

基于流形学习的单字符字体辨别

何秀玲^{1,2}, 杨 扬², 陈增照^{1,2}, 喻 莹^{1,2}, 董才林¹

HE Xiu-ling^{1,2}, YANG Yang², CHEN Zeng-zhao^{1,2}, YU Ying^{1,2}, DONG Cai-lin¹

1. 华中师范大学 数学与统计学学院 最优控制与离散数学重点实验室, 武汉 430079

2. 北京科技大学 信息工程学院, 北京 100083

1. The Center for Optimal Control & Discrete Mathematics, Institute of Mathematics and Statistics, Central China Normal University, Wuhan 430079, China

2. Institute of Information Engineering, University of Science & Technology Beijing, Beijing 100083, China

E-mail: xlhe@mail.ccnu.edu.cn

HE Xiu-ling, YANG Yang, CHEN Zeng-zhao, et al. Type identification of single character based on manifold learning. Computer Engineering and Applications, 2008, 44(6): 206-209.

Abstract: The identification of language and character type has been an active area of research after recognition of machine printed text. Research on identification of handwritten text and printed text is seldom conducted. But it is common used in recognition of form. For character type identification, manifold learning algorithm Locally Linear Embedding (LLE) is imported. A generalizing method and a parameters estimation method are proposed. Experiments in identification printed/handwritten Chinese characters and digits show that its performance is higher than Support Vector Machine (SVM) classification. The combination of dimensionality reduction of LLE and Linear Discriminant Analysis (LDA) classification achieves a similar accurate rate as or higher than the combination of LLE and SVM classification but runs much faster than it.

Key words: identification of character type; manifold learning; Locally Linear Embedding (LLE); parameter estimation

摘 要: 文字种类识别及字体辨别已成为继印刷体文字识别以后新的国内外研究的热点, 关于单字的手写体和印刷体辨别的研究不多, 但在表单中却极为常用。对于字体辨别问题, 引入流形学习算法局部线性嵌入(LLE), 假定数据为存在于嵌入高维空间的一个低维流形。提出了用于单字符体辨别的 LLE 泛化方法及邻域和内在维数的参数估计方法, 基于印刷体/手写体汉字字符及数字的辨别实验表明, 其性能优于直接支持向量机(SVM)分类, 且经过 LLE 降维后的数据直接用线性判别分析方法(LDA)分类可以获得与 LLE 计算后 SVM 分类相近甚至更高的正确率和更快的分类速度。

关键词: 字体辨别; 流形学习; 局部线性嵌入(LLE); 参数估计

文章编号: 1002-8331(2008)06-0206-04 **文献标识码:** A **中图分类号:** TP391.4

1 引言

在文档分析与理解中, 字体信息是版面分析、理解和恢复的重要依据^[1]。实际的文档中往往同时包含了多种字体和语言, 将不同字体分离开, 把多字体识别变成单字体识别非常重要, 这样有助于实现高性能的字符识别系统。近年来, 用于区分不同语言的文字种类识别已成为国内外研究的热点, 字体识别和文字种类识别是一个相近的问题。实际上, 打印体和手写体混排的情况很常见, 尤其是在表单文档中^[2]。一个表单文档通常由两部分组成, 一个是预打印的印刷体, 另外一种手写填入文本。印刷体和手写体的识别方法和机制是完全不同的。为了得到优化的性能, 必须区分这两种不同的文本。一旦文本被判断是印刷体, 它将会被送到印刷体识别核心, 否则就送入手写字符识别核心, 这样可得到最好的识别率。但是, 区分印刷体和手

写体的方法还没有被广泛讨论, 这也是文档分析和识别的一个新的研究点。

文献中基于手写、印刷体文本鉴别的方法, 从处理基元的角度可以分为基于文本行、词、单字(字母)的不同级别。在手写/印刷体辨别中, 大部分技术是文本行级的, 即基于一行或几行文本的。打印文本排列规则, 而手写文本行排版变化大, 可以利用多通道 Gabor 滤波、尺度小波分析等方法得到纹理特征对字体进行辨别。但是, 在一些情形下, 不同的语言是混合在一起的, 有时不能得到整行文本。例如, 在表单识别中, 很多表单区域只有 1-3 个字符, 例如单元内填有姓名、性别、民族等, 无法得到较多的字符用来组成一个字符块或文本行。因此, 在上述情形, 需要发展基于单个字符的字体识别方法。虽然在字符级, 很少的信息可用, 但是考虑到人很容易辨别手写体和印刷体,

