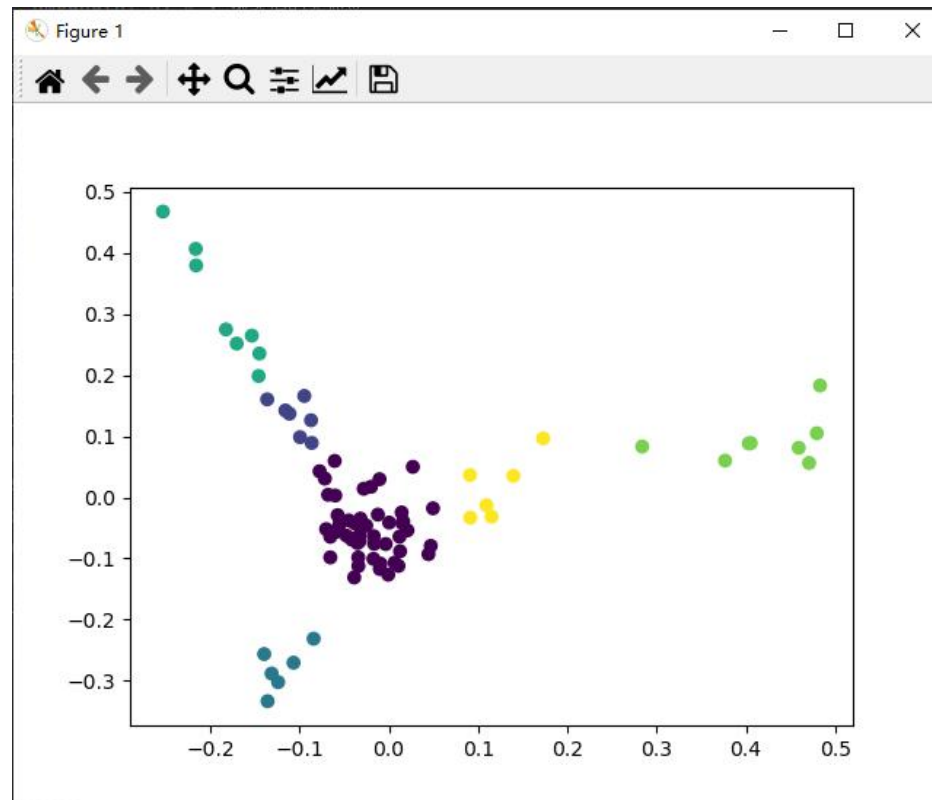
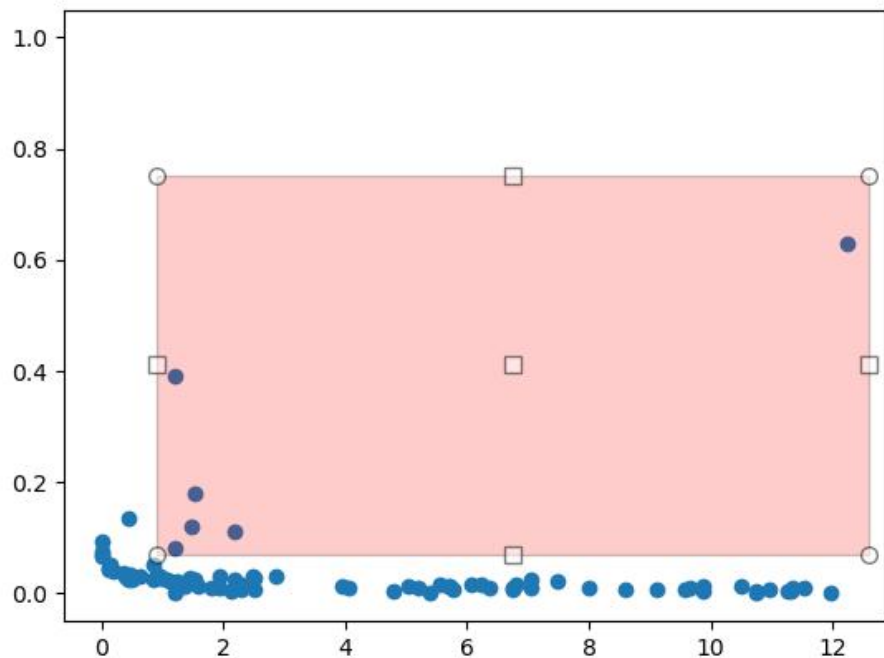
《鱼翅与花椒》', '《中国人超会吃》', '《随园食单》', '《料理图鉴》', '《五味》', '《于谦：人间烟火》', '《厨艺的常识》', '《鸭川食堂》', '《这本书好吃吗》', '《人间滋味》', '《至味在人间》', '《川菜》', '《金牌主厨的面点课》', '《世间味道》

