

◎工程与应用◎

POP 海洋模式在四核至强集群上的并行计算

张理论, 赵 军, 吴建平, 宋君强

ZHANG Li-lun, ZHAO Jun, WU Jian-ping, SONG Jun-qiang

国防科学技术大学 计算机学院, 长沙 410073

School of Computer Science, National University of Defense Technology, Changsha 410073, China

E-mail: zll0434@sina.com

ZHANG Li-lun, ZHAO Jun, WU Jian-ping, et al. Parallel computing of POP ocean model on quad-core intel xeon cluster. Computer Engineering and Applications, 2009, 45(5): 189-192.**Abstract:** Based on analysis of the equations and numerical discretization of POP ocean model, discuss the local blocking technique and load-balance data distribution, focusing on their influence upon parallel performance on quad-core cluster. Aggregation optimization is adopted for the communication bottleneck. Research and experiment show that local blocking technique and way of data distribution have remarkable effect on performance, and the aggregation optimization is effective on quad-core cluster, especially for the case of load-balance data distribution.**Key words:** Parallel Ocean Program (POP) ocean model; quad-core cluster; parallel computing**摘要:** 分析了 POP 海洋模式原理、离散方法。在四核至强集群上, 研究分析 POP 模式中计算局部块技术和平衡并行数据剖分及其对模式性能的影响。针对模式的通信性能瓶颈, 采用聚合通信优化技术。研究结果表明局部块技术和数据剖分方式对于 POP 模式并行性能影响显著; 通过通信聚合优化, POP 模式在四核集群上性能获得明显提升。**关键词:** POP 海洋模式; 四核集群; 并行计算**DOI:** 10.3778/j.issn.1002-8331.2009.05.056 **文章编号:** 1002-8331(2009)05-0189-04 **文献标识码:** A **中图分类号:** TP391

1 引言

POP 海洋数值预报模式由美国洛斯阿拉莫斯国家实验室根据美国能源部气候改变预测计划开发, 利用该模式能够推进 10 年期和大尺度气候预测科学的发展。现在 POP 模式已经成为许多气候模拟器中的标准模块, 例如著名的 NCAR 公共气候系统模型 (CCSM), POP 模式为其重要组件。本文所涉及的 2.0.1 版 POP 并行海洋程序于 2004 年发布, 最大测试算例水平网格为 3 600×2 400, 垂直方向 40 层, 已经达到全球 0.1 度, 相当于赤道上的 10 公里。研究表明该空间分辨率是合理表达海洋中尺度涡(eddy-resolving)动能的尺度门槛^[1]。随着 POP 模式的广泛使用, 他已经成为一个重要的大规模并行 Benchmark 测试程序, 在 NEC 地球模拟器, Cray 红色风暴, IBM 蓝基因等著名的高端并行计算机上通过测试, 其测试数据成为各大高端并行计算机厂商展示其机器性能的重要依据之一(见表 1)。其中模式的数值模拟效率定义为在单位墙钟时间(天)内模式的积分步进时间(年), 也即模式积分时间跨度年数/墙钟时间天数。

目前, 多核 CPU 在高性能计算机上已普遍采用。多核 CPU 在单个硅片上采用多个较低频的紧耦合处理核心, 在控制功耗、处理器成本(成本主要取决于硅片 silicon die)的前提下, 通过 CPU 核间并行达到提高单 CPU 处理性能的目的。当前集群

表 1 POP 模式 0.1 度算例在几个高端并行机的运行性能^[2]

并行计算机 (装机时间)	体系结构特征	使用 CPU 数目及 峰值性能(万亿次)	数值模拟效率
IBM 蓝基因 L (2005 年)	IBM PPC 440 @ 700 MHz 3-D Torus 32×32×64	29 000, 81.1TFlops	8.5
NEC 地球模拟器 (2002 年)	NEC Vector @ 500 MHz 640×640 Cross Bar	500, 3.5TFlops	3
Cray 红色风暴 (2005 年)	AMD Opteron @ 2.0 GHz 3-D Mesh 27×16×24	8 000, 32.1TFlops	6

