

基于 ML 数据清洗算法的下游聚类应用算法的适应性分析

常添

08/23/2024

1 理论分析结论表格

清洗算法	聚类算法					
	K-Means	GMM	AP	HC	OPTICS	BIRCH
Baran	A	B	B	A	A	B
BoostClean	B	A	B	A	B	C
CPClean	A	A	C	B	B	A
ED2	B	A	C	A	B	A
HoloClean	A	B	B	A	B	A
Metadata-Driven	A	B	B	A	A	B
Raha	A	A	B	A	A	A
Scare	B	A	B	C	B	A

表 1: 8 种基于 ML 的数据清洗算法与聚类算法的适应性分析

注: 适应程度用 A、B、C 三个等级来评估, 其中 A 代表高度适应, B 代表中等适应, C 代表低适应。

2 适应性评估的研究思路

2.1 算法机制与模型特性分析

深入分析算法所采用的多种机器学习模型, 探讨算法进行错误检测和错误修复的机制, 同时分析算法如何通过这些模型的协同作用来提高数据清洗的精度和效率。

2.2 数据清洗对特征的影响分析

根据算法的特点, 系统性分析原始数据在清洗过程中的特征变化, 探讨各类数据错误 (如格式错误、语义错误、依赖性错误) 在算法的不同步骤中如何被识别、修正或可能被引入新的误差, 并评估这些变化对数据完整性和一致性的影响, 从而推断其对后续聚类任务可能产生的影响。

2.3 清洗后数据对聚类算法的适应性评估

基于清洗后数据的特征变化，预测清洗后的数据对下游聚类应用的潜在影响，并评估新数据对六种经典聚类算法（K-means、Gaussian Mixture Models、Affinity Propagation、Hierarchical Clustering、OPTICS、BIRCH）的适应性。

3 适应性评估的理论依据

3.1 Baran 算法

3.1.1 算法机制与模型特性分析

Baran 算法采用了一种统一的上下文表示和迁移学习的方法来实现数据错误的检测和修复。算法通过三个主要的模型类型来处理数据错误：基于值的模型（Value-based Models）、基于邻域的模型（Vicinity-based Models）和基于领域的模型（Domain-based Models）。每个模型类型分别利用不同的上下文信息来生成可能的修正候选项。基于值的模型主要利用数据值本身来进行修正，而基于邻域的模型依赖于数据值所在行的其他值来进行修正，基于领域的模型则依靠同一列中的其他值来进行修正。

这些模型通过一个二步流程来完成错误修正：首先，每个模型生成一组初步的修正候选项；接着，这些候选项被整合为最终的修正结果。通过迁移学习，Baran 可以从类似领域的数据中预训练这些模型，使其在数据稀缺的情况下仍能有效工作。

3.1.2 数据清洗对特征的影响分析

在数据清洗过程中，Baran 算法通过不同模型类型识别并修复多种类型的数据错误，例如格式错误、语义错误和依赖性错误。基于值的模型在处理格式错误（如日期格式不一致）时表现出色，能够识别和修正诸如日期格式或拼写错误的简单语法问题。基于邻域的模型则特别擅长处理依赖性错误，例如在同一行中，城市名称和邮政编码不匹配的问题。而基于领域的模型可以利用一列中出现的频繁值来修正异常值。

这种多层次的清洗机制确保了数据的完整性和一致性，虽然有时也可能引入新的误差，例如在语义修正时，由于上下文信息不足导致的错误修正。总体来说，这种系统化的清洗过程会极大地改善数据的质量，但在处理复杂语义错误时，可能会影响下游聚类任务的结果准确性。

3.1.3 清洗后数据对聚类算法的适应性评估

- **K-means**: 适应性评级为 A。对格式错误的修正能够显著提高 K-means 算法的聚类效果，因为该算法对数值一致性要求较高。
- **Gaussian Mixture Models (GMM)**: 适应性评级为 B。由于 GMM 对数据分布的敏感性，清洗后的数据如果修正不当，可能会影响 GMM 的结果。
- **Affinity Propagation**: 适应性评级为 B。该算法对噪声和异常值较为敏感，清洗后的数据如果存在错误修正，可能会影响聚类中心的选择。

