

基于随机傅里叶特征空间的高斯核近似模型选择算法

张 凯¹, 门昌骞¹, 王文剑^{1,2}

(1. 山西大学计算机与信息技术学院, 太原 030006; 2. 山西大学计算智能与中文信息处理教育部重点实验室, 太原 030006)

摘 要:核方法是一种把低维空间的线性不可分问题转化为高维空间中线性可分问题的方法,其广泛应用于多种学习模型。然而现有的核模型选择方法在大规模数据中计算效率较低,时间成本很大。针对这一问题,本文引入随机傅里叶特征变换,将原始核特征空间转换为另一个相对低维的显式随机特征空间,并给出核近似误差上界理论分析以及在核近似的随机特征空间中训练学习模型的误差上界,得到核近似的收敛一致性和误差上界与核近似参数之间的关系。基于随机傅里叶特征空间选择出最优模型参数,避免了对最优原始高斯核模型参数的大规模搜索,从而大幅降低原始高斯核模型选择所需的时间成本。实验表明,本文给出的误差上界确由核近似参数控制,核近似选择的最优模型相较于原始高斯核模型有较高的准确率,并且模型选择时间相对网格搜索法大幅减小。

关键词:核方法;高斯核;傅里叶变换;核近似;模型选择

中图分类号: TP301 **文献标志码:** A

Gaussian Kernel Approximation Model Selection Algorithm Based on Random Fourier Feature Space

ZHANG Kai¹, MEN Changqian¹, WANG Wenjian^{1,2}

(1. College of Computer and Information Technology, Shanxi University, Taiyuan 030006, China; 2. Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education, Shanxi University, Taiyuan 030006, China)

Abstract: Kernel method transforms the linear non-separable problem in low-dimensional space into the linear separable problem in high-dimensional space. It is widely used in a variety of learning models. However, the existing kernel selection methods have low computational efficiency and high time cost in large-scale data. Aiming at above problems, this paper introduces the random Fourier feature to transform the original kernel feature space into another relatively low dimensional explicit random feature space. The theoretical analysis of the upper bound of the kernel approximation error and the upper bound of the error of training the learning model in the kernel approximation random feature space are given. The convergence consistency of kernel approximation and the relationship between error upper bound and kernel approximation parameters are obtained. Moreover, the optimal model parameters are selected based on random Fourier feature space, which can avoid the large-scale search for the optimal original Gaussian kernel model parameters, so as to greatly reduce the time cost required for the selection of the original

基金项目:国家自然科学基金(U21A20513, 62076154);中央引导地方科技发展资金(YDZX20201400001224);山西省国际科技合作计划项目(201903D421050);山西省自然科学基金(201901D1111030)。

收稿日期: 2022-05-12; **修订日期:** 2022-08-28

Gaussian kernel model. Experiments show that the error upper bound proved in this paper is controlled by the kernel approximation parameters. The optimal model selected by the kernel approximation has good performance compared with the original Gaussian kernel function model, and the model selection time is greatly reduced compared with the grid search method.

Key words: kernel method; Gaussian kernel; Fourier transform; kernel approximation; model selection

引 言

核方法为解决机器学习中的许多问题提供了强大、灵活和坚实的理论基础。其核心思想是通过某种非线性映射将原始数据映射到合适的高维特征空间,再利用线性学习器在高维特征空间中分析和处理数据。核方法需要借助核函数实现低维空间向高维空间的转换,不同的核函数对应不同的特征空间。核方法的理论基础 Cover's theorem^[1],对于非线性可分的训练集,可以大概率通过将其映射到一个高维空间来转化成线性可分的训练集。关于核方法对应的线性模型已经有大量的研究。经过核函数进行高维特征空间转化后,原始非线性问题可转化为高维线性可分问题,大量线性模型所对应的算法被提出,包括随机梯度下降算法^[2]、割平面算法^[3]以及对偶坐标下降算法^[4]等。利用具有线性时间复杂度的线性模型算法,普通PC机就可以有效地处理大规模问题^[5]。线性模型虽然在超高维文本和基因分类问题上能取得较好的学习效果,但是对于一般问题,预测精度很难超越非线性核模型,并且在计算机视觉、自然语言处理、系统设计和许多其他领域中的学习问题规模通常在数十万左右,甚至可能超过数百万,这就给运用非线性模型解决实际问题带来巨大的计算困难。然而,线性或亚线性时间复杂度的求解算法,计算开销小,易于扩展到大规模问题,这促进了应用线性学习算法求解非线性问题的研究工作。非线性核模型的计算时间复杂度在理论上不低于 $O(m^3)$,因此研究的重点在于设计快速核近似算法。利用核近似方法逼近原始核是利用线性算法来求解非线性问题的一种方法。随着核近似算法的深入研究,一些关于高斯核的核近似算法被提出,如Nyström方法^[6]。Yang等^[7]在低阶泰勒公式中近似高斯核,从而加速共轭梯度优化算法。Cao等^[8]同样在低阶泰勒公式下近似高斯核,然后加速支持向量机(Support vector machine, SVM)预测。

随着大量研究的开展,出现了许多不同的核近似方法,其中Rahimi等^[9]提出了一种随机傅里叶特征映射方法,该方法主要思想是通过近似平移不变核显式构造特征映射,保证任意两点对应随机特征映射的内积近似其对应的核函数值,从而为显式构造假设空间求解大规模核方法提供了一种有效途径。Feng等^[10]在此基础上构造了一个循环随机假设空间,该方法运行效率较高,但是解释性不足。文献[11]定义了近似高斯核并建立其与高斯核的关系,实现了用线性SVM算法高效求解大规模非线性SVM。另外,还有一些利用核近似来构建特征空间的算法,文献[12]显式描述核函数对应积分算子的特征值与特征函数,并构造其对应的再生核希尔伯特空间和特征映射。文献[13]在大规模样本进行子采样,然后应用随机傅里叶映射显式构造随机空间,从而达到高效求解的目的。已有工作对于核近似的误差上界描述较宽松,并且对于随机傅里叶特征空间维度的关系定义不够清晰。在此基础上,本文进一步分析了核函数中基于随机傅里叶特征变换的核近似方法对支持向量机、核岭回归模型这两种广泛使用的学习模型影响,并且给出这些算法对应的稳定性界。通过对核近似误差上界的分析,可以证明算法的准确率主要和傅里叶特征变换的空间维度以及核函数参数相关。模型选择通常采用含交叉验证的网格搜索法找到最优参数组合,但这种方法需要多次训练模型,在大规模数据中多次训练模型所需的时间成本巨大。针对这一问题,本文在核近似的基础上提出一种在近似空间上的模型选择算法。该算法通过在近似线性空间上高效选择模型参数,将对应参数组合运用到原始核中,在选择出最

