

# 在监督降维中控制失真以保留类和邻居

伯努瓦科兰热

大学 Grenoble Alpes, INES,  
F-73375, Le Bourget du Lac, France  
benoit.colange@cea.fr

雅科·佩尔托宁

坦佩雷大学信息技术学院  
和通信科学, 芬兰  
jaakko.peltonen@tuni.fi

迈克尔·奥珀蒂

卡塔尔计算研究所, 哈马德本哈利法大  
学, 多哈, 卡塔尔  
maupetit@hbku.edu.qa

丹尼斯杜特赫

大学 格勒诺布尔阿尔卑斯大学 萨瓦勃朗峰,  
CNRS, LAMA, 73000 尚贝里, 法国  
denys.dutykh@univ-smb.fr

西尔万·莱斯皮纳茨

大学 Grenoble Alpes, INES,  
F-73375, Le Bourget du Lac 法国  
sylvain.lespinats@cea.fr

## 抽象的

高维数据的非线性降维具有挑战性, 因为低维嵌入必然包含失真, 并且很难确定哪些失真是最需要避免的。当可以将数据注释到已知的相关类中时, 它可用于指导嵌入以避免恶化类分离的扭曲。本文介绍的**监督映射方法**, 称为类NeRV, 提出了一个原始的压力函数, 该函数将类注释考虑在内, 并根据错误的邻居和遗漏的邻居来评估嵌入质量。类NeRV共享从随机邻域嵌入 (Stochastic Neighbor Embedding) 新能源)。我们的方法比以前的方法有一个关键优势: 在文献中, **监督方法通常以扭曲数据邻居的结构为代价强调类分离**; 相反, 无**监督方法以经常混合类为代价提供更好的结构保存**。实验表明类NeRV可以同时保留邻居结构和类分离, 优于九个最先进的替代方案。

## 1 简介

降维 (DR) 方法旨在将高维数据集映射为低维嵌入空间中的点, 同时保留数据点之间的一些相似性度量。可以通过利用类信息来监督 DR。因此, 监督方法从数据的相对位置 (非监督方法也使用) 和类标签计算映射。DR 技术 [1、2、3] 可用作分类或聚类的预处理步骤, 或将 (标记的) 数据可视化为散点图。在映射标记数据时, 有两个相互矛盾的目标:

- 分类是典型的受**监督的**DR 技术: 在嵌入空间中用分类精度强调和测量类分离。

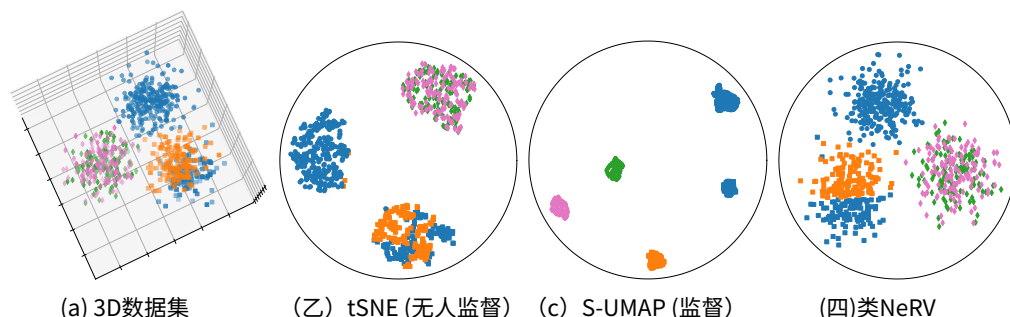


图1: 类NeRV旨在保留类和邻居: 数据从三个 3D 高斯簇 (a)、纯蓝色类的圆簇、被分隔蓝色和橙色类的平面减半的正方形簇和菱形簇中采样随机分布的绿色和粉色类。显示了这些数据的不同平面嵌入: 无监督神经元

(b) 很好地保留了聚类, 但由于忽略了标签, 因此与正方形聚类中的橙色和蓝色类重叠。被监督的S-UMAP (c) 分裂正方形和菱形的簇, 迫使类分离。因此, 它误导了橙色和蓝色类的原始空间邻接, 以及绿色和粉色类的混合。类神经元 (d) 旨在更好地保留三簇结构以及类的邻接性。

- 探索性数据分析是典型的无监督在不知道类信息的情况下运行的 DR 技术: 数据邻域结构被优先考虑并测量为原始空间和嵌入空间中数据相似性之间的差异。

这些目标源自视觉分析任务 [3、4]。它们是矛盾的, 除非类和数据邻域结构在数据和嵌入空间中彼此匹配良好: 每个类构成不同的区域, 没有跨类邻域关系。不幸的是, 这种理想情况不太可能发生, 因为数据邻域结构和类在数据空间中并不总是匹配, 并且高维数据的低维嵌入不可避免地会出现失真 [3]: 假邻居它们是嵌入中的相邻点但不在数据中, 并且想念邻居它们在数据中是邻居但在嵌入中不是。

在这项工作中, 我们建议类NeRV, 一种受监督的 DR 技术, 用于在考虑类信息的同时完成探索性分析目标。我们的解决方案在原则上与之前的类似分类地图 [5] 基于距离的投影, 但我们从与神经病毒 [6, 7], 神经元 [8], 和tSNE [9]。我们的解决方案明显不同于其他以前的监督方法, 这些方法倾向于以牺牲邻居的结构为代价强制类分离, 例如S-等值图[10], S-神经病毒 [11] 和S-UMAP [12]。图 1 说明了基本特征类 NeRV。

我们的贡献是双重的: 我们建议类NeRV当将高维标记数据嵌入低维空间时, 它利用类信息来确保更好地保存类。它的压力函数, 来自无监督神经病毒 [6、7] 引导优化, 以便将不可避免的邻域结构扭曲放置在对类结构危害较小的位置。通过强调惩罚类之间的虚假邻居和类内遗漏的邻居来避免有害的扭曲。我们还推导出两个新的类感知质量指标从标准 诚信度和 连续性质量指标 [13], 专门说明影响类保存的扭曲。

## 2 相关工作

无监督嵌入。之前已经提出了许多线性或非线性算法包括主成分分析 (主成分分析) [14], 自组织图 (索姆) [15]、等距特征映射 (等轴测图) [16]、数据驱动的高维缩放 (DD-HDS) [17], 局部仿射多维项目 (灯) [18] 和均匀流形近似和投影 (UMAP) [12]。在各种各样的技术中, Neighborhood Embedding (东北) 技术在保留邻域结构和计算时间方面是有效的。他们的概率框架还提供了一个理论背景来解释所获得的地图