基金项目: 湖北省重点新产品计划资助项目(No.2003BDST004)。

作者简介: 何秀玲(1971-), 女, 博士生, 主要研究方向为图像处理与模式识别; 杨扬(1955-), 男, 教授, 博士生导师, 主要研究方向为图像处理与模式识别、语音识别、多媒体技术及网络通信; 陈增照(1971-), 男, 博士生, 主要研究方向为图像处理与模式识别; 喻莹(1973-), 女, 博士生, 研究方向为模式识别与图像处理; 董才林(1963-), 男, 博士, 副教授, 研究方向为管理科学与工程。

收稿日期: 2007-06-26

修回日期: 2007-08-29

这也促使研究者们研究字符级别的字体辨别。

事实上, 目前针对单个汉字或字符进行的字体识别研究还很少。文献[3]提出一个用对称性等特征的基于神经网络的方法, 利用字母的结构特点, 得到关于单字母的 78.5% 的正确率。文献[4]用游程投影图特征和笔画密度投影等特征辨别手写和印刷体的中文字符。这些方法应用于本文的样本并未取得可实用的结果。本文针对需要批量处理的表单, 提出一种基于流形学习的单个字符类型辨别方法, 包括表单中常用的手写/打印中文字符和手写/打印数字的辨别。这在银行等行业广泛使用的表单数字化中占有重要的地位。

2 流形学习及算法

如何能从数据中发现和学习其内在的规律性一直是机器学习与多元数据分析的主要目标^[5]。近年来, 神经科学的研究取得很多重大发展。2000 年在《Science》上, Seung 提出感知以流形方式存在, 并在神经生理学上发现整个神经细胞群的触发率可以由少量的变量组成的函数来描述, 如眼的角度和头的方向。这隐含了神经元群体活动性是由其内在的低维结构所控制^[5]。流形是微分几何中的一个基本概念, 20 世纪微分几何得到高速发展, 为研究感知流形的形成及其性质提供了坚实的数学理论基础, 几何和拓扑的研究方法为研究感知流形提供了新的思路。

流形学习旨在发现高维数据集分布的内在规律性, 其基本思想是: 高维观测空间中的点由少数独立变量的共同作用在观测空间张成一个流形, 如果能有效地展开观测空间卷曲的流形或发现内在的主要变量, 就可以对该数据集进行降维。近年来, 在流形学习上形成了许多的算法。局部线性嵌套 (Locally Linear Embedding, LLE)^[6], 等度规映射 (Isometric Mapping, Isomap), 拉普拉斯特征映射算法 (Laplacian Eigenmaps), 局部切空间排列算法 (Local Tangent Space Alignment, LTSA) 方法作为流形学习的突出代表, 以其各自的优势引起了众多理论和应用上的关注。其中, LLE 算法有解析的全局最优解, 不需迭代, 低维嵌入的计算归结为稀疏矩阵特征值的计算, 这样计算复杂度相对较小, 易于执行, 同时又具有学习高度非线性流形的能力, 因此受到了研究者的关注。

局部线性嵌入 (LLE) 算法是非线性降维方法, 认为数据可能不是存在于全局线性的流形上, 但在局部意义下, 数据的结构是线性的, 或者说局部意义下的点在一个超平面上, 这样每个数据点都可以表示成一个其邻域构成的加权线性组合。这种逼近的系数表征了高维空间的局部几何, 随后用来在低维空间构造保持几何的低维嵌入。这样可以把输入数据映射到统一的一个全局低维坐标系, 并保留邻接特性。LLE 算法的学习目标是在低维空间中保持每个邻域中的权值不变, 即假设嵌入映射在局部是线性的条件下, 最小化重构误差。由于关于 LLE 的关键假设是即使流形嵌入到的一个高维空间从总体考虑是非线性的, 如果每个数据点和它的邻域存在于或邻近于流形的一个局部线性片断, 它还是被假设成局部线性的。也就是说, 流形可以被几个局部线性片断覆盖 (可能有交叉), 同时分析时, 可以产生关于流形的全局几何的信息。在用线性超平面替换非线性流形的主要点是这个操作不会带来重大错误, 因为当局部分析时, 流形的曲率不大, 即, 流形可以考虑为局部平坦的。

LLE 的输入为一个由 N 个 D 维的向量集合, 组成一个 $D \times N$ 大小的矩阵 X , 输出是 N 个 d 维向量 ($d \ll D$) 组成的大小为