优参数的同时,得到的最优模型参数在原始核中仍有较好的效果,解决了传统模型选择的多次迭代,大幅降低了模型选择的时间成本,达到模型选择的目的。

1 预备知识

本节简要介绍随机傅里叶特征变换的基本概念。

定理 1 Bochner 定理^[14]。连续函数 $f: \mathbf{R}^d \mapsto S$ 正定,当且仅当 f 为某一有限非负 Borel 测度 μ 的傅里叶变换,则

$$f(x) = \int_{\mathbf{R}^d} e^{-ix^T v} d\mu(v) \quad (1)$$

式中 x^T 为输入矩阵; v 取自概率分布 $p(v)$ 。

选择合适的核函数 $k(\cdot)$, 由 Bochner 定理可以确定存在一个概率分布为 $p(\cdot)$ 的傅里叶变换与 $k(\cdot)$ 对应, 即有如下推论。

现有核函数种类众多, 在众多核函数中高斯核是一种通用核, 使用比较广泛其定义形式为

$$K(x, y) = \langle \phi(x), \phi(y) \rangle = \exp(-\gamma \|x - y\|^2) \quad (2)$$

对于高斯核函数, 通过 $K(x, y)$ 的傅里叶变换计算 $p(w)$, 可得 $w \sim N(0, 2\gamma I)$, 其中 I 表示单位矩阵。由上可得

$$K(x, y) = E_w [e^{-i w^T (x - y)}] = E_w [\cos(w^T (x - y))] = E_{w,b} [\sqrt{2} \cos(w^T x + b) \sqrt{2} \cos(w^T y + b)] \quad (3)$$

式中 w^T 表示权重矩阵。令 $Z_{w,b}(x) = \cos(w^T x + b)$, 有

$$K(x, y) = E_w [\langle Z_{w,b}(x), Z_{w,b}(y) \rangle] \quad (4)$$

则 $\langle Z_{w,b}(x), Z_{w,b}(y) \rangle$ 是高斯核函数的一个无偏估计。运用标准蒙特卡洛近似积分方法逼近高斯核, 构造如下随机特征映射

$$\Psi_{\text{RKS}}: x \mapsto \sqrt{\frac{2}{D}} [\cos(w_1^T x + b_1), \cos(w_2^T x + b_2), \dots, \cos(w_D^T x + b_D)]^T \quad (5)$$

可以看出, 随机映射由随机投影与余弦变换组成。其中 $w_i \in \mathbf{R}^d$ 是一个高斯随机变量, 每个元素均独立同分布采样于 $N(0, 2\gamma I)$, b_i 为均匀随机向量, 服从均匀分布 $U[-\pi, \pi]$, $i = 1, 2, \dots, D$ 。

2 随机特征空间的近似误差分析

2.1 核近似误差分析

本文利用傅里叶特征变换将原始特征空间转换为相对低维随机特征空间。下面分析傅里叶特征变换与原核函数的误差上界, 命题 1 给出了随机傅里叶特征变换的近似一致收敛界^[15]。

命题 1 设 K 是一个连续可微核函数, 其满足式(4)的条件, 并具有相关的测度 P 。此外假设 $\sigma_p^2 = E_{w \sim p} [\|w\|^2] < \infty$, 并且 X 是紧的, d 表示其维度。设 R 表示包含 X 的欧氏球半径。然后对于式(5)中定义的 $\Psi \in \mathbf{R}^D$ 和任意的 $0 < r \leq 2R$ 和 $\epsilon > 0$, 设 \sup 表示 (x, y) 分布在欧氏球中的确界, 以下成立

$$P \left[\sup_{x, y \in X} |\Psi(x) \cdot \Psi(y) - K(x, y)| \geq \epsilon \right] \leq 2N(2R, r) \exp\left(-\frac{D\epsilon^2}{8}\right) + \frac{4r\sigma_p}{\epsilon} \quad (6)$$

根据命题 1 提出的随机傅里叶特征一致收敛界, 核近似的一致收敛界与随机傅里叶特征变换的特征空间维度之间的关系表达如下。

推论 1 对于任意 $\delta \in (0, 1)$, 式(7)以 $1 - \delta$ 的概率成立, 即

$$P \left[\sup_{x, y \in X} |\Psi(x) \cdot \Psi(y) - K(x, y)| \geq \sqrt{\frac{1}{\frac{\delta}{(48R\sigma_p)^2} + \frac{D}{4(d+2)}}} \right] \leq \delta \quad (7)$$

式中: d 表示样本原始特征数; R 为包含 x 的欧几里得半径; $\sigma_p^2 = E_{w \sim p}[\|w\|^2] < \infty$; D 为随机傅里叶变换的特征空间维度。

证明: 定义 $F(x, y) = \Psi(x) \cdot \Psi(y) - K(x, y)$, $z(x, y) = \Psi(x) \cdot \Psi(y)$, L_F 表示 $F(x, y)$ 的利普希茨常数, 对于所有的 i , 当 $|F((x, y)_i)| < \frac{\epsilon}{2}$ 以及 $L_F < \frac{\epsilon}{2r}$, 对于所有 $X \in R^d$, 有 $|F(x, y)| < \epsilon$ 。

根据马尔可夫不等式

$$P \left[L_F \geq \frac{\epsilon}{2r} \right] \leq \left(\frac{2r\sigma_p}{\epsilon} \right)^2 \quad (8)$$

在 ϵ -net 联合霍夫丁不等式

$$P \left[\bigcup_{i=1}^W |F((x, y)_i)| \geq \epsilon/8 \right] \leq 2W \exp(-D\epsilon^2/2) \quad (9)$$

