

流式光谱聚类

Shinjae Yoo
计算科学中心
布鲁克海文国家实验室
厄普顿, 纽约州11973-5000
电子邮件: shinjae@gmail.com

郝黄
机器学习实验室
通用电气全球研究部
San Ramon, CA 94583
电子邮件: haohuanghw@gmail.com

Shiva Prasad Kasiviswanathan
三星研究美国
Mountain View, CA 94043
电子邮件: kasivisw@gmail.com

Abstract-Clustering是一种经典的数据挖掘任务发现数据中相似性的相互关联模式。在许多现代领域,数据正在不断变化作为一条小溪。出于可扩展性的原因,将点聚类在数据流中需要设计单通道,有限的内存流聚类算法。但是,表现已知流式聚类算法,通常使用K-在原始特征空间上的手段(或其变体)往往会受到影响当特征空间是高维的。要克服这一点,我们提出了一种流式光谱聚类算法。我们的算法保持归一化的近似值,随着时间的推移拉普拉斯数据流并有效地更新在流媒体中改变这个拉普拉斯算子的特征向量。它只需要一次传递数据,消耗有限内存,并且对数据流的排序是稳定的。我们为流式光谱聚类提供理论分析,算法和我们的实验结果表明,虽然获得了在可扩展性方面,其性能与其他已知产品相当批处理/流式聚类方法。

1. 引言

聚类是一种重要的无监督学习技术,通常它是现代使用的第一步数据分析。任何良好的可扩展性的基本特征聚类算法是处理大量数据的能力高维特征空间中的数据。最现代的高尺寸数据,如文档,图像和多媒体来自网络自然地以流媒体方式到达。怎么样-永远,检测如此大量和高的集群尺寸数据流也是一个具有挑战性的问题,原因如下: 1) 数据流可能是无限的,所以任何尝试存储整个的离线算法用于分析的流最终会耗尽内存, 2) 由于概念漂移,簇随时间动态演化,因此旧集群可能会合并,新集群可能会合并, 3) 对于许多流行的人来说,有一个维度的诅咒, 4) 各种在线应用,在(近)实时获得聚类很重要。

在本文中,我们研究了流式传输中的聚类设置。虽然以前有相当多的文献流聚类算法[47], [42], [30]和大部分这些和其他流式聚类算法可以有效地完成处理大容量数据流,其性能趋于存在高维数据时降级[4]。要克服这个问题,我们提出了光谱的流式自适应聚类算法。光谱聚类获得了巨大的成功,数据挖掘社区在过去十年中的受欢迎程度因为它能够发现嵌入式结构数据(又称歧管)。在其最流行的形式,光谱

聚类算法包括两个步骤:第一,特征向量使用核函数构造的归一化拉普拉斯算子用于嵌入数据集,第二,用于嵌入数据集聚类算法应用于嵌入数据集[40] [32]。

使频谱聚类适应的第一个挑战流设置在于规范化的构建拉普拉斯。构造精确的归一化拉普拉斯算子是本质上是一个批处理操作,因为它需要访问整体数据和存储空间长度为二次方流。因此,传统的拉普拉斯建筑技术无法适应流媒体设置。在本文中,我们专注于两个最流行的内核,余弦和高斯,对于这两个内核,我们提出了流式拉普拉斯算子近似技术。

下一个挑战在于更新频谱嵌入拉普拉斯的有效和高效,以便我们能够近乎实时地处理数据流。为此,我们利用矩阵草图的最新想法¹保持低级别在每个时间步骤逼近整个观察数据,并且这种近似值作为新数据不断更新到达。对于矩阵草图,我们采用最近的算法自由提出的(称为频繁方向)[27]。运用这个低阶近似,在每个时间步,我们对齐学会了光谱嵌入之间具有相同的基础,每两个相邻的时间步,这样就不会敏感概念漂移。鉴于数据流的嵌入,我们然后应用流式K-means算法[41]处理创建和整合集群。

我们提出了我们的算法和理论分析表明在某些现实的假设下,我们构建了流嵌入是嵌入式的一个很好的近似-由昂贵的批处理技术创建的ding。最好的我们的知识,我们的第一个流光谱聚类在真正的数据流设置中运行的算法(即,使用一次通过任何数据样本)。我们建议的流媒体谱聚类算法是有效的

以下方式:
(a) 空间和时间有效,而只需要一个在内存占用有限的情况下传递数据。让m是数据维度。如果在时间t, n点到达数据流,我们的算法只需要O(m l + m n t)空间和草图矩阵所在的O(max{m n t l, m l 2})时间尺寸m x l。在实践中,设置l很多就足够了

¹ 非正式地,矩阵Z的草图是另一个较小的矩阵Z。大小比Z,但仍然接近它[27]。

令 $SVD(K(Z)) = U \Sigma K V^T$ 。然后

$$\frac{\|Z - X_F\|_F}{\|Z\|_F} \leq \frac{\|Z - U^{(k)} \Sigma^{(k)} V^{(k)T}\|_F}{\|Z\|_F} \leq \frac{\sum_{j=k+1}^m \sigma_j^2}{\sum_{j=1}^m \sigma_j^2} \leq \frac{1}{k}.$$

第2页

英语原文: Google
at time t is an m x n t matrix).

与相似更,因此,无论是时间还是空间要求几乎是输入大小的线性(作为输入在时间t是m x n t矩阵)。
(b) 它很容易适应数据流上看不见的模式。在每次执行步骤t,它都会提供更新的频谱嵌入

到流式光谱聚类中，这种列表分配流中的所有数据点。

实证研究表明我们的方法是有效的与空间和时间相比，效率更高各种流行的批量/流式聚类算法来自文本，图像，网络和协作过滤的数据集域。

II. B 背景

在本节中，我们将简要回顾一下基本概念在谱聚类和流K-means方法背后。我们从正式定义我们的符号开始。我们表示 $[n] = 1: n$ 。向量总是以列为单位的方式用粗体字母表示。对于向量 \mathbf{v} ， v 表示它的转置和 $\|\mathbf{v}\|$ 表示其欧几里德范数。对于向量 $(\alpha_1, \dots, \alpha_M) \in \mathbb{R}^M$ ， $\text{DIAG}(\alpha_1, \dots, \alpha_M) \in \mathbb{R}^{M \times M}$ 表示对角矩阵，其中 $1, \dots, m$ 为对角线条目。我来吧表示维数为 $m \times m$ 的单位矩阵。对于矩阵 $Z \in \mathbb{R}^{M \times N}$ ，行总和 $(Z) \in \mathbb{R}^m$ 是具有1进入一个矢量等于Z的第 i 行中的条目总和。我们使用 $z_{i,j}$ 表示Z的第 (i, j) 个元素。定义光谱范数如 $Z = \sup \{ \|\mathbf{Zv}\| : \|\mathbf{v}\| = 1 \}$ 。我们也使用 $\|\cdot\|_p$ 由 Z 表示的范数，其中 $p = 2$ 给出（Frobenius范数） $\|Z\|_F = \sqrt{\sum_{i,j} z_{i,j}^2}$ 和 $p = \infty$ 给出 $\|Z\|_\infty = \max_{i,j} |z_{i,j}|$ 。给定一组矩阵 Z_1, \dots, Z_t ，我们使用符号 $Z^{(t)}$ 表示通过水平连接获得的矩阵 Z_1, \dots, Z_t ，即 $Z^{(t)} = [Z_1, \dots, Z_t]$ 。

我们使用 $S_{VD}(Z)$ 来表示奇异值分解 - Z的分解，即 $S_{VD}(Z) = U\Sigma V^T$ 。这里 U 是 $m \times m$ 正交矩阵， Σ 是 $m \times n$ 对角矩阵， V 是 $n \times n$ 正交矩阵。 Σ 的对角线条目，其中 $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m$ （给定 $m \leq n$ ），被称为Z的奇异值。我们按照惯例列出奇异值以非递增的顺序排列。对于对称矩阵 $S \in \mathbb{R}^{m \times m}$ ，我们用 $E_{IG}(S)$ 来表示它的特征值分解，即 $U\Lambda U^T = E_{IG}(S)$ 。这里 U 是 $m \times m$ 正交矩阵和 Λ 是 $m \times m$ 对角矩阵（真实的）条目被 $\lambda_1, \dots, \lambda_m$ 的称为的本征值 S （再次以非递增顺序列出）。