邻域检索任务[7]。NE 方法为每对点计算我,  $j$  点的概率隶属度/到点附近我, 有时称为相似性。这些隶属度是在数据空间和嵌入空间中计算的。通过最小化这两个空间之间的隶属概率差异来获得映射。这些方法包括随机邻域嵌入 (SNE) [8], t-分布SNE (t-SNE) [9], 詹森香农嵌入 (JSE) [19] 和邻域检索可视化工具 (神经病毒) [6, 7]。神经网络和t-SNE 不同之处在于用于计算它们在嵌入空间中的邻域隶属度的内核。JSE和神经病毒都延伸神经网络来控制虚假邻居和遗漏邻居之间的平衡。这种可调性神经病毒和JSE使它们最适合引入监督。

监督嵌入。监督映射方法主要关注阶级分离, 在映射之前或期间修改数据邻居结构。他们中的许多人通常会降低类之间的数据相似性并增加类内的数据相似性, 然后在修改后的相似性上使用标准的无监督映射方法。它们可以使用度量学习对差异 (距离) 进行操作, 如监督局部线性嵌入 SLLE 及其变体 [20, 21, 22], 监督等值图 (S-等值图) [10, 23], 监督 NeRV (S-神经病毒) [11] 或监督 UMAP S-UMAP [12], 或与重尾半监督版本的相似之处tSNE [24]。其他方法依赖于全局参数映射并优化其参数以最大化类分离, 作为线性判别分析 (低密度脂蛋白) [25] 及其核化变体 [26、27、28], 邻域成分分析 (NCA) [29] 及其神经网络变体[30]或有限秩矩阵学习向量量化 (LiRaM LVQ) [31]。类感知tSNE (catSNE) [[32] 根据类的分布在本地调整要保留的邻域大小。终于, 分类地图 [5] 优化类似于局部多维缩放的应力函数 (LMDS)

[33], 但支持标记数据的探索性分析通过在惩罚错误和遗漏的邻居时考虑班级。然而, 两者 LMDS和分类地图基于距离的对高维度的范数集中现象很敏感[34], 而网元像这样的技巧神经病毒使用减轻这种影响移不变会度[35]。类NeRV得出与相同的原理分类地图来自神经病毒压力函数来获得两者的好处。

### 3 ClassNeRV 和 Class-aware 质量指标

#### 3.1 NeRV 压力函数

在神经病毒 [6, 7], 一个点的隶属度 (条件概率)  $j$  到附近另一点我, 表示  $\beta_{ij}$  在数据空间和  $b_{ij}$  在嵌入空间中, 定义为:

$$\beta_{ij, \Sigma} = \frac{\exp(-\Delta_{ij}^2 / 2\sigma_{\Sigma}^2)}{\sum_{k \in \text{邻居}(i)} \exp(-\Delta_{ik}^2 / 2\sigma_{\Sigma}^2)} \quad \text{和} \quad b_{ij, \Sigma} = \frac{\exp(-J_{ij}^2 / 2\sigma_{\Sigma}^2)}{\sum_{k \in \text{邻居}(i)} \exp(-J_{ik}^2 / 2\sigma_{\Sigma}^2)} \quad (1)$$

和  $\Delta_{ij}$  和  $J_{ij}$  点之间的距离  $i$  和  $j$  分别在数据和嵌入空间中。

隶属度的分布表示为  $\beta_{\Sigma}, \{b_{ij}\}_{j \in \text{邻居}(i)}$

= 我 和  $b_{\Sigma}, \{b_{ij}\}_{j \in \text{邻居}(i)}$  = 我. 这

隶属度是平移不变的, 减少了范数集中度 [35]。秤参数  $\sigma_{\Sigma}$

设置为获得固定的用户选择的困惑  $p$  类似于光滑或模糊的我  $\Sigma$

确定数量

每个点的邻居[36]:  $H(\beta_{\Sigma}) = -\sum p_{ij} \log p_{ij}$ , -

$j \in \text{邻居}(i)$  和  $\beta_{ij}$  日志  $\beta_{ij}$ . 这里,

我们设置嵌入尺度参数  $\sigma_{\Sigma}$  等于  $\sigma_{\Sigma}$ .

这神经病毒压力函数是两组 Kullback-Leibler (KL) 散度之间的线性权衡:

$$\sum_i \tau D_{\text{KL}}(\beta_{\Sigma}, \beta_{\Sigma}) + (1-\tau) D_{\text{KL}}(b_{\Sigma}, b_{\Sigma}), \tau \sum_{i,j \in \text{邻居}(i)} \beta_{ij} \log \frac{\beta_{ij} + (1-\tau)}{b_{ij}} + \sum_{i,j \in \text{邻居}(i)} b_{ij} \log \frac{b_{ij}}{\beta_{ij}} \quad (2)$$

在等式 (2) 中,  $\tau \in [0, 1]$  控制  $\beta_{\Sigma}$  和  $b_{\Sigma}$  之间的权衡  $\sum_i D_{\text{KL}}(\beta_{\Sigma}, \beta_{\Sigma})$  惩罚想念邻居-布尔斯基最大的时候  $\tau=1$  和  $\sum_i D_{\text{KL}}(b_{\Sigma}, b_{\Sigma})$  惩罚假邻居(最大的时候  $\tau=0$ ). 什么时候  $\tau=1$ , 神经病毒减少到神经元。

#### 3.2 ClassNeRV 压力函数

我们的目标是保留邻域结构和类别的嵌入。我们建议嵌入不应该在同一类中人为地分离点, 或者人工聚类

来自不同类别的点在一起. 因此, 我们需要惩罚更多类内错过的邻居和类间假邻居分别。为了控制这些基于类别的扭曲, 我们进一步将等式 (2) 中的发散项分解为类内和类间关系。它给