式中 $W = (6R/r)^d$, 因为

$$N(R, r) \leq \left(\frac{3R}{r} \right)^d \quad (10)$$

由式(8~10)可得

$$P \left[\sup_{x, y \in X} |F(x, y)| \leq \epsilon \right] \geq 1 - \left[2 \left(\frac{3R}{r} \right)^d \exp(-D\epsilon^2/8) - \left(\frac{2\sigma_p}{\epsilon} \right)^2 \right] \quad (11)$$

上述不等式右式为 $1 - n_1 r^{-d} - n_2 r^2$ 。其中 $n_1 = 2(3R)^d \exp(-D\epsilon^2/8)$, $n_2 = \left(\frac{2\sigma_p}{\epsilon} \right)^2$ 。令 $r =$

$\left(\frac{n_1}{n_2} \right)^{\frac{1}{d+2}}$ 和 $\frac{24R\sigma_p}{\epsilon} \geq 1$, 得

$$P \left[\sup_{x, y \in X} |\Psi(x) \cdot \Psi(y) - K(x, y)| \geq \epsilon \right] \leq \left(\frac{48R\sigma_p}{\epsilon} \right)^2 \exp\left(-\frac{D\epsilon^2}{4(d+2)}\right) \quad (12)$$

令 $\left(\frac{48R\sigma_p}{\epsilon} \right)^2 \exp\left(-\frac{D\epsilon^2}{4(d+2)}\right) = \delta$ 。其中 $\delta \in (0, 1)$, 由泰勒中值定理, 将式(12)中部分变形, 由于 ϵ^2 趋近与 0, 则

$$\exp\left(-\frac{D\epsilon^2}{4(d+2)}\right) \approx 1 - \frac{D}{4(d+2)} \epsilon^2 \quad (13)$$

将式(13)代入式(12)得

$$\left(\frac{48R\sigma_p}{\epsilon} \right)^2 \left(1 - \frac{D}{4(d+2)} \epsilon^2 \right) = \delta \quad (14)$$

将 ϵ 等式代换得

$$\epsilon = \sqrt{\frac{1}{\frac{\delta}{(48R\sigma_p)^2} + \frac{D}{4(d+2)}}}$$

证毕。

由式(7)可以得到,由于 $\Psi(x) \cdot \Psi(y) = K'(x, y)$, $\sqrt{1/\left(\frac{\delta}{(48R\sigma_p)^2} + \frac{D}{4(d+2)}\right)}$ 是 $|\mathbf{K}'x, \mathbf{y} - \mathbf{K}(x, \mathbf{y})|$ 的概率置信上界,并且上界由随机傅里叶变换的特征空间维度 D 和 σ_p^2 影响,对于高斯核, $\mathbf{K}(x, \mathbf{y}) = \exp(-\gamma\|x - \mathbf{y}\|^2)$,有 $\sigma_p^2 = 2d\gamma$ 。

近似特征空间维度 D 越大, γ 越小,对应的 σ_p^2 也越小,近似的核函数与原始核函数误差上界越小。但是参数选择并不是仅要求误差最小,当维度 D 过大时,虽然两者的近似误差很小,但随机傅里叶特征空间中对于内核的计算效率会很低,失去了引入随机特征空间的加速效果和结构简明的优势。

当参数 γ 取值过小时,原高斯核空间数据分布间距趋于相等,分类效果极差。虽然核近似的核函数与原始核函数近似误差一直在减小,但得到的是两个极差的模型,失去了核近似的意义。为了得到一个效果较好,并且近似误差较小的核近似方法,所以 γ 会取在适中位置,而随机特征空间维度 D 要有合适的范围,在降低近似误差的同时要求模型效果良好。

2.2 核近似误差边界分析

本节通过核岭回归和支持向量机模型来进一步分析核近似误差在学习器中的影响因子和误差边界。

2.2.1 核岭回归

核岭回归(Kernel ridge regression, KRR)解决的对偶优化问题^[16]为

$$\max_{\alpha \in \mathbf{R}^m} \lambda \alpha^T \mathbf{a} + \mathbf{a} \mathbf{K} \alpha - 2 \alpha^T \mathbf{y} \quad (15)$$

式中 $\lambda = m\lambda_0 > 0$ 是岭参数。该问题允许封闭形式的解 $\alpha = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$ 。

命题2^[17] 设 h 表示核矩阵 $\mathbf{K} \in \mathbf{R}^{m \times m}$ 核岭回归函数模型, h' 表示使用近似核矩阵 $\mathbf{K}' \in \mathbf{R}^{m \times m}$ 时核岭回归函数模型。定义 $r > 0$,对于所有 $x \in X, \mathbf{K}(x, \mathbf{y}) \leq r$ 和 $\mathbf{K}'(x, \mathbf{y}) \leq r$ 。以下的不等式适用于所有的 $x \in X$,即

$$|h'(x) - h(x)| \leq \frac{rM}{\lambda_0^2 m} \|\mathbf{K}' - \mathbf{K}\|_2 \quad (16)$$

根据式(7)和命题2中的式(16),有以下推论成立。

推论2 定义 $r > 0$,对于所有 $x \in X, \mathbf{K}(x, \mathbf{y}) \leq r$ 和 $\mathbf{K}'(x, \mathbf{y}) \leq r$,以下的不等式适用于所有的 $x \in X$

$$|h'(x) - h(x)| \leq \frac{rM}{\lambda_0^2 m} \|\mathbf{K}' - \mathbf{K}\|_2 \leq \frac{rM}{\lambda_0^2 m} \sqrt{\frac{1}{\frac{\delta}{(48R\sigma_p)^2} + \frac{D}{4(d+2)}}} \quad (17)$$

证明

由于 $\|\cdot\|_2 \leq \|\cdot\|_F$,并且 $\sqrt{1/\left(\frac{\delta}{(48R\sigma_p)^2} + \frac{D}{4(d+2)}\right)}$ 是 $|\mathbf{K}' - \mathbf{K}|$ 的高概率置信上界,所以式(17)不等式成立。证毕。