最佳秩 k 近似（在谱和谱中）Frobenius范数有义）到矩阵 $Z \in \mathbb{R}^{m \times n}$ 是 $Z^{(k)} = \sum_{i=1}^k \sigma_i U_i V_i^T$ ，其中 U_i 和 V_i 是 U 和 V 的列。具有相关的左和右奇异向量 $U_i \in \mathbb{R}^m$ 和 $V_i \in \mathbb{R}^n$ 。我们使用 $S_{VD}^k(Z)$ 来表示 $Z^{(k)}$ 的截断奇异值分解，即， $Z^{(k)} = S_{VD}^k(Z) = U^{(k)} \Sigma^{(k)} V^{(k)}$ 。这里 $\text{DIAG}(\sigma_1, \dots, \sigma_k) \in \mathbb{R}^{K \times K}$ ， $U^{(k)} = [U_1, \dots, U_k] \in \mathbb{R}^{m \times k}$ 和 $V^{(k)} = [V_1, \dots, V_k] \in \mathbb{R}^{n \times k}$ 。以下是众所周知的定理界定了最佳秩- k 的近似误差近似。

定理2.1: [Golub et al. [20]] 设 $Z \in \mathbb{R}^{m \times n}$ $N > m$ ，并且让 $\sigma_1 \geq \dots \geq \sigma_m$ 为Z的奇异值

A. 光谱聚类

算法1: S PECTRAL C LUSTERING [S PECTRAL C LUSTERING] [32]

输入: 输入数据 $Y = [y_1, \dots, y_n] \in \mathbb{R}^{m \times n}$ 和 $k \in \mathbb{R}$ (簇数)
输出: n 个实例的群集分配
1 构造核亲和度矩阵 $W \in \mathbb{R}^{n \times n}$ ，例如，
a) $w_{i,j} \leftarrow \exp(-\frac{\|y_i - y_j\|^2}{2\sigma^2})$ (高斯核)，或
b) $w_{i,j} \leftarrow \frac{y_i \cdot y_j}{\|y_i\| \|y_j\|}$ (余弦核)
2 度矩阵 $D = \text{diag}(d(y_1), \dots, d(y_n))$ 其中 $d(y_i) = \sum_j w_{i,j}$ $\tilde{D} = D - W$
3 归一化拉普拉斯L(符号) $L = D - 1/2 WD - 1/2$
4 $PAP \leftarrow \text{eig}(L)$ (L符号)
5 $V \leftarrow$ 归一化 P ，单位 L 行规范
6 将 V 行聚类成 k 个簇 (使用K-means)

算法S PECTRAL C LUSTERING概述了一个非常受欢迎的

基于嵌入数据集的聚类方法使用核函数，并利用顶部特征向量规范化拉普拉斯算子的发现，以发现潜在的群集TER值。谱聚类与图形有很强的联系分区问题，因为本征空间用于解决放松形式的平衡图分区问题[32]。谱聚类的另一个方面是它可以捕获数据的多种结构，这是很难或不可能的通过K-means风格算法实现。算法S PECTRAL C LUSTERING首先建立亲和力或相似性矩阵（我们显示最流行的两个内核，余弦内核和高斯内核但它不仅限于这两个）。然后我们构建一个拉普拉斯。在算法S P谱 C LUSTERING中，我们提出最流行的对称归一化拉普拉斯算子（ L_{sym} ）。拉普拉斯语的其他流行选择包括非标准化Laplacian $L_{\text{un}} = D - W$ 和随机游走归一化 $L_{\text{rw}} = I - D^{-1}W$ 。一旦我们计算拉普拉斯算子，我们就做了rank- k 特征值分解得到谱嵌入限制在前 k 个特征值。行标准化步骤在一个简单的聚类算法之前使用（例如当K-means）应用于 V 以识别聚类。

B. 流式K-means

在本节中，我们将讨论最近的快速流式Euclidean Shindler等人的K均值聚类算法。[41]（呈现在算法S TRM KM EANS中。由于光谱聚类使用K-在我们提出的方法中，它意味着它的最终聚类步骤在构造的上应用这种流式K-means算法歧管。如果输入数据点的数量是 n 和目标 k 是输出簇，算法维持一组微观表示为C的集群或设施，它们被合并到最后形成 k 簇。算法S TRM KM S TEP，其中形成Algorithm S TRM KM EANS的关键部分，详细说明如何将新数据点 x 添加到当前设施集C。概率 δ/f ，其中 δ 是最小值的平方

算法3: S TRM KM S TEP (重述自[41])
输入: datapoint x ，当前设施集 C ， p 是最大设施数， β 是标量， f 是目前的设施成本
产出: 新设施成本 f ，新设施集 C 和数据点分配

1 测量 $\delta = \min_{y \in C} \|x - y\|^2$ (如果 C 为空则 $\delta = f$)
2 设 r 是0到1之间的均匀随机数
3 如果 $r \leq \delta/f$ 那么
4 设置 $C \leftarrow C \cup \{x\}$
5 其他
6 将 x 分配给 C 中最近的工厂
7 结束
8 而 $C \neq \emptyset$
9 设 $f \leftarrow \beta f$
10 将每个 $z \in C$ 移动到点的质心分配给该群集
11 令 w_z 为分配给 $z \in C$ 的点数
12 $C \leftarrow C$ 的第一个设施
13 对于每个 $z \in C$ 做
14 设 $\delta = \min_{y \in C} \|z - y\|^2$
15 如果发生概率 δ/f 事件则
16 Set $C \leftarrow C \cup \{y\}$
17 其他
18 将 z 分配给其最近的 C 设施
19 结束

算法2: S TRM KM EANS (重述自[41])

输入: 数据流 S ， p 是最大数量设施， k 是簇的数量， β 是标量。
输出: 设施集 C 和数据点分配
1 初始化 $f = 1 / (k (1 + \log n))$ 和空集 C
2 而流 S 没完成呢

x 与 C 和 C 中任何设施之间的欧氏距离 $f = 1 / (k (1 + \log n))$ (记为设施成本)，新设施初始化只有 x 。剩余概率 x 被分配给最近的现有设施。算法试图确保不超过 $p = O(k \log n)$ 设备。如果设施的数量达到 p ，那么算法重新组织（合并）设施以获得更小的设施数量并相应更新设施中心。最后的“球K-means”步骤（这让人联想到执行类似Lloyd的重新聚类）以获得聚类中心。Shindler等人。[41]从理论上分析了算法S TRM KM EANS的性能也提供了对其性能的实验支持。但是，自从算法S TRM KM EANS在输入维度中运行，它会受到高维度的诅咒维度数据集。但是，这不是我们的问题流光谱聚类方法（在下文中介绍）因为Streaming K-means算法在 a 上运行低维流形。

3 从流中读取下一个点
4 在流中更新 $\text{KMS}_{\text{TEP}}(x, C, \rho, \beta, f)$
5 结束
6 在加权点C上运行批次K-means
7 执行球K-means [9] 在结果集群上
获得最终的聚类中心C。

III. STREAMING SPALC LUSTERING ALGORITHM

有两个主要困难需要克服
同时设计光谱聚类算法
流媒体环境：1) 第一个挑战是如何实现
在流上方便地构造归一化拉普拉斯矩阵
(算法SPECTRAL CLUSTERING的前三个步骤)。
拉普拉斯结构本质上是非流媒体任务
因为亲和度矩阵W需要访问整体
数据集，度矩阵D和拉普拉斯算子L。
因此，我们需要新颖的想法来近似拉普拉斯算子
流媒体设置中的矩阵。2) 给定拉普拉斯矩阵，
第二个挑战是构建流式流形
V，这很难，因为概念漂移(主题变化)
在流中导致这些嵌入变化很多
时间。因此，我们需要用于鲁棒流式传输的新技术
嵌入式结构既有效又能够
适应固有的概念 - 漂移。

