

基于非参数估计与随机模拟的不确定数据流相似性度量方法

迟荣华 黄少滨 李熔盛

哈尔滨工程大学 信息与通信工程学院 哈尔滨

摘要: 对不确定数据流相似性度量方法，提出一种基于非参数估计与随机模拟的相似性度量方法。该方法首先对不确定数据流进行非参数估计，然后利用随机模拟技术生成具有相似性的数据流。仿实验结果表明，该方法不仅可以准确地度量不确定数据流的相似性，而且在数据流规模较大、分布复杂的情况下，依然具有良好的稳定性和鲁棒性。

关键词: 不确定数据流 相似性度量 非参数估计 随机模拟

中图分类号:

文献标识码:

文章编号:

An Uncertain Data Stream Similarity Measurement Method Based on Nonparametric Estimation and Stochastic Simulation

College of Computer Science and Technology/Harbin Engineering University, Harbin

China

Abstract

Keywords

随着信息技术的飞速发展，信息数据量急剧增加，信息存储与处理工作变得日益繁重。在信息科学、通信工程、生物医学、金融工程等领域，数据流的相似性度量问题日益突出。例如，在信号处理、模式识别、数据挖掘、网络流量分析、系统故障诊断等方面，都需要对数据流的相似性进行度量。然而，由于数据流的不确定性、非平稳性和非线性等特点，使得相似性度量的难度大大增加。

为了对不确定数据流的相似性进行度量，本文提出了一种基于非参数估计与随机模拟的相似性度量方法。该方法首先对不确定数据流进行非参数估计，然后利用随机模拟技术生成具有相似性的数据流。仿实验结果表明，该方法不仅可以准确地度量不确定数据流的相似性，而且在数据流规模较大、分布复杂的情况下，依然具有良好的稳定性和鲁棒性。

本文首先介绍了不确定数据流相似性度量的背景和意义，然后详细阐述了基于非参数估计与随机模拟的相似性度量方法的原理和实现步骤。接着，通过仿实验验证了该方法的有效性。最后，对全文进行了总结和展望。

1 基于非参数估计的不确定数据流对象建模以及距离度量

1.1 不确定数据流对象建模

非参数估计在不假设数据流服从某种特定分布的前提下，基于数据流本身的信息进行建模。对于不确定数据流，其分布往往是未知的，因此非参数估计提供了一种有效的建模方法。本文主要研究基于非参数估计的不确定数据流对象建模问题。

收稿日期:

修回日期:

作者简介: 迟荣华，男，博士，主要从事不确定数据流及人工智能方面的研究工作。

字优先出

合。其中密估一基于，于估密函参，不先假变函关，其估如（）。

$$f_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$$

其中K为一个函，如函、三函、均匀函。密估密函够在与够大，到任一密函，因在分况下，为了取个不定对分，利参估取不定对分。

假不定U，包含n个不定对，任一个不定对u_i含l个，以不定如U={u_{1,l},...,u_{i,l}}, 0 ≤ i ≤ n。合上参估密估，

U中不定对u_i密可 为 $\hat{f}_{u_i}(x) = \frac{1}{nh} \sum_{j=1}^l K\left(\frac{x-u_{ij}}{h}\right)$ ，公一个

列中一个不定值密估函，个不定对包含了多少个，密函就多少个函。其中h函宽，一均定。

1.2 相似性度量

假两个不定对XY，其中U_x代不定对X，U_x={u_{1,l},...,u_{i,l}}, 0 ≤ i ≤ N，u_i于不定产到分合，大小为l，因存在不定|u_i| > 1，即元取值不唯一。在不定一中，假不定XY定义 dists(X,Y) = {L(x,y)|x ∈ U_x, y ∈ U_y}，L一，么XY定

值为: $PR(\text{distance}(X,Y) \leq \varepsilon) = \frac{|\{d \in \text{dists}(X,Y) | d \leq \varepsilon\}|}{|\text{dists}(X,Y)|}$

其中，d X与Y之任元之，假可，为任元之，即在固定匹中也可NI。从中不出，一够合不定似，一合似，但元中大历，则低。

借了，合参估，将不定似化为密似。密取值区域不同，但密中似，也好匹。假两个不定对U={u_{i,l}}和V={v_{i,l}}，可各密函F_U={f_{u_i}}和F_V={f_{v_i}}，0 ≤ i ≤ N，因任两个不定中元似可为:

$$\text{distance}(X,Y) = \frac{1}{N} \sum (|\int_{I_{u_i} \cup I_{v_i}} (\hat{f}_{u_i}(x) - \hat{f}_{v_i}(x)) dx \leq \varepsilon \cap \int_{I_{u_i} \cup I_{v_i}} |(\hat{f}_{u_i}(x) - \hat{f}_{v_i}(x))| dx \leq \delta|)$$