根据以上推论可以得出,模型的近似误差上界仍由核函数近似误差所影响,其近似误差上界是由随机傅里叶变换的特征空间维度 D 和 σ_p 所影响。近似特征空间维度 D 越大,核近似对应的模型与原始核模型之间的误差上界越小; γ 越小, σ_p^2 越小,两者误差上界越小,但要在降低近似误差的同时,要求模型效果良好。

2.2.2 支持向量机

下面分析支持向量机在 M 个点的样本 S 上训练时返回的假设 h 与使用核 K' 在相同样本上训练时获得的假设 h' 之间的差异。命题3度量了SVM应用核矩阵得到的假设与应用近似核矩阵得到的假设之间的误差界。

命题3 设 h' 表示使用近似核矩阵 $K' \in \mathbf{R}^{M \times M}$ 时支持向量机返回的假设。定义 $p > 0$,对于所有 $x \in X, K(x, y) \leq p$ 和 $K'(x, y) \leq p$ 。那么,以下不等式适用于所有的 $x \in X$

$$|h'(x) - h(x)| \leq \sqrt{2} p^{\frac{3}{4}} C_0 \|K' - K\|_2^{\frac{1}{4}} + p^{\frac{1}{2}} C_0 \|K' - K\|_2^{\frac{1}{2}} \quad (18)$$

式中 C_0 为一常量,根据式(7)和命题3中的式(18),有以下推论成立。

推论3 定义 $p > 0$,对于所有 $x \in X, K(x, y) \leq p$ 和 $K'(x, y) \leq p$,那么,以下不等式适用于所有的 $x \in X$,有

$$|h'(x) - h(x)| \leq \sqrt{2} p^{\frac{3}{4}} C_0 \|K' - K\|_2^{\frac{1}{4}} + p^{\frac{1}{2}} C_0 \|K' - K\|_2^{\frac{1}{2}} \leq \sqrt{2} p^{\frac{3}{4}} C_0 \left(\sqrt{\frac{1}{\frac{\delta}{(48R\sigma_p)^2} + \frac{D}{4(d+2)}}}} \right)^{\frac{1}{4}} + p^{\frac{1}{2}} C_0 \left(\sqrt{\frac{1}{\frac{\delta}{(48R\sigma_p)^2} + \frac{D}{4(d+2)}}}} \right)^{\frac{1}{2}} \quad (19)$$

证明: 由于 $\|\cdot\|_2 \leq \|\cdot\|_F$,并且 $\sqrt{1/\left(\frac{\delta}{(48R\sigma_p)^2} + \frac{D}{4(d+2)}\right)}$ 是 $|K' - K|$ 的高概率置信上界,所以式

(19)不等式成立。证毕。

通过式(17)和式(19)可以得出,随机傅里叶特征变换后的核近似在核岭回归和支持向量机上的应用,模型中原始核矩阵和核近似矩阵对应的两种假设 h' 和 h 之间存在误差上界,并且上界由具体的随机特征空间维度 D 以及 σ_p^2 所影响,维度 D 越大,对应模型近似误差上界越紧,即相应误差越小。在具体应用时应注意控制计算代价,随机特征空间 D 的取值以及 σ_p^2 所对应 γ 的取值应尽量控制在合理范围内。

在理论验证核近似收敛一致性的情况下,要保持近似核与原始核之间的误差在可接受范围内,同时通过核近似算法在近似空间上选择一个性能较好的模型,并且利用核近似所得到的模型参数组合选出一个表现良好的原始核模型。

2.3 算法设计

本文提出算法RFFDG-RBF目的在于选出一个核近似模型的最优参数组合。由于核近似函数与原始核函数有误差上界并且很小,所以在近似空间上选择出最优模型参数并运用到原始核空间,在保证性能的前提下,原始核搜索参数所需时间将大幅降低,所以既可以保证使用核近似参数的原始核准确度较高,又比传统网格搜索法选择模型速度快很多。选定合适大小的参数 D ,输入需要搜索的 σ_p^2 值,其中 D 为随机特征空间维度,对于高斯核 $K(x, y) = \exp(-\gamma\|x - y\|^2)$,有 $\sigma_p^2 = 2d\gamma$,所以需要输入 γ 的值。下面给出了RFFDG-RBF算法的主要步骤。

RFFDG-RBF 算法

(1) 输入训练样本 $N = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ 、参数空间 $R(D_*, \gamma_*, C_*)$;

(2) 输出 最优参数组合 (D_o, γ_o) ;

(3) 原始核准确率 rbf-bestacc;

(4) for i in range(1, D) and each $(\mathbf{x}, \mathbf{y}) \in N$ do

$$\mathbf{w}_i \in \mathbf{R}^d: \mathbf{w}_i \sim \Lambda(0, 2\gamma \mathbf{I}) \mathbf{b}_i \in T[-\pi, \pi]$$

计算 (\mathbf{x}, \mathbf{y}) 在随机特征空间中的映射 $(\Psi(\mathbf{x}_i), \mathbf{y})$

$$\Psi(\mathbf{x}_i) = \sqrt{\frac{2}{D}} [\cos(\mathbf{w}_1^\top \mathbf{x} + b_1), \cos(\mathbf{w}_2^\top \mathbf{x} + b_2), \dots, \cos(\mathbf{w}_D^\top \mathbf{x} + b_D)]^\top;$$

(5) end for

(6) 在训练集 $N = \{(\Psi(\mathbf{x}_i), \mathbf{y}_i)\}_{i=1}^n$ 上运用线性核模型训练

(7) for each $(D, \gamma)_i \in (D, \gamma)_*$;

(8) D in range $(10^1, 10^4, \text{step}=10)$;

(9) γ in range $(10^{-2}, 10^1, \text{step}=10)$ do;

(10) result = cross_val_score $(\Psi(\mathbf{x}_i), \mathbf{y}_i, \text{cv}=10)$;

(11) if result > rff-btacc;

(12) rff-btacc = result;

(13) $(D_o, \gamma_o) = (D, \gamma)_i$;

(14) Return (D_o, γ_o) ;

(15) end for;

(16) for C in range $(10^{-2}, 10^2, \text{step}=10)$ do;

(17) rbf-acc = accuracy (D_o, γ_o, C_i) ;

(18) if rbf-acc > rbf-bestacc;

(19) rbf-bestacc = rbf-acc;

(20) Return rbf-bestacc;

(21) end for.

