

ZHICHI XIANGLIANGJI YU YOUHUA SUANFA

ZAI LINXIASHEN GUANGHUANJING PINGJIA XITONG ZHONG DE YANJIU

支持向量机与优化算法 在林下参光环境评价系统中的研究

武海巍 著



東北大學出版社
Northeastern University Press

支持向量机与优化算法 在林下参光环境评价系统中的研究

武海巍 著



东北大学出版社
·沈阳·

© 武海巍 著 2013

图书在版编目 (CIP) 数据

支持向量机与优化算法在林下参光环境评价系统中的研究/武海巍著. —沈阳: 东北大学出版社, 2013. 12

ISBN 978-7-5517-0519-6

I. ①支… II. ①武… III. ①向量计算机—应用—人参—自然光—环境质量评价②最优化算法—应用—人参—自然光—环境质量评价 IV. ①S567. 5

中国版本图书馆 CIP 数据核字 (2013) 第 300766 号

出版者: 东北大学出版社

地址: 沈阳市和平区文化路 3 号巷 11 号

邮编: 110819

电话: 024—83680267 (社务室) 83687331 (市场部)

传真: 024—83680265 (办公室) 83680178 (出版部)

网址: <http://www.neupress.com>

E-mail: neuph@neupress.com

印刷者: 沈阳市第二市政建设工程公司印刷厂

发行者: 东北大学出版社

幅面尺寸: 165mm × 235mm

印 张: 7.25

字 数: 108 千字

出版时间: 2013 年 12 月第 1 版

印刷时间: 2013 年 12 月第 1 次印刷

组稿编辑: 郭爱民

责任编辑: 潘佳宁

责任校对: 叶 子

封面设计: 刘江旸

责任出版: 唐敏志

ISBN 978-7-5517-0519-6

定 价: 20.00 元

摘要

支持向量机是一种基于统计学习理论的机器学习方法，在解决小样本、非线性及高维模式识别中表现出许多特有的优势，已在各个领域得到广泛应用，但在林下参光环境领域中的应用研究尚未见报道。作为在支持向量机中起着决定性作用的核函数，已引起国内外学者的关注，构建新型核函数成为其研究热点之一。通过研究自然界中一些自然现象而总结出的仿生智能优化算法，能够可靠解决全局最优化问题，且这些优化算法具有普遍适应性。支持向量机中的参数优化程度影响着构建模型的预测精度和泛化能力，将仿生智能算法用于参数优化过程，可寻找出最佳的预测模型。

林下资源是人类宝贵的物质资源，人参是最重要的林下资源之一，其作为阴性植物，对生态环境提出了很高的要求，尤其是对光环境极具敏感性。本研究瞄准本领域的科技发展前沿，选取林下种植生态环境中重要的自然光环境问题，设计了一种人性化、便于日后扩展其他功能的光照强度监控系统；以生态系统的物理和生物学原理为基础，利用系统分析和机器学习方法，建立林下参种植的自然光环境动态模型。

本书的主要研究内容和结论如下。

(1) 构造了新型核函数。研究核函数原理，分析常见核函数中高斯核函数 $K(\mathbf{x}, \mathbf{x}_i) = \exp(-\gamma \|\mathbf{x} - \mathbf{x}_i\|^2)$ 、多项式核函数 $K(\mathbf{x}, \mathbf{x}_i) = (\langle \mathbf{x}, \mathbf{x}_i \rangle + c)^d$ 和感知器核函数 $K(\mathbf{x}, \mathbf{x}_i) = \tanh(p\langle \mathbf{x}, \mathbf{x}_i \rangle + c)$ 的性能，采用高斯核函数 $K_1(\mathbf{x}, \mathbf{x}_i) = \exp(-\gamma \|\mathbf{x} - \mathbf{x}_i\|^2)$ 和多项式核函数 $K_2(\mathbf{x}, \mathbf{x}_i) = (\langle \mathbf{x}, \mathbf{x}_i \rangle + c)$ 为基准核函数，构造核函数 $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}')$ ， $K(\mathbf{x}, \mathbf{x}') = \lambda K_*(\mathbf{x}, \mathbf{x}')$ ， $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}')K_2(\mathbf{x}, \mathbf{x}')$ ， $K(\mathbf{x}, \mathbf{x}') = \exp(K_1(\mathbf{x}, \mathbf{x}'))$ 和 $K(\mathbf{x}, \mathbf{x}_i) = \sum_{i=1}^2 a_i K_i$ ，并对各构造核