系统上较常见的多核 CPU 主要有 AMD Opteron 系列和 Intel Xeon 系列。IBM Power 系列不同于 AMD 和 Intel 的货架产品, 一般专用于 IBM 高端并行系统。较早的 Power4 采用多芯片模块(MCM), 在单芯片上集成两个处理核心, 4 个芯片构成一个 8 核的处理模块。随着多核并行机特别是多核集群的涌现, 对其性能进行评测很有必要, 特别是对众多传统大型计算问题在多核并行计算机上的实测性能和性能提升。针对 POP 海洋模式, 研究分析其在四核至强集群上的 MPI 并行性能和优化问题, 以期对多核并行应用开发和多核并行机性能测试提供参考。POP 模式运行需要 MPI、NETCDF、GNU makefile 等软件环境支持, 本文涉及的运行平台为四核至强集群, 其基本配置见表 2。

基金项目: 国家自然科学基金(the National Natural Science Foundation of China under Grant No.40505023, No.10505030)。**作者简介:** 张理论(1975-), 男, 博士, 副研究员, 主要从事大规模并行计算和数值并行软件研究。**收稿日期:** 2008-06-18 **修回日期:** 2008-09-10

表2 四核至强集群测试平台

处理器	八路四核 Intel® Xeon™ Clovertown/2.33 GHz L1 cache 为 4x(32 KB 数据+32 KB 指令), L2 cache 为 8 MB 单核浮点性能 9.32 Gflops
互联	互联采用 Voltaire Infiniband DDR, 单链路双向带宽 40 Gb/s
内存和 IO	计算结点节点内存 4 GB IO 节点内存 8 GB
系统平台环境	操作系统版本 Redhat AS4 2.6.18, Intel 编译器版本 10.0. 023

2 POP 海洋模式数值离散

2.1 模式方程

POP 海洋模式可以用薄层流体的三维原始动力方程来描述^[3]。基于流体静力平衡与 Boussinusq 近似, 在水平球面坐标与垂直 z 坐标下, 动量方程为

$$\frac{\partial u}{\partial t} + L(u) - (uv \tan \phi) / a - fv = -\frac{1}{\rho_0 \cos \phi} \frac{\partial p}{\partial \lambda} + F_{Hx}(u, v) + F_V(u)$$

$$\frac{\partial v}{\partial t} + L(v) - (u^2 \tan \phi) / a + fu = -\frac{1}{\rho_0 a} \frac{\partial p}{\partial \phi} + F_{Hy}(u, v) + F_V(v)$$

$$L(a) = \frac{1}{a \cos \phi} \left(\frac{\partial}{\partial \lambda} (u\alpha) + \frac{\partial}{\partial \phi} (\cos \phi v\alpha) \right) + \frac{\partial}{\partial z} (u\alpha)$$

$$F_{Hx}(u, v) = A_M \left(\nabla^2 u + u(1 - \tan^2 \phi) / a^2 - \frac{2 \sin \phi}{a^2 \cos^2 \phi} \frac{\partial v}{\partial \lambda} \right)$$

$$F_{Hy}(u, v) = A_M \left(\nabla^2 v + v(1 - \tan^2 \phi) / a^2 + \frac{2 \sin \phi}{a^2 \cos^2 \phi} \frac{\partial u}{\partial \lambda} \right)$$

$$\nabla^2 \alpha = \frac{1}{a^2 \cos^2 \phi} \frac{\partial^2 \alpha}{\partial \lambda^2} + \frac{1}{a^2 \cos \phi} \frac{\partial}{\partial \phi} \left(\cos \phi \frac{\partial \alpha}{\partial \phi} \right)$$

$$F_V(\alpha) = \frac{\partial}{\partial z} \mu \frac{\partial}{\partial z} \alpha$$

连续性方程为 $L(1) = 0$

$$\text{流体静力学方程为 } \frac{\partial p}{\partial z} = -\rho g$$

$$\text{状态方程为 } \rho = \rho(\Theta, S, p) = \rho(\Theta, S, z)$$

轨迹输运方程为

$$\frac{\partial}{\partial t} \varphi + L(\varphi) = D_H(\varphi) + D_V(\varphi)$$

$$D_H(\varphi) = A_H \nabla^2 \varphi$$

$$D_V(\varphi) = \frac{\partial}{\partial z} \kappa \frac{\partial}{\partial z} \varphi$$

其中 $\lambda, \phi, z, r = a$ 是经度、纬度、与相对于平均海平面 $r = a$ 的深度, g 是重力加速度, $f = 2\Omega \cos \phi$ 是 Coriolis 参数, ρ_0 是海水的背景密度。这些方程中的预报变量是东西方向的速度分量 u 与 v 、垂直速度 w 、压力 p 、密度 ρ 、势温 Θ 与盐度 S 。状态方程中的压力通常用一个只与深度有关的函数来近似, A_H 与 A_M 分别是水平散度与黏度系数, κ 与 μ 分别表示对应的垂直混合系数, 依赖于局部 Richardson 数。 $F_V(u)$ 与 $F_V(v)$ 是对流导数所引起的度量项, 分别作用在 λ, ϕ 方向的单位向量上。由风应力、热、淡水通量引起的强迫项用作摩擦与扩散项 F_V 和 D_V 的表面边界条件。