在介绍我们提出的算法之前，我们先来看看
描述我们的流设置。我们假设数据到了
在流中，每个数据点都有一个指示的时间戳
当它到来时。令 $S = \{Y_{t_n} \in \mathbb{R}^{m \times n} \mid t_n = 1, 2, \dots\}$ 表示
一系列流数据矩阵，其中 Y_{t_n} 代表
数据点在时间 t_n 到达。这里 m 是特征的大小
空间，并且 $n_{t_n} \geq 1$ 是到达时间的数据点的数量

20 结束
22 结束

2 令 $Y_{[t]} = [Y_1, \dots, Y_t] \in \mathbb{R}^{m \times n_{[t]}}$ 表示所有流
数据点到达时间 t 。

在本节的其余部分，我们将介绍和分析
我们的流式光谱聚类方法，在Algo-中概述
rithm SPALC。我们首先描述各种建筑
我们算法的块。

A. 归一化拉普拉斯结构

亲和矩阵。亲和基质的尺寸增长
因为我们不断观察流中的数据，
我们不需要明确地构造亲和矩阵
用于谱聚类。这里我们关注两个流行的内核
亲和力结构。

对于余弦内核的情况，在时间 t ，给定流
 $Y_{[t]}$ 与每列(点)具有单位 L_2 -范数，亲和力
矩阵 $W_{\cos} = Y_{[t]} Y_{[t]}^T$ 。令 $SVD(Y_{[t]}) = U \Sigma C^T V^T$ 。我们可以
在不构造亲和力的情况下得到特征值分解
矩阵因为：

$$\text{EIG}(W_{\cos}) = \text{EIG}(Y_{[t]} Y_{[t]}^T) = V C^2 V^T$$

换句话说，只要我们可以做流式奇异值
分解 $Y_{[t]}$ ，我们可以得到光谱嵌入 $V C$
频谱聚类所需。

2 许多先前的流式算法假设只有一个点到达
每个时间步长，即， $n_{t_n} = 1$ 。通过允许 $n_{t_n} \geq 1$ ，我们允许更
灵活的设置。

第4页

算法4: SPALC

输入：数据流 S ， $l \in \mathbb{R}$ ， $\kappa (\leq l) \in \mathbb{R}$ ， $\rho \in \mathbb{R}$ ，
 $k \in \mathbb{R}$ ， $\beta \in \mathbb{R}$
输出：S中所有实例的群集分配
1 初始化 $f = 1 / (k(1 + \log n))$ 和空集C
2 $s_0 \leftarrow$ 全零向量 $\in \mathbb{R}^m$
3 $B_0 \leftarrow$ 全零矩阵 $\in \mathbb{R}^{m \times l}$
4 $U_0 \leftarrow U_1$ (跳过第一次旋转)
5 当流S没完成的时候
6 令 $Y_{t_n} \in \mathbb{R}^{m \times n_{t_n}}$ 是批次与时间戳 t_n 在
流S。
7 $[B_{t_n}, U_{t_n}(\kappa), V_{t_n}, S_{t_n}] \leftarrow$ 小号 $\text{TRME MB}(Y_{t_n}, B_{t-1},$
8 $V_{t-1}, \text{归一化 } V_{t_n}, \text{与单元中的 } L_2 \text{ 行规范}$
9 $R_{t_n} \leftarrow U_{t_n}^T U_{t-1}(\beta)$
10 用 R_{t_n} 代替每个 $z \in C$
11 对于每列 x in V_{t_n} do
12 $[C, f] \leftarrow \text{SPALC}_{\text{KMS}_{\text{TEP}}}(x, C, \rho, \beta, f)$
13 结束
14 结束
15 在加权点C上运行批量K-means以形成簇
集群

算法5: SPALC MB

输入： $Y_{t_n} \in \mathbb{R}^{m \times n_{t_n}}$ ， $B_{t-1} \in \mathbb{R}^{m \times l}$ ， $S_{t-1} \in \mathbb{R}^m$ 和
 $\kappa \in \mathbb{R}$
输出：乙 t_n ， $U_{t_n}(l)$ ， V_{t_n} ， S_{t_n}
1 $s_t \leftarrow s_{t-1} + \text{行和}(Y_{t_n})$
2 $c_t \leftarrow s_t / s_t$
3 $D_t \leftarrow \text{diag}(\langle y_1, c_t \rangle, \dots, \langle y_{n_{t_n}}, c_t \rangle)$
其中 $Y_t = [y_1, \dots, y_{n_{t_n}}]$
4 $Y_{t_n} \leftarrow Y_{t_n} D^{-1/2}$
5 $c_t \leftarrow [B_{t-1}, Y_{t_n}]$
6 $U_t(l) \leftarrow \Sigma_t(l) \leftarrow V_t(l) \leftarrow SVD(c_t)$
7 $\Sigma_t \leftarrow \text{诊断} \left(\sqrt{\rho_1^2 - \lambda_{\alpha 2}}, \dots, \sqrt{\rho_{l-1}^2 - \lambda_{\alpha 2}}, 0 \right)$
其中 $\Sigma_t(l) = \text{DIAG}(\sigma_{T1}, \dots, \sigma_{Tl})$
8 $B_t \leftarrow U_t(l) \Sigma_t(l)$
9 $U_t(\kappa) \leftarrow [u_1, \dots, u_{\kappa}]$ 其中 $U_t(l) = [u_1, \dots, u_l]$
10 $\Sigma_t(\kappa) \leftarrow \text{DIAG}(\sigma_{T1}, \dots, \sigma_{T\kappa})$
11 $V_t \leftarrow \sum_{i=1}^{\kappa} U_t(i) \Sigma_t(i) \sim Y_t$

对于高斯核的情况，给定流 $Y_{[t]}$ ，
(i, j) 亲和度矩阵的条目 W_{gau} 等于 $\exp(-y_i -$

来自 $[0, \pi]$ 的均匀分布和余弦函数
入门应用。根据[中的分析] 38]，作为
样本数 d 增加，这个随机傅立叶的误差
基数近似值为零。以上投影即可
制定如下步骤：

- 1) 从 $p(\omega)$ 中绘制 d 样本 $\omega(1), \dots, \omega(d)$
 $\omega \sim \mathcal{N}(0, 1)$ 和 $p(\omega)$ 是快速傅立叶变换
在 $[0, \pi]$ 上的分布；
- 2) 从均匀分布中绘制 d 样本 $b(1), \dots, b(d)$
在 $[0, \pi]$ 上的分布；
- 3) 计算投影数据，其中 $h(x) = \cos(\omega x + b)$ ；

由于每个点可以独立于其他点进行投影，
它非常适合流媒体应用。这个想法是要取代
流中的每个点的投影为 $h(y)$ 。设 $H_{[t]} =$
并 $[h(Y_1), \dots, h(Y_{N_{[t]}})] \in \mathbb{R}^{d \times N_{[t]}}$ 的带 $SVD(H_{[t]}) = U \Sigma V^T$ ，

$$\text{EIG}(W_{\text{GAU}}) \approx \text{EIG}(H_{[t]} H_{[t]}^T) = V \Sigma^2 V^T$$