《中国历代政治得失》', '《明朝那些事儿》', '《红星照耀中国》', '《枪炮、病菌与钢铁》', '《巨人的陨落》', '《第三帝国的到来》', '《东晋门阀政治》', '《毛泽东选集》', '《全球通史》', '《中世纪之美》', '《建筑师》', '《我还未读懂漫山白雪》', '《穿墙透壁》', '《制造东京》', '《生活艺术家的手作私宅》', '《空间的诗学》', '《城市发展史》', '《图像中国建筑史》', '《神雕侠侣》', '《射雕英雄传》', '《倚天屠龙记》', '《鹿鼎记》', '《小李飞刀：多情剑客无情剑》', '《侠客行》', '《逝去的武林：一代形意拳大师口述历史》', '《武林：一代形意拳大师口述历史》', '《蜀山剑侠传》', '《笑傲江湖》', '《献给阿尔吉侬的花束》', '《克莱因壶》', '《平面国》', '《小镇奇谈》', '《仿生人会梦见电子羊吗？》', '《醉步男》', '《球状闪电》', '《机器人大师》', '《雪崩》', '《环界1》', '《雅舍谈吃》', '《行走的柠檬》', '《四口吃遍江户》', '《美食与文明》', '《肚子饿万岁》', '《流动的餐桌》

DPC



(DPC) Click and drag to rectangle-select cluster centers.





DPC

```

开始聚类:          DPC
[1 1 1 1 1 1 1 1 1 5 5 5 5 6 6 6 1 5 5 5 6 6 5 1 6 3 3 3 3 3 1 1 1 1 1
 1 1 1 1 1 2 4 1 1 4 4 1 2 1 2 4 2 4 2 4 2 1 2 1 4 1 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1]
['《中国历代政治得失》', '《明朝那些事儿》', '《红星照耀中国》', '《枪炮、病菌与钢铁》', '《巨人的陨落》', '《第三帝国的到来》', '《东晋门阀政治》', '《毛泽东选集》', '《全球通史》', '《中世纪之美》', '《制造东京》', '《城市发
['《索拉里斯星》', '《永恒的终结》', '《小镇奇谈》', '《仿生人会梦见电子羊吗?》', '《你一生的故事》', '《其主之声》', '《球状闪电》'] 7
['《心理健康论文一》', '《心理健康论文二》', '《心理健康论文三》', '《心理健康论文四》', '《心理健康论文五》', '《心理健康论文六》'] 6
['《三体》', '《齐马蓝》', '《星之继承者》', '《银河帝国：基地七部曲》', '《三体II》', '《三体III》', '《沙丘》', '《惨败》'] 8
['《认识建筑》', '《造房子》', '《建筑的故事》', '《贝聿铭全集》', '《建筑也可以很好玩：欧洲篇》', '《建筑家安藤忠雄》', '《安藤忠雄》', '《手绘紫禁城：遗失在日本的北京皇城建筑艺术》'] 8
['《建筑师》', '《我还未读懂漫山白雪》', '《穿墙透壁》', '《生活艺术家的手作私宅》', '《空间的诗学》', '《图像中国建筑史》'] 6
计算轮廓系数：
0.5408256899760746
It's over!!