2.2 模式数值离散分析

模式空间离散采用水平交错 B 网格(标量取空间块中间值, 矢量取空间块角点值), 有限差分的导数与平均值取:

$$\delta_x \psi = \frac{\psi(x + \Delta_x/2) - \psi(x - \Delta_x/2)}{\Delta_x}$$

$$\bar{\psi}^x = \frac{\psi(x + \Delta_x/2) + \psi(x - \Delta_x/2)}{2}$$

$$\text{梯度算子离散为 } \nabla \psi = \bar{x} \delta_x \bar{\psi}^y + \bar{y} \delta_y \bar{\psi}^x$$

\bar{x}, \bar{y} 为方向向量。水平散度算子离散为

$$\nabla \cdot u = \frac{1}{\Delta_y} \delta_x \bar{\Delta}_y u_x + \frac{1}{\Delta_x} \delta_y \bar{\Delta}_x u_y$$

旋度算子的垂直部分为

$$\bar{z} \cdot \nabla \times u = \frac{1}{\Delta_y} \delta_x \bar{\Delta}_y u_z + \frac{1}{\Delta_x} \delta_y \bar{\Delta}_x u_z$$

粘性扩散项中的拉普拉斯算子离散为

$$\nabla \cdot G \nabla \psi = \frac{1}{\Delta_y} \delta_x [\bar{\Delta}_y G \delta_y \bar{\psi}^y] + \frac{1}{\Delta_x} \delta_y [\bar{\Delta}_x G \delta_x \bar{\psi}^x]$$

海洋模式接近海表层部分和深层的空间运动尺度差别很大, 海表层的重力波速往往达到 200 m/s, 而深层运动速度通常比表层低两个数量级, 因而海洋表层波速对时间积分步长要求严格。在时间离散中, 需要将其模式方程分为分为两维正压的外模态(Barotropic)和三维斜压的内模态(Baroclinic)进行处理, 其中外模态可以用动量方程和连续方程的垂直平均来近似。内模态时间离散采用三个时间层的改进蛙跳格式, 具有两阶精度。经典蛙跳在时间积分中一般采用奇偶步分离, 往往会导致较明显的结果扰动, POP 模式中对蛙跳格式采用时间滤波。两维正压外模态部分需要求解一个关于采用表面压力的椭圆方程, POP 模式中采用对角预条件 PCG 求解表面压力方程, 其通信部分以标量归约主导, 通信密集。三维斜压外模态部分采用显式差分离散, 计算密集, 通信以边界更新通信为主, 通信量少和频度低, 并行可扩展性较好。

3 影响模式性能的几项因素

在没有特别说明的情况下, 文中均为满核运行 MPI 任务, 即每个相关 CPU 核上只运行一个 MPI 任务; 编译选项取“-O3 -ip”。

3.1 编译开关对性能的影响

编译开关的影响简单直接。intel 的 Fortran 编译器常可选用以下性能优化参数进行编译优化:

表3 intel 典型的性能优化开关

编译选项	功能描述
-O3	除了激活在-O2 中包含的全局指令调度、软件流水、预测、与数据预取, 同时还包括内部函数的内嵌、常量的复制、僵码的删除、全局寄存器分配、循环展开、优化代码选择、部分冗余的消除、溢出处理优化等之外, 还包括标量重分布、循环变换等
-ip	激活在当前源程序内部的其他优化过程, 包括对定义于当前源程序文件内部的函数调用的内联展开
-pad	通过补边对数组变量的内存分布进行调整, 可能利于发挥 cache 性能
-ipo	通过内联展开不同文件间的函数来提高性能
-no-prec-div	通过降低除法的浮点精度来提高性能

在至强集群系统考察编译选项对性能的影响, 采用 128 任务满核计算, 通过对编译选项进行调整, 发现 POP 编译选项采用 -O3 -ip -pad -no-prec-div -ipo 比采用 -O3 -ip 约快 2% 以上, 比单纯采用 -O3 快约 6%。