- **Hierarchical Clustering:** 适应性评级为 A。层次聚类对数据修正后的变化较为鲁棒，清洗后的数据适应性较好。
- **OPTICS:** 适应性评级为 A。由于 OPTICS 算法主要用于处理噪声较多的数据，Baran 清洗后的数据通常会减少噪声。
- **BIRCH:** 适应性评级为 B。对于大规模数据集，BIRCH 依赖于输入数据的特征。清洗后的数据如果特征变化较大，可能会影响其树结构的构建。

3.2 BoostClean 算法

3.2.1 算法机制与模型特性分析

BoostClean 算法的核心在于利用机器学习模型来自动检测和修复数据错误，尤其是领域值违规问题。该算法采用了一种基于“Boosting”技术的集合方法，将多个弱学习器组合成一个强学习器，以此来选择最优的错误检测和修复操作组合。BoostClean 首先通过一个检测器生成器库（如 Isolation Forests 和 Word2Vec 神经网络模型）自动生成检测数据错误的规则，然后利用修复函数库对检测到的错误进行修复。通过在训练数据和测试数据上应用这些检测和修复操作，BoostClean 不断优化下游预测模型的准确性。

3.2.2 数据清洗对特征的影响分析

在 BoostClean 的清洗过程中，数据的特征会发生明显变化。算法中的检测器可以识别诸如格式错误（例如数据类型不匹配）、语义错误（如类别属性中出现的异常值）、以及依赖性错误（如某些属性组合出现的逻辑错误）等多种数据错误。修复操作则可能包括填补缺失值、纠正异常值或者删除有问题的记录。

然而，这些修复操作在提升数据一致性的同时，也有可能引入新的误差，例如在填充缺失值时可能破坏数据中的原有统计特征，从而对下游任务造成负面影响。此外，BoostClean 的修复策略是基于提高模型预测准确性，而非完全恢复数据的逻辑一致性，因此在某些场景下，可能会因为过度修复导致特征的偏差加大，进而影响后续的聚类任务。

3.2.3 清洗后数据对聚类算法的适应性评估

- **K-means:** 适应性评级为 B。由于 K-means 依赖于数据的均值和方差，BoostClean 在清洗过程中如果过度平滑数据，可能会降低 K-means 对数据中心的识别能力。
- **Gaussian Mixture Models (GMM):** 适应性评级为 A。GMM 能够很好地适应 BoostClean 生成的相对干净且分布合理的数据，尤其是在处理高斯分布的数据时表现更好。
- **Affinity Propagation:** 适应性评级为 B。此算法对数据的特征分布较为敏感，BoostClean 可能会引入一些偏差，影响代表点的选择。
- **Hierarchical Clustering:** 适应性评级为 A。层次聚类对数据分布的要求较低，清洗后的数据可以更好地反映样本间的层次关系。

- **OPTICS**: 适应性评级为 B。此算法关注密度聚类，清洗可能会影响数据的局部密度，导致部分样本在聚类过程中被误识别。
- **BIRCH**: 适应性评级为 C。BIRCH 算法对输入数据的结构信息非常敏感，清洗过程中的任何特征扭曲都有可能对簇划分结果的明显变化。

3.3 CPClean 算法

3.3.1 算法机制与模型特性分析

CPClean 算法的核心思想是基于“确定预测”（Certain Predictions, CP）的概念。该算法通过引入可能世界（Possible Worlds）的概念，扩展了机器学习中对不完整数据的处理方法。其主要机制包括：

- **错误检测**: CPClean 使用 KNN 分类器来分析不完整数据集的可能世界，检查数据集中缺失或不准确值对分类结果的影响。算法通过构建“l-极端世界”（l-extreme world）来推断是否存在某个可能世界可以得出不同的预测结果。
- **错误修复**: 通过排序和扫描候选值的相似度，CPClean 能够动态计算每个可能世界下的边界集，并基于此进行修复。该过程通过逐步减少数据的不确定性，最终达到提高模型在验证集上的分类准确性。