```



DPC

《心理健康论文一》', '《心理健康论文二》', '《心理健康论文三》',
'《心理健康论文四》', '《心理健康论文五》', '《心理健康论文六》

《三体》', '《齐马蓝》', '《星之继承者》', '《银河帝国：基地七部
曲》', '《三体 II》', '《三体 III》', '《沙丘》', '《惨败》

《认识建筑》', '《造房子》', '《建筑的故事》', '《贝聿铭全集》',
'《建筑也可以很好玩：欧洲篇》', '《建筑家安藤忠雄》', '《安藤忠雄》',
'《手绘紫禁城：遗失在日本的北京皇城建筑艺术》

《建筑师》', '《我还未读懂漫山白雪》', '《穿墙透壁》',
'《生活艺术家的手作私宅》', '《空间的诗学》', '《图像中
国建筑史》

《索拉里斯星》', '《永恒的终结》', '《小镇奇谈》', '《仿生人会
梦见电子羊吗？》', '《你一生的故事》', '《其主之声》', '《球状
闪电》

《中国历代政治得失》', '《明朝那些事儿》',
'《红星照耀中国》', '《枪炮、病菌与钢铁》',
'《巨人的陨落》', '《第三帝国的到来》', '《东
晋门阀政治》', '《毛泽东选集》', '《全球通史》',
'《中世纪之美》', '《制造东京》', '《城市发
展史》', '《神雕侠侣》', '《射雕英雄传》',
'《倚天屠龙记》', '《鹿鼎记》', '《小李飞刀:多
情剑客无情剑》', '《侠客行》', '《逝去的武林:
一代形意拳大师口述历史》', '《武林:一代形意
拳大师口述历史》', '《蜀山剑侠传》', '《笑傲
江湖》', '《献给阿尔吉侬的花束》', '《克莱因
壶》', '《未来学大会》', '《平面国》', '《醉步
男》', '《机器人大师》', '《雪崩》', '《环界1》',
'《鱼翅与花椒》', '《中国人超会吃》', '《雅
舍谈吃》', '《行走的柠檬》', '《四口吃遍江户》',
'《随园食单》', '《料理图鉴》', '《美食与文
明》', '《五味》', '《于谦：人间烟火》', '《肚
子饿万岁》', '《厨艺的常识》', '《鸭川食堂》',
'《这本书好吃吗》', '《人间滋味》', '《至味在
人间》', '《流动的餐桌》', '《川菜》', '《金牌
主厨的面点课》', '《世间味道》



结果分析

K-means	时间复杂度低，算法简单 需要手动输入参数K 由于随机性，算法不稳定
Birch	速度快、对异常值影响不大 对高斯簇效果不好
DPC	对任意形状均可聚类，可通过决策图直观推断出簇个数 dc等参数需要提前输入，决策图的选择具有主观性 对密度差异大的难以区分

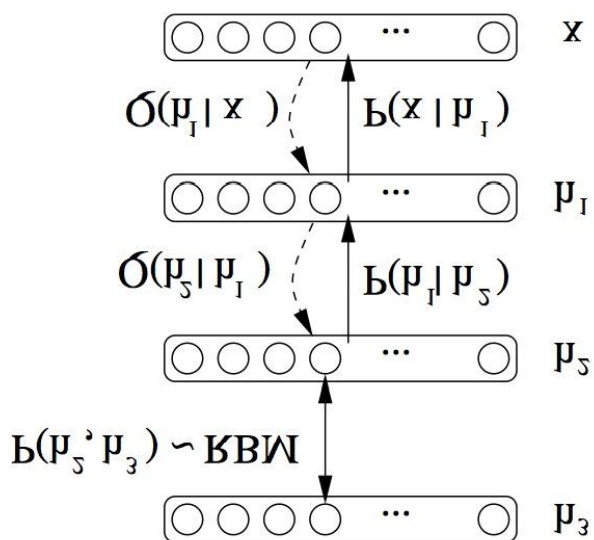
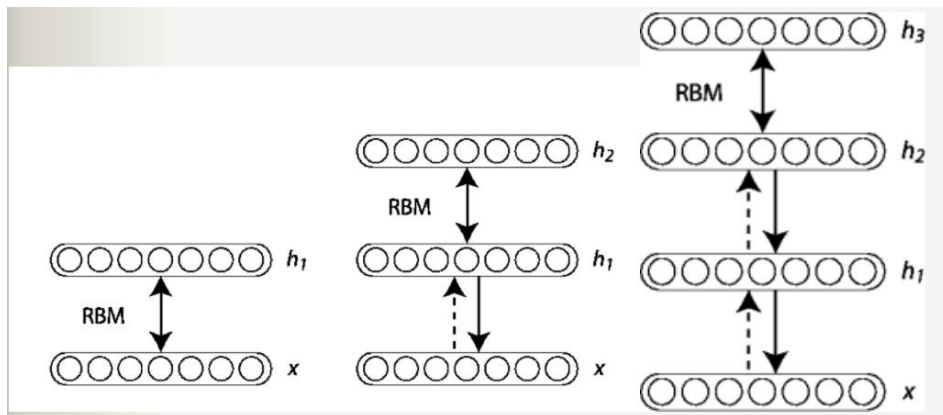