函数进行性能分析。通过特征曲线分析得知，高斯核函数具有较强的局部学习能力，并且参数 γ 影响决策函数的判断能力： γ 值变大，则训练集中的支持向量增多，模型泛化能力下降， γ 过大会导致过学习现象； γ 值变小，训练集中样本被错误分类的几率增大， γ 过小会导致欠学习现象。多项式核函数具有较强的全局学习能力， d 值的增加会增强该能力。感知器核函数兼顾局部学习能力和全局学习能力，从其特征曲线得知，该核函数的局部学习能力较强，而全局学习能力较弱，可以通过减小 p 值增加其全局学习能力，但同时局部学习能力相应减弱。采用式 $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}')$ 方式构造的核函数 $K(\mathbf{x}, \mathbf{x}_i)$ ，兼备局部学习能力和全局学习能力，其局部学习能力完全由高斯核函数 $K_1(\mathbf{x}, \mathbf{x}_i)$ 决定，全局学习能力则完全由多项式核函数 $K_2(\mathbf{x}, \mathbf{x}_i)$ 决定。采用式 $K(\mathbf{x}, \mathbf{x}') = \lambda K_*(\mathbf{x}, \mathbf{x}')$ 方式构造的核函数 $K(\mathbf{x}, \mathbf{x}_i)$ ，其性能完全取决于基准函数的特性，若基准核函数采用 $K_1(\mathbf{x}, \mathbf{x}_i)$ ，则 $K(\mathbf{x}, \mathbf{x}_i)$ 仅改善局部学习能力；若基准核函数采用 $K_2(\mathbf{x}, \mathbf{x}_i)$ ，则 $K(\mathbf{x}, \mathbf{x}_i)$ 仅改善全局学习能力；采用式 $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}')K_2(\mathbf{x}, \mathbf{x}')$ 方式构造的核函数 $K(\mathbf{x}, \mathbf{x}_i)$ ，特征曲线变化规律与 $K_1(\mathbf{x}, \mathbf{x}_i)$ 相类似，呈现出局部学习能力强，全局学习能力弱的特点。通过调节 $K_1(\mathbf{x}, \mathbf{x}_i)$ 和 $K_2(\mathbf{x}, \mathbf{x}_i)$ 的参数，仅对 $K(\mathbf{x}, \mathbf{x}_i)$ 的局部学习能力有所改善，而对全局学习能力影响甚微。采用式 $K(\mathbf{x}, 0.2) = a_1 \exp(-2 \|\mathbf{x} - 0.2\|^2) + a_2 (\langle \mathbf{x}, 0.2 \rangle + 1)$ 构造的核函数，可以通过调节 $a_i (i = 1, 2)$ 值达到调节其局部学习能力和全局学习能力的目的。

(2) 构建了新型仿生智能算法——追踪算法。利用 Needle-in-a-haystack 函数和 Schaffer 函数检验遗传算法、粒子群算法和追踪算法的全局寻优能力。Needle-in-a-haystack 函数，在区间 $-5 \leq x \leq 5, -5 \leq y \leq 5$ 范围内有全局最小值 $f(0,0) = -3600$ ，对于 Schaffer 函数在区间 $-20 \leq x \leq 20, -20 \leq y \leq 20, i = 1, 2$ 范围内有全局最小值 $f(0,0) = 0$ ，遗传算法取个体数目为 20，最大遗传代数为 200，变量的二进制位数为 25，交叉概率为 0.9，变异概率为 0.08，对于 Needle-in-a-haystack 函数，经过遗传代数为 40 代以后，最佳适应度值趋于全局最优解 -3600，