$d \times N$ 的矩阵 Y , 向量可看成是在 R^d 或 R^d 内的点。LLE 算法主要包含 3 个步骤^[6]。

步骤 1 在高维空间计算构成局部片断的点集。对每个点 $X_i, i=1, \dots, N$, 计算并排序另外 $N-1$ 个点与 X_i 的欧式距离, 构成 $N \times N$ 邻域矩阵 A , 取 X_i 的 K 最近邻。

步骤 2 计算相邻点对之间的权值矩阵 W_{ij} 。对每个数据点找到其 K 最近邻以后, 下一步就要为每个相邻的点对赋一个权值。这个权值 W_{ij} 表征了两点间的接近程度, 即第 j 个数据点对重构第 i 个数据点的所做的贡献。数据点的重构误差的用的成本函数来衡量:

$$\varepsilon(W) = \sum_{i=1}^N \left\| X_i - \sum_{j=1}^N W_{ij} X_{j \in A_i} \right\|^2 \quad (1)$$

为了得到最优权值, 对成本函数进行最小值计算满足两个约束: 如果 X_i 与 X_j 不相邻, $W_{ij}=0$; $\sum_{j=1}^N W_{ij}=1$ 。最优权值 W_{ij} 可以通过解线性方程得到:

$$\sum_{j \in A_i} C_{jk} W_{jk} = 1, \forall j \in A_i \quad (2)$$

其中 C 是局部协方差矩阵, 其元素 C_{jk} 定义为:

$$C_{jk} = (\phi(X_i) - \phi(X_j)) \cdot (\phi(X_i) - \phi(X_k)), \forall j, k \in A_i \quad (3)$$

当 $D < K$ 时, 需要对 C 进行正则化, 即将 C 的对角元素加入一个小的正数。

步骤 3 计算低维嵌入 Y_i 。由于目标是尽量正确地在低维空间保持高维空间的局部线性结构, 高维观察值 X_i 被映射为低维向量 Y_i , 权值 W_{ij} 保持不变, d 维向量 Y_i 通过最小化成本函数得到:

$$\delta(Y) = \sum_{i=1}^N \left\| Y_i - \sum_{j=1}^N W_{ij} Y_j \right\|^2 \quad (4)$$

满足约束: 正规化单位协方差 $\frac{1}{N} \sum_{i=1}^N Y_i Y_i' = I$; 平移不变嵌入 $\sum_{i=1}^N Y_i = 0$, 提供唯一解。

3 单字字体辨别

取自不同采集设备的图像数据通常是多维的, 因此它们不是很适合一般能正确对小集合的相关特征分类方法。因此, 需要维数约减以减少或去除处于次要的信息并保留或突出有意义的信息。由于真实世界数据的本性常常是非线性的, 线性维数约减技术在数据映射到低维时无法保持数据的结构和它们之间的关系。这就意味着此时需要非线性维数约减方法。LLE 方法具有以下优势: 只有两个自由参数需要设置; LLE 是一个非迭代方案, 避免了收敛到局部极小问题; 一个高维数据在嵌入空间的局部几何的好的保持 (邻域保持); 一个嵌入空间的全局坐标系。因为 LLE 出现的时间不长, 在应用中存在一些问题, 本文从泛化, 参数估计方面进行了探讨, 其方法用于对真实数据包括手写/印刷体数字和中文单个字符字体辨别测试中取得了很好的效果。

3.1 LLE 的泛化

原始的 LLE 对数据是固定的, 也就是说, 为了将其映射到嵌入空间需要一个整个的点集作为输入, 是按批处理的方式进行操作。当新的点到来, 映射它们的唯一方法是将旧的和新的点集中起来, 并对这个聚集再次返回 LLE。换句话说, 原始的 LLE 缺乏对新数据点的泛化能力。这意味着它不适合变化, 动

态的环境。

文献[7]基于假定来自高维空间的部分点已经明确地被LLE投影这样一个事实,令点集 $X_i(i=1, \dots, N)$ 作为LLE的输入,组成一个集合 γ 。对一个不可见的点 X_{N+1} ,在所有 $X_i \in \gamma$ 中寻找距离点 X_{N+1} 最近的点 X_j ,令 Y_j 为 X_j 在嵌入空间的投影。以下等式近似成立: $Y_j = ZX_j$,其中 Z 是一个未知的线性转换矩阵,大小为 $d \times D$,可以直接定义为 $Z = Y_j X_j^{-1}$ 。因为 X_j 和 X_{N+1} 彼此很接近(这里集合 γ 必须很充分地代表了潜在流形或者说流形被很好的采样了), Y_{N+1} 可以用 ZX_{N+1} 计算,即,投影两点时使用的是同样的转换矩阵。这种方法对噪声非常敏感,因此本文采用改进的泛化方法依赖于一个LLE低维和高维空间的自然映射,计算新输入 X_{N+1} 的输出 Y_{N+1} 方法如下:

- (1)在训练输入集合中搜索 X_{N+1} 的 K 最近邻;
- (2)计算线性权重 W_j 使得能将 X_{N+1} 从其邻域最好重构,并遵从和为1的约束, $\sum \mu_j = 1$;
- (3)输出 $Y_{N+1} = \sum \mu_j Y_j$,其中,求和操作是对 X_{N+1} 邻域的相应输出进行。

从嵌入空间到输入空间的映射可以用同样的方式推导:为计算一个新输出 Y_{N+1} 的输入 X_{N+1} ,在训练输出中找到 Y_{N+1} 的最近邻,计算重构权重 $X_{N+1} = \sum \mu_j X_j$ 。这种方法无需重新计算特征向量,因此计算简单,速度快。

3.2 参数估计

原始LLE有两个参数需要调整, K 和 d ,分别为每个点的最近邻的数目和嵌入空间的维数,即流形的内在维数。正确选择参数 K 至关重要,原因是一个大的最近邻会导致将平滑或流形中小尺度结构排除在外;相反,太小的邻域会错误地将连续流形分解成了不连接的子流形。另外,在分类问题中数据样本的内在维数无法预知,而错误地估计 d 会导致LLE行为病态。

在理想的情况下,数据 K 邻域严格局部平坦。可以证明:令 z 为矩阵 M 的零特征值数目,那么 $d < z^{[8]}$ 。由于实际数据可能不是局部邻域内严格平坦的,因此 M 的特征值个数 z 取为接近0的个数 \hat{z} , d 在区间内,这样得到了 d 的取值范围。

原始LLE数据点的重构误差的的成本函数(1)衡量的。 K 的取值应使其最小,即

$$K^* = \underset{K}{\operatorname{argmin}} \mathcal{E}(W) = \underset{K}{\operatorname{argmin}} \sum_{i=1}^N \left\| X_i - \sum_{j=1}^K W_{ij} X_j \right\|^2 \quad (5)$$

这里的 \mathcal{E} 对 K 通常有多于一个最小值,因此得到 K 的取值范围 $[1, K_{\max}]$ 。

得到 K, d 的取值范围后,进一步采用迭代算法,计算类内离散度和类间离散度,以期获得尽量大的类间散度和尽量小的类内散度,从而使得可分性最好。

为了定义类内离散度和类间离散度,先定义几个基本的参量。

在 d 维 Y 空间, N 个样本 $\{y_1, \dots, y_N\}$,分为 L 类 $\{c_1, \dots, c_L\}$,每类样本的个数分别为 N_1, \dots, N_L ,第 c 类的样本集为 $Y_c, \{y_{c,1}, \dots, y_{c,N_c}\}$ 。

- (1)各类样本均值向量 m_c 及总均值 m

$$m_c = \frac{1}{N_c} \sum_{y \in Y_c} y, c=1, 2, \dots, L \quad (6)$$

$$m = \frac{1}{N} \sum_{y \in Y} y \quad (7)$$

- (2)第 c 类样本类内离散度 S_c 和总类内离散度矩阵 S_w

$$S_c = \sum_{y \in Y_c} (y - m_c)(y - m_c)^T, c=1, \dots, L \quad (8)$$

$$S_w = \sum_{c=1}^L S_c \quad (9)$$

- (3)样本类 c_1 与 c_2 类间离散度矩阵 S_{bcc} ,总类间离散度矩阵 S_b

$$\begin{cases} S_{bcc} = (m_{c_1} - m_{c_2})(m_{c_1} - m_{c_2})^T, & c_1, c_2 = 1, \dots, L \\ S_b = \sum_{c=1}^L (m_c - m)(m_c - m)^T \end{cases} \quad (10)$$