下面对算法时间复杂度进行分析。记 n 为样本规模, d 为样本特征数, D 为随机特征空间维度, R_γ 、 R_C 为参数 γ 和 C 对应的搜索长度。进行随机特征转换的时间复杂度为 $O(ndDR_\gamma)$, 在 k -折交叉验证中, 核近似算法选择参数 γ 和 C 的时间复杂度为 $O(knDR_\gamma)$, 最后选择最佳参数 C 的时间复杂度为 $O(knR_c)$, 故算法总时间复杂度为 $O(ndDR_\gamma + knDR_\gamma + knR_c)$ 。

3 实验与结果

本节首先对上文所提出的误差上界受参数影响的理论在回归和分类数据集上分别进行验证, 然后利用 2.3 节提出的算法分别在核岭回归模型和 SVM 模型找到最优参数组合, 将最优参数组合运用到原始核上, 从而达到通过近似核选出一个较优原始核模型的效果, 并且参数选择的计算效率更高, 最后将几种核方法进行对比。实验在 Intel(R) Core(TM) i7-4700 CPU 3.60 GHz, 内存 16 GB DDR3 的机器上运行, 采用公开数据集^[18]如表 1 所示, 对于每个数据集根据其自身特征数, 选择相应合适的核近似特征空间维度 D 。

表 1 实验所实用数据集
Table 1 Datasets used in the experiment

Dataset	Dimension	Train	Test	Task
ijcnn1	22	49 990	91 701	classification
a8a	123	22 696	9 865	classification
letter	16	15 000	5 000	classification
w7a	300	24 692	25 057	classification
shuttle	9	43 500	14 500	classification
phishing	68	8 844	2 211	classification
mg	6	1 039	346	regression
abalone	8	3 133	1 044	regression
housing	13	380	126	regression
space_ga	6	2 263	754	regression
cpusmall	12	6 144	2 048	regression
mpg	7	294	98	regression

3.1 随机特征空间维度 D 对近似误差的影响

本节将验证核近似的误差上界受到近似特征空间维度 D 影响。随机傅里叶变换后的核近似与原始核的均方误差 MSE 为

$$MSE = \frac{1}{n} \sum_{i=1}^n (\mathcal{K}(x, x') - \mathcal{K}(x, x'))^2 \tag{20}$$

图 1 和图 2 分别为按照表 1 顺序在回归和分类数据集上 MSE 随特征空间维度 D 变化的情况。由于本节主要讨论 D 的变化与误差的关系,所以参数 γ 统一设定为 1.0。从图中可以看出,随着维度 D 增大,均方误差整体趋势在减小,当 $D \in (200, 800)$ 之间时,误差上界已经收敛在较小范围内。

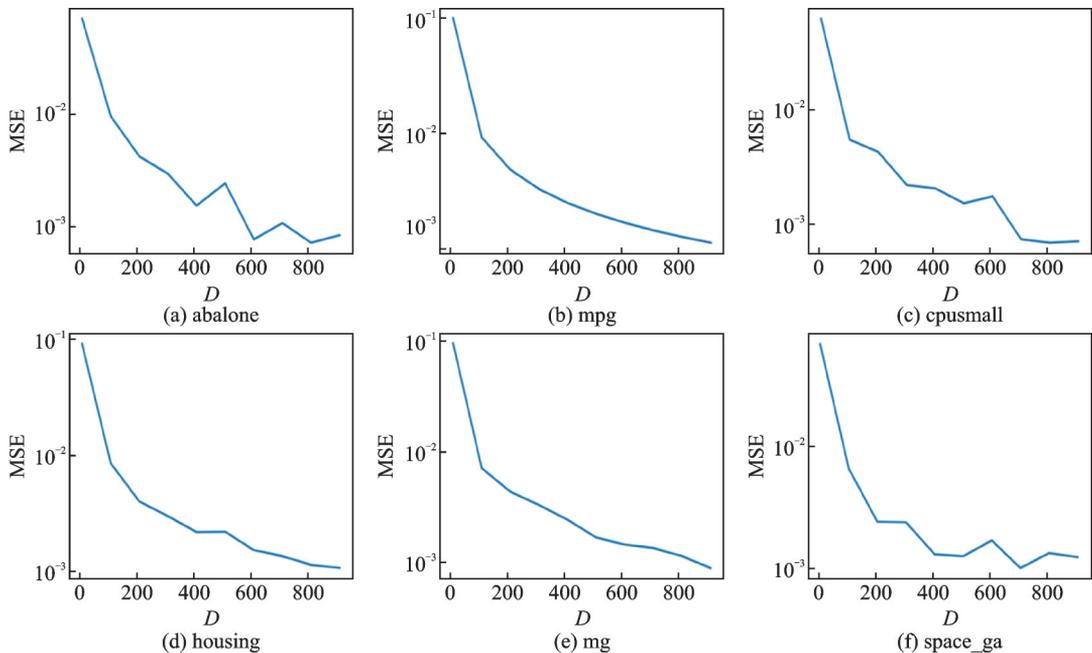
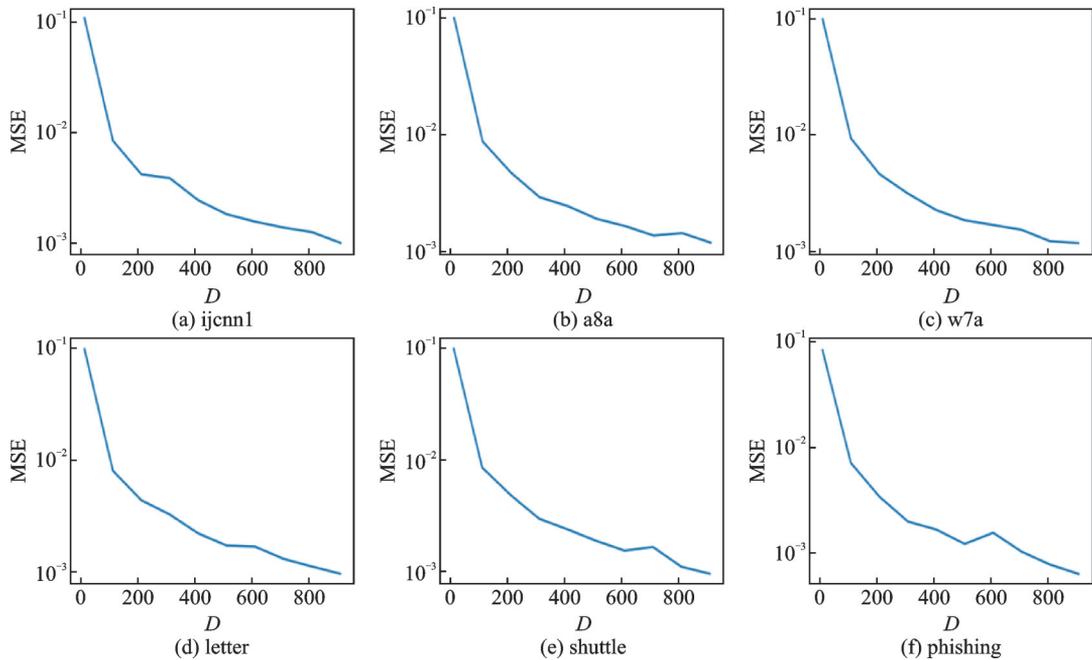