在公中，I_{u_i} ∪ I_{v_i}元u_i与v_i大取值，∫_{I_{u_i} ∪ I_{v_i} (f_{u_i}(x) - f_{v_i}(x)) dx密}

分，在I_{u_i} ∪ I_{v_i}区域中，一义上∫_{I_{u_i} ∪ I_{v_i} (f_{u_i}(x)) dx = 1、∫_{I_{u_i} ∪ I_{v_i} (f_{v_i}(x)) dx = 1，为了，公可化为:}}

$$\text{distance}(X,Y) = \frac{1}{N} \sum (|\int_{I_{u_i} \cap I_{v_i}} (\hat{f}_{u_i}(x) - \hat{f}_{v_i}(x)) dx \leq \varepsilon \cap \int_{I_{u_i} \cap I_{v_i}} |(\hat{f}_{u_i}(x) - \hat{f}_{v_i}(x))| dx \leq \delta|)$$

公可以出，似包含三个主因，取值、对以及对。I_{u_i} ∩ I_{v_i}代了取值，交为Φ似，∫_{I_{u_i} ∩ I_{v_i} (f_{u_i}(x) - f_{v_i}(x)) dx ≤ ε对，}

f_{u_i}(x)与f_{v_i}(x)取值完全同，对似为。∫_{I_{u_i} ∩ I_{v_i} |(\hat{f}_{u_i}(x) - \hat{f}_{v_i}(x))| dx ≤ δ|对，了分密几何上不同，可以互不。取值决定了在允}

一个不定中元，假取值I=[a,b]，ε，因：

$$1 - \left[\frac{1}{nh} \sum \int_{-\infty}^a K(\cdot) - \frac{1}{nh} \sum \int_{-\infty}^b K(\cdot) \right] = \varepsilon$$

含义参密函信区。一况下，假密函对，因可 $\frac{\varepsilon}{2} =$

$\frac{1}{nh} \sum_{-\infty}^a K(\cdot)$, $\frac{\varepsilon}{2} = 1 - \frac{1}{nh} \sum_{-\infty}^b K(\cdot)$ 。于, $x_1 \leq x_2 \leq \dots \leq x_n$, 可 小值与 大值 取值 区域 似值, 令 $\frac{\varepsilon}{2} = \frac{1}{nh} \sum_{-\infty}^a K\left(\frac{x-x_1}{h}\right)$, $\varepsilon = \frac{2}{h} \int_{-\infty}^a K\left(\frac{x-x_1}{h}\right)$ 以及 $\frac{\varepsilon}{2} = 1 - \frac{1}{nh} \sum_{-\infty}^b K\left(\frac{x-x_n}{h}\right)$, $\varepsilon = 2 - \frac{2}{h} \int_{-\infty}^b K\left(\frac{x-x_n}{h}\right)$

复合函 导 , $\varepsilon = \frac{2}{\sqrt{2\pi}h\sigma} \int_{-\infty}^a e^{-\frac{1}{2\sigma^2}\left(\frac{x-x_1}{h}\right)^2} dx$

令, $Z = \frac{x-x_1}{h}$, 则 $\varepsilon = \frac{2}{\sqrt{2\pi}h\sigma} \int_{-\infty}^{\frac{a-x_1}{h}} e^{-\frac{z^2}{2\sigma^2}} h dx$, 可 $\varepsilon = 2\phi\left(\frac{a-x_1}{h\sigma}\right)$

, 可 到 a , 同 可 b 。且, $\varepsilon_a = \frac{2}{nh} \sum_{-\infty}^a K(\cdot)$, $\varepsilon_b = 2 - \frac{2}{nh} \sum_{-\infty}^b K(\cdot)$

可 $\varepsilon_a, \varepsilon_b \leq \varepsilon$, 因 可 少 估 密 函 取值区 I 。假 不 定对 元 u 与 v , 其 实 密 函 与 估 密 函 分别 为 f_u, f_v 与 $\widehat{f}_u, \widehat{f}_v$, 在 ε 件 可 U 与 V 信 区 分别 为 I_u 和 I_v 。可 对 与 对 分别 为:

$$\int_{I_u \cap I_v} (f_u - \widehat{f}_v) dx$$

$$\int_{I_u \cap I_v} |f_u - \widehat{f}_v| dx$$

于 $f_u f_v > 0$, 且 $\int f_u = 1, \int f_v = 1$, 则可 $0 \leq \int_{I_u \cap I_v} (f_u - \widehat{f}_v) dx \leq 1$