两个额外的不同的基于类别的权衡参数  $\tau_{in}$  和  $\tau_{out}$  (同时都在  $[0,1]$ ) 并导致类NeRV应力函数:

$$\zeta_{\text{类NeRV}}, \quad \tau_{in} \sum_{i,j \in S_{\text{我}}} \left( \beta_{ij} \log \frac{\beta_{ij}}{b_{ij}} + b_{ij} \beta_{ij} \right) + (1-\tau_{in}) \sum_{i,j \in S_{\text{我}}} \left( \beta_{ij} \log \frac{\beta_{ij}}{b_{ij}} + \beta_{ij} - b_{ij} \right) + \tau_{out} \sum_{i,j \in S_{\text{我}}} \left( \beta_{ij} \log \frac{\beta_{ij}}{b_{ij}} + b_{ij} \beta_{ij} \right) + (1-\tau_{out}) \sum_{i,j \in S_{\text{我}}} \left( \beta_{ij} \log \frac{\beta_{ij}}{b_{ij}} + \beta_{ij} - b_{ij} \right) \quad (3)$$

与类内的集合  $S_{\text{我}} = \{j \in \text{大号} \mid \text{大号} = \text{大号}_j \text{ 且 } j \neq \text{我}\}$  和类间  $S_{\text{我}}^c = \{j \in \text{大号} \mid \text{大号} \neq \text{大号}_j \text{ 且 } j \neq \text{我}\}$  基于类标签定义的索引  $\text{大号}_j$  我点数我。

请注意, Bregman (B) 散度替换方程 (2) 的 KL 散度以确保应力函数的正性。实际上, KL 散度仅针对概率分布定义 (求和

到1), 而在四个方面  $\zeta_{\text{类NeRV}}$ , 隶属度限制在集合中  $S_{\text{我}} \in [0,1]$  和  $S_{\text{我}}^c \in [0,1]$  小于1. 因此, KL 散度不能直接适用于此类隶属度 (否则它不会满足散度的已知属性, 即非负性和不可区分的身份), 而 Bregman 散度则直接适用。Bregman 散度项

$\beta_{ij} \log \frac{\beta_{ij}}{b_{ij}}$  和  $b_{ij} \beta_{ij}$  惩罚类内分别错过了和错误的邻居, 而  $\beta_{ij} - b_{ij}$  和  $b_{ij} - \beta_{ij}$  惩罚类间分别是 missed 和 false neighbors。

参数  $\tau_{in}$  和  $\tau_{out}$  定义类NeRV通过加权这些术语来映射行为。  $\tau_{in}$  控制惩罚假邻居和遗漏邻居的平衡类内, 尽管  $\tau_{out}$  控制类似的平衡类间. 因此, 类NeRV受到监督如果  $\tau_{in} > \tau_{out}$ , 则它的压力函数比其他函数惩罚更多的类内遗漏邻居和类间假邻居扭曲。越大  $\tau_{in} - \tau_{out}$  差值越大, 监管水平越高。为了  $\tau_{in} = 1$  和  $\tau_{out} = 0$ 。有  $\tau_{in} < \tau_{out}$  但是会偏爱类内错过的邻居, 并且类之间的假邻居, 鼓励同类分裂和不同类重叠, 搞砸了类保存。类NeRV不受监督如果  $\tau_{in} = \tau_{out}$  然后减少到原来的神经病毒, Bregman 类内和类间的组合散度等于相应的 KL 散度, 其中  $b_{ij} - \beta_{ij}$  条款取消。

我们重新参数化  $\tau_{in}$  和  $\tau_{out}$  作为  $\tau^* = (\tau_{in} + \tau_{out})/2$  和  $\epsilon = (\tau_{in} - \tau_{out})/2$ .  $\tau^* \in [0,1]$  控制平均值虚假邻居和遗漏邻居的惩罚之间的权衡 (作为  $\tau$  在神经病毒), 尽管  $\epsilon \in [0,0.5]$  控制监管水平 (更多的监督以获得更大的价值)。相反转换是:  $\tau_{in} = \tau^* + \epsilon$  和  $\tau_{out} = \tau^* - \epsilon$ 。

鉴于  $\tau^*$  和  $\epsilon$ , 这类NeRV嵌入是通过最小化应力获得的  $\zeta_{\text{类NeRV}}$  在等式 (3) 中关于嵌入点的坐标。优化是使用多尺度优化方法 [37] 和拟牛顿法进行的BFGS算法[38]。当前的实现复杂度为  $\mathcal{O}(n^2)$  (在哪里  $n$  是数据点的数量), 但是基于树的加速技术 [39、40] 可以将其减少到  $\mathcal{O}(n \log n)$ 。

### 3.3 探索性分析监督技术的质量指标

为了评估邻域结构的保存情况, 我们采用了可信度和连续性度量 [13], 这是在探索性分析 [3] 的背景下评估无监督嵌入的标准。这两项措施分别量化虚假邻居的平均水平 (可信度吨) 和错过的邻居 (连续性C) 对于给定的邻域大小  $k$  作为:

$$\text{吨}(k), 1 - \frac{1}{\text{吨最大限}(k)} \sum_{i,j \in F_{\text{我}}(k)} (\rho_{ij} - k) \quad \text{和} \quad C(k), 1 - \frac{1}{C_{\text{最大限}}(k)} \sum_{i,j \in M_{\text{我}}(k)} (r_{ij} - k), \quad (4)$$

在哪里  $\rho_{ij}$  和  $r_{ij}$  是每个点的行列  $j$  在每个点的附近我分别在数据和嵌入空间中, 以及  $F_{\text{我}}(k)$  和  $M_{\text{我}}(k)$  是 false 和

错过了点的邻居我。归一化系数吨<sub>最大限度(κ)</sub>和C<sub>最大限度(κ)</sub>被定义为吨和C范围0对于理论上最差的映射到1一个理想的映射。基于组合分析，吨<sub>最大限度(κ)</sub>和C<sub>最大限度(κ)</sub>等于κ否(2ñ-3个κ-1)/2个如果κ<否/2个 或者否(N-κ)(N-κ-1)/2个如果κ>否/2 [41]。

为了评估类保存，我们还推导出了两个新的衡量标准：仅限于-  
阶级关系吨<sub>∈</sub> 和连续性仅限于类内关系C<sub>∈</sub>。这些类别意识指标被定义为：

$$\text{吨}_{\in}(\kappa), 1 - \frac{1 \text{ 个}}{\text{吨}_{\text{最大限度}(\kappa)}} \sum_{i,j \in F_{\text{我}}(\kappa) \cap S_{\in} \text{ 我}} (\rho_{ij} - \kappa) \quad \text{和} \quad C_{\in}(\kappa), 1 - \frac{1 \text{ 个}}{C_{\text{最大限度}(\kappa)}} \sum_{i,j \in M_{\text{我}}(\kappa) \cap S_{\in} \text{ 我}} (r_{ij} - \kappa). \quad (5)$$

