

面向大数据流的多任务加速在线学习算法*

李志杰, 李元香, 王峰, 匡立

武汉大学软件工程国家重点实验室, 武汉市 430072

Big Data Stream Oriented Multi-Task Accelerated Online Learning Algorithm

Li Zhijie, Li Yuanxiang, Wang Feng, Kuang Li

(State Key Laboratory of Software Engineering, Wuhan University, Wuhan 430072, China)

Abstract Multi-task online learning with stream computing mode is a promising tool for big data stream analyses. However, the convergence rate of current multi-task online algorithm is only $O(1/\sqrt{T})$ up to T -th iteration, and its low convergence rate has become a bottleneck of algorithm performance. This paper propose a novel multi-task accelerated online learning algorithm which obtain low computational time complexity and optimal convergence rate $O(1/T^2)$. The proof of a closed-form solution theorem which efficiently updates the weight matrix W_t at each iteration is provided, and the detailed theoretical analysis for algorithm convergence is conducted. Experimental results show that proposed multi-task accelerated online learning algorithm can improve real-time performance and scalability of big data stream processing, and it is a realistic method for big data stream analyses.

Key words big data stream; multi-task; accelerated; online learning; convergence analysis

摘要 多任务在线学习框架采用直接数据处理的流式计算模式, 是大数据流分析很有前途的一种工具。然而目前的多任务在线学习算法收敛率低, 仅为 $O(1/\sqrt{T})$, T 为算法迭代次数。提出一种新颖的多任务加速在线学习算法, 在保持多任务在线学习快捷计算优势的基础上, 达到最优收敛率 $O(1/T^2)$ 。对多任务权重学习矩阵 W_t 的迭代邻近解表达式进行了推导, 对提出算法的收敛性进行了详细的理论分析。实验表明, 提出的多任务加速在线学习算法能够更好地保障大数据流处理的实时性和可伸缩性, 有较广泛的实际应用价值。

关键词 大数据流; 多任务; 加速; 在线学习; 收敛分析

中图法分类号 TP311

0 引言

*本课题得到国家自然科学基金(Nos. 61070009, 61103125); 国家高技术研究发展计划(863)(No. 2007AA01Z290)资助

云计算、物联网、社交网络等新兴信息技术和应用模式的快速发展,越来越多的软件系统部署在以互联网为代表的开放、动态、难控和不确定的变化环境,产生了急剧增加的大数据流。大数据流以其实时性、无序性、无限性、易失性、突发性等显著特征,使得其与传统批量大数据在数据计算的要求、方式等方面有着明显的不同^[1]。大数据的批量计算首先进行数据的存储,然后再对存储的静态数据进行集中计算。Hadoop 是典型的大数据批量计算框架,由 HDFS 分布式文件系统负责静态数据的存储,并通过 MapReduce 将计算逻辑分配到各数据节点进行数据计算和价值发现。而在流式计算中,无法确定数据的到来时刻和到来顺序,也无法将全部数据存储起来。因此,不再进行流式数据的存储,而是当流动的数据到来后在内存中直接进行数据的实时计算。如 Twitter 的 Storm、Yahoo 的 S4 就是典型的流式数据计算框架。因此,如何通过创新性的数据分析方法实现对大数据流的快速、高效、及时的分析与计算,是大数据分析技术领域面临的新挑战^{[1][2]}。

在围绕实时大数据流处理这一需求展开的研究中,在线机器学习算法采用数据流直接处理的模式,每次迭代处理一个随机流数据,学习变量的迭代更新只经过简单的计算,从而在实时性和准确率之间取得一个平衡,是解决该难题很有前途的方案,特别适合于训练流式大数据^[3]。

研究表明^[4-12],许多大数据实际应用领域,如金融时间序列预测,自然语言处理,生物信息学多任务数据集分析等,同时学习多个相关的任务获取共享信息模式要优于单个任务分别学习的方式。还有一些应用,如协同过滤的用户特征矩阵优化与物品特征矩阵优化,矩阵分类等,优化对象本身是矩阵变量,直接类同于多任务学习算法的权重矩阵变量的求解。因此,近年来许多相关研究侧重于多任务学习框架,如文献[3~14]。不过这些工作大多是批训练模式,在线多任务学习框架的研究工作很少。

文献[15]提出的正则化对偶平均方法(DA-STL, regularized dual averaging method for single task learning),以凸函数优化为基础,通过迭代计算过去所有损失函数梯度值的平均值,降低优化目标表达式的求解复杂度,达到高效更新权重向量的目的。正则化对偶平均方法是单任务在线学习领域的标志性研究工作。基于正则化对偶平均思想, Yang 提出了一种多任务的在线学习框架^[3,14]并应用于多任务特征选择(DA-MTFS, dual averaging based multi-task feature selection)。DA-MTFS 算法收敛率为 $O(1/\sqrt{T})$,每次

迭代时间与空间复杂度均为 $O(dQ)$ 。其中, T 为算法迭代次数, d 为数据特征维度, Q 是任务个数。显然,该算法依然存在在线算法面临的普遍问题:低计算复杂性常常与算法的低收敛率联系在一起。

加速在线学习方面,文献[15]提出了正则化对偶平均在线学习框架并讨论了加速分支的情况。文献[16]提出了随机优化与在线学习的加速梯度方法。不过,这些算法的加速效果取决于参数的满足条件,极端情况下才能获得最优收敛率 $O(1/T^2)$,并且讨论的都是单任务在线算法的加速效果。

文献[17]提出了一种迹范数最小化加速梯度方法,在理论上证明了多任务学习算法的最优收敛率可以达到 $O(1/T^2)$ 。不过,该算法是基于批处理模式。

并且,由于应用了矩阵的奇异值分解 SVD,每次迭代的昂贵计算代价使得该方法并不适用于大规模数据场景。

本文结合使用加速技术和正则化对偶平均方法,首次提出一种新颖的多任务加速在线学习算法,算法计算高效且加速收敛。迭代更新权重学习矩阵的邻近解计算表达式,以及算法的收敛率分析,均有完整的数学推导过程。多任务加速在线算法可以很方便地应用于金融时间序列预测、在线推荐系统协同过滤等动态多变的大数据流场景中,实验结果表明它能显著提升大规模数据流处理的实时性和可伸缩性,是有效学习大数据流的一种很有前途的工具。

1 多任务在线学习框架

1.1 问题构建

定义 1. 多任务学习 (MTL, multi-task learning). 假设有 Q 个任务,它们的数据来自于同一空间 $X \times Y$, 其中, $X \subset R^d, Y \subset R$ 。每个任务有 N_q 个数据点 $D_q = \{z_i^q = (x_i^q, y_i^q)\}_{i=1}^{N_q}$, 它们组成多任务数据集 $D = \cup_{q=1}^Q D_q$ 。 D_q 取样于 $X \times Y$ 空间的分布 P_q , 每个任务的 P_q 各不相同,但这些不同任务的 P_q 是相关的。多任务学习的目标是学习 Q 个函数 $f_q: R^d \rightarrow R, q=1, \dots, Q$, 使得 $f_q(x_i^q)$ 估计 y_i^q 。当 $Q=1$, 它是标准的单任务学习问题。