3.2 局部块规模对性能的影响

POP 模式默认在水平方向上, 基于笛卡儿坐标分布进行数据剖分, 并对每个计算节点的剖分数据进一步分解成为若干局部数据块, 其规模作为参数可调。这一点类似于 GRAPES 和

表4 128 并行的不同性能优化开关的6小时积分墙上时间

编译性能优化开关					时间
-O3	-ip	-pad	-no-prec-div	-ipo	458.82
	-ip	-pad	-no-prec-div	-ipo	473.27
	-O3	-ip			468.58
	-O3				487.16

WRF 模式中并行分区(patch)和数据瓦片(tile)的两层数据剖分方式^[4]。其区别在于 POP 模式中各个处理器所对应的局部数据块数目可能不止一个(可调),而后两者每个处理器只对应一个瓦片(tile)数据块。

局部块技术具有以下优势^[5]:

(1)Cache 友好。采用自动性能调整参数,计算局部数据块(cache)。POP 有 2 个参数 block_size_x 和 block_size_y 控制运行时数据的局部化,从而影响 Cache 数据的重用率和消息传递的数据通信量。

(2)便于陆地点剔除和实现调度平衡。海洋模式中往往需要考虑陆地岛屿的影响,陆地地点的海洋状态变量为空,需要剔除掉,传统的数据剖分中,各个计算任务的数据块往往较大,无法将其分辨出来,因而带来计算不便和并行计算的负载不平衡。采用局部块技术,各个局部块足够小时能够分辨出陆地点(陆地点例如超过局部块的 50%),便于将其剔除,因而基于局部块的数据分布相对而言更易实现负载平衡。

(3)更好的混合同行选择。局部数据块技术也为混合同行提供一种便捷的实现途径,在各个计算节点完成数据剖分后,对节点内的多个局部数据块可以进行基于 OpenMP 的多线程并行,节点间并行采用 MPI。各个线程对局部数据块的操作相对独立,这实质上是一种基于 SPMD 方式的 OpenMP 并行,其实现性能往往好于循环级的 OpenMP 并行。

我们发现不同的局部块规模选择对于程序可执行码的内存占用量也存在较大影响,不同分块规模对应的各个并行度可执行码静态规模(size 命令查看)见图 1(单位 MB)。

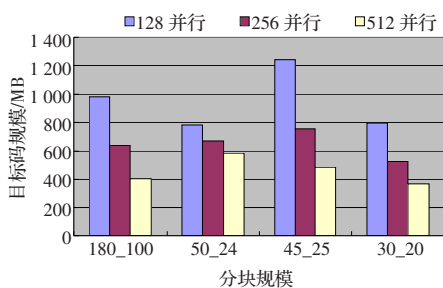


图1 不同局部分块所产生的目标执行码规模

不同局部块规模对模式程序性能的影响非常明显,以下仅以全部采用笛卡儿并行数据剖分为例进行说明。

表5 128 并行任务6小时积分采用不同局部块(笛卡儿分布)

块规模	180_100	90_50	88_48	36_24	30_20	15_10	18_12	20_16
时间	841.49	561.93	625.88	496.86	468.58	732.84	619.66	509.87

不难发现在并行度 128 时不同计算块规模选择对性能影响比较显著,各种分块规模对应的程序执行时间存在较大时间落差。其中 180×100 是 POP 测试程序中默认的分块大小,在该集群系统上效果较差,针对 128 并行时采用 30×20 较为合适。通过进一步对更大并行度的试验分析,采用同样的 30×20 分块,对其他分块规模仍然具有相对性能优势。

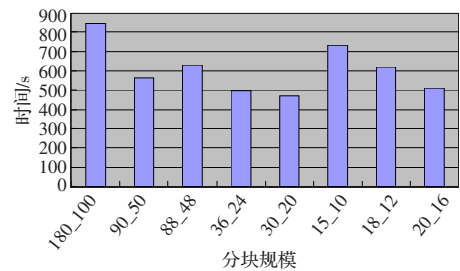


图2 128 并行任务6小时积分不同局部块规模对性能的影响

3.3 并行数据剖分方式对性能的影响

POP 模式并行计算的数据剖分可以分为笛卡儿坐标和平衡分布两种方式,不同的数据剖分形式不影响计算结果,但对程序性能可能产生较大影响。笛卡儿数据剖分最为常见:对并行任务数进行二维拓扑分解,分别将 x, y 方向的任务数 N_x, N_y 和相应方向上的物理网格数作简单线性映射。MPI 平台提供笛卡儿剖分相应的程序接口,实现起来简洁自然。一般而言,各个物理格点上计算量比较一致,则这种数据剖分在适当调整 N_x, N_y 后可以较好地实现负载平衡。