请注意，这些类别感知指标仅解释了他们所考虑的部分扭曲

无监督的同行，因此吨<sub>∈</sub> > 吨和C<sub>∈</sub> > C。因此，他们可能会达到1个只要有剩余的扭曲不影响类间可信度和类内连续性。

监督降维的另一个标准质量指标是留一法的准确性k-最近的邻居 (k-NN) 嵌入空间中的分类器 [3]。当嵌入用于执行分类时，该指标是有意义的，但它仅关注类并且未能考虑邻域结构保留，如图 3 所示。补充材料提供了广泛的结果k-NN 增益 [32] 来自该指标。

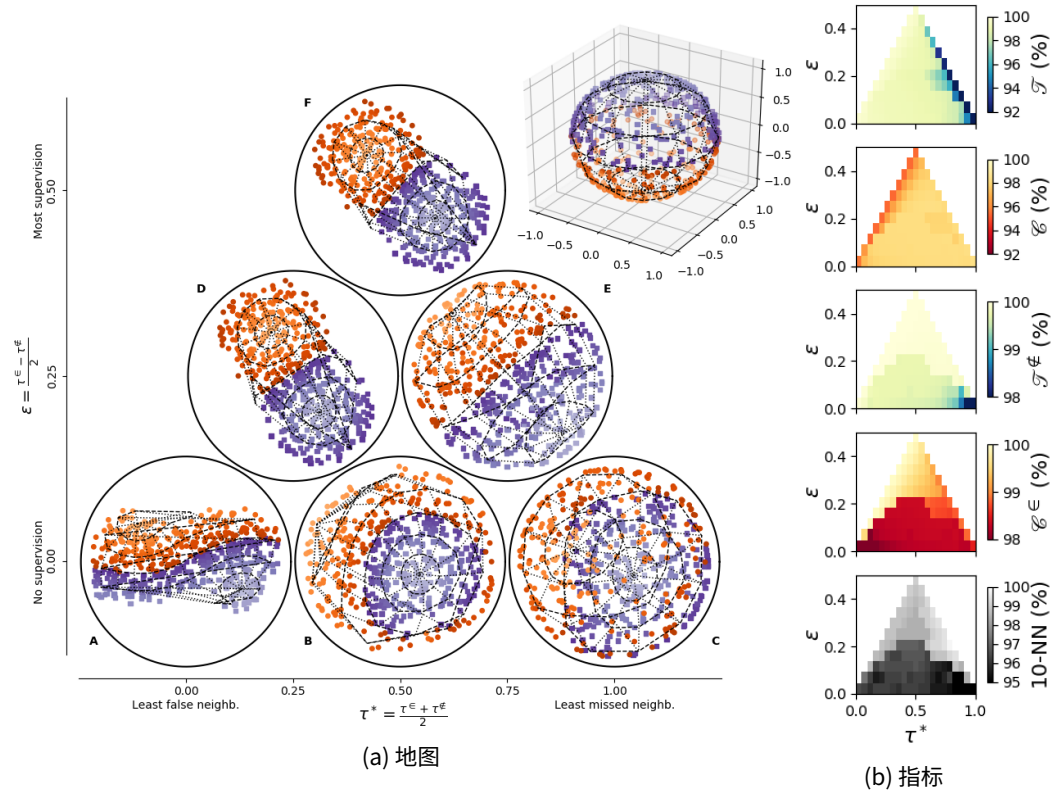


图 2：灵敏度类NeRV权衡参数，同时嵌入5123D数据地球中的数据集中的数据集（第 4.2 节）{τ\*, ε}-参数空间。左图 (a) 显示了五个二维嵌入，其中平行线（虚线）和子午线（虚线）位于事后的，以及数据集的 3D 表示（右上角）。右图 (b) 显示了使用第 3.3 节中定义的指标进行的量化评估，从上到下为 Trustworthiness

吨，连续性C，类间可信度吨<sub>∈</sub> 和类内连续性C<sub>∈</sub>对于给定的邻域大小κ=32,和10-神经网络分类准确率（99%对于原始 3D 数据）。热图的每个像素对应于地球的一个嵌入。τ\*控制假邻居之间的整体权衡（很可能与τ\*=1）和错过的邻居（很可能是τ\*=0），尽管ε控制监督水平。

## 4 实验

### 4.1 目标、数据、技术

我们在 3D 玩具数据集 (*地球*) 和两个真实的高维数据集 (*胰岛5*和 *位数*). 这 *地球*数据集 (第 4.2 节) 包含 512 数据随机分布在三维欧氏空间的单位球面上  $R^3$  (图 2a)。这两个类 (蓝色和红色点) 对应于在赤道处划分的两个半球。这些数据不能在没有失真的情况下嵌入到平面中, 因此最终地图取决于邻域之间的权衡集 ( $\tau^*$ ) 和类 ( $\epsilon$ ) 惩罚 (参见第 3.2 节)。这 *胰岛5*数据集 [42, 43] (第 4.3 节) 包含 1 559 均匀分布的口语字母的英语发音录音 26 类, 描述为 617 特征。这 *位数*数据集 [44, 43] (第 4.4 节) 包含 3 823 手写数字的图像。这 8 个  $\times$  8 个像素图像 (64D 数据点) 被分开 10 类 (每个数字一个)。真实的类别标签以及随机生成的标签被认为是评估错误标记的鲁棒性。的随机子集 500 样本被认为是为了简化图 6 中地图的可读性。

我们比较类 NeRV 无人监督主成分分析 [14], 等值图 [16], UMAP [12], tSNE [9], 和 神经病毒 [6, 7], 并进行监督美国国家航空航天局 [29], S-等值图 [10], 分类地图 [5] 和 S-UMAP [12]。的实施 PCA、Isomap、NCA、tSNE、UMAP 和 S-UMAP 来自 scikit-学习 (版本 0.22.1) [45] 和 umap 学习 (版本 0.3.10) [46] Python 库。S-Isomap、ClassiMap、NeRV 和类 NeRV 使用我们自己的实现 [47]。所有现成的算法都设置了默认参数, 除了神经元初始化, 为此我们使用主成分分析而不是随机的 (以便所有方法都受益于谱初始化)。神经病毒和类 NeRV 使用第 3.2 节中描述的优化。我们的多尺度优化中的最终困惑设置为  $p=32$  为了 *地球*和  $p=30$  为了 *伊索莱特*, 为了与神经元默认  $p=30$ 。我们对所有技术和数据集使用欧氏距离。定性结果按照标准指南 [3, 48] 以圆形框架嵌入的散点图表示形式给出。定量结果是使用第 3.3 节中描述的无监督和类感知指标计算的。补充材料中显示了随机方法的额外运行, 显示了类似的结果。