各类样本之所以可以分开是因为它们处于特征空间中的不同区域,这些区域之间的距离越大可分性越大,越有利于分类,即希望类间离散度 S_b 越大越好;同一类别各样本相互间的距离越小可分性越大,同类样本的越密集越有利于分类,即类内离散度越小越好 S_w 。因此,定义适应度函数

$$J = S_b / S_w \quad (11)$$

这里的 J 需要计算已知矩阵 M 的特征值个数变化的情况下的适应度,因此 J 与 d 有关,适应度越大越好。这样

$$(d^*, k^*) = \underset{d, k}{\operatorname{argmax}} (S_b / S_w) \quad (12)$$

显然应该寻找使得分子尽可能大,分母尽可能小,也就是使 J 尽可能大的 d 取值。算法如下:

- (1) $K=1$;
- (2)执行LLE的步骤1、步骤2,得到稀疏矩阵 M ;
- (3)求 M 的所有特征值并排序,取相邻特征值间距产生突变或小于阈值 ε 的特征值个数(近似0) z ,令 d 的候选集合为 $[1, z]$;
- (4)对 d 的每个取值,根据LLE步骤3计算 Y ,并且根据式计算适应度函数 J ;
- (5) $K=K+1$,若 $K < K_{\max}$ 则转第(2)步,否则继续;
- (6)根据式(12)求取最大的 J 值对应的维数 d 和 K 。

3.3 分类方法

基于流形学习的特征降维过程可以看成是特征提取的过程,用于分类时采用基于监督的LLE算法(Supervised LLE, SLLE)^[9]。在完成特征提取的过程以后,接下来是要选择分类器对提取的特征进行分类。本文采用两类常用的分类器,即支持向量机(SVM)及线性判别分析(LDA),其中SVM参数采用文献[10]中的优化方法得到。

4 实验方法及结果分析

分别搜集了金融票据中常用的手写体汉字4000个,印刷体汉字2075个,手写体数字4000个,印刷体数字样本4000个,总共14075个样本。取其中每类样本中每个字符一半作为训练样本集,一半作为测试样本集。将每个字符规格化为二值 24×24 点阵图像,分别用SVM直接分类,SLLE+SVM,SLLE+LDA几种分类方法分别对手写体汉字与印刷体汉字,手写数字与印刷体数字进行辨别。

将规格化后的 24 点阵图像转换成一个 D 维特征矢量($D=24 \times 24=576$),像素值用作特征值。对原数据图像采用SLLE进行降维后再SVM分类,根据邻域数估计和内在维数估计,汉字辨别时 K, d 分别取11,12,数字辨别时 K, d 分别取13,18。实验数据均在配置为Intel P4 CPU 2.6 GHz,1 G内存的PC机上得出,见表1。

表1 辨别方法比较

任务	训练样 测试样		方法	正确率/%	支持向 量个数	总训练 时间/s	总测试 时间/s	训练集 LLE 时间/s	测试集 LLE 时间/s
	本数	本数							
手写体汉字与 印刷体汉字	3 066	3 009	SVM	96.84	783	40	30	/	/
			SLLE+SVM	98.60	473	0.8	0.3	43.4	13.5
			SLLE+LDA	98.54	/	0.046	0.016	43.4	13.5
手写数字与 印刷体数字	4 000	4 000	SVM	89.03	736	86	85	/	/
			SLLE+SVM	93.65	240	1.0	0.5	52.6	16.8
			SLLE+LDA	94.6	/	0.047	0.015	52.6	16.8

较之直接 SVM 分类的情况, SLLE+SVM 用于样本训练和测试的时间大为减少, 仅为其 0.5%~2%; 加上 SLLE 计算特征所需的时间, 总的识别时间也只有原来的 46%~56%, 节省了近一半。

经 SLLE 降维后的数据用 LDA 分类得到的识别率与 SVM 相差仅 0.06%, 识别速度进一步降低, 平均每个字符仅需要不到 0.005 ms。

用线性 LDA 分类, 在识别率没有显著降低的情况下, 更进一步提高了速度。这意味着避免了 SVM 中核函数的选择及其参数调整问题。

因此, 为了当约减维数时得到一个紧的类的表示, SLLE 可以用来作为分类前的预处理步骤。

参考文献:

- [1] 陈力, 丁晓青. 基于小波特征的单字符汉字字体识别[J]. 电子学报, 2004, 32(2): 177-180.
- [2] Ma J K, Guo M Y. Separating handwritten material from machine printed text using hidden Markov models[C]//Proceedings of International Conference on Document Analysis and Recognition, 2001: 439-443.