对于 Schaffer 函数，经过遗传代数为 20 代以后，最佳适应度值陷入局部最优解 0.085；粒子群算法取进化代数为 200，种群规模为 20，对于 Needle-in-a-haystack 函数，经过进化代数为 10 代以后，最佳适应度陷入局部最优解 -2500，对于 Schaffer 函数，经过进化代数为 60 代以后，最佳适应度趋于全局最优解 0；追踪算法取种群规模为 200，个体追踪路程为 20，对于 Needle-in-a-haystack 函数，经过搜索代数为 20 代以后，最佳适应度趋于全局最优解 -3600，对于 Schaffer 函数，经过搜索代数为 20 代以后，最佳适应度趋于全局最优解 0。

(3) 利用 TSL2561 对可见光敏感特性，ATMega16L 具备 I²C 和 SPI 总线功能，采用主机 - 从机架构，结合上位机监控软件，设计了林下参光照强度实时监控系统。本系统采用单个主机、多个从机的 SPI 总线连接方式，使得主从机传输距离达到 1000 m，为今后更方便增加试验单位监控点提供保障。该方法简单易行且光照强度传输数据受外界因素干扰较小，适用于林下参光环境中所要测量的光照强度区域较大的特点。将各从机与 TNHY - 9 监测仪放入标准光照环境中，测试光照强度。在数据传输过程中，加入校验码以确保数据传输的可靠性，加入从机编号以区分林下基地不同试验单位。主 - 从机传输距离为 500 m 时，除 5 号从机所测得光照强度值为 341 lux 外，其余从机所测值均为 340 lux，其方差为 8.5。主 - 从机传输距离增加为 1000 m 时，各从机所得光照强度开始变化，其方差为 11.5。可见，主 - 从机传输距离的增加使得该系统性能有所下降，但从整体看，系统主 - 从机传输距离不超过 1000 m 时，本系统对光照强度数据传输具备相当可靠性。本研究构建的实时监测软件可以根据林下光环境及对光照强度测量精度的要求，调整采样频率，为林下光环境数据的测量提供了新方法。

(4) 本研究利用支持向量机建立预测模型，通过可见光光谱组成成分分配比关系，预测个体净光合速率 (Pn)，通过直射辐射 (PFDdir) 和散射辐射 (PFDdif) 预测光合有效辐射 (PAR)，为林下光环境的预测和评价提供了新方法。采用 epsilon-SVR 公式，nu-SVR 公式，linear 核函数 (K_1)，polynomial 核函数 (K_2)，radial basis function 核函数

(K_3), $K(\mathbf{x}, \mathbf{x}_i) = \sum_{i=1}^3 a_i K_i$ 核函数 ($a_i \geq 0$ 且 $\sum_{i=1}^3 a_i = 1$), $K_1 K_2$ 核函数, $K_1 K_3$ 核函数, $K_2 K_3$ 核函数, $K_1 K_2 K_3$ 核函数, 惩罚参数 c 和 gamma 值采用 grid-search, 遗传算法, 粒子群算法和追踪算法进行参数寻优, 以上多种组合建立不同的支持向量模型, 在加入其他影响因素的 ϵ 粒子后, 经进行交叉试验, NRTA 模型为预测 Pn 的最优模型, 对 2011 年 8 月 14 日—8 月 28 日的 Pn 拟合程度为 90.903%; EGSK (0.1, 0, 0.9) 模型为预测 PAR 的最优模型, 对 2010 年 7 月 21 日—7 月 30 日的 PAR 拟合程度为 86.897%。

关键词: 构造核函数; 仿生智能算法; 实时监控; 光环境

Abstract

Support vector machine (SVM) is a kind of machine learning method based on statistical learning theory, and shows many unique advantages in solving small sample, nonlinear and high dimension data. It is applied to many fields, except in ginseng under forest light environment. Kernel is the key in SVM and many scholars in the world are interested in it. It is a hot topic on constructing a new kernel function. Bionic intelligent optimization algorithm is sourced from the study about some natural phenomena. It can solve global optimization problems. Moreover, it has adaptability. Parameters optimization of SVM affects the prediction accuracy and generalization ability. If bionic intelligent algorithm is used to optimize the parameters, the best predictive model will be established.