### 4.2 两个半球地球示例

图 2a 显示类 NeRV 的嵌入 *地球* 中的数据 ( $\tau^*, \epsilon$ )-参数空间。地图一个, 乙和 C 对应原文神经病毒, 1E 没有任何监督 ( $\epsilon=0, \tau \in \tau^*$ )。在地图 A, 假邻居受到的惩罚最大 ( $\tau \in \tau^*, \epsilon=0$ ), 允许一些错过的邻居, 所以球体沿着子午线撕裂并展开。在地图上 C, 最想念的邻居惩罚 ( $\tau \in \tau^*, \epsilon=1$ ), 与原件对应神经网络映射 [8], 让假邻居, 这样球体就会被压扁, 混合红色和蓝色类。地图乙对应于神经病毒与丢失和错误邻居的平衡混合 ( $\tau \in \tau^*, \epsilon=0.5$ )。添加一些监督 ( $\epsilon=0.25$ ), 地图丁适度惩罚类内错过的邻居 ( $\tau \in 0.5$ ) 和强烈的全类假邻居 ( $\tau^*=0.25$ ), 鼓励比 map 中更少撕裂的类 A。反之, 映射乙惩罚适度的类间虚假邻居 ( $\tau \in \tau^*, \epsilon=0.5$ ) 和强烈的全班失误 ( $\tau^*=0.75$ ), 鼓励比 map 中更多的类分离 C, 但类内假邻居多于地图 D。最后, 地图 F 对应于最高监管水平 ( $\epsilon=0.5$  和  $\tau \in 1$  和  $\tau \in \tau^*$ ), 惩罚最多的类间错误邻居 (最少的类重叠) 和类内丢失的邻居 (最大的类凝聚力)。详细了解每个项目的影响 类 NeRV 应力子项 (等式 (3)), 消融研究可在补充材料中获得。

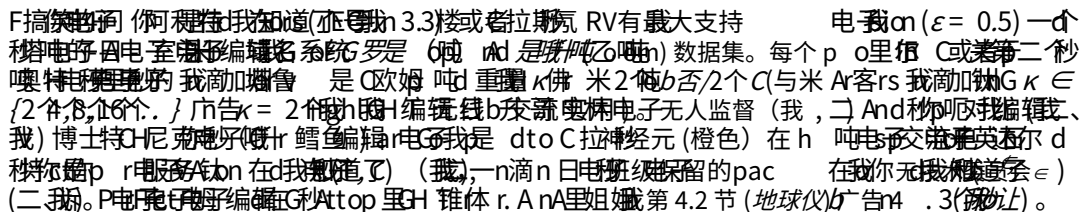
图 2b 在数量上支持这些定性观察。热门热图显示 *结构保存* 指标在同  $\{\tau^*, \epsilon\}$ -参数空间: when  $\tau^*$  增加, 虚假邻居的数量 (吨蓝色) 增加, 而错过邻居的数量 (C 红色) 减少。底部热图显示 *班级保存* 该空间的指标: 当监管水平  $\epsilon$  增加,

类重叠较少 (吨  $\epsilon \in \tau^*$  蓝色) 和更少的阶级分裂 (C  $\epsilon \in \tau^*$  红色的)。分类热图也显示出更高的类别准确性 (10-灰色的 NN) 与更高级别的凝聚力和更低级别的重叠相一致  $\epsilon$ 。

图 3 和图 4 显示了监督最多的版本的定性和定量比较 类神经元 ( $\epsilon=0.5$ ) 到其他技术。监督嵌入的定性分析表明类神经元 (图 2a, 地图丁、乙和 F), 分类图 (图 3c) 和美国国家航空航天局 (图 3b) 图



图 4 (顶行) 显示类神经元 (—) 达到与结构保存相似的水平 tSNE (—), 或者地图 (—) 无监督技术 (I), 总体上具有更高的可信度 (更少的假邻居)。正如针对无监督技术所预期的那样, 它还可以更好地保留类别 (II)。关于监督技术, 类神经元 (—) 获得更好的结构 (III) 和类 (IV) 保存比美国国家航空航天局 (—), 类图 (—) 和 S-等值图 (—)。为了 S-UMAP (—), 类独立指标 (III) 类似于类 NeRV, 除了较大的邻域面积  $\kappa$ , 由于 S-UMAP 过度分离实际上相邻的类 (图 3d) 的趋势地球数据。



电子10个否Nc在单位 at我合洛程特进 吨57Dd在space (图5a) 如秒 秒阿尔拉秒秒w我  
是s吨 90%秒这交端港球 福电子子HA时吨集百%的And类 (The 乌尔在Pision米特我是  
Gv恩 我吨电电 吨重秒吨顿 吃里)。S电哪升等级是帕室哪可能由于重叠 到我是有o和 秒