通常,在多任务学习模型中,第 q 个任务的决策函数 f_q 假定是以模型权重向量 w_q 为参数的一个超平面,即, $f_q(x) = w_q^T x, q=1, \dots, Q$ 。 Q 个权重向量 w^q 组

成大小 $d \times Q$ 的矩阵 W ,

$$W = (w^1, w^2, \dots, w^Q) = (W_{\bullet 1}, \dots, W_{\bullet Q}) = (W_{1\bullet}^T, \dots, W_{d\bullet}^T)^T.$$

这里, d 表示数据特征维度。

权重矩阵 W 是多任务学习的目标。多任务学习是通过最小化经验风险以及权重的正则化项来完成学习权重矩阵 W 的目标。批训练方式的多任务学习目标 W 优化表达式如下,

$$\min_W \phi(W) = \sum_{q=1}^Q \frac{1}{N_q} \sum_{i=1}^{N_q} l^q(W_{\bullet q}, z_i^q) + \Omega_\lambda(W). \quad (1)$$

其中, $\lambda \geq 0$ 是平衡损失与正则化项的一个常数,

$l^q(W_{\bullet q}, z_i^q)$ 定义第 q 个任务的样本点 $z_i^q = (x_i^q, y_i^q)$

的损失。本文损失函数设定为平滑可导的强凸函数,

$$l(w, z) = \frac{1}{2} [y - w^T x]^2 \quad (2)$$

基于对偶平均的在线算法, 采用直接数据处理的流式计算模式, 多任务学习目标 W 优化表达式如下,

$$\min_W \phi(W) = \left\{ \langle \bar{G}_t, W \rangle + \Omega_\lambda(W) + \frac{\beta_t}{t} h(W) \right\}. \quad (3)$$

其中, 对偶变量 $G_t = \partial l_t$ 表示第 t 次迭代的损失

函数 l_t 对原变量 W 的偏梯度。 \bar{G}_t 则表示第 1 到第 t 步迭代的偏梯度的平均值。 $h(W)$ 是附加的强凸函数, $\{\beta_t\}_{t \geq 1}$ 是非递减的非负输入序列。

因此, 对偶平均在线算法的每次迭代过程主要由三步组成: 1. 计算损失函数关于权重的偏梯度; 2. 计算偏梯度的平均值; 3. 基于平均偏梯度计算更新权重变量。

对于损失函数(2), 单个任务偏梯度计算式如下,

$$\partial l_t(w) = \partial l(w, z_t) = (w^T x_t - y_t) x_t \quad (4)$$

其中, w 表示这个任务的权重向量, (x_t, y_t) 代表该任务在第 t 次迭代时的特征向量和响应值。

1.2 基于对偶平均的多任务特征选择在线算法

Yang 在文献[3]提出的基于对偶平均的多任务特征选择在线算法 DA-MTFS 是典型的多任务在线学习框架, 为便于比较与启发思路, 我们在算法 1 中概述其基本思想。

算法 1 中, $\|\bullet\|_F$ 表示矩阵的 Frobenius 范数, 正

则化项 $\Omega_\lambda(W)$ 取权重矩阵 W 混合范数 $L_{1/2,1}$, 由 W

的 $L_{1,1}$ 范数 $\|W_{j\bullet}^T\|_1$ 和 $L_{2,1}$ 范数 $\|W_{j\bullet}^T\|_2$ 线性组合而成,

$$\Omega_\lambda(W) = \lambda \sum_{j=1}^d (c \|W_{j\bullet}^T\|_1 + \|W_{j\bullet}^T\|_2) \quad (5)$$

c 是控制任务选择稀疏性的一个常数。 $\Omega_\lambda(W)$ 中的

$L_{1/2,1}$ 混合范数具有特征与任务选择的功能。算法 1 的权重矩阵 W 更新优化表达式(6), 由多任务目标优化

式(3)的输入序列取 $\beta_t = \gamma \sqrt{t}$ 得到。算法采用流式计算模式, 数据在内存直接处理, 这就要求优化式(6)的求解简单化, 从而获得在线算法急需的计算效率。

由于根据优化式(6), 推导出了权重矩阵 W 更新的邻近解(closed-form solution)计算表达式, 这样, 优化式(6)每次更新计算 W 的时间变为常数。算法 1 的一次迭代时间复杂度是 $O(dQ)$, 收敛率 $O(1/\sqrt{T})$, 其基本特征是计算高效但收敛率低。

算法 1: 多任务特征选择在线学习框架 DA-MTFS

输入:

$$W_0 = \operatorname{argmin}_W h(W), \quad h(W) = \frac{1}{2} \|W\|_F^2;$$

给定常数 $\lambda > 0, \gamma > 0$ 。

初始化: 令 $W_t = W_0, \bar{G}_0 = 0$ 。

for $t = 1, 2, 3, \dots$, do

1) Q 个任务依次出现一个实例 $z_t^q, q = 1, \dots, Q$,

根据(4)式分别计算 $g_t^q = \partial l(w^q, z_t^q)$ 值, 构成

一个 $d \times Q$ 矩阵, $G_t = (g_t^1, \dots, g_t^Q)$ 。

2) 计算偏梯度的平均值,

$$\bar{G}_t = \frac{t-1}{t} \bar{G}_{t-1} + \frac{1}{t} G_t.$$

3) 基于 \bar{G}_t 计算更新下次迭代的权重矩阵 W_{t+1} ,

$$W_{t+1} = \operatorname{arg min}_W \phi(W),$$

$$\phi(W) = \langle \bar{G}_t, W \rangle + \Omega_\lambda(W) + \frac{\gamma}{\sqrt{t}} h(W) \quad (6)$$

end for

2 加速的多任务在线学习算法

2.1 加速的多任务在线学习框架

我们在正则化对偶平均方法框架内, 结合使用加速技术的微批量方法(mini-batch approach), 提出一种新颖的多任务加速在线学习框架 ADA-MTL(accelerated dual averaging method for multi-task learning).

算法 2: 加速的多任务在线学习框架 ADA-MTL

输入:

$$W_0 = \operatorname{argmin}_W h(W), \quad h(W) = \frac{1}{2} \|W\|_F^2;$$

给定常数 $\lambda > 0$.

初始化: 令 $U_1 = W_0, \bar{G}_0 = 0, \alpha_1 = 1$.

for $t = 1, 2, 3, \dots$, do

1) Q 个任务依次出现一个实例 $z_t^q, q = 1, \dots, Q$, 计算损失函数偏梯度在查询点 U_t 的值 G_t .