图 1 各回归数据集上维度 D 与 MSE 的关系

Fig.1 Relationship between dimension D and MSE under different regression datasets

图2 各分类数据集上维度 D 与 MSE 的关系Fig.2 Relationship between dimension D and MSE under different classification datasets

3.2 各种方法的模型选择效果对比

实验选取本文提出的 RFFDG-RBF 随机傅里叶核近似算法, RBF 为原始高斯核, Nyström 为 Nyström 近似方法, Mid-KSRC 为中值法^[19]。Nyström 近似方法是一种基于列的子集求解较大核矩阵的低秩逼近的有效方法。Nyström 近似的质量在很大程度上取决于所使用列的子集, 用 W 表示由所选 c 列与相应的 c 行的相交组成的 $c \times c$ 矩阵。Nyström 近似形式如下

$$\tilde{K} = CW_k^+C^T \approx K \quad (21)$$

式中: W_k 为 W 的最优秩 k 近似; W_k^+ 为 W_k 的广义逆。

KSRC 可以通过使用 RBF 内核来解决分类问题, 但必须确定相应的参数。在文献[19]的实验中, 用中值法 Mid-KSRC 来设置核参数。实验采用高斯核作为原始核函数实现核岭回归和 SVM 问题, 在 6 个回归数据集和 6 个分类数据集中用 RFFDG-RBF 算法进行随机傅里叶特征变换, 找出最优的参数组合 (γ, D) , 将最优参数组合运用到原始 RBF 核, 利用核近似选择出的核参数 γ , 再对 SVM 中参数 C 以及核岭回归参数 λ 进行搜索。利用核近似选择出的参数 γ , 达到和原函数相近甚至相等的准确率。表 2 给出了各数据集上核近似选择最优模型对应的测试均方误差和训练时间, 表 3 给出了各数据集上核近似选择最优模型对应的测试准确率 Acc 和训练时间, 所有随机算法在每个数据集上重复 10 次, 并且记录平均值。

从表 2 可以得出: RFFDG-RBF 算法最优模型的测试均方误差在 mg、abalone、space_ga、cpusmall、和 mpg 数据集上比中值法 Mid-KSRC 和 Nyström 方法最优模型的均方误差小, 并且与原始 RBF 核对应的最优模型均方误差相近或相等。在 housing 数据集上均方误差较大是由于数据集本身样本量较小, 再通过核近似会一定程度上增大均方误差的值, 当在充足样本量的数据集中, 误差就会相对较小。并且

表2 各种方法的最优模型对应的测试MSE以及训练时间

Table 2 Test MSE and training time corresponding to the optimal model of various methods

Dataset	RFFDG-RBF		RBF		Nyström		Mid-KSRC	
	Time/s	MSE	Time/s	MSE	Time/s	MSE	Time/s	MSE
mg	3.925 5	0.124 45	4.391 1	0.124 47	16.435 5	0.138 99	4.245 8	0.129 8
abalone	19.197 2	2.089 55	50.154 7	2.089 54	39.628 4	2.089 56	23.904 0	2.135 4
housing	0.884 4	3.957 15	0.933 1	3.790 92	9.369 9	5.261 12	0.367 8	3.962 3
space_ga	12.675 8	0.108 80	25.607 6	0.108 80	29.395 3	0.108 81	15.125 5	0.122 3
cpusmall	61.328 8	2.968 58	201.899 6	3.005 52	93.486 6	3.036 40	99.086 4	4.196 5
mpg	0.854 9	2.707 37	0.600 4	2.855 80	7.819 1	3.424 80	0.176 6	2.539 7

表3 各种方法的最优模型对应的测试准确率以及训练时间

Table 3 Test accuracy and training time corresponding to the optimal model of various methods

Dataset	RFFDG-RBF		RBF		Nyström		Mid-KSRC	
	Time/s	Acc/%	Time/s	Acc/%	Time/s	Acc/%	Time/s	Acc/%
ijcnn1	285.526	96.68	1 537.922	96.69	427.597	98.24	395.680 0	92.33
a8a	547.090	85.26	8 506.023	85.26	632.908	82.90	765.428 0	85.20
letter	209.395	96.74	656.080	96.74	675.339	96.74	98.014 5	89.30
w7a	5 144.605	97.23	9 938.427	97.35	5 106.867	97.05	294.816 0	97.05
shuttle	219.801	99.86	1 490.305	99.86	862.808	99.77	591.522 0	98.01
phishing	31.422	96.88	326.978	97.06	58.124	95.86	98.972 0	94.90