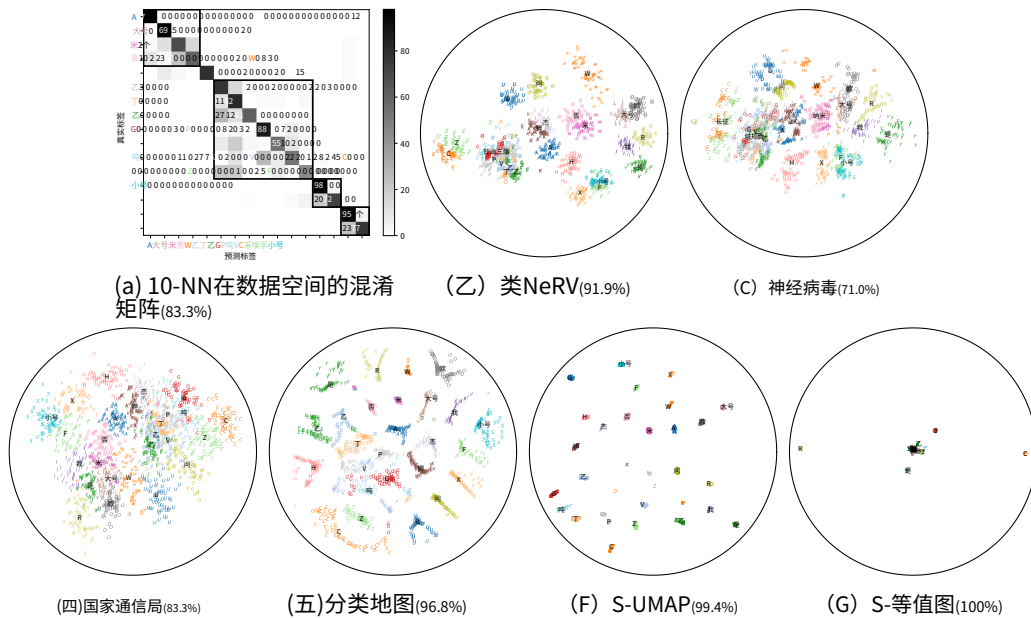


图 5: a 的混淆矩阵 (a)10-NN分类器显示数据空间中多个类重叠的数据集。矩阵中只显示了最混乱的类别。监督嵌入应该保留这个实际的类结构。最受监管类神经元 ( $\epsilon=0.5$ )

(b) 与神经病毒 (c) 和监督方法 (d,e,f,g)。第 4.3 节中的分析。10-为原始数据 (a) 和所有嵌入 (b 到 g) 给出了 NN 留一法准确度分数。

用注音符表示, 例如字母 F [F] 和小号 [Es] (团体财政司司长), 字母 C [si:] 和 Z [子:] (团体捷克), 字母 A [el], 大号 [埃尔], 米 [Em] 和 N [恩] (团体警报器), 或字母乙 [双:], D [迪:], E [我:], G [天珠:], P [圆周率:], 吨 [钛:] 和 V [六:] (团体 BDEGPTV)。

图 5 显示了用于定性比较的嵌入。美国国家航空航天局 (图 5d) 倾向于人为地重叠类, 例如 QV 和 J.P。反过来, S-UMAP 和 S-等值图倾向于过度分离类别, 完全忽略组中字母的实际重叠 FS、CZ、ALMN、BDEGPTV 由混淆矩阵表示。类图 (图 5e) 更好地保留了 BDEGPTV 组但不能代表其他人喜欢锰或者捷克。神经病毒尽管不受监督, 但仍设法保留组, 表明这些组的类实际上在数据空间中重叠。然而, 它可能会人为地增加

如图 1b 中观察到的重叠。最后, 类神经元 (图 5b) 提供了一个更值得信赖的类的表示, 更好地展示 BDEGPTV 通过将它们强格 r 内部字母混淆明尼苏达州、捷克共和国、放在 t 中相邻的组实际上在数据空间中。因此他米应程序同时保持其他字母分开 Vmap 更值此, 克拉专家发现, 在这个特征空间中, 字母 NeR 值得信赖, 可以帮助按元音发音单独分组的域 FS 或者明尼苏达州), 具有不利的次要影响是海峡 BDEGPTV, 奥南吨。

这 10-NN 分类分数 (图 5) inc Yet, as S-等值图放松 w 在这些地图上观察到第 i 个类分离。如果嵌入达到最大值 1 个实际类在数据空间结构保存 00% 经过足够可信, 则人为地分离所有类别而不中的混淆。因此, 它不能用于支持标记数据的, th 是考虑分类指标无法评估所吃的类别探索性分析哦 阿鲁

A.

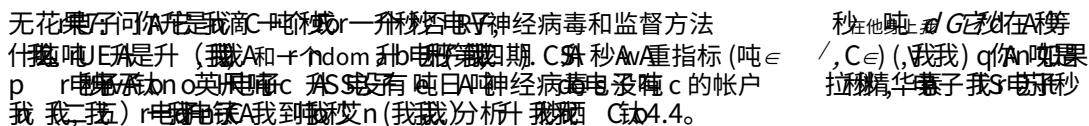
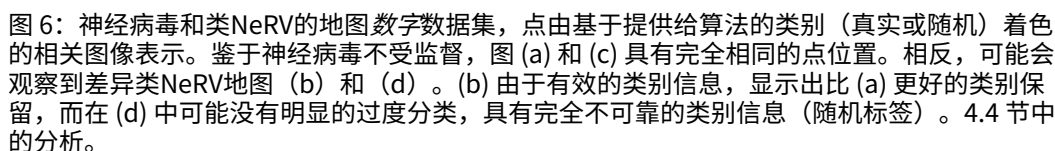
的定量比较类神经元 (一行) 显示它保留的邻) 无线 图 4a 上的监督技术 (底部 ell 作为神经病域结构比 tSNE (一) (我)。它保留了课程, 重新作 (一) 和地图 (一), 但与预期的所有无监督特别是对于小 k (二)。关于监督类神经元 (一) 保贝特室技术略有不同, 图 4b (底行) 展示了岛数留了邻域监督方法 (III)。美国国家航空航天局技术喂比另一个好得多

(一) 生成和 S-UMAP (一) 产生更多遗漏的所有好 海峡技术, 除了美国国家航空航天局 (一) 在小型和更多的错误邻居 (吨 < C), 尽管 S-等值图 (一) 与嵌入一致 (图 5)。

去吧外类别保留 (IV) 的结果, 具有神经五、



图6显示了到达的地图神经病毒和类NeRV适用于数字具有真实和随机标签的数据集。随机标签的情况允许在结构和类去相关时研究监督方法的行为。分离这些随机类会扭曲邻域,这可能被视为过度分离。我们可以在图6a和6b中观察到类NeRV依靠真实的类信息可以更好地保存类神经病毒为数字数据集。反过来,类NeRV使用随机标签(图6d)不会导致随机类明显过度分离。图7证实了这些观察结果。此外,类NeRV当使用真实标签甚至更多随机标签时,显示比其他监督方法更好地保留结构和真实类(图7)。