Forest resources are valuable resources, and ginseng is the most important one of them. As a kind of plants living in the cold and humid environment, it is highly sensitive to optical environment. This paper researches of the field of ginseng under the forest planting, and study the question about natural light environment. In the paper, we design a kind of humanity, to facilitate future expansion to other features of the light intensity monitoring system. Based on physical and biological principle in ecological system, system analysis and machine learning method, this paper establishes the light environment impacting ginseng growing dynamic model.

The main research contents and conclusions are as follows:

- (1) Constructing some new type of kernel functions. Research nuclear function principle, and be analysis of common kernel function, such as

gauss kernel function $K(\mathbf{x}, \mathbf{x}_i) = \exp(-\gamma \|\mathbf{x} - \mathbf{x}_i\|^2)$, polynomial kernel function $K(\mathbf{x}, \mathbf{x}_i) = (\langle \mathbf{x}, \mathbf{x}_i \rangle + c)^d$ and perceptron kernel function $K(\mathbf{x}, \mathbf{x}_i) = \tanh(p\langle \mathbf{x}, \mathbf{x}_i \rangle + c)$. Construct these kernels and research their function, which are $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}')$, $K(\mathbf{x}, \mathbf{x}') = \lambda K_*(\mathbf{x}, \mathbf{x}')$, $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}')K_2(\mathbf{x}, \mathbf{x}')$, $K(\mathbf{x}, \mathbf{x}') = \exp(K_1(\mathbf{x}, \mathbf{x}'))$ and $K(\mathbf{x}, \mathbf{x}_i) = \sum_{i=1}^2 a_i K_i$ based on gauss kernel function $K_1(\mathbf{x}, \mathbf{x}_i) = \exp(-\gamma \|\mathbf{x} - \mathbf{x}_i\|^2)$, polynomial kernel function $K_2(\mathbf{x}, \mathbf{x}_i) = (\langle \mathbf{x}, \mathbf{x}_i \rangle + c)$. Through analysising the characteristic curve, gauss kernel function has a strong local learning ability, and the parameters γ influencing the decision function judgement: If γ value is larger, the support vector number is more and the model generalization ability is more. But it can lead to overfit phenomenon while γ value is over a certain range. If γ value becomes smaller, the training set samples were misclassified probability. But it will cause less learning phenomenon if γ value is too small. Polynomial kernel function has a strong global learning ability. The larger d value will lead to enhance the ability of global learning ability. Perceptron nuclear balances local learning ability and global learning ability. Analysising of the special curve, the kernel has the stronger local learning ability, and the global learning ability is weaker. It can increase its global learning ability by decreasing the p value, but the local learning ability is weakened accordingly. The kernel function $K(\mathbf{x}, \mathbf{x}_i)$ sourced from the $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}') + K_2(\mathbf{x}, \mathbf{x}')$ mode structure has local and global learning ability, the local learning ability is impacted completely by the gauss kernel function and global learning ability is impacted completely by the polynomial kernel function. If the kernel function $K(\mathbf{x}, \mathbf{x}_i)$ is sourced from the $K(\mathbf{x}, \mathbf{x}') = \lambda K_*(\mathbf{x}, \mathbf{x}')$ mode structure, its performance depends solely on the baseline function property. If the baseline kernel function used $K_1(\mathbf{x}, \mathbf{x}_i)$, $K(\mathbf{x}, \mathbf{x}_i)$ is only to improve local learning ability; If the baseline kernel function used $K_2(\mathbf{x}, \mathbf{x}_i)$, $K(\mathbf{x}, \mathbf{x}_i)$ is only to improve global learning ability.