因此，高斯核可以被认为是应用余弦
 H 上的内核(忽略规范化)在下面，我们
假设每个流数据点都已使用转换
上述投影操作，因此，只能集中注意力
关于余弦核的情况。

度矩阵(算法SPALC MB中的步骤1-3)。该
下一个问题是，当我们从流中观察到更多数据时，
对角度矩阵不仅在其大小上增加
但也在其条目的值(因为亲和力矩阵
 W 是非负的)。我们使用聪明的人来克服这个问题
特技。在时间 t ，流 $Y_{[t]}$ 中的数据点 y 的程度
对于余弦内核，可以按如下方式计算：

$$d_t(y) \Sigma_z = y \sum_{z \in Y_{[t]}} z = y s_t \quad (1)$$

其中 $s_t = \sum_{z \in Y_{[t]}} z$ 是行和 $(Y_{[t]})$ 是数据集总和
向量在时间 t 。换句话说，没有构建亲和力
矩阵，只要我们知道，我们就可以计算度矩阵
整个数据集和向量。在流设置中，给定 s_t ，
我们可以计算出 $Y_{[t]}$ 中数据点的程度
流 $Y_{[t]} = [Y_1, \dots, Y_t]$ 。但是，正如我们观察到的那样
更多的数据点， s_t 的范数不断增加，因而在
时间 t ，我们必须重新计算所有数据点的度数
在 t 之前到达。为了克服这个问题，我们建议使用
数据集质心 $c_t = s_t / s_t$ 对于程度近似为
如下。限定，

$$\sim d_t(y) \Sigma_{s_t}^{-1} = y c_t. \quad (2)$$

请注意，在流式设置中，可以合理地假设
数据集质心 $c_t = s_t / s_t$ (对于足够大的 t) 是

$v_i^2/(2\sigma^2)$ 。因此, W_{gau} 的构造要多一点, 因为内核函数不是线性的, 所以很复杂。然而, 我们仍然可以通过首先使用线性化内核来进行近似 [38]。这个想法是明确地将数据映射到欧几里德内部使用随机特征图 h 的产品空间: $R_m \rightarrow R_d$ 这样内核评估可以用变换对之间的欧几里德内积。Concretely, 对于任何两个数据点的 $x, y \in R_m$, 高斯内核可以表示为 $h(x) \cdot h(y)$ 的期望, 其中 $h(x) = \cos(\omega x + b)$, $\omega \in R_m \times d$ 来自适当地缩放高斯分布, $B \in R_d$ 被拉

稳定 (即, c 随时间变化缓慢), 因为大多数话题可以观察到分布。以下引理限制任何两个连续时间之间的质心转换脚步。

引理3.1: 令 $y \in R_m$ 为一个单位矢量。转变了以上定义的时间步长 t 之间的归一化 y 度和 $t+1$ 满足:

$$\sim dt+1(y) - \sim dt(y) / \sqrt{n[t+1] - n[t]}$$

其中 $n[t+1]$ 是在时间 $t+1$ 到达的点数
 $n[t+1]$ 是直到时间 $t+1$ 观察到的总点数。

第5页

证明:

$$\begin{aligned} \sim dt+1(y) - \sim dt(y) &= YC_{t+1} - YC_t \\ &\leq yC_{t+1} - yC_t \\ &= \left\| \begin{matrix} S_{t+1} \\ S_t \end{matrix} \right\| \\ &\leq \left\| \begin{matrix} S_{t+1} - S_t \\ S_t \end{matrix} \right\| \\ &= \sqrt{n[t+1] - n[t]} \end{aligned}$$

在这里我们使用 $y = 1$ 。第二个不等式使用事实上, S_t 中的条目 ≥ 0 并且以入口为主由 S_{t+1} 。

上面的引理表明了一个点的归一化程度一旦数据足够充足, 流中的变化不大已观察到, 通常为 $n[t] \gg n[t]$ 。这表明我们可以设置数据点 y 到达的归一化程度在时间 t 处 $\sim dt(y)$ 并且不在随后的时间调整它脚步。从上面的论证中, 我们知道 $\sim dt(y)$ 仍在继续保持与实际标准化的良好近似即使新数据到达, 流中的 y 度也是如此。我们用这个归一化程度近似在其余部分纸。

我们现在在实践中测试这种近似的质量。图 1 显示了两个不同的程度近似误差

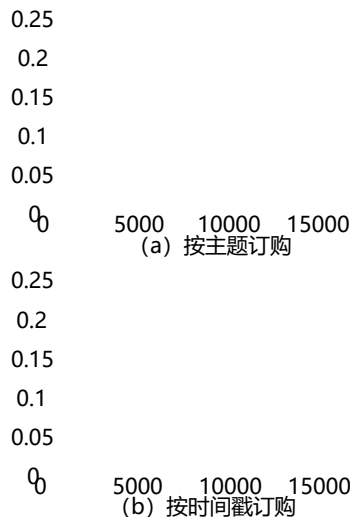


图. 1. 20 上每个文档的度近似误差图新闻组 (20NG) 数据集。Y 轴是两者之间的绝对误差之和, 使用整个数据集近似度和真实归一化度。图 1(a) 显示了主题排序的流数据, 这是最多的难以接近的情况下, 图 1(b) 显示了订购的流数据通过时间戳, 这是更现实的情况。

20 新闻组中的不同流订单 (主题或时间戳) (20NG) 数据集 (有关数据集的详细信息, 请参阅重刑 [4])。batchsize (n_t) 设置为 1000。图 1(a) 显示只需翻转特征向量和特征值的符号。

根据主题所有文档排序的结果属于同一主题 (群集) 的点) 在一起。这是最难的情况是近似因为每当一个新的主题介绍, 它将创造一个重要的概念漂移。在尽管观察了流中明确的主题变化, 程度随着观察越来越多, 近似误差减小文档。最终平均真实标准化程度为 0.1388 但平均近似度数为 0.1441, 平均而言度近似误差为 0.0053 (<4%)。图 1(b) 显示更实际的案例, 即订购文件的情况通过他们真正的时间戳。在这种情况下, 程度近似误差要小得多, 平均为 0.0009 (<2%)。有在流的开头仍然有一些尖峰, 但错误随时间减少。

对于高斯核来说, 通过使用来克服这些问题我们可以使用前面描述的随机傅立叶投影在投影空间上的相同程度近似。

归一化拉普拉斯算子 (算法 S-TRM E-MB 中的步骤 4)。构建亲和力和矩阵和正常程度后 - 我们必须构造的特征向量对称归一化拉普拉斯算子。在这里, 我们声称没有明确地构建拉普拉斯算子, 我们可以确定 L_{sym} 的特征向量。在批处理设置中, 给定亲和力和矩阵 $W_{\text{cos}} = YY$ 和相应的度矩阵 D , 我们可以构造一个对称的归一化图拉普拉斯算子如下:

$$L_{\text{sym}} = I - W_{\text{sym}} = I - D^{-1/2} W_{\text{cos}} D^{-1/2}$$

L_{sym} 和的特征向量之间存在等价性我们在下面的引理中建立了 W_{sym} 。

引理3.2: L_{sym} 的第 i 个最小特征值是 $(1 - \lambda_i)$, 其中 λ_i 为 W 符号的第 i 个最大特征值。

证明: 从特征分解的定义, 我们知道:

$$L_{\text{sym}} v = \lambda v.$$

我们可以将条款重新排列为以下形式:

$$(L_{\text{sym}} - \lambda I) v = 0$$

如果我们用 L_{sym} 替换 W_{sym} , 那么我们有以下内容这两者之间的关系:

$$\begin{aligned} (L_{\text{sym}} - \lambda I) v &= 0 \\ (I - W_{\text{sym}} - \lambda I) v &= 0 \\ (-W_{\text{sym}} + I - \lambda I) v &= 0 \\ (W_{\text{sym}} - I + \lambda I) v &= 0 \\ (W_{\text{sym}} - (1 - \lambda) I) v &= 0 \\ W_{\text{sym}} v &= (1 - \lambda) v \end{aligned}$$

由于 W_{sym} 的特征值总是 ≥ 0 并且因为对于归一化, 特征值从 1 到 0, 即 L_{sym} 的特征值范围从 0 到 1 [3]

我们使用 W_{sym} 的前 k 个有效特征向量是 L_{sym} 的 k 个最不重要的特征向量。仅限因此, 光谱聚类需要特征向量

3 在实践中, 特征值的范围从 -1 到 1。如果发生这种情况, 我们可以只需翻转特征向量和特征值的符号。

第6页

标准化拉普拉斯L symm的谱嵌入，可以很容易从 $YD^{-1/2}$ 上的SVD获得，没有构造拉普拉斯显式。

但是，在流设置中，构建精确度矩阵是不可能的（没有存储所有观察到的），因此，我们使用流程度近似等式的想法（2）。设 $Y[t] = [Y_1, \dots, Y_t]$ 表示原始输入流。而不是 $Y[t]$ ，我们专注于修改流：

$$\sim Y[t] = [\sim Y_1, \dots, \sim Y_t] \text{ 其中 } \sim Y_t \text{ 对于 } t \in \mathbb{N}, D^{-1/2}$$