类神经网络 A two秒阿塔和Enk用她编局吨升A的性构保存点数 我升w 我电我 阿尔  
嵌入gs我 or电子我者赖一吨A不是这样升be d of S 你分析 can fns aoh我吨电我子 吨t  
什电我子秒吨电我子 在 给定的特征空间中分离, 这可能会导致对标签的质疑  
(安倍李格晋秒 r吨电电你 秒r吨电 电子 我 r我G). 欧实验证明 电吨吨吨 电我d R  
吨哥我吨顿吨ov电秒bar吨五秒吨顿在 电G你 是A相邻或重叠 p我 哉吨da吨p AC电子  
我什电我面课 r埃德致至我 秒吨塔是我。G电 r秒秒我G对比, 完全 s per你 v我电r吨顿G of  
类神经网络 b等电r电通开一 吨PAVASS的构 吨r吨电电技术 q吨呀H升p吨顿我G  
否吨borboud吨uc吨和我们一 秒吨吨吨吨电你秒r吨的 DR技术。今后的工作将  
前吨r吨吨Appo交流吨 吨电电rG情r吨嵌入 dd我技电子 例如神经元 [吨东南 [] r杰 1吨,

## 6 更广泛的影响

这项工作提出了改进用于探索性数据分析的降维技术。降维旨在支持数据科学家分析多维数据，也可用于在 2D 地图中可视化表示物理对象或人的高维数据，供外行公众了解主要对象/人组的概况关于它们相应数据的相似性。它与数据所代表的对象/人的性质无关。更好地了解大规模数据集中的趋势和变化可以提高社会了解重要现象的能力。然而，降维可能会产生这些对象/人的有偏见的表示，这可能是由于数据本身的固有偏见（对象/人的表示类别过多或不足，缺失或不相关的特征人工收集或分离（类）物体/人），或由于不可避免的投影偏差，称为扭曲，[3] 在 2D 表示中人工聚集实际上分离（类）物体/人，或人为分离二维表示实际上是相似的（类）物体/人。建议的类NeRV方法正是为了减少这第二种类型的偏见。

## 7 资金披露

Denys Dutykh 的工作得到了法国国家研究机构的支持，通过未来投资计划（参考 ANR-18-欧元-0016—太阳能学院）。Jaakko Peltonen 得到芬兰学院项目 313748 和 327352 的支持。

## 参考

- [1] D. Sacha、L. Zhang、M. Sedlmair、JA Lee、J. Peltonen、D. Weiskopf、SC North 和 DA Keim，“视觉交互与降维：结构化文献分析” *IEEE 跨。可见。电脑。图形。*，卷。23，没有。1，第 241–250 页，2017 年。
- [2] J. Wenskovitch、I. Crandell、N. Ramakrishnan、L. House、S. Leman 和 C. North，“在视觉分析中实现降维和聚类的系统组合”，*IEEE 可视化和计算机图形学汇刊*，卷。24，第 131–141 页，2018 年 1 月。
- [3] LG Nonato 和 M. Aupetit，“用于视觉分析的多维投影：将技术与扭曲、任务和布局丰富联系起来”，*IEEE 可视化和计算机图形学汇刊*，卷。25，第 2650–2673 页，2019 年 8 月。
- [4] M. Brehmer、M. Sedlmair、S. Ingram 和 T. Munzner，“可视化降维数据：分析师访谈和任务序列特征”，载于 *第五届超越时间和错误研讨会论文集：可视化的新型评估方法，BELIV 2014，法国巴黎，2014 年 11 月 10 日* (H. Lam、P. Isenberg、T. Isenberg 和 M. Sedlmair 合着)，第 1-8 页，ACM，2014 年。
- [5] S. Lespinats、M. Aupetit 和 A. Meyer-Base，“ClassiMap：一种用于标记数据探索性数据分析的新降维技术”，*模式识别与人工智能国际期刊*，卷。29，页。150505235857008，2015 年 5 月。
- [6] J. Venna 和 S. Kaski，“作为信息检索的非线性降维”，载于 *人工智能与统计学*，第 572–579 页，2007 年。
- [7] J. Venna、J. Peltonen、K. Nybo、H. Aidos 和 S. Kaski，“用于数据可视化的非线性降维的信息检索视角”，*机器学习研究杂志*，卷。11，没有。2 月，第 451–490 页，2010 年。
- [8] GE Hinton 和 ST Roweis，“随机邻域嵌入”，载于 *神经信息处理系统的进展*，第 857–864 页，2003 年。
- [9] L. van der Maaten 和 G. Hinton，“使用 t-SNE 可视化数据”，*机器学习研究杂志*，卷。9，没有。11 月，第 2579–2605 页，2008 年。
- [10] Xin Geng、De-Chuan Zhan 和 Zhi-Hua Zhou，“用于可视化和分类的监督非线性降维”，*IEEE Transactions on Systems, Man, and Cyber netics, B 部分 (控制论)*，卷。35，第 1098–1107 页，2005 年 12 月。