2) 计算偏梯度的平均值, $\bar{G}_t = \frac{t-1}{t} \bar{G}_{t-1} + \frac{1}{t} G_t$.

3) 输出,

$$W_t = \operatorname{arg min}_W \phi(W, U_t) = \left\{ \langle \bar{G}_t, W \rangle + \Omega_\lambda(W) + \frac{1}{2} \|W - U_t\|_F^2 \right\} \quad (7)$$

4) 更新输入序列, .

$$\alpha_{t+1} = \frac{1 + \sqrt{1 + 4\alpha_t^2}}{2}. \quad (8)$$

5) 计算查询点 U_{t+1} ,

$$U_{t+1} = W_t + \left(\frac{\alpha_t - 1}{\alpha_{t+1}} \right) (W_t - W_{t-1}). \quad (9)$$

end for

比较算法 1 和算法 2, 不难看出, 算法 2 增加了一个查询点 U_t , 每次迭代计算的 G_t 是损失函数偏梯度在查询点 U_t 的值。查询点由(9)式调整, 其中 α_t 是由(8)

产生的输入序列, 其作用是使查询点 U_t 尽量靠近权重矩阵 W_t 。必须指出, 与算法 1、文献[15]的加速算法等其它算法不同的是, 我们在算法 2 的优化式(7)中, 附加的 $h(W)$ 强凸函数形式已由 $\frac{1}{2} \|W\|_F^2$ 改为

$\frac{1}{2} \|W - U_t\|_F^2$ 。因为, 这时不仅需要增加优化函数的

凸性, 更重要的是, 通过把 $\frac{1}{2} \|W - U_t\|_F^2$ 融入算法的优化目标, 使得优化变量 W 与查询点 U_t 不断靠近, 从而获得算法持续加速的必要条件。因此, 我们在算法 2 中采用的加速技术是一种改进的微批量方法, 加速效果更好。

2.2 求解 W_t 的邻近解

算法 2 的计算复杂度主要取决于优化式(7)权重矩阵 W_t 的求解。我们通过推导出优化式(7)的 W_t 邻近解, 高效迭代更新权重矩阵 W_t 。

定理 1. 已知迭代 t 对偶平均梯度 \bar{G}_t , 正则化项

$$\Omega_\lambda(W) = \lambda \sum_{j=1}^d (c \|W_{j\cdot}^T\|_1 + \|W_{j\cdot}^T\|_2), \quad \text{如(5)所示。令}$$

$(\bar{R}_{j,q})_t = [(\bar{G}_{j,q})_t - (u_q)_t | -\lambda c]_+ \cdot \operatorname{sign}((\bar{G}_{j,q})_t - (u_q)_t),$
 $j=1, \dots, d, \quad q=1, \dots, Q$ 。则算法 2 可用下面(10)式计算的邻近解来高效更新优化式(7)的目标 W_t , 其中, 运算 $[x]_+ = x, \text{ if } x > 0; [x]_+ = 0, \text{ if } x \leq 0$ 。

$$(W_{j\cdot})_t = \left[1 - \frac{\lambda}{\|(\bar{R}_{j\cdot})_t\|_2} \right]_+ \cdot (\bar{R}_{j\cdot})_t. \quad (10)$$

定理 1 的证明见附录 A。

由于定理 1 能够为算法 2 的优化式(7)提供 W_t 每次迭代更新的邻近解(10), 这样, 优化式(7)每次更新计算 W_t 的时间变为常数。同时, 对偶偏梯度平均值 \bar{G}_t 充分利用上次迭代计算与存储过的平均值

\bar{G}_{t-1} , 计算简单。算法 2 每次迭代只需 $O(dQ)$ 空间来存储要求的信息, 平均偏梯度和权重矩阵。算法每次迭代的时间复杂度也是 $O(dQ)$, d 为数据特征维度, Q 是任务个数。因此, 算法 2 与算法 1 一样, 计算高效。与算法 1 不同的是, 算法 2 能够加速收敛, 通过下节的算法收敛理论分析, 我们证明算法 2 的收敛率可以达到最优值 $O(1/T^2)$ 。

2.3 算法收敛理论分析

正则化随机多任务学习是对随机多任务样本 Z 的期望损失最小化, 加上正则化项 Ω_λ ,

$$\min_W \phi(W) = E_Z l(W, Z) + \Omega_\lambda(W) \quad (11)$$

$\varphi(W)$ 称为多任务学习的目标值。现实中,由于 Z 的分布一般难以确定,期望损失 $E_Z l(W, Z)$ 常由经验平均值 $\sum_{q=1}^Q \frac{1}{N_q} \sum_{i=1}^{N_q} l^q(W_{\bullet,q}, z_i^q)$ 代替,如(1)所示。而

对于多任务对偶平均方法,期望损失则用 $\langle \bar{G}_t, W \rangle$ 代替如下,

$$\varphi(W) = \langle \bar{G}_t, W \rangle + \Omega_\lambda(W) \quad (12)$$

定义 2. 多任务在线学习收敛率(CR-MTOL, convergence rate for multi-task online learning). 多任务在线学习目标值 $\varphi(W)$ 定义如(12), 如果算法优化问题 $\min_W \varphi(W)$ 存在最优解 W^* , 那么, 该算法运行到第 T 步的收敛率为 v_T ,

$$v_T = \varphi(W_T) - \varphi(W^*) \quad (13)$$

显然, 如果算法的收敛率大, 则算法效率高, 迭代次数少。

定理 2. 设算法 2 优化问题 $\min_W \varphi(W)$ 的最优解是 W^* , 则算法 2 运行到第 T 次迭代的收敛率为 v_T ($T \geq 1$),

$$v_T = \varphi(W_T) - \varphi(W^*) \leq \frac{2 \|W_0 - W^*\|_F^2}{(T+1)^2},$$

$$\text{即, } v_T = O\left(\frac{1}{T^2}\right).$$

定理 2 的证明见附录 B。

3 评价与应用

本节我们分析比较的方法包括本文提出的多任务加速在线学习算法 ADA-MTL、文献[3]的多任务在线特征选择算法 DA-MTFS、文献[15]的正则化对偶平均单任务在线学习方法 DA-STL、和文献[9]的多任务批训练学习算法 BT-MTL(batch-training method for multi-task learning)。对提出的多任务加速在线学习算法 ADA-MTL, 研究大数据流背景下, 在金融时间序列预测(FTSP, financial time series prediction)场景的应用 ADA-MTL for FTSP, 和在协同过滤的概率矩阵分解(PMF, probabilistic matrix factorization)在线学习场

景的应用 ADA-MTL for PMF。分析评价相应的大数据流挖掘系统的实时性、可伸缩性、数据吞吐量等性能指标。数据流抽取(data stream sampling) 算法采用环形循环滑动窗口方法^[22]。实验环境为 2.80GHz CPU, 3.93GB 内存 PC。所有算法在 Matlab 上运行。