其中 $\sim D$ 是数据点的度近似值 \hat{y} 。注意that $\sim Y[t]$ 可以从 $Y[t]$ 的以流来获得时间。

B.流媒体流形构造

流形构造（算法中的步骤5-11） $S_{TRM E MB}$ 。现在剩下的主要挑战在于有效地构建流的频谱嵌入。给定 $\sim Y[t]$ ，呈现流式歧管结构在算法 $S_{TRM E MB}$ 步骤5-11中。

我们使用基于矩阵草图的方法进行构建以流式方式嵌入 $\sim Y[t]$ 的频谱。在他的重新分纸[27]，Liberty引入了一种优雅的算法（称为矩阵草图的频繁方向）。频繁的指令-tions算法在流模型中运行并构造一个草图矩阵使用（令人惊讶的）简单的“收缩”概念一些正交向量。在原始的频繁方向算法设置[27]，输入是矩阵 $Z \in \mathbb{R}^{p \times d}$ 。在每一步， Z 行的一行由算法处理，并且算法迭代地更新矩阵 $Q \in \mathbb{R}^{l \times d}$ ($l \ll p$)。对于任何单元向量 $x \in \mathbb{R}^d$ ， $z \cdot x^2 - q \cdot x^2 \leq 2Z^2$ 这里参数 l 调整了大小之间的权衡草图矩阵和近似误差界限（较大的 l 增加了计算和存储的要求算法，同时给出更好的结果）。最近，Ghashami和飞利浦[18]，重新分析了频繁方向算法表明它为低秩矩阵提供了相对误差界限近似。

我们构建谱嵌入的方法流 $\sim Y[t]$ （在Algo-步骤5-9之间概述）rithm $S_{TRM E MB}$ ）基于扩展频繁指令-tions算法[27]在更广泛的环境中。每一步，我们添加 $n_t \gg 1$ 个新列。4和Frequent一样方向，我们的算法只需要传递一次数据流。在算法 $S_{TRM E MB}$ 中， B_t 是维护的草图stream $Y[t]$ 的流程，并在每个时间步骤更新为新的数据到了。参数 l （如前所述）定义草图矩阵的大小和参数 k ($\leq l$) 定义嵌入数据的维度。我们讨论设置之后这两个参数。

嵌入（重新）对齐（算法中的步骤9和10） $S_{TRM SC}$ ）。虽然我们已经解决了这些挑战构建embedding $\sim V$ 吨对于每个流批次中Algo-rithm $S_{TRM E MB}$ ，在时间 t 的低秩空间基础， $\sim U_t$ (k)，趋向

4 基于类似草图的低秩矩阵近似方法是最近在完全不同的特征选择环境中使用[23]。

由于概念漂移而随时间变化。因此，我们不能只需使用发现embedding $\sim V$ 吨集群。每时每刻步骤有必要重新调整过去的设施或微集群质心到新的基础（算法 $S_{TRM SC}$ 中的步骤9-10）。为此，我们构造对齐矩阵 R_t ，如下所示：

$$R_t = \sim U_{t-1}^{-1} \sim U_t^{(k)}.$$

然后我们先前通过 $R_t z$ 旋转找到每个设施 z 。

C.计算复杂性

在任何时间 t ，Algo-的运行时间rithm $S_{TRM E MB}$ 是 $O(\max\{mn_t l, ml\})$ ，使用电源-SVD的迭代或等级显示QR分解[20]。空间复杂度为 $O(mn_t + ml)$ 。对于批量光谱聚类算法，在时间 t 使用所有数据直到时间 t ，空间复杂度为 $O(mn_t)$ ，时间复杂度为 $O(mn_t k)$ 。一个人注意到算法 $S_{TRM SC}$ 很多因为 n_t 非常增长，所以比它的批处理对应物更有效迅速。

D.理论分析

我们现在证明Algorithm $S_{TRM E MB}$ 的有效性，通过在合理的假设下显示光谱由算法 $S_{TRM E MB}$ 构造的嵌入接近利用整个数据流构建的嵌入。由于空间限制，我们在此省略了详细的证明。

在算法 $S_{TRM E MB}$ 的时间 t ，考虑秩- k 近似 $C_t = [B_{t-1}, \sim Y_t]$ 。让

$$\text{小号 } VD_k(C_t) = \sim U_t^{(k)} \Sigma_t^{(k)} \sim V_t^{(k)}$$

算法 $S_{TRM E MB}$ 构造一个embedding $\sim V_t$ 吨 of $\sim Y_t$ （原样总结 $\sim \Sigma_t^{-1}(k)$ 存在）定义为，

$$V_t := \sim \Sigma_t^{-1}(k) \sim U_t^{(k)} \sim Y_t.$$

现在 B_{t-1} 是 $\sim Y_{[t-1]}$ 的草图。考虑秩- k 近似 $\sim Y_t = [\sim Y_{[t-1]}, \sim Y_t]$ 。让

$$\text{小号 } VD_k(\sim Y_t) = U_t^{(k)} \Sigma_t^{(k)} V_t^{(k)}$$

我们比较算法所构建的embedding $\sim V_t$ 使用实际流 $\sim Y_t$ 构造的嵌入 V_t ，其中 V_t 定义为（假设 $\Sigma_t^{-1}(k)$ 存在）

$$V_t = \Sigma_t^{(k)} U_t^{(k)} Y_t^{(k)}.$$

我们claim $\sim V$ 牛逼是第 V 的良好近似牛逼。

$$\begin{aligned} V_t - \sim V_{五六} &= V_t - \sim V_{五六} \\ &= \Sigma_t^{-1}(k) U_t^{(k)} Y_t^{(k)} - \sim \Sigma_t^{-1}(k) \sim U_t^{(k)} \sim Y_t^{(k)} \\ &\leq \Sigma_t^{-1}(k) U_t^{(k)} \left(Y_t^{(k)} - \sim Y_t^{(k)} \right) \end{aligned} \quad (3)$$

第7页

让我们专注于限制 $\Sigma_t^{-1}(k) U_t^{(k)} \left(Y_t^{(k)} - \sim Y_t^{(k)} \right) U_t^{(k)} F_{五六}(k)$ 命题3.4:

$$\begin{aligned} &\Sigma_t^{-1}(k) U_t^{(k)} \left(Y_t^{(k)} - \sim Y_t^{(k)} \right) U_t^{(k)} F_{五六}(k) \\ &= U_t^{(k)} \left(Y_t^{(k)} - \sim Y_t^{(k)} \right) U_t^{(k)} F_{五六}(k) \\ &= U_t^{(k)} \left(Y_t^{(k)} - \sim Y_t^{(k)} \right) U_t^{(k)} F_{五六}(k) \\ &\leq U_t^{(k)} \left(Y_t^{(k)} - \sim Y_t^{(k)} \right) U_t^{(k)} F_{五六}(k) \\ &\leq U_t^{(k)} \left(Y_t^{(k)} - \sim Y_t^{(k)} \right) U_t^{(k)} F_{五六}(k) \\ &\leq U_t^{(k)} \left(Y_t^{(k)} - \sim Y_t^{(k)} \right) U_t^{(k)} F_{五六}(k). \end{aligned} \quad (4)$$

第二个不等式使用 Σ 的事实 $\Sigma_t^{-1}(k)$ 和 $\sim \Sigma_t^{-1}(k)$ 是对角矩阵， $\sim U_t$ (k)的列是正交的。对于最后的不等式，我们使用了自 Σ 以来的事实 $\Sigma_t^{-1}(k)$ 是一个对角矩阵，它的Frobenius范数至多是正方形其维度的根乘以最大对角线值。

因此，一个结合开 $V_t - \sim V_{五六}$ 从各自边界如下关于 $U_t^{(k)} - \sim U_t^{(k)}$ 和 $\Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k)$ 。对于前者，

$$\begin{aligned} &\Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k) \\ &\leq \Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k) \\ &\leq \Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k) \end{aligned} \quad (6)$$

证明：请注意 $\Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k)$ 可以绑定为

$$\Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k) \leq \Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k) + \Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k).$$

通过在Algorithm $S_{TRM E MB}$ 中构建，

$$\begin{aligned} &\Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k) \\ &= \Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k) \\ &= \Sigma_t^{-1}(k) - \sim \Sigma_t^{-1}(k). \end{aligned}$$