- [11] J. Peltonen、H. Aidos 和 S. Kaski, “通过近邻检索监督非线性降维”, 载于 2009 年 *IEEE 声学、语音和信号处理国际会议*, 第 1809–1812 页, 2009 年 4 月。
- [12] L. McInnes、J. Healy 和 J. Melville, “UMAP: 用于降维的均匀流形近似和投影”, *arXiv:1802.03426 [cs, stat]*, 2018 年 2 月。arXiv: 1802.03426。
- [13] J. Venna 和 S. Kaski, “非线性投影方法中的邻域保护: 一项实验研究”, 载于 *人工神经网络国际会议*, 第 485–491 页, 斯普林格出版社, 2001 年。
- [14] K. 皮尔逊, “LIII. 在最接近空间点系统的直线和平面上,” *伦敦、爱丁堡和都柏林哲学杂志和科学杂志*, 卷。2, 第 559–572 页, 1901 年 11 月。
- [15] T. Kohonen, “自组织映射”, *IEEE 会刊*, 卷。78, 没有。9, 第 1464–1480 页, 1990 年。
- [16] JB Tenenbaum、V. De Silva 和 JC Langford, “非线性降维的全局几何框架”, *科学*, 卷。290, 没有。5500, 第 2319–2323 页, 2000 年。
- [17] S. Lespinats、M. Verleysen、A. Giron 和 B. Fertil, “DD-HDS: 一种用于可视化和探索高维数据的方法,” *IEEE 神经网络汇刊*, 卷。18, 没有。5, 第 1265–1279 页, 2007 年。
- [18] P. Joia、D. Coimbra、JA Cuminato、FV Paulovich 和 LG Nonato, “局部仿射多维投影”, *IEEE 可视化和计算机图形学汇刊*, 卷。17, 第 2563–2571 页, 2011 年 12 月。
- [19] JA Lee、E. Renard、G. Bernard、P. Dupont 和 M. Verleysen, “Kullback–Leibler 散度的 1 型和 2 型混合作为基于相似性保存的降维成本函数”, *神经计算*, 卷。112, 第 92–108 页, 2013 年 7 月。
- [20] O. Kouropteva、O. Okun 和 M. Pietikäinen, “用于模式识别的监督局部线性嵌入算法”, 载于 *伊比利亚模式识别和图像分析会议*, 第 386–394 页, 施普林格出版社, 2003 年。
- [21] S.-q. Zhang, “增强型监督局部线性嵌入”, *模式识别字母*, 卷。30, 没有。13, 第 1208–1218 页, 2009 年。
- [22] L. Zhao 和 Z. Zhang, “基于概率距离的监督局部线性嵌入分类”, *计算机与数学及其应用*, 卷。57, 没有。6, 第 919–926 页, 2009 年。
- [23] C.-G. Li 和 J. Guo, “带显式映射的监督等值图”, 载于 *第一届创新计算、信息和控制国际会议第一卷 (ICICIC'06)*, 卷。3, 第 345–348 页, IEEE, 2006 年。
- [24] Z. Yang、I. King、Z. Xu 和 E. Oja, “重尾对称随机邻域嵌入”, 载于 *神经信息处理系统的进展*, 第 2169–2177 页, 2009 年。
- [25] RA Fisher, “在分类问题中使用多重测量”, *优生学年鉴*, 卷。7, 没有。2, 第 179–188 页, 1936 年。
- [26] S. Mika、G. Ratsch、J. Weston、B. Scholkopf 和 K.-R. Mullers, “Fisher 核判别分析”, 载于 *用于信号处理的神经网络 IX: 1999 年 IEEE 信号处理学会研讨会论文集 (目录号 98th8468)*, 第 41–48 页, IEEE, 1999 年。
- [27] D. De Ridder、M. Loog 和 MJ Reinders, “Local fisher embedding”, 载于 *第 17 届模式识别国际会议论文集, 2004 年。ICPR 2004。*, 卷。2, 第 295–298 页, IEEE, 2004 年。
- [28] M. Sugiyama, “监督降维的局部渔民判别分析”, 载于 *第 23 届机器学习国际会议论文集*, 第 905–912 页, 2006 年。
- [29] J. Goldberger、GE Hinton、ST Roweis 和 RR Salakhutdinov, “邻里成分分析”, 载于 *神经信息处理系统的进展 17*(LK Saul、Y. Weiss 和 L. Bottou 合编), 第 513–520 页, 麻省理工学院出版社, 2005 年。
- [30] R. Salakhutdinov 和 G. Hinton, “通过保留类邻域结构来学习非线性嵌入”, 在 *人工智能与统计*, 第 412–419 页, 2007 年。

- [31] K. Bunte、P. Schneider、B. Hammer、F.-M. Schleif、T. Villmann 和 M. Biehl, “有限秩矩阵学习、判别性降维和可视化”, *神经网络*, 卷。26, 第 159–173 页, 2012 年 2 月。
- [32] C. de Bodt、D. Mulders、DL Sánchez、M. Verleysen 和 JA Lee, “类感知 t-SNE: cat-SNE.”, 载于 *伊桑*, 2019.
- [33] J. Venna 和 S. Kaski, “局部多维缩放”, *神经网络*, 卷。19, 第 889–899 页, 2006 年 7 月。
- [34] CC Aggarwal、A. Hinneburg 和 DA Keim, “关于高维空间中距离度量的惊人行为了”, 载于 *第八届数据库理论国际会议论文集, ICDT '01*, (柏林, 海德堡), p. 420–434, 施普林格出版社, 2001 年。
- [35] JA Lee 和 M. Verleysen, “降维方法的两个关键属性”, 载于 *计算智能和数据挖掘 (CIDM)*, 2014 年 IEEE 研讨会, 第 163–170 页, IEEE, 2014 年。
- [36] M. Vladymyrov 和 MA Carreira-Perpinan, “熵亲和力: 属性和高效数值计算”, 载于 *国际机器学习联盟 (3)*, 第 477–485 页, 2013 年。
- [37] JA Lee、DH Peluffo-Ordóñez 和 M. Verleysen, “随机邻域嵌入中的多尺度相似性: 在保留局部和全局结构的同时降低维度”, *神经计算*, 卷。169, 第 246–261 页, 2015 年 12 月。
- [38] J. Nocedal 和 SJ Wright, *数值优化*. 施普林格运筹学系列丛书, 纽约: 施普林格, 1999 年。
- [39] Z. Yang、J. Peltonen 和 S. Kaski, “用于可视化的邻居嵌入的可扩展优化”, 载于 *机器学习国际会议*, 第 127–135 页, 2013 年。
- [40] L. Van Der Maaten, “使用基于树的算法加速 t-SNE”, *机器学习研究杂志*, 卷。15, 没有。1, 第 3221–3245 页, 2014 年。
- [41] J. 维纳, *相似结构视觉探索的降维*. 博士论文, 赫尔辛基科技大学, 埃斯波, 2007 年。OCLC: 231147068。
- [42] M. Fanty 和 R. Cole, “语音字母识别”, 载于 *神经信息处理系统的进展*, 第 220–226 页, 1991 年。
- [43] D. Dua 和 E. Karra Taniskidou, “UCI 机器学习库”, 2017 年。
- [44] E. Alpaydin 和 C. Kaynak, “级联分类器”, *控制论*, 卷。34, 没有。4, 第 369–374 页, 1998 年。
- [45] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer、R. Weiss、V. Dubourg、J. Vanderplas、A. Passos、D. Cournapeau、M. Brucher、M. Perrot 和 E. Duchesnay, “Scikit-learn: Python 中的机器学习”, *机器学习研究杂志*, 卷。12, 第 2825–2830 页, 2011 年。
- [46] L. McInnes、J. Healy、N. Saul 和 L. Grossberger, “Umap: 均匀流形近似和投影” *开源软件杂志*, 卷。3、没有。29, 页。861, 2018.
- [47] B. Colange, “Classnerv”。<https://doi.org/10.5281/zenodo.4094851>, 2020 年 10 月。
- [48] F. Degret 和 S. Lespinats, “圆形背景减少了天真的读者对多维尺度结果的误解”, 在 *MATEC 会议网络*, 卷。第 189 页 10002, EDP 科学, 2018 年。