RFFDG-RBF核近似算法在mg、abalone、space_ga和cpusmall数据集上的训练所需时间都比另外3种方法少,在样本规模较大的数据abalone、space_ga和cpusmall中的加速效果更加明显。通过表2实验结果说明5个数据集上的RFFDG-RBF算法都得到了与原始高斯核相近的均方误差,从而证明本文提出的方法可以在提升模型选择计算效率的基础上,并且运用核近似方法选择出表现较好的原始RBF核模型。

从表3可以得出:RFFDG-RBF算法的最优模型的测试准确率在a8a、letter、shuttle、w7a和phishing数据集上比Nyström方法的准确率高,在全部数据集上比中值法Mid-KSRC的准确率高,并且与原始高斯核准确度相近或相等。在ijcnn1数据集上RFFDG-RBF算法比Nyström方法的准确率略低,但RFFDG-RBF算法的测试准确率仍然与原始RBF核准确率相近,并且计算效率大幅提升。另外,在w7a数据集上由于数据集本身的特征数量较大,导致RFFDG-RBF算法选择特征空间维度 D 也相应增大,但仍比Nyström方法准确度高,并且测试准确率与原始RBF核准确率相近。虽然在w7a数据集RFFDG-RBF算法需要时间比中值法Mid-KSRC要多,但是RFFDG-RBF算法比中值法Mid-KSRC的测试准确率高。同时在ijcnn1、a8a、letter、shuttle和phishing数据集上RFFDG-RBF算法的模型选择计算效率相比原始RBF核以及Nyström近似方法和中值法Mid-KSRC的计算效率显著提高,在w7a数据集上RFFDG-RBF算法模型搜索时间仍然比原始RBF核模型搜索时间短。

结合表2和表3,本文提出的算法在高斯核模型选择上大幅提高计算效率的同时,保证核近似在回归数据集上的均方误差与原始RBF核相近,在分类数据集上的准确率与原始RBF核相近,对原始高斯核模型的选择进行了优化。

3.3 各种方法在样本规模不同时的训练时间以及测试结果对比

随机傅里叶特征变换将数据映射到一个相对低维的特征空间中,其维度是 D 。因为映射后的核近似算法在计算上可以避免其本身数据量 N ,而原始高斯核计算仍需考虑数据量 N ,所以本文提出的算法相对于原始RBF核,在规模越大的数据集上计算效率的提升幅度越大。本节分别在回归数据集 cpusmall 和分类数据集 phishing 上进行实验,将RFFDG-RBF算法与原始RBF核,中值法 Mid-KSRC 以及 Nyström 方法的模型选择结果进行对比,如图3、4所示。随着样本规模的增大,得出其训练时间增长速率和最优模型选择的测试均方误差和测试准确率的变化规律。

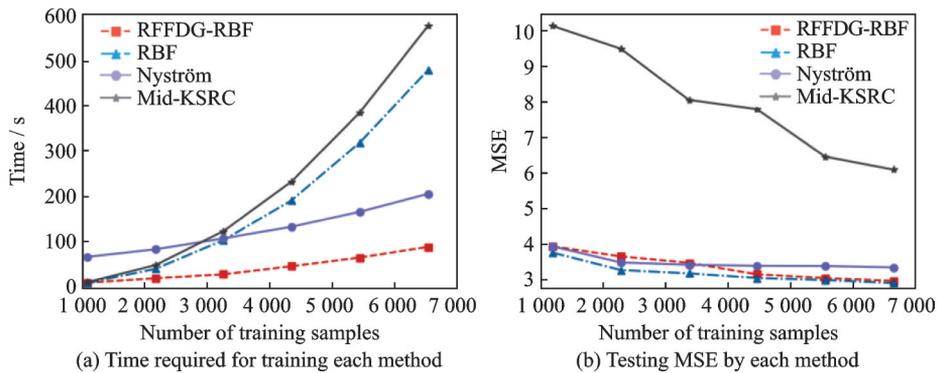


图3 4种模型选择算法的训练时间和测试MSE

Fig.3 Training time and test MSE of three model selection algorithms

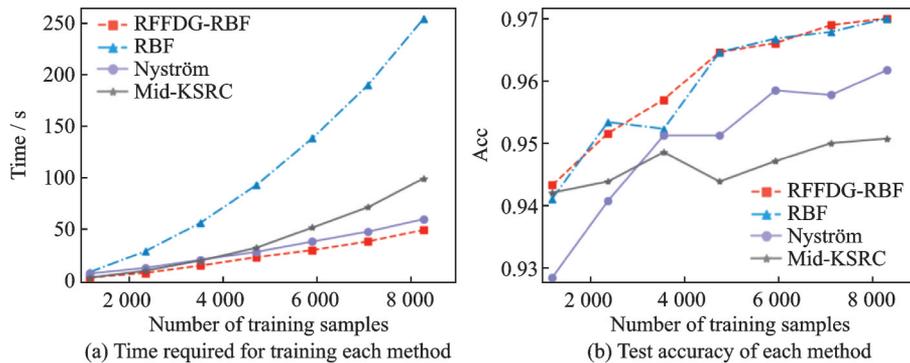


图4 4种模型选择算法的训练时间和测试准确率

Fig.4 Training time and test accuracy of three model selection algorithms

结合图3和图4可以看出,随着样本规模的增大,RFFDG-RBF算法在模型选择所需时间上相对于原始RBF核的模型选择时间大幅降低,同时比Nyström方法和中值法Mid-KSRC的模型选择所需的时间少;并且随着数据规模越大,四者训练所需时间差距越大。在cpusmall数据集上RFFDG-RBF算法比原始高斯核模型选择训练所需时间缩短400s,仅为原始RBF核训练所需时间的五分之一,比Nyström方法训练所需时间缩短100s,是Nyström方法所需时间的二分之一。在phishing数据集上RFFDG-RBF算法比原始高斯核模型选择训练所需时间缩短200s,比中值法Mid-KSRC所需时间缩短60s,比Nyström方法所需时间缩短20s。

在cpusmall和phishing两个数据集上,原始RBF核计算所需时间增长速率都远大于RFFDG-RBF算法对应的所需时间增长速率,并且随着样本规模的增大,两者的时间成本差距越来越大,原始RBF核