3.3.2 数据清洗对特征的影响分析

CPClean 算法在数据清洗过程中对特征的影响主要体现在两个方面：

- **特征的改变**: 由于 CPClean 考虑了多个可能的修复选项，在实际清洗过程中可能会对原始数据特征产生显著影响。例如，算法会评估每个候选修复值对分类结果的影响，并选择那些能够提高分类精度的修复选项。这一过程虽然有效减少了原始数据中的错误，但当原始数据的候选修复值质量较差时，也可能引入新的误差。
- **数据完整性与一致性**: CPClean 通过动态计算不同候选修复值对最终分类的影响，确保在最小化数据不确定性的同时，最大程度地保持数据的一致性。然而，这种修复过程可能会在多个特征之间存在复杂关联时引入数据之间新的依赖性错误。

3.3.3 清洗后数据对聚类算法的适应性评估

- **K-means**: 适应性评级为 A。CPClean 清洗后的数据能够很好地保留特征间的相似性和差异性，K-means 依赖于欧几里得距离或相似性度量，因此对清洗后的数据适应性较高。
- **Gaussian Mixture Models (GMM)**: 适应性评级为 A。GMM 同样依赖于数据的分布和相似性度量，CPClean 能够很好地保留这些特征，因此 GMM 对清洗后的数据适应性较高。
- **Affinity Propagation**: 适应性评级为 C。Affinity Propagation 对数据结构要求较高，CPClean 可能会引入一些特征间新的关联性，从而影响算法对数据的适应性，因此适应性较弱。

- **Hierarchical Clustering**: 适应性评级为 B。层次聚类依赖数据的结构特性，CPClean 清洗后的数据能够在一定程度上保持原始数据的结构特性，但由于清洗过程可能引入新的依赖性错误，因此适应性稍差。
- **OPTICS**: 适应性评级为 B。OPTICS 算法依赖于数据的局部结构和密度特性，CPClean 清洗后的数据在保持这些特性上表现一般，因此适应性评级为 B。
- **BIRCH**: 适应性评级为 A。BIRCH 算法对输入数据的结构信息和相似性度量非常敏感，CPClean 清洗后的数据保留了这些特征，因此 BIRCH 对清洗后的数据适应性较高。

3.4 ED2 算法

3.4.1 算法机制与模型特性分析

ED2 算法的核心机制是基于主动学习的错误检测方法。该算法通过两阶段采样策略来优化用户标注的效率。ED2 使用多个分类器（每个列一个分类器）对数据进行分类，先从所有列中选择一个最需要标注的列，然后在该列中选择最有可能含有错误的单元格进行标注。这个策略显著减少了用户所需标注的单元格数量，并且通过逐步优化分类器模型，使得错误检测的准确率达到最优。

模型特性方面，ED2 算法主要依赖于 XGBoost 作为分类器，由于 XGBoost 对不相关特征的鲁棒性以及在处理不平衡数据集上的表现优异，ED2 能够在大多数情况下快速收敛。此外，ED2 的设计还包括了对各列分类模型的交叉验证和超参数调优，进一步增强了模型的精度和稳定性。

3.4.2 数据清洗对特征的影响分析

在数据清洗过程中，ED2 的特征提取器会从数据的属性、元组和数据集级别提取信息，形成特征向量。该过程能够识别和纠正不同类型的错误，如格式错误（例如错别字）、语义错误（例如值交换错误）和依赖性错误（例如违反约束的错误）。在处理这些错误时，ED2 通过频率排序和离群点检测等方法预过滤可能的错误数据点，随后通过主动学习模型进行进一步确认。

然而，尽管 ED2 在识别和修正错误方面表现出色，但在某些情况下，特征提取和主动学习策略可能引入新的误差。例如，当模型在一个数据列上过早地收敛，其他列的数据错误可能未被充分捕捉，导致清洗后的数据存在隐性错误。此外，ED2 的模型依赖于用户标注的数据，若初始标注数据中存在偏差，可能导致后续清洗过程中的误差累积。