5 最新研究进展



近年来的文本聚类研究与传统的基于密度，基于划分，基于层次，基于网格等聚类算法相比，更多的结合了深度学习，这些方法按照时间顺序可以分为4个大类：

1. Deep Belief Networks
2. Convolutional neural network
3. Recurrent neural network
4. Autoencoder



Deep Belief Network 是Hinton等在2006年提出了一种新的方法来求得这种比较接近最优解的初始权重。

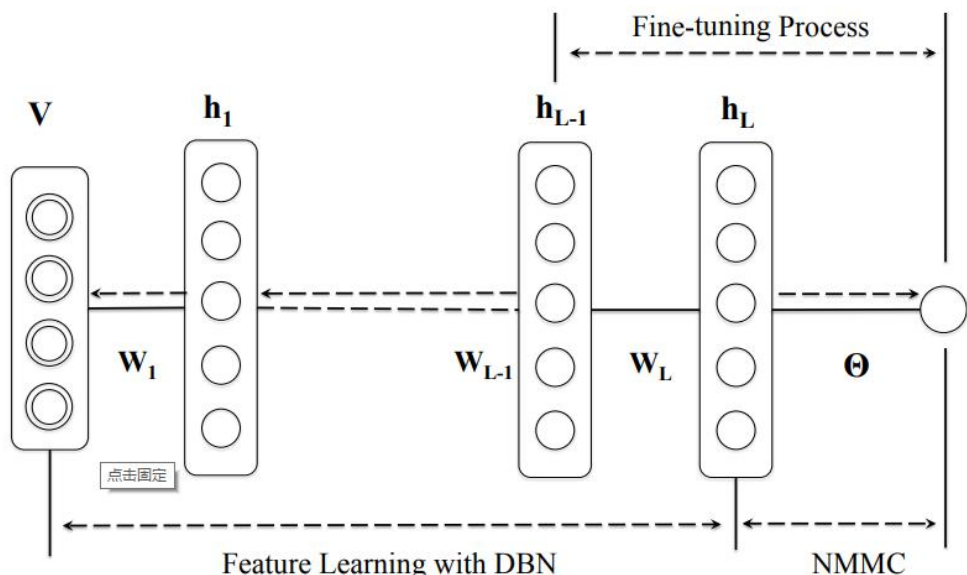
在使用上述的逐层无监督方法学得节点之间的权重以及节点的偏置之后(亦即初始化)。

使用逐层无监督方法来初始化了权重值, 使其比较接近最优值, 解决了之前多层神经网络训练时存在的问题, 能够得到很好的效果。

年份	DBN
2015	Deep learning with nonparametric clustering
2016	A personalized Markov clustering and deep learning approach for Arabic text categorization

Deep Learning with Nonparametric Clustering ——Gang Chen

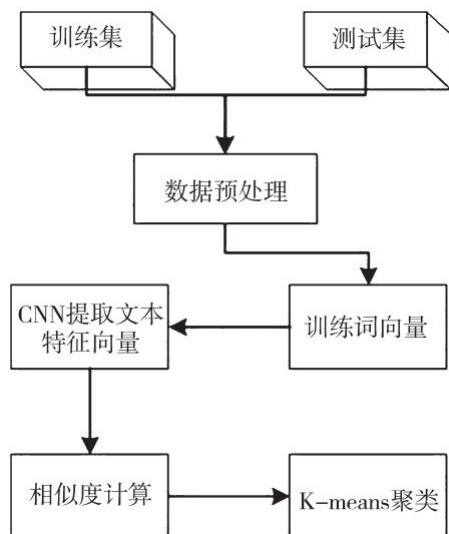
最大边际聚类(maximum margin clustering, MMC)是一种判别方法，将每个实例的标签作为潜在变量，使用支持向量机(SVM)进行大边际聚类。