和自适应索引结构，用于维持流和新婚夫妇。阿克曼等人。[1]构造了一个小加权通过使用称为的新数据结构来获取数据流的样本核心树。Shindler等人。[41]最近提出了一个有效的流的高效算法 (Algorithm S_{TRM} KM EANS) 使用在线设施位置文献中的想法进行聚类。他们的算法与现有算法相比都是有利的理论和实验上。最近, Liberty等人。[28]提出了一种在线K-means聚类算法空间和时间要求只是多对数的流的长度。所有上述技术都在运行整个功能空间。但是, 当数据集很高时维度, 这些技术依赖于邻域搜索面部稳定性问题并且可能表现不佳, 这是通常被称为维度的诅咒。Compa-我们的方法使用低维嵌入派生来自数据流的奇异值分解。那里-它具有类似的稳定性和性能优势

设置, 据我们所知, 还没有人知道这里的泛化是对包含的数据流进行操作的高维数据点。因此, 我们不包括他们在我们的实验比较中。

高效的静态聚类。我们提出的算法主要 - 将一组集群中心候选者作为新的候选者更新并更新数据流来了。已经有许多有效的方法建议[39], [46], [44], [6], [48], [29] 试图恢复高维数据的稀疏聚类中心。但是全部它们需要访问整个数据集。因此不适合用于流式设置。因此, 我们也不包括它们在我们的以下实验比较中。

V. EXPERIMENTAL ANALYSIS

在本节中, 我们通过实验测试我们提出的ap-在效率和效率方面有所作为。所有的经验 -

第9页

表I. 小号 TATISTICS的实验数据集。

数据集	#instance	#特征	#clusters
1 扇形	9619	55197	105
2 ohscal	11162	10493	10
3 20NG的	3707	29476	20
4 RCV1	193844	47236	103
5 MNIST	70000	784	10
6 小	10000	3072	75062
7 SenVec (SensIT车辆)	9528	100	3
8 滑稽演员	73421	100	86

在英特尔 (R) Xeon (R) CPU X5650 2.67GHz上运行处理器, 128GB内存。

A. 实验设置

数据集。我们在四个文本数据集上进行了实验 (扇区, ohscal, 20Newsgroup和RCV1), 两个图像数据集 (MNIST和Tiny), 一个传感器网络数据集 (SenVec) 一个协同过滤评级数据集 (jester), 所有摘要 - 表 I 列出了。对于前四个文本数据集, 我们使用余弦核亲和矩阵, 和高斯核亲和矩阵用于剩下的四个。四个文本数据集是L₂为每个文档标准化。

行业数据集由公司网页分类组成在工业部门的层次结构中。ohscal数据集是来自包含文件的OHSUMED集合来自医疗类别。20NG (20Newsgroup) 是一个平衡的涵盖20个新闻主题的数据集。RCV1数据集包含制作手动分类的新闻专题报道的档案可由路透社[26]。我们使用了RCV1数据集的一个子集, 只选择那些只有一个标签的文件。MNIST数据集有10类手写数字和784图像功能。所有上述数据集都可以在[10]或[16]。Tiny是非参数对象的大型Web图像集合和场景识别 (从[下载 45])。从80年代开始百万个小数据集图像, 我们创建了一个百万的数据集图像由60,000个标记图像组成, 覆盖100个课程 (来自[25]和其他随机抽样的图像。该此处的评估仅限于标记的60,000幅图像。我们使用了Tiny图像以及所有3072个原始特征, 这是3色通道 (RGB) 中的32×32彩色图像, 用图像来演示图像的聚类性能高维特征向量。SenVec (SensIT车辆) 数据集包含车辆的分布式传感器网络数据分类。Jester数据集包含匿名评级数据用来推荐笑话。

基线。我们测试了我们提出的算法的两个版本: SSC-1, 它使用Algorithm S_{TRM} E MB来构造流的低维嵌入, 但我们执行收集低维后的最终K均值聚类嵌入整个流; 和SSC-2, 这是一样的算法S_{TRM} SC (流K-means聚类是用于每一步)。SSC-1的结果证明了这一点光谱嵌入的质量, 而SSC-2的结果表示完整的流谱聚类。

选择五种其他算法进行比较: 经典算法K-means (简称KM), 经典谱聚类算法 S (简称NJW), BIRCH [49], 流K-

5 对于谱聚类, 我们使用算法S PECTRAL C LUSTERING。

手段 6 (简称SKM) 和HDDStr [35]。首先两种算法是最广泛使用的批处理解决方案集群。BIRCH构建了一个分层的数据结构递增地传入数据流。讨论SKM早期是最近提出的快速准确的流媒体K-意味着聚类方法。HDDStr是一种有效的预测高维数据的流聚类方法。

评估指标和参数设置。我们用两个流行的评估指标: 规范化的互信息 (NMI) 和纯度。为了显示每种算法的稳定性, 我们还报告了这两个指标的标准差每个算法和每个数据集。

假设数据集簇k的数量为N_{cl}为1000 (N的稳定性试验将在第显示VD) 数据集谱嵌入的维数m如第 III-D部分所述。算法S_{TRM} SC (K) 被设置为k的倍数。的大小矩阵草图设置为m如第 III-D部分所述。对于使用高斯内核的随机特征映射, 我们设置带宽比例为5000, 投影尺寸 (d) 为2000。微集群 (设施) 的最大数量p在算法S中, _{TRM} KM EANS设置为k log n, 其中n为数据样本的数量 (假设是先验已知的如[41])。对于BIRCH, SKM和HDDStr, 我们遵循建议在[49], [41], [35]各自的参数设置。对于所有比较算法, 我们执行了WCSS K均值 (最小化群内平方和, 用30个内环和30个外环) 以获得稳定的簇分配。

B. 一般业绩比较

表 II 显示了比较的一般性能算法。对于每个数据集, 我们随机调整数据流顺序30次并报告平均NMI (上) 和纯度 (底部)。我们可以从中得出以下观察结果这些结果:

- (1) 与NJW相比, 我们提出的SSC-1达到了92%NMI, 平均纯度超过99%。这个结论我们的流光谱嵌入近似的公司接近计算上昂贵的高质量批量版本。SSC-1具有非常相似的结果稳定性 (使用标准偏差的测量) 与NJW一样, 平均而言, 这比KM要好得多。
- (2) SSC-1优于那些运行的算法高处的完整特征空间 (KM, BIRCH和SKM) 维数据集 (四个文本数据集和Tiny图像数据集)。
- (3) 毫不奇怪, 所有流式算法都能执行比批次K-means (KM) 差。BIRCH, SKM和HDDStr平均只有30%的NMI和60%纯度与KM相比。但我们提出的SSC-2, 其中在真正的流媒体设置中运行, 平均达到58% NMI和97%纯度的KM, 这是更有效的比其他流媒体竞争对手 (在NMI中≈两倍) 和1.6倍的纯度)。这证明了有效性SSC-2比其他流媒体竞争对手 (BIRCH, SKM, 和HDDStr)。

6 对于流K-means, 我们使用Algorithm S_{TRM} KM EANS。

第10页

(4) 流媒体算法也更敏感
比批处理算法 (NJW, KM等) 不稳定
更改数据顺序。但是, 通过利用低
维度嵌入, 我们提出的SSC-2算法
实现与KM相当的稳定性。和...相比
无论是BIRCH, SKM还是HDDStr, SSC-2都是平均水平
在NMI中稳定三倍以上, 大约十倍
纯度更稳定。

C. 订单敏感度

众所周知, 流式算法对此很敏感
数据的排序或概念漂移。测试性能
在这种情况下我们的方法, 我们按数据流排序
以下两个附加条件, 时间和主题顺序。
时间顺序反映了现实情景, 而
主题顺序是模拟最困难的条件
由于大概概念漂移而进行度近似
集群之间。对于后者, 除了分组点
属于一个集群, 我们也尽力保证
类似的主题彼此相邻 (例如, 在20NG数据集中,
comp.sys.ibm.pc.hardware主题中的文档集
和comp.sys.mac.hardware主题放在每个主题旁边
其他)。