对 POP 全球海洋模式而言,由于存在大量陆地点和岛屿,其对应物理格点上的计算量可以忽略不计,因而简单地采用笛卡儿数据剖分往往会导致负载不平衡。平衡数据剖分实质上是对笛卡儿剖分的一种采用改进,其主要思想是:在简单笛卡儿分布完成后,对并行剖分中的各个局部块的计算量进行统计,并在 x, y 方向按照各个局部块的计算量大小作优先级重排序,尔后依据优先级对以局部块为单位进行数据重分布,并尽可能保证局部块分布在基于笛卡儿分布的拓扑任务上或者只迁移到邻近的拓扑任务上。下面给出在至强集群上 6 小时积分的不同并行数据剖分形式的计算时间对比。

表6 POP 积分6小时不同数据剖分方式的墙钟时间

并行剖分	128		256		512	
	笛卡儿	平衡	笛卡儿	平衡	笛卡儿	平衡
时间	468.58	405.50	228.06	216.72	108.82	106.44
差异程度	13.46%		4.97%		2.18%	

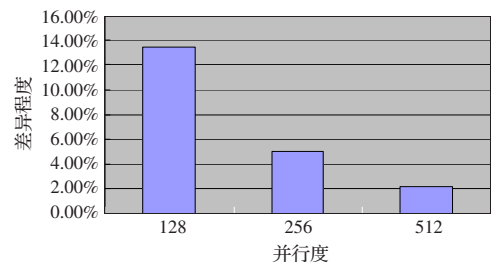


图3 不同剖分方式的性能差异程度

其中剖分方式导致的性能差异程度定义为:|(笛卡时间-平衡时间)|/笛卡时间。并行数据剖分方式在很大程度上决定了 MPI 进程的计算和通信平衡程度,因而对程序性能产生较大影响。显然由简单笛卡儿分布造成的负载不平衡,是程序的一个重要性能瓶颈,该瓶颈效应随着并行度的增大相对减弱。

4 POP 模式的通信优化

(1)通信瓶颈分析

采用 prof 性能分析工具发现 POP 并行程序耗费大量时间

进行进程等待,在MPI函数MPI_wait占用时间高达10%以上。除MPI_wait外,比较耗时的程序包括boundary_2d_dbl:二维场量的边界通信;impvmixt:垂直混合计算;impvmixt_correct:垂直混合计算修正;global_sum_dbl:计算全局和。由于POP模式计算量主要集中在显式时间积分三维斜压部分,本文没有考虑优化二维正压部分的PCG求解过程。通过分析源码发现POP模式中所有高维数组的边界通信均先转化为二维场数据,再分别在各个垂直方向上调用子程序boundary_2d_dbl来实现的。整个POP模式在6小时积分过程中,共调用二维边界更新子程序25350次,其中包括273次三维物理变量数据边界交换(40层,每个含40次边界通信),55次四维迹(tracer)变量边界交换(每个含80次二维边界通信)。

(2)通信聚合

①将三维物理变量和四维迹变量的边界交换的小消息通信转化为大消息包通信。

②将相邻的3个二维变量边界交换聚合后一次实现。

经过优化将三维数据和四维数据中的二维边界通信进行聚合,聚合后总的边界通信次数由25350降低至10358次。各个并行度的两种数据剖分方式的通信优化效果见表8和图3。

表7 POP积分6小时通信优化效果

并行剖分	128		256		512	
	笛卡儿	平衡	笛卡儿	平衡	笛卡儿	平衡
优化前时间	468.58	405.50	228.06	216.72	108.82	106.44
优化后时间	463.47	394.63	220.94	187.29	98.09	77.30
优化效率	1.3%	2.68%	3.12%	13.58%	9.86%	27.37%

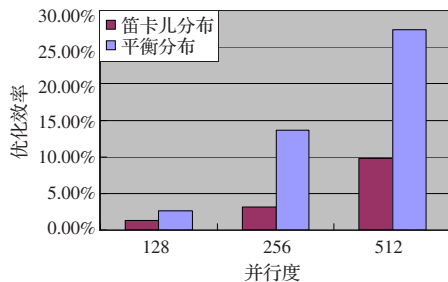


图4 POP通信优化在至强集群上的性能提升

优化效率定义为:(优化前时间-优化后时间)/优化前时间。可见通信聚合优化对于平衡数据分布方式来说效果更为显著,这显示了负载均衡和降低通信次数对于多核并行程序性能均有重要意义。