3.1 时空代价分析

大数据流的数据规模大, 到达速率非常快, 要求数据流挖掘算法在有限的内存空间中实时处理, 这就要求面向大数据流分析算法的时间、空间复杂度低。

假设有 Q 个任务, 每个任务有 k 个 d 维的样本, 文献[9]的多任务批处理算法 BT-MTL 的浮点操作数(flops)的界为: $O(dkQ+dQ)=O(dkQ)$ 。由于我们这里考虑的训练样本是流数据, 因此, 当一个新样本出现时, 批处理算法必须重新训练。当每个任务的样本数达到 N 时, 则总浮点操作数的界是,

$$O\left(\sum_{k=2}^N dkQ\right) = O(dQN^2). \quad (14)$$

该多任务批训练算法的存储代价界是,

$$O(dQN). \quad (15)$$

不同的是, 在线算法 DA-MTFS 和 ADA-MTL 每次迭代的最坏时间复杂度为 $O(dQ)$, 并且当样本一个个顺次出现时它能相应地更新, 不必重新训练。这样, 当每个任务的样本数是 N , 总浮点操作数是,

$$O(dQN). \quad (16)$$

多任务在线学习算法的存储代价界是,

$$O(dQ). \quad (17)$$

比较式(14)、(15)和(16)、(17), 多任务的在线学习算法相对于多任务批训练算法耗费的时空代价分别下降了一个数量等级。图 1 画出了四个算法的运行时间—训练样本数的 Log-log 对数关系图。

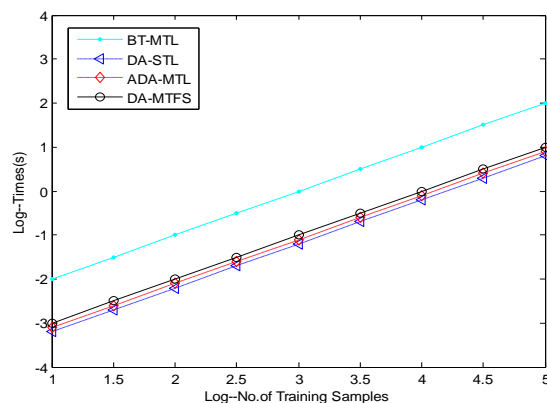


Fig. 1 Log-log plot of computation time on training samples.

图 1 运行时间-训练样本数的 Log-log 对数关系图

3.2 收敛性

为了验证算法 2(ADA-MTL) 相对于算法 1(DA-MTFs)的加速收敛效果, 我们选择四个多任务数据集: yeast, letters, digits, 和 DMOZ。yeast 数据集由一个酵母基因分类问题得出^[18], 有 14 个任务; letters、digits 分别是手写体单词(8 个任务)和数字数据集^[18](10 个任务); DMOZ 是文本分类数据集(10 个任务, <http://www.dmoz.org/>)。我们随机抽取 10% 的各个任务的数据做训练样本。抽取数据流采用环形循环滑动窗口方法。

为了评估算法的效率, 所有算法当目标值的相对变化小于 10^{-8} 时终止运行, 所耗费的时间我们称之为收敛时间(convergence time)。图 2 比较了两个多任务在线算法 DA-MTFs 和 ADA-MTL 在 yeast, letters, digits, 和 DMOZ 四个数据集上分别进行十次随机试验的平均收敛时间。我们观察到加速算法 ADA-MTL 在四个数据集所需的收敛时间均大大小于 DA-MTFs 算法, 这是由于 ADA-MTL 收敛快, 趋于稳定值的时间大大节省所致。

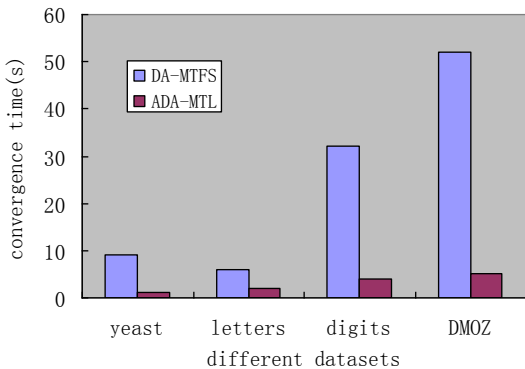


Fig. 2 Comparison of the convergence times(in seconds).

图 2 算法 2(ADA-MTL)与算法 1(DA-MTFs)收敛时间比较。

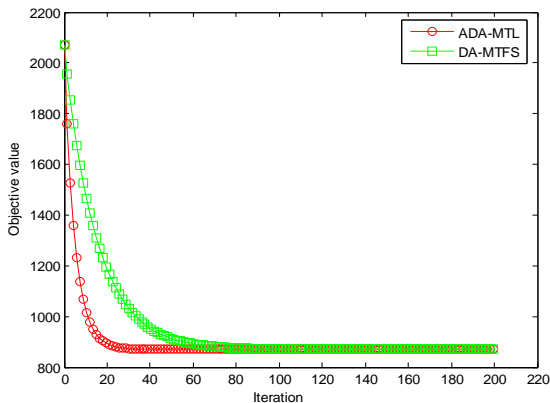


Fig.3 The convergence of algorithms on yeast data set.

图 3 算法 1 与算法 2 在 yeast 数据集的收敛性比较。

为了进一步观察比较算法 1 和算法 2 的收敛特性, 图 3 画出了这两个多任务在线学习算法在 yeast 数据集上运行时的目标值迭代曲线。我们观察到 ADA-MTL 仅需 30 次左右迭代趋于稳定值, 而 DA-MTFs 需 120 次左右的迭代才能收敛至稳定值。

ADA-MTL 的加速收敛与低时间复杂度特性, 带来优异的实时性能。假设数据流挖掘系统要求新产生的数据流在 T 时间内处理完毕并反馈处理结果, 我们称之为实时度 T 。显然, 实时度 T 越小, 对算法的实时处理能力要求越高。如果数据流处理能力不能满足实时度 T 的要求, 那么某些流数据会因为不能得到及时处理而丢失。我们观察到当实时度 T 低于某个阈值, 例如 $T=0.5s$, DA-MTFs, 尤其是 DA-STL, 流数据实际平均处理率 P 开始显著下降。而这时 ADA-MTL 基本不变, 维持近 100% 的流数据处理率。这说明 ADA-MTL 的实时性能明显好于 DA-MTFs, 尤其是 DA-STL。ADA-MTL 优良的实时性能确保了它适应随机动态的大数据流场景的能力。

3.3 金融数据流回归分析

金融数据分析是一种典型大数据流分析。股票价格波动预测, 属于大数据回归分析的范畴, 是金融领域研究的焦点问题之一。我们将多任务在线回归分析算法应用于股市预测当中, 并对 ADA-MTL、DA-MTFs、DA-STL、BT-MTL 等算法在金融时间序列预测中的性能进行比较, 以验证我们提出的模型与在线算法 ADA-MTL 在金融数据流回归分析中的特性与优势。