时间成本增长速率越来越快,而 RFFDG-RBF 算法计算时间成本增长速率则较为缓慢。实验结果表明:RFFDG-RBF 算法在不同规模数据集上的模型选择中都实现了计算效率的大幅提升,并且数据集规模越大,计算效率提升越明显。最优模型测试结果方面,随着样本规模的增大,在 cpusmall 数据集上 RFFDG-RBF 算法对应的最优模型测试 MSE 逐渐减小,并且最后与原始 RBF 核对应的核岭回归最优模型的测试 MSE 很接近。在 phishing 数据集上,RFFDG-RBF 算法的最优模型测试准确率逐渐增加并且高于 Nyström 方法和中值法 Mid-KSRC 最优模型对应的测试准确率。在样本规模最大时,RFFDG-RBF 算法的最优模型测试准确率与原始 RBF 核支持向量机最优模型的测试准确率相近,仅相差 0.002%。上述实验表明本文提出的算法相对于原始 RBF 核模型,在规模越大的数据集上计算效率提升越大,提升计算效率的同时仍可以保证核近似选择的最优模型拥有良好的准确率。

4 结束语

本文首先提出基于随机傅里叶特征变换生成的核函数与原始核函数的近似误差存在上界,并且证明误差上界会受到随机傅里叶特征空间维度 D 以及核函数参数 γ 的影响,从理论上论证核近似的收敛一致性。然后,构建一个低维显式的特征空间,在这个特征空间中进行模型选择相对透明,增强了模型的可解释性。另一方面,选择出核近似相应的最优模型相对于原始核模型选择的网格搜索法,不仅实现了计算效率的提高,并且与原始高斯核的测试准确率相近。在规模越大的数据集上,计算效率的提升幅度越大,与其他核近似方法进行对比,仍有较好的准确率以及较少的训练时间,所以本文提出的算法提供了在模型选择中提升计算效率和保持准确率之间的优化。下一步工作将进一步探索在其他核函数中运用核近似的合理性,并进行深入的理论依据研究。

参考文献:

- [1] THOMAS H, BERNHARD S, ALEXANDER J. Kernel methods in machine learning[J]. *The Annals of Statistics*, 2008, 36(3): 1171-1220.
- [2] ZHANG T. Solving large scale linear prediction problem using stochastic gradient descent algorithms[C]//*Proceedings of the 25th International Conference on Machine Learning*, New York: ACM, 2004: 919-926.
- [3] JOACHIMS T. Training linear SVMs in linear time[C]//*Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. New York: ACM, 2006: 217-226.
- [4] HSIEH C J, CHANG K W, LIN C J, et al. A dual coordinate descent method for large-scale linear SVM[C]//*Proceedings of the 25th International Conference on Machine Learning*. New York: ACM, 2008: 408-415.
- [5] FAN R E, CHANG K W, HSIEH C J, et al. LIBLINEAR: A library for large linear classification[J]. *Mach Learn Res*, 2008, 9(2): 1871-1874.
- [6] ASSIA O, ABDELAZIZ O, MOKHTAR K. Low complexity method for DOA estimation based on nyström method[J]. *International Journal on Communications Antenna and Propagation*, 2017, 1(6): 239-245.
- [7] YANG C J, DURAISWAMI R, DAVIS L. Efficient kernel machines using the improved fast Gauss transform[C]//*Proceedings of Advances in Neural Information Processing Systems*. Cambridge, MA: MIT Press, 2004: 1561-1568.
- [8] CAO H, NAITO T, NINOMIYA Y. Approximate RBF kernel SVM and its applications in pedestrian classification[C]//*Proceedings of the 1st International Workshop on Machine Learning for Vision-based Motion Analysis*. Berlin: Springer, 2008: 120-129.
- [9] RAHIMI A, RECHT B. Random features for large-scale kernel machines[C]//*Proceedings of the 21st Annual Conference on Neural Information Processing Systems*, Vancouver. Red Hook: Curran Associates, 2008: 1177-1184.
- [10] FENG C, LIAO S Z. Large-scale kernel methods via random hypothesis spaces[J]. *Frontiers of Computer Science and Technology*, 2018, 12(5): 785-793.
- [11] LIU Y, JIANG L S, LIAO S Z. Approximate Gaussian kernel for large-scale SVM[J]. *Journal of Computer Research and*

- Development, 2014, 51(10): 2171-2177.
- [12] LIU Y, LIAO S Z. Framework core selection method based on an explicit description of the integral operator space[J]. *Scientia Sinica Informationis*, 2016, 46(2): 757-765.
- [13] FENG C, LI Z D, LIAO S Z. Efficient algorithm for large-scale support vector machine[J]. *Computer Science*, 2015, 42(9): 195-198.
- [14] RUDIM W. *Fourier analysis on groups*[M]. San Francisco: John Wiley & Sons, 2011: 1-285.
- [15] BHARATH K.S, ZOLTÁN S. Optimal rates for random Fourier features[J]. *NIPS*, 2015, 28(1): 1144-1152.
- [16] MEHRYAR M, AFSHIN R, AMEET T. *Foundations of machine learning*[M]. Cambridge, Massachusetts: The MIT Press, 2018: 22-812.
- [17] CORTES C, MOHRI M, TALWALKAR A. On the impact of kernel approximation on learning accuracy[C]//*Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*. Brookline, MA: Microtome Publishing, 2010: 113-120.
- [18] CHANG C C, LIN C J. LIBSVM : A library for support vector machines[J]. *ACM Transactions on Intelligent Systems and Technology*, 2011, 2(3): 27: 1-27: 27.
- [19] ZHANG L, ZHOU W D, CHANG P C, et al. Kernel sparse representation-based classifier[J]. *IEEE Trans on Signal Process*, 2012, 60 (4): 1684-1695.

作者简介:



张凯(1997-),男,硕士研究生,研究方向:机器学习、数据挖掘, E-mail: 573437827@qq.com。



门昌骞(1982-),男,博士,讲师,研究方向:支持向量机、机器学习理论和核方法。



王文剑(1968-),通信作者,女,博士,教授,研究方向:机器学习、计算智能、图像处理。

(编辑:刘彦东)