图2显示了按时间顺序排列的稳定性结果
输入流顺序。对于扇区数据集, 只有少数几个集群
因此, 在一开始就出现了较早的表现
比后者好。总的来说, 所有三个数据集都显示出来
表现相当稳定。

图3显示了按主题排序的输入流。外加
有条不紊地, 通过分类的主题, SSC结果是均匀的
高于表2的平均表现。我们推测
虽然有大量的主题变化, 但数据流
草图可能更好地代表了因为
属于同一主题的文件在一起。这个
导致更丰富的光谱嵌入, 可以捕获更多
来自流的信息。一般来说, 两个数字都表明了
对于我们提出的SSC-2结果非常稳定
1和SSC-2算法, 它们优于其他三种算法
流媒体算法, 包括时间顺序和主题数据
排序。

D. 批量大小灵敏度

不同批次尺寸 (n_t) 的SSC-1和SSC-2测试是
如图4所示。实验设置与那些相同
在表11中, 批量大小在500,800,1000,1200,1500,
1800和2000。性能表明稳定的行为
SSC-1和SSC-2具有不同的批量大小。

E. 可扩展性比较

图5显示了之间的可伸缩性比较
Tiny数据集上的各种聚类算法 (
采样和下采样)。所有的流媒体算法
只需要几分钟即可集群超过一百万
数据点。特别是, 我们观察到了我们的两个版本
SSC与其他高效流媒体相比毫不逊色
算法虽然具有更好的效果 (表2)。COM的
与KM和NJW相比, 我们提出的算法表明
卓越的可扩展性, 同时具有可比性。

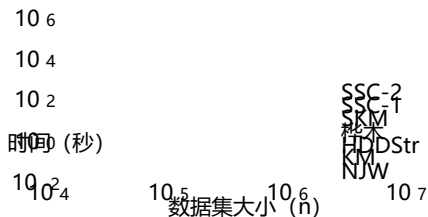


图5. 可扩展性实验。KM和NJW不会扩展到大型数据集
(由于其极高的内存开销而失败)。

VI. C 结论

我们提出了一种流式光谱的新方法
集群。我们的算法构建了一个简洁的近似值
拉普拉斯算子, 并从矩阵草图到适应思想
有效地构建了一个低维谱嵌入
流。我们的算法只需要一次通过
数据, 利用有限的存储, 很好地适应概念漂移
并近乎实时地运作。实验结果表明
我们提出的方法可以胜过其有效性
高维流行的流媒体聚类算法
数据流同时实现良好的可扩展性。

REFERENCES

- [1] M. R Ackermann, M. Märtens, C. Raupach, K. Swierkot, C. Lam-
梅森和C. 索勒. Streamkm ++: 数据的聚类算法
流. ACM JEA, 2012.
- [2] CC Aggarwal. 流聚类算法调查, 2013.
- [3] CC Aggarwal, J. Han, J. Wang和PS Yu. 一个框架
聚集不断发展的数据流。在VLDB, 2003年。
- [4] CC Aggarwal, J. Han, J. Wang和PS Yu. 一个框架
投影的高维数据流聚类。在VLDB, 2004年。
- [5] CC Aggarwal和PS Yu. 用于群集大量文本的框架
和分类数据流。在SIAM SDM, 2006年。
- [6] M. Azizyan, A. Singh和L. Wasserman. Minimax理论的高 -
具有稀疏平均分离的尺寸高斯混合物。在NIPS中,
2013.
- [7] D. 巴巴拉. 集群数据流的要求。ACM SIGKDD
探索通讯, 2002年。
- [8] C. Böhm, K. 栏杆, HP Kriegel和P. Kroger. 密度连接
使用本地子空间首选项进行聚类。在IEEE ICDM, 2004年。
- [9] Vladimir Braverman, Adam Meyerson, Rafail Ostrovsky, Alan Royt-
男子, Michael Shindler和Brian Tagiku. 良好的流式k-means
可群集数据。在SODA, 第26-40页, 2011年。
- [10] D. 蔡. 用于特征学习的Matlab代码和数据集。 [HTTP://www.cad.zju.edu.cn/home/dengcai/Data/data.html](http://www.cad.zju.edu.cn/home/dengcai/Data/data.html).
- [11] F. Cao, M. Ester, W. Qian和A. Zhou. 基于密度的聚类
带有噪音的不断发展的数据流。在SIAM SDM, 2006年。
- [12] M. Charikar, L. O'Callaghan和R. Panigrahy. 更好的流媒体
聚类问题的算法。在ACM STOC, 2003年。
- [13] Y. Chen和L. Tu. 用于实时流数据的基于密度的聚类。
在ACM SIGKDD, 2007年。
- [14] Y. Chi, X. Song, D. Zhou, K. Hino和BL Tseng. 发展的
通过结合时间平滑度进行谱聚类。在ACM
SIGKDD, 2007.
- [15] M. Ester, HP Kriegel, J. Sander和X. Xu. 基于密度的算法
用于在具有噪音的大空间数据库中发现聚类。在ACM
SIGKDD, 1996.
- [16] RE Fan和CJ Lin. Libsvm数据: 分类, 回归和
多标签。 <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/数据集/>.

第11页

表二. P ERFORMANCE比较 NMI (TOP) 及 P URITY (BOTTOM). † HE NUMBERS中括号里的标准偏差
分类指标. w ^ e COULDN ' 牛逼运行NJW ON THE RCv1及牛逼 INY数据集DUE TO OUT MEMORY运行. † HE号码NJW ON
THE RCv1 DATASET从 [43].

	NJW	SSC-1	KM	桦木	SKM	HDDStr	SSC-2
扇形	0.4039 (0.0794)	(0.00833961) (0.00008790)	(0.003354 (0.003376)	(0.00872 (0.008372)	(0.00921 (0.00623)	(0.00821 (0.00617)	(0.00361 (0.00316)
ohscal	0.3408 (0.0897)	(0.0183217 (0.0048833)	(0.01641 (0.00299)	(0.00876 (0.00902)	(0.00860 (0.007405)	(0.00781 (0.005893)	(0.003273 (0.0018564)
20ng的	0.4836 (0.0643)	(0.0084699 (0.00008532)	(0.00886 (0.00376)	(0.00942 (0.00793)	(0.003601 (0.00402)	(0.04282 (0.007611)	(0.00812 (0.00315)
RCV1	[43] 0.2875 (0.0203)	(0.00833 (0.00432)	(0.00376 (0.00721)	(0.00872 (0.00611)	(0.00921 (0.006274)	(0.00821 (0.00557)	(0.00361 (0.0017557)
MNIST	0.4965 (0.0868)	(0.0103792 (0.00008788)	(0.00473 (0.00809)	(0.00272 (0.00492)	(0.00389 (0.00301)	(0.00281 (0.00288)	(0.00253 (0.0012)
滑稽演员	0.2099 (0.0444)	(0.00008422 (0.00002810)	(0.00446 (0.00831)	(0.00438 (0.00689)	(0.00911 (0.00975)	(0.006857 (0.00654)	(0.00996 (0.00221)
SenVe	0.3007 (0.06769)	(0.00006921 (0.0001744)	(0.00376 (0.00744)	(0.00872 (0.00817)	(0.00921 (0.00321)	(0.00821 (0.00588)	(0.00361 (0.001203)
小	- (-)	0.9569 (0.00008790)	(0.003354 (0.003376)	(0.00872 (0.008372)	(0.00921 (0.00623)	(0.00821 (0.00617)	(0.00361 (0.00316)