中国股市自 1990 年 12 月 19 日正式营业至今已有二十多年的历史, 交易数据量巨大。我们选取上证综合指数和四川长虹、民生银行等个股进行实例研究, 每个实例视为一个任务, 共 100 个任务。股票历史交易数据和实时交易数据从新浪财经网下载。训练数据集选取 2009 年 4 月 13 日到 2010 年 3 月 30 日, 一共 240 个交易日的数据, 测试集选取 2010 年第二、三季度的交易日里的数据, 一共 120 个交易日的数据。使用证券交易信息作为任务的数据特征, 具体信息见表 1。

Table 1 Data feature

表 1 数据特征信息表

1.transdate	2.openprice	3.lastprice	4.topprice	5.lowprice
6.DIFF	7.DEA	8.MACD	9.K	10.D
11.J	12.WR1	13.WR2	14.WR3	15.volume
16.EMA12	17.EMA26	18.avg5	19.avg10	20.avg20
21.avg30	22.BBI	23.UPR	24.DWN	25.No
26.EMPMA	27.EMPMA	28.EMPMA	29.EMPMA	30.EMPMA

下载的原始数据集各特征值先进行预处理,即归一化处理,将其缩放到[-1,1]之间,再抽取数据流。算法里的参数可以利用 LIBSVM 工具的 gridregression.py 函数进行参数寻优过程。实验的方案是多任务流数据在线学习,训练样本顺次到达,算法迭代时每个任务都有一个待处理的样本。表 2 列出各比较算法在选定数据集上回归分析的性能平均值,这里使用预测误差 RMSEs(root mean square errors)和权重矩阵的 NNZs(number of non-zeros)值,以及数据吞吐量(throughout, /s)来衡量模型与算法性能。其中,预测误差 RMSEs 是真实值与预测值误差的平均平方根。两个无数据的单元格表示这些项的数据意义不明显。

Table 2 RMSEs and NNZs of the regression analysis
表 2 选定数据集上回归分析的 RMSEs 和 NNZs 值

Methods	RMSEs	NNZs	Throughout	Parameters
BT-MTL	2.15	2981	—	$\lambda=4.2, c=0.01$
DA-STL	2.16	—	151	$\lambda=0.8, c=0.01$
DA-MTFS	1.81	2059	185	$\lambda=6.9, c=0.01$
ADA-MTL	1.81	2129	204	$\lambda=6.9, c=0.01$

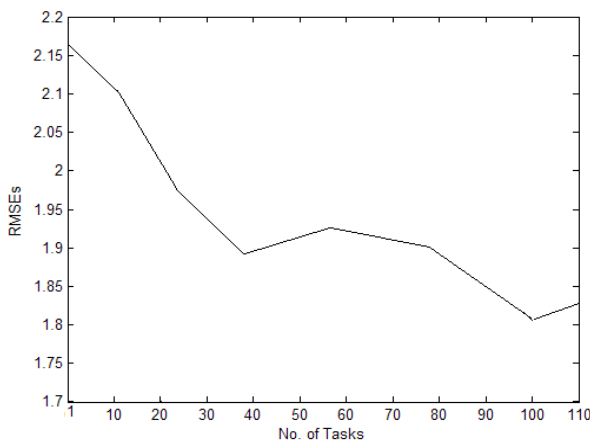


Fig. 4 RMSEs vs. number of tasks

图 4 ADA-MTL 的预测性能 RMSEs 值与任务个数的关系

图 4 画出了 ADA-MTL 的预测性能 RMSEs 值与任务个数的关系。

从表 2 和图 4 我们观察到:

- 1) 多任务学习方式($Q=100$)比各个任务单独学习方式的回归预测性能要好 20%左右;
- 2) 在线学习算法和相应批处理算法性能相近;
- 3) ADA-MTL 回归预测性能取得最好值,并且 NNZs 值仅高于 DA-MTFS 方法,权重矩阵解的稀疏性较强。
- 4) 表 2 第 4 列是在数据流速 210/s、实时度 $T=0.3s$ 条件下,三个在线算法的流数据吞吐量(throughout)。

ADA-MTL 吞吐量最高,接近于数据流速。DA-MTFS 次之,DA-STL 最低。

3.4 协同过滤在线学习

协同过滤 (CF, collaborative filtering) 将用户评价物品作为一种偏好指示,根据已知的某些用户的偏好预测用户的未知偏好,是构建推荐系统的一种主要方式。随着 MovieLens, Yahoo!Music, Amazon 等大规模在线用户贡献网站和在线购物网站的出现,用户和物品库急剧扩大,系统始终处于随机获取新的三元组(u, i, r),即新的用户、新的物品和新的评价的动态场景中,呈现大数据流式计算的特征。这里,我们使用的协同过滤技术是概率矩阵分解^[18-20] (PMF, probabilistic matrix factorization)。PMF 的目标是学习两个低秩矩阵,用户特征矩阵 U 和物品特征矩阵 V ,使用 $U^T V$ 拟合评价矩阵 R 。因此,PMF 的优化矩阵是用户特征矩阵 U 和物品特征矩阵 V ,对应于算法 2 的权重矩阵 W ,可以分别使用算法 2 学习优化。我们选择数据集 Yahoo!Music(<http://kddcup.yahoo.com>)和 MovieLens(<http://www.cs.umn.edu/Research/Group-Lens>) 研究所比较算法的实验性能。Yahoo!Music 是一个很大的数据集,包含超过 250 百万的评分值。但这个数据集非常稀疏,只有 0.4%的单元是已知的。表 3 给出了这两个数据集的基本统计。

Table 3 Statistics of datasets

表 3 数据集的基本统计

	MovieLens	Yahoo!Music
No. ratings	1,000,209	252,800,275
No. users	6,040	1,000,990
No. items	3,952	624,961
Rating range	[1,5]	[0,100]

$$\text{我们采用 } RMSE = \sqrt{\sum_{(u,i,r) \in Z} (\hat{r}_{u,i} - r)^2 / |Z|} \text{ 来评估}$$

法的性能。预测误差 RMSE (root mean square error) 是真实评分值与预测评分值误差的平均平方根。

为了评估本文的流数据在线挖掘算法的可伸缩性,我们的实验有下列三种设置: **1) T1:** 随机抽取所有(u, i, r)三元组的 10%训练,使用剩余的 90%作评估使用; **2) T5:** 随机抽取所有(u, i, r)三元组的 50%训练,使用剩余的 50%作评估使用; **3) T9:** 随机抽取所有(u, i, r)三元组的 90%训练,使用剩余的 10%作评估使用。

表 4 记录了各种不同设置下在线与批训练 PMF 算法的 RMSEs 性能值。

从实验中,我们观察到下列现象:

1) 总的来说, 在线 PMF 算法在不同数据集上的性能与批训练算法相近。

2) 在 T1 设置下, 在线算法的性能甚至要超出批训练算法一点点, 这可能是由于下面事实引起: 在很少训练样本的场景下, 相对于批训练算法, 在线学习算法陷入局部最优解的可能性要小一些。

3) 由于 Yahoo!Music 数据集的巨大数据, 我们不能使用 T9 设置执行批训练算法 BT-MTL for PMF。在 T5 设置下, BT-MTL for PMF 完成 120 次迭代收敛用了 6 个多小时, 而使用我们提出的多任务加速在线算法的 ADA-MTL for PMF 只用大约 1 分钟完成所有 180 百万个评分的处理并达到相似的性能指标。时间的节省十分惊人。因此, ADA-MTL for PMF 十分轻松地适应了 T1、T5、T9 等不同的环境设置, 具有很好的可伸缩性, 可以满足大规模数据流的需求。

Table 4 RMSEs of PMF algorithm on different settings

表 4 各种设置下算法的 RMSEs 比较

	MovieLens			Yahoo!Music	
	T1	T5	T9	T1	T5
BT-MTL	1.002	0.902	0.868	28.88	23.54
DA-STL	0.896	0.895	0.869	29.24	23.57
DA-MIFS	0.992	0.901	0.902	28.41	22.84
ADA-MTL	0.898	0.894	0.901	28.35	22.82

4 结束语

本文首次提出一种面向大数据流的、新颖的多任务加速在线算法, 每次迭代时空复杂度 $O(dQ)$, 最优收敛率 $O(1/T^2)$ 。与当前的在线算法和批训练算法相比, 提出的算法在各种应用实验中均取得相当或更好的分析预测性能。更重要的是, 该算法显著提升了大规模数据流处理的实时性和可伸缩性, 在大数据流领域有较广泛的实际应用价值。

本研究工作尚存需要进一步完善之处。1) p -范数 RDA 方法能够更好地适应学习问题的几何结构, 如何与加速机制结合在特定情况下获得更好的稀疏解与收敛率值得研究; 2) 大数据环境下的多核学习问题是大数据机器学习的基础性问题之一, 如何设计出高效的多任务加速在线多核学习算法, 更好地满足大数据流应用的各种需求, 是我们下一步研究要做的工作。

参考文献

[1] Sun DW, Zhang GY, Zheng WM. Big data stream computing: technologies and instances. Journal of Software, 2014-1-23 (in chinese)

<http://cnki.net/kcms/doi/10.13328/j.cnki.jos.000015.html>.

(孙大为, 张广艳, 郑纬民. 大数据流式计算: 关键技术及系统实例. 软件学报, 2014-1-23 jos在线出版)

[2] Meng XF, Ci X. Big data management: concepts, techniques and challenges. Journal of Computer Research and Development, 2013, 50(1): 146-169 (in Chinese)
(孟小峰, 慈祥. 大数据管理: 概念, 技术与挑战. 计算机研究与发展, 2013, 50(1): 146-169)

[3] Yang HQ, Lyu MR, King I. Efficient online learning for multi-task feature selection. ACM Transactions on Knowledge Discovery from Data, 2013, 1(1): 1-28

[4] Ando RK, Zhang T. A framework for learning predictive structures from multiple tasks and unlabeled data. Journal of Machine Learning Research, 2005, 6: 1817-1853.

[5] Argyriou A, Evgeniou T, Pontil M. Multi-task feature learning. In NIPS, 2006, 41-48

[6] Ben DS, Borbely RS. A notion of task relatedness yielding provable multiple-task learning guarantees. Machine Learning, 2008, 73(3): 273-287.

[7] Caruana R. Multitask learning. Machine Learning, 1997, 28(1): 41-75.

[8] Evgeniou T, Pontil M. Learning multiple tasks with kernel methods. Journal of Machine Learning Research, 2005, 6: 615-637

[9] Liu J, Ji S, Ye J. Multi-task feature learning via efficient $l_{2,1}$ norm minimization. In UAI, 2009.

[10] Pong TK, Tseng P, Pi S, Ye J. Trace norm regularization: reformulations, algorithms, and multi-task learning. SIAM Journal on Optimization, 2010

[11] Yang HQ, King I, Lyu MR. Multi-task learning for one-class classification. In IJCNN. Barcelona, Spain, 2010, 1-8.

[12] Zhang Y. Multi-task active learning with output constraints. In AAAI, 2010

[13] Yang HQ, Xu ZL, King I, Lyu MR. Online learning for group lasso. In: Proc. of the 38th International Conference on Machine Learning. Haifa, Israel, 2010, 1191-1198

[14] Yang HQ, King I, Lyu MR. Online learning for multi-task feature selection. In: Proc. of the 19th ACM Conference on Information and Knowledge Management. Toronto, Canada, 2010, 1693-1696

[15] Xiao L. Dual averaging method for regularized stochastic learning and online optimization. Journal of Machine Learning Research, 2010, 11: 2543-2596.

[16] Hu CH, Kwok JT, Pan WK. Accelerated Gradient Methods for Stochastic Optimization and Online Learning. In Advances in Neural Information Processing Systems 22, Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, Eds. 781-789

[17] Ji SW, Ye JP. An Accelerated Gradient Method for Trace Norm Minimization. In Proceedings of the 26th International Conference on Machine Learning, Montreal, Canada, 2009.

[18] Mairal, J., Bach, F., Ponce, J., and Sapiro, G. 2010. Online learning for matrix factorization and sparse coding. *Machine Learning Research*, vol. 11, pp. 19-60.

[19] Liu, NN and Yang, Q. 2008. Eigenrank: a ranking-oriented approach to

collaborative filtering In SIGIR, pp. 83–90.

- [20] Liu, NN, Zhao, M., and Yang, Q. 2009. Probabilistic latent preference analysis for collaborative filtering. In *CIKM*, pp. 759–766.
- [21] Argyriou, A., Evgeniou, T., Pontil, M. 2008. Convex multi-task feature learning. *Machine Learning* (2008) 73: 243–272. DOI: 10.1007/s10994-007-5040-8.
- [22] Zhan, Y, Wu, CM, Wang BJ. An algorithm for data stream sampling based on ring circular sliding window tightly coupled with buffer. *Acta Electronica Sinica*, 2011, 39(4): 894–898 (in chinese)
(詹英, 吴春明, 王宝军. 一种与缓冲区紧耦合的环形循环滑动窗口的数据流抽取算法. 电子学报, 2011, 39(4): 894–898).