(上接150页)

由于半监督的SVM分类算法还是一个较新的领域,还有许多方面的课题留待进一步研究。算法的性能取决于有标签样本的分布情况,在某些样本分布特性下,该算法的性能不如PTSVM算法理想。如何使SVM算法更好地适应各种分布特性,如何更进一步地降低算法的时间复杂度和提高样本的精度,这些都是今后所要进一步研究的课题。

参考文献:

- [1] Joachims T. Transductive inference for text classification using support vector machines[C]//Proceedings of the 16th International Conference on Machine Learning (ICML). San Francisco: Morgan Kaufmann Publishers, 1999: 200-209.
- [2] Altun Y, McAllester D, Belkin M. Maximum margin semi-supervised learning for structured variables[C]//Advances in Neural Information Processing, 2005.

5 结论和展望

在四核至强集群上,POP模式中所采用的局部块技术和平衡数据剖分对并行性能有着显著的影响,这两项技术可为从事大型并行程序开发人员提供提高性能的技术参考。通过对POP模式的高维数组采用聚合通信优化,降低通信次数,有效提升了程序性能,这一点特别对平衡数据剖分尤为显著。

从提高POP模式性能的角度分析,目前外模态中的共轭梯度迭代,只采用简单的对角预条件子,未来可以实现更为高效的预条件子,也可以采用通用并行数值软件平台(PETSc)和高性能预条件软件平台(HYPRE)进行性能改进;对并行数据剖分,未来可采用权重空间填充曲线技术,进一步改善负载均衡。对集群多核特点,也需要进一步分析优化POP模式的混合并行性能。

致谢:感谢朱敏副研究员在并行计算环境方面提供的帮助。

参考文献:

- [1] 陈显尧,宋振亚,王永刚.并行计算在海洋环流数值模式中的应用[J].高性能计算发展与应用,2005(4).
- [2] Dennis J M. Expedition computing: exploring the petascale frontier [EB/OL]. (2007). http://www.cisl.ucar.edu/dir/CAS2K7/final_agenda-2007.html.
- [3] Smith R D, Gent P. Reference manual for the Parallel Ocean Program (POP), Los Alamos Unclassified Report LA-UR-02-2484[R]. 2002.
- [4] 伍湘君,金之雁,陈德辉,等.新一代数值预报模式GRAPES的并行方案设计与实现[J].计算机研究与发展,2007(3).
- [5] Kerbyson D J, Jones P W. A performance model of the parallel ocean program[J]. The International Journal of High Performance Computing Applications, 2005, 19(3): 261-276.
- [6] Kim Dong-Hoon, Nakashiki N, Yoshida Y. Computation of super high-resolution global ocean model using earth simulator[C]//Proceedings of Coastal and Ocean Engineering in Korea, 2003.
- [7] Zaki T, Moulton D, Nadig B, et al. Multigrid preconditioning in fully-implicit evolution of the ocean[R/OL]. T-7, MS B284, Theoretical Division, Los Alamos National Laboratory, Los Alamos, NM 87545. <http://math.lanl.gov/>.
- [8] Jones P W, Worley P H, Yoshida Y, et al. Practical performance portability in the Parallel Ocean Program (POP)[J]. Concurrency and Computation: Practice and Experience, 2005, 17: 1317-1327.
- [3] Chapelle O, Zien A. Semi-supervised classification by low density separation[C]//Proceedings of the 10th International Workshop on Artificial Intelligence and Statistics (AISTAT 2005), 2005.
- [4] 陈毅松,汪国平,董士海.基于支持向量机的渐进直推式分类学习算法[J].软件学报,2003,14(3):451-460.
- [5] Burges C J C. A tutorial on support vector machines for pattern recognition[Z]//Data Mining and Knowledge Discovery, 1998: 121-167.
- [6] Cristianini N, Shawe-Taylor J. 支持向量机导论[M].北京:电子工业出版社,2004: 82-90.
- [7] 许建华,张学工.一种基于核函数的非线性感知器算法[J].计算机学报,2002, 25.
- [8] 张召,黄国兴.一种改进的SMO算法[J].计算机科学,2003, 30.
- [9] 骆世广,杨晓伟.一种改进的序贯最小优化算法[J].计算机科学,2006, 33.
- [10] 边肇祺,张学工.模式识别[M].2版.北京:清华大学出版社,2000: 86-90.