$\int_{I_u \cap I_v} |f_u - \widehat{f}_v| dx \leq \int_{I_u \cap I_v} (f_u + \widehat{f}_v) dx \leq 2$, 则可 $0 \leq \int_{I_u \cap I_v} |f_u - \widehat{f}_v| dx \leq 2$

公 可, 两个 不 定对 取值区域 完全 不同, 对 和 对 分别 为 和, 取值区域 完全 同, 对 为, 对 可 取 任何 值。因, 可以 即使 两个 对 在 取值 完全 同, 仍可 于 分 不同 产

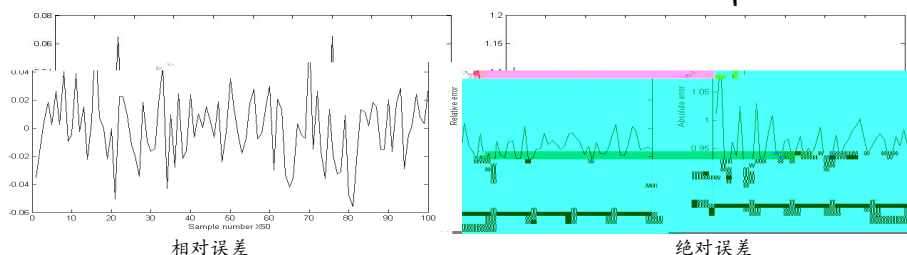


图 1 不同样本规模下的相对误差 (a) 与绝对误差 (b) 比较

Fig.1 Comparison between relative error (a) and absolute error (b) at different sample size

为了分 不 定对 密 对 与 对, 了两个 同取值, 但分 完全, 且不同 不 定对 对 与 对 变化, 从图 可以 出, 在 下, 对 基 定在 密, 但 对 可 体 出不同 变 化。在 可以, 不 定对 可以 密 对 与 对 准, 但如何 同 一个复。

2 仿真实验分析

为了从实 可, 仿 实。先, 分别 三 仿, 对 分 分别为 分、均匀分 和 分, 个 值, 值 取值 为, 分 均值为 为, 分 均值为。分别对 三 不 定对。密 函 宽 参 估 一个, 宽 合 与 否 决定 了 密 函。前 为 基于 均

假 不 定对 实 密 函 $f(x)$ 与 估 密 函 $\widehat{f}(x)$ 均

$$MSE(x;h) = E\{\widehat{f}(x) - f(x)\}^2$$

可 $h_{opt} = 1.06 \times \sigma \times n^{-\frac{1}{5}}$ (前 假 实 密 函)。

先, 为了分 在 不 定对 属于 同区域 分 密 不同, 了 属于 两个 不 同分 型 不 定对 分 属于 不同 别, 对 与 对 变化 况。如图 可, 两个 对 在 同区域 内 密 几乎 完全 同, 分 型 不同, 对

仍可导 两个对 具 大 别。 后， 不同 代 ， 察分 三分 变化 况如图 。

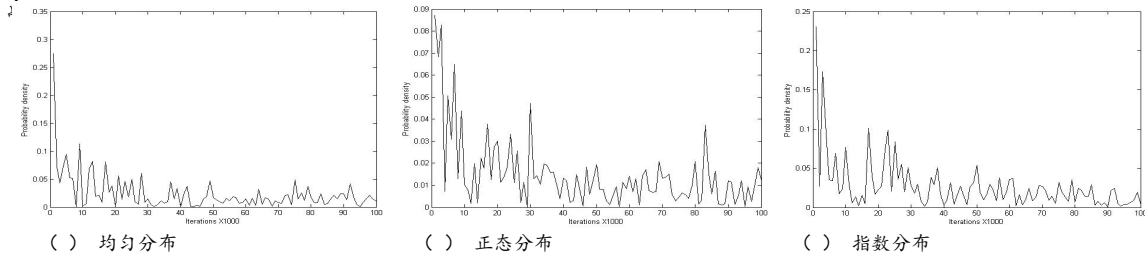


图 2 三种典型分布下的模拟效果示意图

Fig.2 Simulation performances of three typical distributions

仿 实 可 均匀分 ， 分 之， 分 小， 原 因可 与 函 假 关，但 可 以 制在 以下，实 ， 且 多 实 ， 依 定。

3 总 结

向 属 存在 不 定 ， 对不 定对 以及不 定对 似 中 临 主 ， 出基于 参估 及 合 似 。 利 参 估 对不 定对 ， 可以 因为属 以及存在 导 不 定 ， 同 又 决了 参 估 函 复 于 似 。 后， 仿 实 对 了 ， 实 不仅 够 地 对 ， 同 又 在 下 地 不 定对 似 ， 和实 了 与准 。

参 考 文 献

et al

et al

et al

et al