针对无监督聚类问题的深度学习，提出一种无参数聚类的深度置信网络。利用了深度学习在特征表示和降维方面的优势。然后在最大边距框架下进行无参数聚类，该框架是一种判别聚类模型。最后，在深度置信网络中对模型参数进行了优化。

深度学习如CNN在图像领域的成功应用，给文本聚类也带了新的思考。

Convolutional neural network

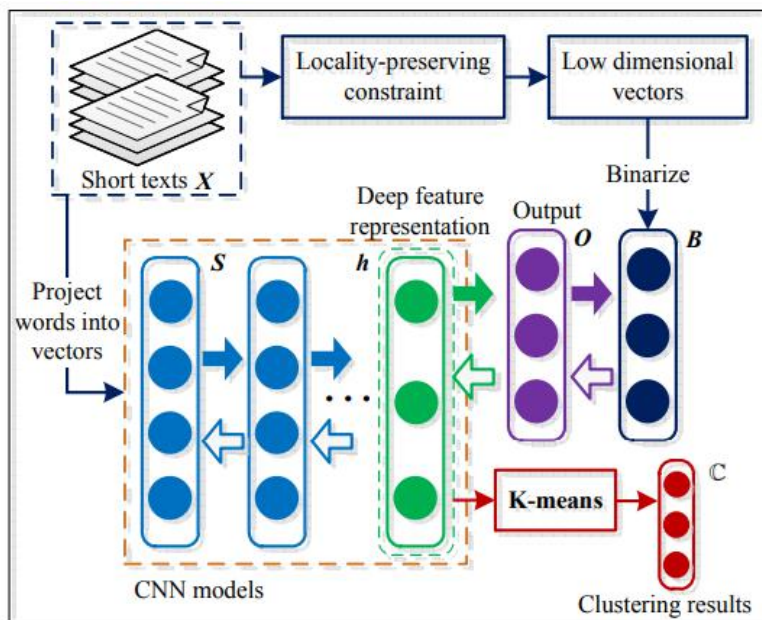


应用CNN的基础方法如左图所示。首先通过word2vec或其他模型来获得语料库中所有词的词向量；然后将文本转换为词向量矩阵，再经过CNN对文本的词向量矩阵进行特征提取，获得文本的特征向量表示；选取合适的相似度计算方法计算文本向量的相似度；最后使用K-means方法进行聚类。

年份	CNN
2015	Short Text Clustering via Convolutional neural networks (STCC)
2016	Semi-supervised clustering for short text via deep representation learning
2017	Self-Taught Convolutional Neural Networks for Short Text Clustering (dubbed STC2)

Short Text Clustering via Convolutional Neural Networks ——Jiaming Xu, Peng Wang

由于短文本的稀疏性，使得聚类成为一个具有挑战性的问题。本文提出了一种基于卷积神经网络的短文本聚类算法(STCC)，通过不使用任何外部标签/标签的自学学习框架，考虑学习特征的一个约束，更有利于聚类。



首先，我们将原始关键字特征嵌入到具有局部保持约束的紧凑二进制码中。然后，研究单词嵌入并将其输入卷积神经网络以学习深度特征表示。在获得学习过的表示之后，我们使用K-means对它们进行聚类。



年份	RNN
2018	Comparison of deep learning based concept representations for biomedical document clustering
2018	Text clustering algorithm based on deep representation learning
2019	Sequential embedding induced text clustering, a non-parametric Bayesian approach
2019	ADC: Advanced document clustering using contextualized representations