3.4.3 清洗后数据对聚类算法的适应性评估

- **K-means**: 适应性评级为 B。ED2 清洗后的数据可能仍存在一些未识别的离群点或噪声数据，可能影响 K-means 的中心点选择和结果一致性。
- **Gaussian Mixture Models (GMM)**: 适应性评级为 A。GMM 可以较好地处理由 ED2 引入的潜在误差，尤其是在数据清洗后仍然存在的一些小规模的不平衡分布。
- **Affinity Propagation**: 适应性评级为 C。由于该算法对噪声和异常值敏感，ED2 清洗后的数据如果仍存在未检测的噪声，可能导致 Affinity Propagation 的收敛性和结果准确性下降。

- **Hierarchical Clustering**: 适应性评级为 A。层次聚类算法能够较好地处理经过 ED2 清洗后的数据，因为清洗过程提高了数据的结构化程度，有助于层次聚类更好地划分数据层次。
- **OPTICS**: 适应性评级为 B。OPTICS 适合处理噪声较多的数据集，ED2 清洗后的数据如果仍保留部分噪声，可能会对 OPTICS 的聚类顺序产生一定影响。
- **BIRCH**: 适应性评级为 A。BIRCH 对于大规模数据集的聚类表现良好，且能够处理由 ED2 清洗后产生的聚类内误差和噪声。

3.5 HoloClean 算法

3.5.1 算法机制与模型特性分析

HoloClean 算法结合了多种机器学习模型，通过概率推断来实现数据清洗。它采用了三种主要信号：完整性约束、外部数据和数据的统计特性。首先，HoloClean 将输入的数据视为一个噪声数据集，通过构建一个概率图模型，对每个数据单元的正确性进行评估。然后，算法通过统计学习和概率推断来确定每个数据单元的最佳修复值。HoloClean 通过将各种信号整合进统一的概率模型中，从而能够在数据修复过程中实现高精度和高召回率的平衡。

3.5.2 数据清洗对特征的影响分析

在数据清洗过程中，HoloClean 算法通过多种方法识别和修复错误。这些方法包括完整性约束的违反、异常检测、以及与外部数据的匹配。对于不同类型的数据错误，例如格式错误、语义错误和依赖性错误，HoloClean 分别利用统计分析、外部知识库和完整性约束进行识别和修复。例如，格式错误可以通过匹配外部数据源进行修复，而语义错误则可能通过完整性约束来识别。但是，在修复过程中，有可能引入新的误差，尤其是在不同信号之间存在冲突时。此外，HoloClean 在减少错误的同时，也可能对数据的分布特性产生影响，从而影响下游任务的数据完整性和一致性。

3.5.3 清洗后数据对聚类算法的适应性评估

- **K-means**: 适应性评级为 A。因为 K-means 依赖于整体数据分布，清洗后的数据通常符合其假设。
- **Gaussian Mixture Models (GMM)**: 适应性评级为 B。GMM 对数据分布假设较强，HoloClean 的修复可能引入分布偏差。
- **Affinity Propagation**: 适应性评级为 B。该算法对数据的局部相似性较敏感，清洗后的数据在局部区域可能产生误差。
- **Hierarchical Clustering**: 适应性评级为 A。该算法对全局数据结构敏感，清洗后的数据整体性较好。
- **OPTICS**: 适应性评级为 B。由于其对噪声较敏感，清洗后的数据可能残留少量误差。
- **BIRCH**: 适应性评级为 A。适用于大规模数据，清洗后的数据一致性有利于其树状结构构建。

3.6 Metadata-driven 算法

3.6.1 算法机制与模型特性分析

Metadata-driven 算法通过结合多种机器学习模型来进行错误检测和修复。主要采用了两种核心策略：Bagging 和 Stacking。这些方法利用了多个错误检测系统的结果，并通过元数据的支持，提升了错误检测的效果。Bagging 策略通过组合多个决策树模型来预测错误，Stacking 则将不同模型的输出作为输入，进一步结合形成更强大的预测模型。该算法的机制通过引入元数据来增强错误检测的准确性，从而更好地识别数据中的错误。