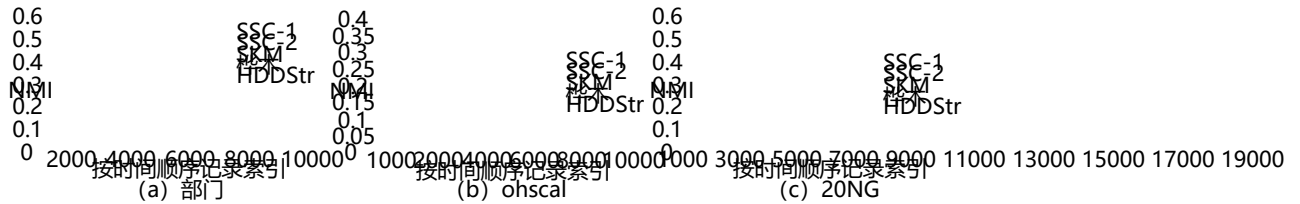


图2.概念随时间的漂移。我们的SSC-1和SSC-2都胜过其他三种流媒体方法。

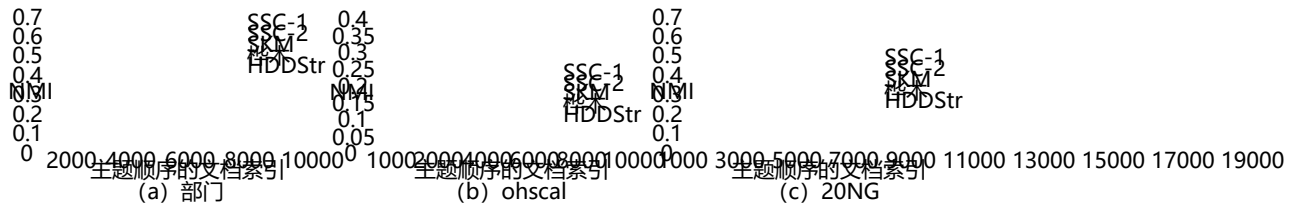


图3.跨排序主题的概念漂移。我们的SSC-1和SSC-2都胜过其他三种流媒体方法。

- [17] MM Gaber, A. Zaslavsky和S. Krishnaswamy. 挖掘数据流中的Liberty. 简单而确定的矩阵草图。在ACM回顾。ACM Sigmod Record, 2005.
- [18] M. Ghashami和JM Phillips. 确定性低的相对误差 [28] E. Liberty, R. Sriharsha和M. Sviridenko. 一种在线算法秩矩阵近似。在SODA, 第707-717页, 2014年.
- [19] GH Golub, M. Heath和G. Wahba. 广义交叉验证 [29] J. Liu, C. Wang, M. Danilevsky和J. Han. 大规模光谱选择良好曲线参数的方法。技术计量学, 21 (2) : 245. 在图上的聚类。在第一十三届国际会议论文集集中, 人工智能联合会议, 第1486-1492页. AAAI, 2013年.
- [20] GH Golub和C. F Van Loan. 矩阵计算, 第3卷. JHU 按, 2012.
- [21] S. Guha, A. Meyerson, N. Mishra, R. Motwani和选择. Calzadas, 2013年.
- [22] S. Guha, N. Mishra, R. Motwani和L. O'Callaghan. 聚类数据流: 理论与实践. IEEE TKDE, 2003.
- [23] H. Huang, S. Yoo和S. Kasiviswanathan. 无监督功能 [31] O. Nasraoui和C. Rojas. 强大的聚类, 用于跟踪噪声演变流。在计算机科学基础. IEEE, 2000.
- [24] P. Kranen, I. Assent, C. Baldauf和T. Seidl. clustree. 索引 [32] A. Ng, M. Jordan和Y. Weiss. 关于谱聚类: 分析与分析选择数据流。在CIKM, 2015年.
- [25] A. Krizhevsky, V. Nair和G. Hinton. cifar-100数据集 [33] H. Ning, W. Xu, Y. Chi, Y. Gong和TS Huang. 增量光谱集群应用于监控不断发展的博客社区。在SIAM SDM, 2007年.
- [26] DD Lewis, Y. Yang, TG Rose和F. Li. Rcv1: 一个新的 [34] H. Ning, W. Xu, Y. Chi, Y. Gong和TS Huang. 增量光谱文本分类研究的集合. ACM SIGKDD探索通讯, 5 (1) : 113-127, 2010.
- [35] I. Ntoutsi, A. Zimek, T. Palpanas, P. Kröger和HP Kriegel. 密度-

第12页

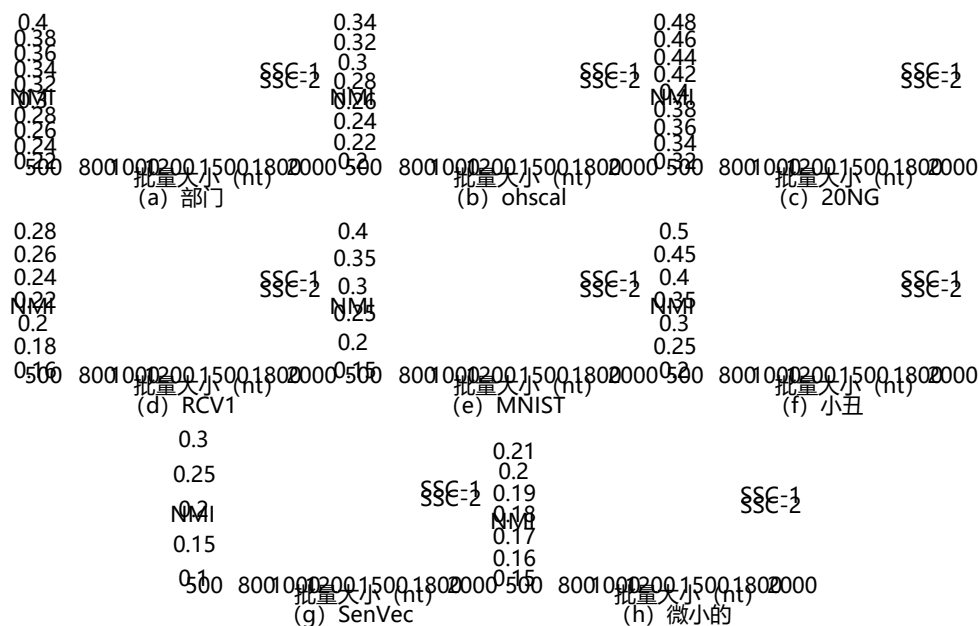


图4.不同批量大小的方法的稳定性。SSC-1和SSC-2都可以通过改变批量大小来表现出稳定的性能。

- 基于高维数据流的预测聚类。在SIAM SDM, 2012.
- [36] L. O'callaghan, A. Meyerson, R. Motwani, N. [49] T. Zhang, R. Ramakrishnan和M. Livny. Birch: 一个有效的数据库非常大型数据库的聚类方法。在ACM SIGMOD, 1996年.
- [37] ME Orlowska, X. Sun和X. Li. 可以独占聚类 [50] Shuang Gu. 高效的流式文本聚类。神经网络, 18 (5) : 790-798, 2005.
- 流数据实现? ACM SIGKDD探索通讯,

- 2006年。
- [38] A. Rahimi和B. Recht。大规模内核的随机特性机器。NIPS, 2007年。
- [39] D. Sculley。网络规模的k均值聚类。在ACM WWW, 2010年。
- [40] J. Shi和J. Malik。标准化剪切和图像分割。 IEEE TPAMI, 22 (8) : 888-905,2000。
- [41] M. Shindler, A. Wong和AW Meyerson。快速准确的k-means对于大型数据集。在NIPS, 2011年。
- [42] Q. Song, J. Ni和G. Wang。基于快速聚类的特征子集高维数据的选择算法。IEEE TKDE, 2013年。
- [43] Y. Song, W. Chen, H. Bai, C. Jin和EY Chang。平行光谱集群。数据库中的机器学习和知识发现, 2008年。
- [44] W. Sun, J. Wang, Y. Fang, et al。正规化的k均值聚类维数据及其渐近一致性。电子期刊统计, 2012年。
- [45] A. Torralba, R. Fergus和WT Freeman。8000万张小图片：用于非参数对象和场景识别的大型数据集。 IEEE TPAMI, 2008年。
- [46] J. Wang, J. Wang, Q. Ke, G. Zeng和S. Li。快速近似k均值通过群集闭包。在IEEE CVPR, 2012年。
- [47] H. Yang, MR Lyu和J. King。多任务的高效在线学习功能选择。ACM TKDD, 2013年。
- [48] J. Yi, L. Zhang, J. Wang, R. Jin和A. Jain。单通道算法有效地恢复高维数据的稀疏聚类中心。在ICML, 2014年。