Recurrent neural network

应用LSTM或Bi-LSTM等的基础做法和CNN类似，也是用于文本的特征提取和表示。一个包含n个单词的文本的概率表示可以是如下两种形式：

$$p(t_1, t_2, \dots, t_n) = \prod_{k=1}^N p(t_k | t_1, t_2, \dots, t_{k-1})$$

$$p(t_1, t_2, \dots, t_n) = \prod_{k=1}^N p(t_k | t_n, t_{n-1}, \dots, t_{k+1})$$

$$p(t_k | t_1, t_2, \dots, t_{k-1})$$

表示给定前这k-1个单词后，第k个单词为 t_k 的概率。

$$p(t_k | t_n, t_{n-1}, \dots, t_{k+1})$$

类似地，表示给定第k+1到n个单词后，第k个单词为 t_k 的概率。



Recurrent neural network

一个基础的应用LSTM来获取文本的特征的方法为：将包含n个单词的文本输入LSTM，得到每个时间步下的隐状态 h_k ：

$$h_k = [\vec{h}_k, \overleftarrow{h}_k]$$

其中， $\vec{h}_k = f(\vec{h}_{k-1}, x_k)$ ， $\overleftarrow{h}_k = f(\overleftarrow{h}_{k+1}, x_k)$ ，这里用函数 f 来表示LSTM的操作， x_k 表示文本的第 k 个单词。

假设隐状态的维度为 d ，可以采取三种方式计算文本的最终特征表示

- ① Max-pooling 对于 n 个隐状态，选取每个维度上的最大值构成文本特征
- ② Mean-pooling 对于 n 个隐状态，计算每个维度上的均值构成文本特征
- ③ last-time 直接选取第1个和第 n 个隐状态构成文本特征



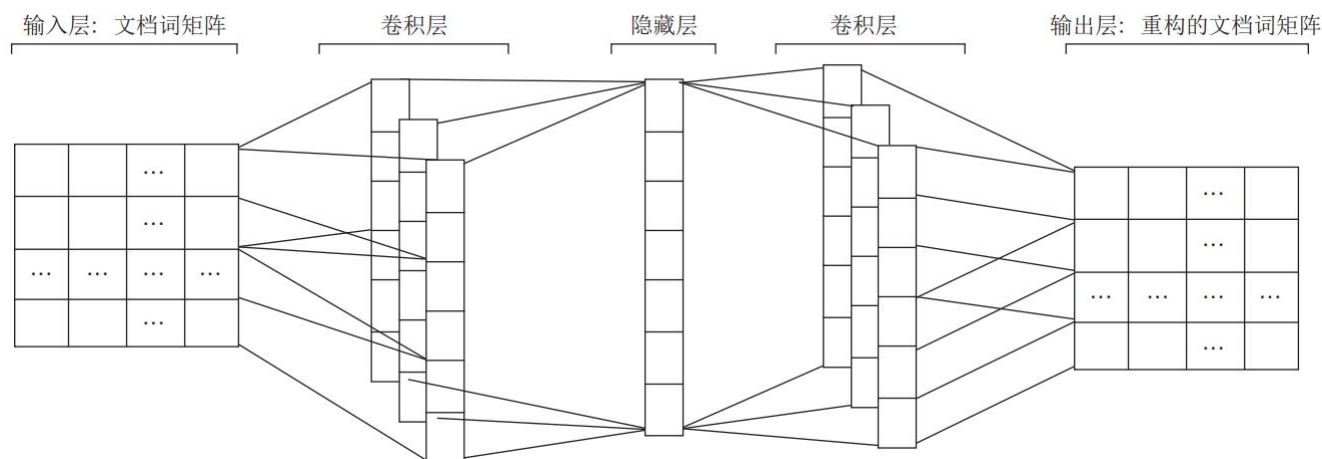
Autoencoder

自编码器被广泛应用于文本聚类前的特征降维。一类方法是先使用自编码器学习从数据域到低维潜在空间的非线性映射，随后在学到的低维空间进行聚类。这些方法将自编码器作为一个预处理阶段，与后续的聚类阶段分开设计。也就是说，特征学习和聚类是按顺序进行的。

当前应用自编码器到文本聚类的一个趋势是希望模型能同时进行特征学习和聚类。这类方法认为，将特征学习和聚类分开设计时，自编码器在低维特征学习过程中不是以促进聚类效果为目标的，学到的数据的潜在低维不一定适合于聚类，因此考虑联合特征降维和聚类。