3.6.2 数据清洗对特征的影响分析

在数据清洗过程中，Metadata-driven 算法利用元数据特征（如数据完整性、数据类型隶属关系、属性域、常见值分布等）来识别和修复不同类型的错误，如格式错误、语义错误和依赖性错误。通过这些特征的关联性，算法能够有效地检测出模式违反、规则违反、离群点和重复冲突。然而，在清洗过程中，也可能引入新的误差，特别是在处理复杂的多列依赖关系时，误差的累积可能会影响数据的完整性和一致性。

3.6.3 清洗后数据对聚类算法的适应性评估

- **K-means**: 适应性评级为 A。由于 K-means 对噪声和异常值敏感，经过清洗后的数据适应性可能会显著提高，尤其是在清除了模式违反和离群点后。
- **Gaussian Mixture Models (GMM)**: 适应性评级为 B。GMM 对数据的假设比较严格，清洗后的数据在消除了多模态错误后适应性较好。
- **Affinity Propagation**: 适应性评级为 B。此算法对数据分布不敏感，清洗后的数据一致性增强，但整体影响有限。
- **Hierarchical Clustering**: 适应性评级为 A。此算法对数据的完整性要求较高，清洗后数据的层次结构更清晰。
- **OPTICS**: 适应性评级为 A。由于 OPTICS 能处理噪声，清洗后的数据会提升其在噪声处理上的效率。
- **BIRCH**: 适应性评级为 B。BIRCH 适用于大规模数据，清洗后的数据可减少簇的重叠。

3.7 RAHA 算法

3.7.1 算法机制与模型特性分析

RAHA 算法是一种无需配置的错误检测系统，旨在通过多种机器学习模型自动检测数据中的错误。其核心机制是生成特征向量，其中每个数据单元格（data cell）的特征向量通过一组配置的错误检测算法来编码。RAHA 使用的主要错误检测算法包括异常值检测、模式违规检测、规则违规检测和知识库违规检测。这些算法分别针对不同类型的错误：异常值检测算法主要用于识别不符合常规分布的数据；模式违规检测算法用于检查数据是否符合预定义的模式；规则违规检测算法评估数据值是否符合完整

性约束；知识库违规检测算法则通过与知识库数据对比来发现错误。RAHA 通过将这些检测算法的输出合并为一个特征向量，并使用无监督聚类方法将相似特征向量的单元格聚类，从而实现错误检测。

3.7.2 数据清洗对特征的影响分析

在数据清洗过程中，RAHA 通过不同的算法识别和修正各种数据错误。对于格式错误（如日期格式不一致），RAHA 利用模式违规检测算法将其识别并标记为错误。对于语义错误（如地名与实际位置不符），知识库违规检测算法通过对比知识库中的正确关系来识别这些错误。在处理依赖性错误时（如两个列之间的函数依赖关系被破坏），规则违规检测算法能够有效检测此类错误。

然而，RAHA 的多种检测算法在数据清洗过程中可能会引入新的误差。例如，当使用异常值检测时，某些合法但不常见的值可能被错误地标记为异常值。此外，如果训练数据存在噪声或者聚类不够精确，可能导致标签传播过程中引入新的错误。

3.7.3 清洗后数据对聚类算法的适应性评估

- **K-means**: 适应性评级为 A。清洗后的数据在分布上更趋向于均匀，因此在 K-means 中表现为高度适应。然而，如果异常值被误识为合法值，可能导致聚类中心的偏移。
- **Gaussian Mixture Models (GMM)**: 适应性评级为 A。与 K-means 类似，清洗后的数据对 GMM 具有高度适应，但在处理复杂分布时，可能会出现适应性问题。
- **Affinity Propagation**: 适应性评级为 B。由于该算法对噪声较为敏感，因此清洗过程中的特征变化可能导致适应性略有下降。
- **Hierarchical Clustering**: 适应性评级为 A。清洗后的数据由于层次关系清晰，通常表现为高度适应。
- **OPTICS**: 适应性评级为 A。对于 OPTICS 而言，清洗后的数据可以更好地反映密度变化，因此适应性较高。
- **BIRCH**: 适应性评级为 A。BIRCH 算法在处理大规模数据时表现较好，清洗后的数据对其有高度适应，但需注意大规模清洗可能带来的层次结构变化。