\mathbf{x}_i) is only improve global learning ability. If the kernel function $K(\mathbf{x}, \mathbf{x}_i)$ is sourced from the $K(\mathbf{x}, \mathbf{x}') = K_1(\mathbf{x}, \mathbf{x}')K_2(\mathbf{x}, \mathbf{x}')$ mode structure, its feature curve is similar with $K_1(\mathbf{x}, \mathbf{x}_i)$'s, which shows the strong local learning ability and the weak global learning ability. By adjusting the $K_1(\mathbf{x}, \mathbf{x}_i)$ and $K_2(\mathbf{x}, \mathbf{x}_i)$ parameters, it is improved to the local learning ability, but have little effect on the global learning ability. If the kernel function $K(\mathbf{x}, 0.2)$ sourced from the $K(\mathbf{x}, 0.2) = a_1 \exp(-2 \|\mathbf{x} - 0.2\|^2) + a_2 (\langle \mathbf{x}, 0.2 \rangle + 1)$ model structure, it can adjust its local learning ability and learning ability of the global by adjusting the $a_i (i = 1, 2)$ value.

(2) Constructing a new type of intelligent bionic algorithm—Tracing Target Algorithm. By using of Needle-in-a-haystack function and Schaffer function, it test the global optimization ability about genetic algorithm, particle swarm optimization algorithm and tracking algorithm. The Needle-in-a-haystack function has a global minimum $f(0, 0) = -3600$ in the range of $-5 \leq x \leq 5, -5 \leq y \leq 5$, and the Schaffer function has a global minimum $f(0, 0) = 0$ in the range of $-20 \leq x \leq 20, -20 \leq y \leq 20, i = 1, 2$. When individual numbers are 20, the maximum algebras are 200, variable binary digits are 25, the probability of crossover is 0.9 and mutation probability is 0.08, genetic algorithm finds the best fitness value to global optimal solution -3600 of Needle-in-a-haystack function through genetic algebra of 40 generations later, and finds the best fitness value into a local optimal solution 0.085 of Schaffer functions through genetic algebra of 20 generations later. When the evolution algebras are 200, population sizes are 20, particle swarm algorithm finds the best fitness into local optimal solution -2500 of Needle-in-a-haystack function after a number of generation is 10, and finds the best fitness tend to global optimal solution 0 of the Schaffer function after a number of generation is 60. When population sizes are 200, tracing the distance are 20, Tracking algorithm finds the best fitness tend to global optimal solution -3600 of Needle-in-a-haystack functions

after searching for algebra of 20 generations, and finds the best fitness tend to global optimal solution 0 of schaffer functions after searching for algebra of 20 generations later.

(3) By using of TSL2561 being sensitive to visible light characteristics and ATMega16L with I²C and SPI bus function, this paper designs the ginseng under forest light intensity real-time monitoring system, which is the host-from machine architecture and combined with the monitoring software of the host computer. The transmission distance is 1000 m between the master and slave in the system using a single host and multiple from machine SPI bus connection method, which is more convenient to increase the test unit monitoring point in the future. The method is simple and the light intensity transmission data have little interference with these outside factors, so the system is suitable for ginseng under forest light environment in which the light intensity is measured. It measures light intensity by putting these slaves and TNHY-9 monitor in the standard illumination environment. In the process of data transmission, it ensures the reliability of data transmission by joining the check codes, and it distinguishes different test unit of forest bases by the serial slaves numbers. When the distance is 500 m between master and slave transmission, except for the measured light intensity value of 341 lux by the number 5 slaver, the measured light intensity value is 340 lux by the other slavers, and the variance is 8. 5. When the distance is 1000 m between master and slaver, the light intensity are different among all slavers, its variance is 11. 5. The performance of the system is decreased for the distance increasing between Master and slaver, but in the whole, the system of the light intensity data transmission is considerably reliability when the distance is not more than 1000 m between master and slaver. The real time monitoring system adjusts the sampling frequency based on requirements of measurement accuracy of light intensity in the forest light environment. It provides a new method for the measurement of light environ-

mental data under the forest.