附录 A. 定理 1 的证明

证明 : 对于在(5)式定义的 $L_{1/2,1}$ -范数正则化, (7)式是

基于向量 $\bar{W}_{j\bullet}^T$ 和 $\bar{G}_{j\bullet}^T$ 的。我们使用 w 表示 $W_{j\bullet}^T$, 则(7)

式变为:

$$\phi(w_t) = (\bar{G}_{j\bullet})_t w_t + \lambda(c \|w_t\|_1 + \|w_t\|_2) + \frac{1}{2} \|w_t - u_t\|_2^2. \quad (\text{A1})$$

该式中的目标向量是基于任务属性的, 因此我们可以考虑一个任务 $q \in [1, Q]$, (A1)式变成:

$$\begin{aligned} \phi((w_q)_t) &= (\bar{G}_{j,q})_t \cdot (w_q)_t + \lambda c |(w_q)_t| + \lambda \|(w_q)_t\|_2 \\ &+ \frac{1}{2} ((w_q)_t - (u_q)_t)^2. \end{aligned}$$

显然, 1) 当 $\bar{G}_{j,q} = (u_q)_t$, $(w_q)_t = 0$;

2) 当 $\bar{G}_{j,q} > (u_q)_t$, $(w_q)_t < 0$;

3) 当 $\bar{G}_{j,q} < (u_q)_t$, $(w_q)_t > 0$ 。

对于 2), 上式变成:

$$\begin{aligned} \phi((w_q)_t) &= ((\bar{G}_{j,q})_t - \lambda c) \cdot (w_q)_t + \lambda \|(w_q)_t\|_2 \\ &+ \frac{1}{2} ((w_q)_t - (u_q)_t)^2. \end{aligned}$$

显然, 当 $\bar{G}_{j,q} - (u_q)_t \leq \lambda c$, $(w_q)_t = 0$;

类似, 对于 3), 当 $-\lambda c \leq \bar{G}_{j,q} - (u_q)_t$, $(w_q)_t = 0$ 。

综上, 当 $|(\bar{G}_{j,q})_t - (u_q)_t| - \lambda c \leq 0$ 时, $(w_q)_t = 0$ 。因此, 令

$$(\bar{R}_{j\bullet})_t = [|(\bar{G}_{j,q})_t - (u_q)_t| - \lambda c]_+ \cdot \text{sign}((\bar{G}_{j,q})_t - (u_q)_t)$$

对所有的 $q=1, \dots, Q$, (A1)式就变成:

$$\phi(w_t) = (\bar{R}_{j\bullet})_t w_t + \lambda \|w_t\|_2 + \frac{1}{2} (\|w_t\|_2^2 + \|u_t\|_2^2) \quad (\text{A2})$$

不难验证, 上式的最优解应该是 $w_t = k(\bar{R}_{j\bullet})_t$, $k \leq 0$ 。

因此, (A2)式的目标变为:

$$\min_{k \leq 0} k \|(\bar{R}_{j\bullet})_t\|_2^2 - \lambda k \|(\bar{R}_{j\bullet})_t\|_2 + \frac{1}{2} k^2 \|(\bar{R}_{j\bullet})_t\|_2^2 \quad (\text{A3})$$

构建上述最优问题的拉格朗日式, 并应用 KKT 条件, 最优解必须满足:

$$\frac{\partial \phi}{\partial k} = \|(\bar{R}_{j\bullet})_t\|_2^2 - \lambda \|(\bar{R}_{j\bullet})_t\|_2 + k \|(\bar{R}_{j\bullet})_t\|_2^2 + v = 0,$$

$$v k = 0, v > 0.$$

由此解得:

$$k = - \left[1 - \frac{\lambda}{\|(\bar{R}_{j\bullet})_t\|_2} + \frac{v}{\|(\bar{R}_{j\bullet})_t\|_2^2} \right]. \quad (\text{A4})$$

根据 KKT 条件, $k < 0$ 当且仅当: $v = 0, \lambda < \|(\bar{R}_{j\bullet})_t\|_2$ 。因此得到:

$$(W_{j\bullet})_t = k(\bar{R}_{j\bullet})_t = - \left[1 - \frac{\lambda}{\|(\bar{R}_{j\bullet})_t\|_2} \right]_+ \cdot (\bar{R}_{j\bullet})_t.$$

定理 1 得证。

附录 B. 定理 2 的证明

证明 : 令算法 2 的学习权重矩阵 W 最优值:

$$W^* = \underset{W}{\text{Argmin}} \{ \phi(W) = \langle \bar{G}_t, W \rangle + \Omega_\lambda(W) \}, \quad (\text{B1})$$

那么, $s_T = \phi(W_T) - \phi(W^*)$ 代表算法迭代 T 次收敛率。

根据算法 2 的(7)式, 令 $p(Y) = \arg \min_X \phi(X, Y)$,

则 $W_t = \arg \min_W \phi(W, U_t) = p(U_t)$ 。

由于算法 2 损失函数 l 强凸, Ω_λ 是泛凸函数, 有:

$$l(X) \geq l(Y) + \langle X - Y, \nabla l(Y) \rangle,$$

$$\Omega_\lambda(X) \geq \Omega_\lambda(p(Y)) + \langle X - p(Y), g(p(Y)) \rangle.$$

这里, $g(p(Y)) \in \partial \Omega_\lambda(p(Y))$ 。上面两不等式相加,

得:

$$\begin{aligned} \varphi(X) &\geq l(Y) + \langle X - Y, \nabla l(Y) \rangle \\ &+ \Omega_\lambda(p(Y)) + \langle X - p(Y), g(p(Y)) \rangle \end{aligned} \quad (\text{B2})$$

因为,

$$\begin{aligned} \phi(p(Y), Y) &= \varphi(p(Y)) + \frac{1}{2} \|p(Y) - Y\|_F^2. \\ \phi(p(Y), Y) &= l(Y) + \langle p(Y) - Y, \nabla l(Y) \rangle \\ &+ \Omega_\lambda(p(Y)) + \frac{1}{2} \|p(Y) - Y\|_F^2. \\ \partial\phi &= \nabla l(Y) + (p(Y) - Y) + g(p(Y)) = 0. \end{aligned} \quad (\text{B3})$$

由(B2)、(B3)关系式, 我们得到:

$$\begin{aligned} \varphi(X) - \varphi(p(Y)) &\geq \varphi(X) - \phi(p(Y), Y) \\ &\geq \langle X - p(Y), \nabla l(Y) + g(p(Y)) \rangle - \frac{1}{2} \|p(Y) - Y\|_F^2 \\ &= \langle Y - X, p(Y) - Y \rangle + \frac{1}{2} \|p(Y) - Y\|_F^2. \end{aligned}$$

$$\begin{aligned} \varphi(X) - \varphi(p(Y)) &\geq \\ \text{即: } &\langle Y - X, p(Y) - Y \rangle + \frac{1}{2} \|p(Y) - Y\|_F^2 \end{aligned} \quad (\text{B4})$$

$$\begin{aligned} \text{令: } v_t &= \varphi(W_t) - \varphi(W^*), \\ V_t &= \alpha_t W_t - (\alpha_t - 1)W_{t-1} - W^*. \end{aligned} \quad (\text{B5})$$