Autoencoder

一个应用自编码器的基础方法如下图所示，首先通过word2vec获得词向量表示，文本矩阵则是由词向量堆叠而成，通过自编码器学习一个文本矩阵的潜在低维表示，使用k-means或谱聚类等聚类方法对学到的低维向量进行聚类。



年份	Autoencoder
2019	A self-training approach for short text clustering
2020	Patent document clustering with deep embeddings(DEC)
2020	Semi-supervised Network Embedding with Text Information

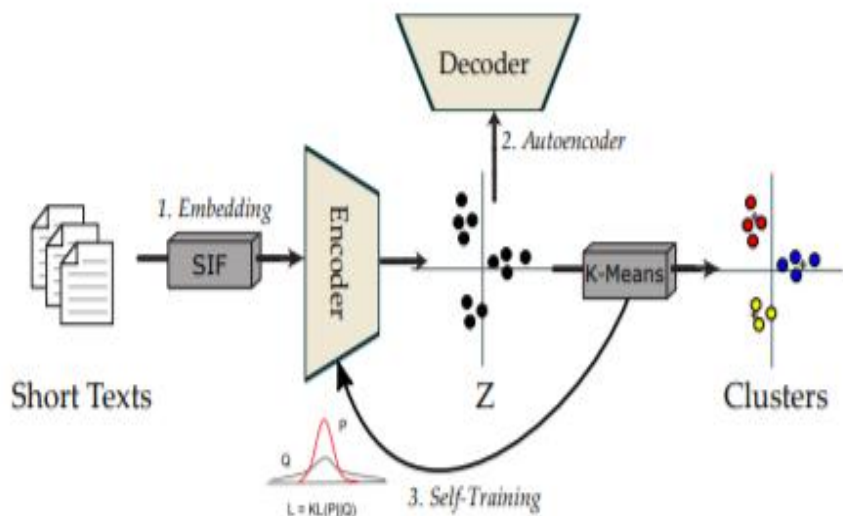
A Self-Training Approach for Short Text Clustering

——Amir Hadifar Lucas Sterckx

短文本聚类是一个具有挑战性的问题。

当采用传统的bag-of-words或TF-IDF表示时，这将导致短文本的稀疏向量表示。

当采用传统的bag-of-words或TF-IDF表示时，这将导致短文本的稀疏向量表示。本文利用深度聚类，提出从自动编码器和句子嵌入中学习鉴别特征，然后利用聚类算法的分配作为监督来更新编码器网络的权值。