(4) This paper uses support vector machine to establish the forecasting model, and predicts individual net photosynthetic rate (Pn) by using of the visible light spectral composition proportion relation, and predicts Photosynthetic active radiation (PAR) by using of Direct radiation (PFDdir) and scattering radiation (PFDdif). It provides a kind of new method for light environmental forecasting and evaluation under the forest. This paper uses the epsilon-SVR formula, the formula of nu-SVR, linear kernel function (K_1), polynomial kernel function (K_2), radial basis function kernel function (K_3), kernel function $K(\mathbf{x}, \mathbf{x}_i) = \sum_{i=1}^3 a_i K_i$ ($a_i \geq 0$ and $\sum_{i=1}^3 a_i = 1$), kernel function $K_1 K_2$, kernel function $K_1 K_3$, kernel function $K_2 K_3$, kernel function $K_1 K_2 K_3$, penalty parameters c and gamma optimized by using grid-search, genetic algorithm, particle swarm algorithm and tracking algorithm parameters optimization, and establishes different support vector model by the above combination. Mixed the other influencing factors called the ϵ particle, NRTA model for predicting Pn is the best optimal model. The fitting degree is 90.903% when NRTA model predicts Pn sourced from August 14, to August 28, 2011. EGSK (0.1, 0, 0.9) model for predicting PAR is the best optimal model and the fitting degree is 86.897% when this model predicts PAR sourced from July 21 to July 30, 2010.

Keywords: Constructing Kernel Function, Bionic Intelligent Algorithm, Real Time Monitoring, Light Environment

目 录

第1章 绪 论	1
1.1 机器学习背景	1
1.2 支持向量机简介	2
1.3 林下参光环境研究现状	5
1.4 相关研究存在的问题	6
1.5 本书研究目标和研究内容	7
1.5.1 研究目标	7
1.5.2 研究内容	7
1.6 本章小结	8
第2章 支持向量机	9
2.1 引 言	9
2.2 最优分类超平面理论	9
2.2.1 线性可分情况.....	10
2.2.2 非线性可分情况.....	13
2.3 相似程度与内积	15
2.3.1 相似程度概述.....	15
2.3.2 两点相似程度与内积关系.....	15
2.3.3 三点相似程度与内积关系.....	16
2.3.4 线性分类机与内积关系.....	18
2.4 核函数的引入	19
2.4.1 二次分划问题.....	19

2.4.2 核函数原理	20
2.5 常见核函数与性能分析	22
2.5.1 常见核函数	22
2.5.2 常见核函数的性能分析	23
2.6 核函数的构造与性能分析	28
2.6.1 核函数的构造	28
2.6.2 构造核函数的性能分析	30
2.7 本章小结	42
第3章 优化算法	43
3.1 引言	43
3.2 进化类算法	43
3.2.1 进化型算法介绍	44
3.2.2 遗传算法	44
3.3 群智能算法	48
3.3.1 群智能算法介绍	48
3.3.2 粒子群算法	48
3.3.3 追踪算法	50
3.4 三种算法性能比较	52
3.4.1 寻找 Needle-in-a-haystack 函数的全局最优解	52
3.4.2 寻找 Schaffer 函数的全局最优解	55
3.5 本章小结	58
第4章 林下参光照强度实时监控系统构建	59
4.1 引言	59
4.2 林下光照强度实时监控系统构建	60
4.2.1 光照强度测定方法	60
4.2.2 主机与从机接口	62
4.2.3 主机与上位机接口	62