应用(B4)式, 分别设 $X=W_t, Y=U_{t+1}$ 和 $X=W^*, Y=U_{t+1}$, 得到如下两个不等式:

$$2(v_t - v_{t+1}) \geq \|W_{t+1} - U_{t+1}\|_F^2 + 2\langle W_{t+1} - U_{t+1}, U_{t+1} - W_t \rangle \quad (\text{B6})$$

$$-2v_{t+1} \geq \|W_{t+1} - U_{t+1}\|_F^2 + 2\langle W_{t+1} - U_{t+1}, U_{t+1} - W^* \rangle \quad (\text{B7})$$

(B6)式两边乘以 $\alpha_{t+1}-1$ 并加(B7)式:

$$\begin{aligned} 2((\alpha_{t+1} - 1)v_t - \alpha_{t+1}v_{t+1}) &\geq \alpha_{t+1} \|W_{t+1} - U_{t+1}\|_F^2 \\ &+ 2\langle W_{t+1} - U_{t+1}, \alpha_{t+1}U_{t+1} - (\alpha_{t+1} - 1)W_t - W^* \rangle \end{aligned}$$

上式两边乘以 α_{t+1} 并利用算法 2 的(8)式导出的关系

式: $\alpha_k^2 = \alpha_{k+1}^2 - \alpha_{k+1}$, 我们得到:

$$\begin{aligned} 2(\alpha_t^2 v_t - \alpha_{t+1}^2 v_{t+1}) &\geq \alpha_{t+1} (W_{t+1} - U_{t+1})\|_F^2 \\ &+ 2\alpha_{t+1} \langle W_{t+1} - U_{t+1}, \alpha_{t+1}U_{t+1} - (\alpha_{t+1} - 1)W_t - W^* \rangle \end{aligned} \quad (\text{B8})$$

因为, $\|B - A\|_F^2 + 2\langle B - A, A - C \rangle = \|B - C\|_F^2 - \|A - C\|_F^2$

对任何大小相同的矩阵 A、B、C 均成立。B8 可变为:

$$\begin{aligned} 2(\alpha_t^2 v_t - \alpha_{t+1}^2 v_{t+1}) &\geq \alpha_{t+1} \|W_{t+1} - (\alpha_{t+1} - 1)W_t - W^*\|_F^2 \\ &- \alpha_{t+1} \|U_{t+1} - (\alpha_{t+1} - 1)W_t - W^*\|_F^2. \end{aligned}$$

利用算法 2 的(9)式以及(B5)式 V_t 的定义, 得:

$$2\alpha_t^2 v_t - 2\alpha_{t+1}^2 v_{t+1} \geq \|V_{t+1}\|_F^2 - \|V_t\|_F^2. \quad (\text{B9})$$

应用上式, t 分别取 $1, 2, \dots, t-1$, 并将各式相加, 得:

$$2v_1 - 2\alpha_t^2 v_t \geq \|V_t\|_F^2 - \|V_1\|_F^2. \quad (\text{B10})$$

应用(B4)式, 设 $X=W^*, Y=U_t$, 得:

$$\begin{aligned} \varphi(W^*) - \varphi(W_1) &= \varphi(W^*) - \varphi(p(U_1)) \\ &\geq \langle U_1 - W^*, p(U_1) - U_1 \rangle + \frac{1}{2} \|p(U_1) - U_1\|_F^2 \\ &= \frac{1}{2} (\|W_1 - U_1\|_F^2 + 2\langle U_1 - W^*, W_1 - U_1 \rangle) \\ &= \frac{1}{2} \|W_1 - W^*\|_F^2 - \frac{1}{2} \|U_1 - W^*\|_F^2. \end{aligned}$$

$$\text{即: } 2v_1 \leq \|U_1 - W^*\|_F^2 - \|W_1 - W^*\|_F^2. \quad (\text{B11})$$

把(B11)式代入(B10)式:

$$\|U_1 - W^*\|_F^2 - 2\alpha_t^2 v_t \geq \|V_t\|_F^2 \geq 0.$$

由于 $\alpha_t \geq (t+1)/2$, 于是:

$$v_t = \varphi(W_t) - \varphi(W^*) \leq \frac{2\|W_0 - W^*\|_F^2}{(t+1)^2}.$$

若算法 1 运行到第 T 步, 则收敛率:

$$v_T = \varphi(W_T) - \varphi(W^*) \leq \frac{2\|W_0 - W^*\|_F^2}{(T+1)^2} = O\left(\frac{1}{T^2}\right).$$

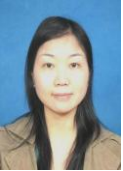
定理 2 得证。



Li Zhijie, born in 1964. PhD candidate from Wuhan University. He is an associate professor at Hunan Institute of Science and Technology, and the member of CCF. His current research interests include machine learning and data mining.



Li Yuanxiang, born in 1962. Professor and PhD supervisor in Wuhan University. His current research interests include intelligent computation and parallel computing etc.



Wang Feng, born in 1981. Associate professor in Wuhan University. Her main research interest include data mining and machine learning.



Kuang Li, born in 1986. PhD candidate from Wuhan University. His current research interest include complex network and evolutionary algorithm.

12. 2.2 节 w_t 邻近解是相对 $\min_w \phi(W)$ 而言, 对于算法优

化式 $\min_w \phi(W, U_t)$ 是最优解。

13. 实验数据集进行了尽量详细的描述。

修改说明:

全文按照“计算机研究与发展”的格式要求进行了认真的修改。写作上, 根据评审意见, 进行了相应的修改。具体如下:

1. 在 2.3 节定义 2 后, 简单说明收敛率与算法效率、迭代次数之间关系。
2. 对原 3.1 节, “通过精心设计由(8)产生的输入序列”, 改变了表达方式。
3. 进一步解释了表 3 中数据集的含义; 表 2 中两个无数据的单元格, 解释了其原因; 图 4, 具体指出是 ADA-MTL 算法取得的结果。
4. 3.2 节用文字进一步描述了三种在线算法实时性能上的差异。
5. 公式(5)等于号左边改为 $\Omega_i(W)$ 。
6. 算法 1 和算法 2 都是在线学习算法, 处理对象为流数据, 所以没有明确算法的执行次数、循环的终止条件。
7. 数据吞吐量(throughout, /s)。
8. 3.3 节明确给出预测误差 RMSE 是真实值与预测值误差的平均平方根。
9. 3.1 节用图 1 画出了四个算法的运行时间—训练样本数的 Log—log 对数关系图。
10. 文中流式大数据的说法大部分已改为大数据流。
11. DA-MTFS 虽为已有工作, 但为了增加论文的可读性, 以及便于比较, 对其作了较详细的描述。