3.8 Scare 算法

3.8.1 算法机制与模型特性分析

Scare 算法通过最大化数据分布下的数据修复似然性来修复数据库中的错误。它利用机器学习（ML）技术来预测更高质量的更新，以修复错误数据。Scare 采用了多种机器学习模型（例如，决策树、贝叶斯网络）来捕捉数据集中的依赖性、相关性和异常值。这些模型主要用于预测数据中的可能更新。为了确保预测的准确性，Scare 首先在数据集上进行水平分区处理，然后在每个分区上应用机器学习方法，最终结合局部预测来生成准确的最终预测。

Scare 算法的核心机制包括两个主要步骤：第一是生成可能的修复更新，第二是选择最合适的修复更新。在更新生成过程中，算法通过学习分类模型来预测数据中的灵活属性（即可能存在错误的属性）

的值，并将这些预测结果存储在临时存储中。然后，Scare 使用图优化方法来整合这些局部预测，最终选择最合适的修复更新。

3.8.2 数据清洗对特征的影响分析

在数据清洗过程中，Scare 算法通过多次水平分区和局部视图的组合，有效地识别并修复数据中的错误。这种方法能够捕捉到数据局部相关性，从而提高修复的准确性。然而，数据清洗过程中可能会引入新的误差，尤其是在局部模型中的弱相关性未能被正确捕捉时。此外，在没有预定义约束的情况下，Scare 依赖于概率模型来进行修复，这意味着其预测结果可能会受到模型假设的影响，尤其是在面对复杂的依赖性错误时。

数据清洗对特征的影响主要表现为：

- **格式错误**：Scare 能够有效地检测并修复简单的格式错误，例如地址中的邮政编码或城市名称的拼写错误。
- **语义错误**：由于 Scare 依赖于机器学习模型来捕捉数据中的相关性，语义错误的识别和修复效果取决于模型的学习能力。例如，在同一领域内，不同城市可能有相似的地理信息，Scare 在修复这些数据时可能会引入轻微的语义误差。
- **依赖性错误**：Scare 通过依赖性网络模型来处理属性间的依赖性，因此在修复复杂的依赖性错误时表现出色。然而，如果数据中存在显著的局部相关性，而全局模型无法有效捕捉到这些相关性，则可能会引入新的依赖性误差。

3.8.3 清洗后数据对聚类算法的适应性评估

- **K-means**：适应性评级为 B。清洗后的数据通过 Scare 的处理，格式错误和部分语义错误得到了有效修复，但依赖性错误的残留可能导致 K-means 在中心点选择上的不稳定性。
- **Gaussian Mixture Models (GMM)**：适应性评级为 A。由于 GMM 能够处理数据中的噪声和不确定性，Scare 的修复策略能够有效增强 GMM 在处理高维数据聚类时的性能。
- **Affinity Propagation**：适应性评级为 B。尽管 Scare 修复了大部分错误，但由于依赖性错误的潜在影响，Affinity Propagation 的性能可能会在高相关性的属性之间有所下降。
- **Hierarchical Clustering**：适应性评级为 C。层次聚类对数据质量要求较高，Scare 修复后的数据在层次结构的构建上可能引入偏差，导致聚类效果不佳。
- **OPTICS**：适应性评级为 B。OPTICS 能够处理噪声和异常点，Scare 修复后的数据在处理过程中表现良好，但在复杂的语义错误情况下可能适应性不强。
- **BIRCH**：适应性评级为 A。BIRCH 算法对大规模数据和局部聚类效果良好，Scare 修复策略能够显著提升 BIRCH 在处理大数据集时的聚类效果。