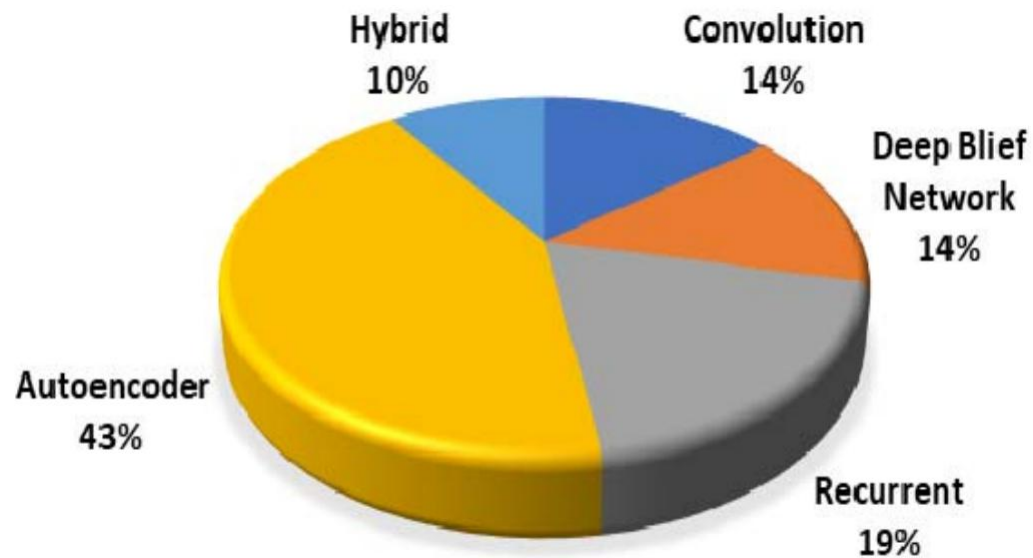


从近5年比较有成效的研究成果来看，自编码器的应用是一个主要趋势

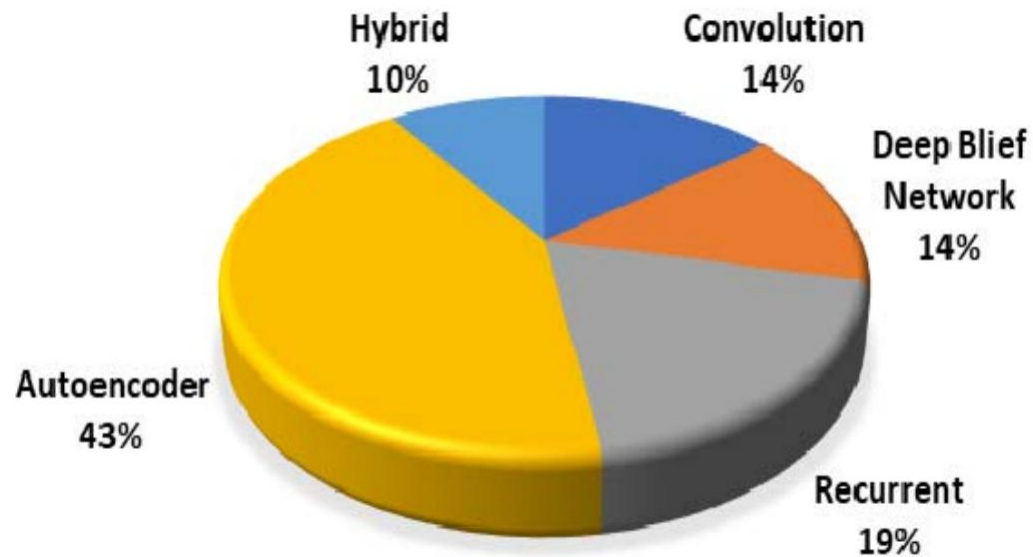
年份	CNN
2015	Short Text Clustering via Convolutional neural networks (STCC)
2016	Semi-supervised clustering for short text via deep representation learning
2017	Self-Taught Convolutional Neural Networks for Short Text Clustering (dubbed STC2)
2019	Deep divergence-based approach to clustering
年份	RNN
2018	Comparison of deep learning based concept representations for biomedical document clustering
2018	Text clustering algorithm based on deep representation learning
2019	Sequential embedding induced text clustering, a non-parametric Bayesian approach
2019	ADC: Advanced document clustering using contextualized representations
年份	DBN
2015	Deep learning with nonparametric clustering
2016	A personalized Markov clustering and deep learning approach for Arabic text categorization

年份	Autoencoder
2016	Unsupervised deep embedding for clustering analysis
2016	Variational deep embedding: An unsupervised and generative approach to clustering(VaDE)
2017	Deep clustering with convolutional autoencoders
2017	Improved deep embedded clustering with local structure preservation
2017	Towards k-means-friendly spaces: Simultaneous deep learning and clustering
2017	Deepcluster: A general clustering framework based on deep learning(DC-GMM)
2018	Spectralnet: Spectral clustering using deep neural networks
2019	Deep co-clustering
2019	A self-training approach for short text clustering
2020	Patent document clustering with deep embeddings(DEC)
2020	Semi-supervised Network Embedding with Text Information

混合和集成深度学习算法：在大多数研究中，混合深度学习算法利用了每种算法的优势，通常会获得比较出色的效果，胜过其他单独的深度学习算法；下图展示了当前深度学习方法的比例，可以看到，混合算法的研究仍数少数，未来有望出现更多新的成果。



填补空白：还有很多深度学习算法没有被应用到文本聚类这个领域，可以尝试将这些算法扩展到文本聚类领域。例如，深度强化学习（Deep Reinforcement Learning）已经在文本情感摘要等方面有所应用，但是还没有任何方法将其用于文本聚类；此外，目前还很少有工作将预训练的深度模型应用到文本聚类上，并对深度学习方法的有效性作解释。



感谢各位专家老师
请您批评指正

德以明理 学以精工