4.2.4 系统数据传输可靠性分析	62
4.2.5 从机不同分布对光照强度测试结果影响分析	64
4.2.6 多套单个主机、10个从机组成的系统组合分析	65
4.2.7 上位机控制光照强度测定分析	65
4.3 本章小结	66
第5章 林下参净光合速率预测模型	69
5.1 引言	69
5.2 光谱测定	70
5.3 净光合速率影响因素分析与数据处理	70
5.4 支持向量机建模	72
5.5 参数寻优	73
5.6 试验结果分析	73
5.7 本章小结	78
第6章 核函数组合优化光合有效辐射预测模型	79
6.1 引言	79
6.2 数据测定	80
6.3 数据处理	80
6.4 支持向量机建模	82
6.5 试验结果分析	83
6.6 本章小结	88
第7章 结论	89
参考文献	92

第1章 绪论

【提要】本章简要介绍机器学习背景和支持向量机在国内外各领域的应用研究成果，阐述林下参光环境研究现状，对建立林下参种植的自然光环境动态模型和实现选定拟种植区域光环境的计算机预测系统的必要性进行分析，指出在支持向量机和林下光环境研究中存在的主要问题，并提出本书研究目标和主要研究内容。

1.1 机器学习背景

构造从经验中学习的机器是科技界的研究目标之一，机器学习从技术角度来看，是从电子计算机的发明中获得了强大原动力，其讨论的是如何让计算机程序进行学习，从已有经验中进行学习，来提高机器的性能^[1,2]。机器学习是一种学习领域，赋予计算机学习的能力，但没有经过显式编码，其最早的一个非正式描述是由 Arthur Samuel 在 1955 年提出的：“Field of study that gives computers the ability to learn without being explicitly programmed^[3]”。更现代的定义是 1998 年由卡内基·梅隆大学的 Tom Mitchell 提出的：“A computer program is said to learn from experience E, with respect to some task T, and some performance measure P, if its performance on T, as measured by P, improves with experience E”^[3]。使用机器学习方法来解决某个任务，首先对这个任务选取合适的原型，如线性回归，Logistic 回归，朴素贝叶斯等，然后通过经验来优化性能度量。原型确定后，需要根据一定的方法来调整原型参数，从而达到优化性能度量的目的。机器学习理论很大程度上都是在讲优化方法，如最小训练误差，最小均方差，最大似然率，凸集优化

等，其学习过程，就是利用经验来对性能度量最优化的过程。在实际应用中，我们并不知道最优是个什么样子，机器学习的结果是对最优值一个估计，这个估计以大概率收敛于最优值。

根据经验不同，机器学习可以分为以下三类：监督学习（Supervised Learning），无监督学习（Unsupervised Learning）及增强学习（Reinforcement Learning）^[3]。监督学习是指在训练经验中明确告诉正确结果，期望在学习后能正确识别出所学习的种类。若能识别出来，则称为收敛；若识别不出来，则说明还需要继续训练。监督学习算法的输出如果是连续的，称为回归（Regression）；如果是离散的，称为分类（Classification）。大部分的机器学习任务都是监督学习。在无监督学习过程中，只有训练样本而没有正确结果，其常用方法是聚类。增强学习只对程序的行为做出评价，程序就会做出更有可能得到正面评价的行为，这在机器人领域中已得到了广泛应用。

1.2 支持向量机简介

支持向量机（Support Vector Machine, SVM）是一种基于统计学习理论的机器学习方法，由 Cortes 和 Vapnik 于 1995 年首先提出，它在解决小样本、非线性及高维模式识别中表现出许多特有的优势，并能够推广应用到函数拟合等其他机器学习问题中^[3-10]。

SVM 建立在统计学习理论的 VC 维理论和结构风险最小原理基础上，根据有限样本信息，在模型的复杂性（即对特定训练样本的学习精度，Accuracy）和学习能力（即无错误地识别任意样本的能力）之间寻求最佳折中，以期获得最好的推广能力^[1,3,11-13]（或称泛化能力），已经在许多智能信息获取与处理领域的应用中获得成功。但是 SVM 也并非完美，在实际应用中，经常出现在训练时模型的精度很好，到预测时精度下降很多直至不符合精度要求的情况，这是 SVM 中参数比较难选择特点的体现。