- [3] Kuhnke K, Simoncini L, Kovacs-V Z M. A system for machine-written and hand-written character distinction[C]//Proceedings of International Conference on Document Analysis and Recognition, 1995: 811-814.
- [4] Zheng Y, Liu C, Ding X. Single character type identification[C]//Proc SPIE Conf Document Recognition and Retrieval, 2002: 49-56.
- [5] 周志华, 曹存根. 神经网络及其应用[M]. 北京: 清华大学出版社, 2004: 172-207.
- [6] Roweis S T, Saul L K. Nonlinear dimensionality reduction by locally linear embedding[J]. Science, 2000, 290(5500): 2323-2326.
- [7] Kouropteva O, Okun O, Hadid A, et al. Beyond locally linear embedding algorithm, MVC-01-2002[R]. Machine Vision Group, University of Oulu, 2002.
- [8] Polito M, Perona P. Grouping and dimensionality reduction by locally linear embedding[M]. [S.l.]: MIT Press, 2002.
- [9] Kouropteva O, Okun O, Pietik'ainen M. Classification of handwritten digits using supervised locally linear embedding algorithm[C]//11th European Symposium on Artificial Neural Networks (ESANN), Bruges, Belgium, 2003: 229-234.
- [10] 陈增照. 基于支持向量机的脱机手写体金融汉字识别研究[D]. 北京: 北京科技大学, 2007.

(上接 135 页)

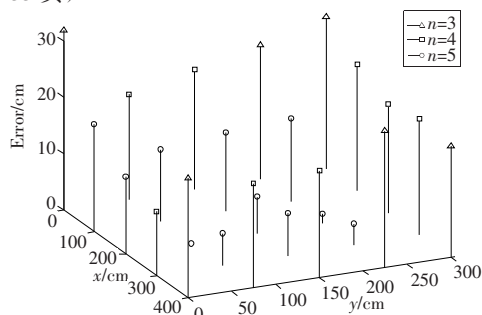


图5 基于在线校正的定位误差

结合线性校正模型, 提出基于在线校正的定位算法; 算法针对无线传感器网络布撒的不均匀性, 充分利用密集布撒下节点定位的冗余信息, 以提高网络测距和定位精度。在线校正适合于传感器节点不可接近或较复杂的网络。实验表明, 算法在实际应用中有效的减小了网络的定位误差。

参考文献:

- [1] Holger Karl, Andreas Willig. Protocols and architectures for wireless sensor networks[M]. [S.l.]: John Wiley & Sons, 2005-06-24.
- [2] Ramanathan N, Balzano L, Burt M, et al. Rapid deployment with confidence: calibration and fault detection in environmental sensor networks[R]. Center for Embedded Networked Sensing, 2006.
- [3] Tolle G, Polastre J, Szewczyk R, et al. A macroscope in the red-

woods[C]//Proceedings of Sensys, 2005.

- [4] Taylor C, Rahimi A, Bachrach J, et al. Simultaneous localization, calibration, and tracking in an ad hoc sensor network[C]//Proceedings of the Fifth International Conference on Information Processing in Sensor Networks, 2006: 27-33.
- [5] Balzano L, Nowak R. Blind calibration in sensor networks [C]//Proceedings of Information Processing in Sensor Networks (IPSN), 2007.
- [6] Whitehouse K, Culler D. Calibration as parameter estimation in sensor networks [C]//Proceedings of the 1st ACM International Workshop on Wireless Sensor Networks and Applications, 2002: 59-67.
- [7] Laura Kathryn Balzano. Addressing fault and calibration in wireless sensor networks[D]. Los Angeles: University of California, 2007.
- [8] Feng J, Megerian S, Potkonjak M. Model-based calibration for sensor networks[J]. IEEE, 2003.
- [9] Burden R L, Faires J D. Numerical analysis[M]. 7th ed. Beijing: Higher Education Press, 2005.
- [10] Chapra S C, Canale R P. Numerical methods for engineers[M]. 3rd ed. [S.l.]: McGraw-Hill Companies Inc, 2000.
- [11] Heath M T. Scientific computing: an introductory survey[M]. 2nd ed. [S.l.]: McGraw-Hill Companies Inc, 2002.
- [12] Bazaraa M, Shetty C M, Sherali H. Nonlinear programming: theory & applications [M]. [S.l.]: Wiley, 1994.
- [13] Press W, Flannery B, Teukolsky S, et al. Numerical recipes[M]. [S.l.]: Cambridge, 1986.
- [14] 王福豹, 史龙, 任丰原. 无线传感器网络中的自身定位系统和算法[J]. 软件学报, 2005.