

文章编号:1001-9081(2009)02-0416-03

## 基于 KDA 和 SVM 的文档分类算法

王自强, 钱 旭

(中国矿业大学(北京)机电与信息工程学院, 北京 100083)

(wzqbox@yahoo.cn)

**摘 要:**为了高效地解决 Web 文档分类问题,提出了一种基于核鉴别分析方法 KDA 和 SVM 的文档分类算法。该算法首先利用 KDA 对训练集中的高维 Web 文档空间进行降维,然后在降维后的低维特征空间中利用乘性更新规则优化的 SVM 进行分类预测。采用了文档分类领域两个著名的数据集 Reuters-21578 和 20-Newsgroup 进行实验,实验结果表明该算法不仅获得了更高的分类准确率,而且具有较少的运行时间。

**关键词:**文档分类;核鉴别分析;支持向量机;数据挖掘

**中图分类号:** TP181 **文献标志码:** A

## Document classification algorithm based on KDA and SVM

WANG Zi-qiang, QIAN Xu

(College of Mechanical Electronic and Information Engineering, China University of Mining and Technology(Beijing), Beijing 100083, China)

**Abstract:** To efficiently solve Web document classification problem, a novel document classification algorithm based on kernel discriminant analysis (KDA) and SVM was proposed. The proposed algorithm firstly reduced the high dimensional Web document space in the training sets to the lower dimensional space with KDA algorithm, then the classification and predication in the lower dimensional feature space were implemented with the multiplicative update-based optimal SVM. The experimental evaluations were performed on the Reuters-21578 and 20-Newsgroup which are two well-known data sets in the field of document classification. Experimental results show that the proposed algorithm not only achieves higher classification accuracy, but also has lower running time.

**Key words:** document classification; Kernel Discriminant Analysis (KDA); Support Vector Machine (SVM); data mining

### 0 引言

为了高效地管理和组织海量信息<sup>[1]</sup>,以帮助用户快速准确地查询到所需要的相关文档,基于内容的 Web 文档自动分类技术日益成为数据挖掘领域备受关注的热点研究问题。

Web 文档分类的主要任务是在预先给定的类别标记集合中,根据文档内容判定它的类别。有关学者已提出了多种文档分类技术,如:决策树、神经网络、k-近邻、贝叶斯方法和支持向量机(SVM)<sup>[2]</sup>等。虽然这些方法在进行文档分类时取得了一定的效果,但是他们大多都是直接在原始文档空间分类,而 Web 文档具有维数高、样本稀疏和特征不太明显的特点。考虑到高维文档空间引起的“维数灾难”问题,我们应当首先把高维文档投影到低维特征子空间,然后在降维后的低维特征空间中利用分类器进行分类预测。基于上述考虑,本文提出了基于核鉴别分析方法(Kernel Discriminant Analysis, KDA)和 SVM 的 Web 文档分类算法,该算法首先使用核鉴别分析方法把高维文档投影到低维特征空间,然后在降维后的低维特征空间利用乘性更新规则优化<sup>[3]</sup>的 SVM 进行分类。实验结果表明该方法有效地提高了文档分类准确率,而且减少了分类器的训练时间。

### 1 基于 KDA 和 SVM 的文档分类

#### 1.1 基于 KDA 的文档低维空间获取

设给定  $n$  个文档样本集  $X = \{x_1, x_2, \dots, x_n\}$  属于  $c$  个不同的类别,其中第  $i(1 \leq i \leq c)$  个类别中有  $n_i$  个样本。核鉴

别分析 KDA<sup>[4-5]</sup>的主要思想是通过构造非线性映射  $\varphi$  把输入空间  $X$  映射到特征空间  $F$ :

$$\varphi: x_i \in X \rightarrow \varphi(x_i) \in F \quad (1)$$

以便在特征空间  $F$  中实现线性鉴别分析 LDA<sup>[6]</sup>的优化目标:最小化类内距离的同时最大化类间的距离,以达到最大限度的类鉴别分析。其形式化描述为:

$$J(w) = \max_w \frac{w^T S_B^o w}{w^T S_W^o w} \quad (2)$$

其中,  $S_B^o$  和  $S_W^o$  分别表示在特征空间  $F$  中的类间散度和类内散度,其定义如下:

$$S_B^o = \sum_{i=1}^c n_i (\mu_i^o - \mu^o) (\mu_i^o - \mu^o)^T \quad (3)$$

$$S_W^o = \sum_{i=1}^c \sum_{j=1}^{n_i} (\varphi(x_j^i) - \mu_i^o) (\varphi(x_j^i) - \mu_i^o)^T \quad (4)$$

$$\mu_i^o = \frac{1}{n_i} \sum_{j=1}^{n_i} \varphi(x_j^i) \quad (5)$$

$$\mu^o = \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \quad (6)$$

其中,  $\mu_i^o$  和  $\mu^o$  分别表示在特征空间  $F$  中的第  $i$  个类别中的样本均值和所有样本的均值。

根据再生核空间理论,任何解向量  $w \in F$  都是由空间  $F$  中的所有样本张成的。可得:

$$w = \sum_{i=1}^n \alpha_i \varphi(x_i) \quad (7)$$

将式(7)代入式(2)可得:

收稿日期:2008-08-12;修回日期:2008-10-15。 基金项目:教育部科学技术研究重点资助项目(107021)。

作者简介:王自强(1973-),男,河南郑州人,博士研究生,主要研究方向:数据挖掘、模式识别; 钱旭(1962-),男,江苏无锡人,教授,博士生导师,主要研究方向:数据挖掘、信息融合。

$$J(\alpha) = \max_{\alpha} \frac{\alpha^T K_B \alpha}{\alpha^T K_W \alpha} \quad (8)$$

其中:

$$K_B = \sum_{i=1}^c n_i (m_i - m) (m_i - m)^T \quad (9)$$

$$K_W = \sum_{i=1}^c \sum_{j=1}^{n_i} (\eta_j - m_i) (\eta_j - m_i)^T \quad (10)$$

$$m_i = \left( \frac{1}{n_i} \sum_{j=1}^{n_i} k(x_1, x_j), \frac{1}{n_i} \sum_{j=1}^{n_i} k(x_2, x_j), \dots, \frac{1}{n_i} \sum_{j=1}^{n_i} k(x_n, x_j) \right)^T \quad (11)$$

$$m = \left( \frac{1}{n} \sum_{j=1}^n k(x_1, x_j), \frac{1}{n} \sum_{j=1}^n k(x_2, x_j), \dots, \frac{1}{n} \sum_{j=1}^n k(x_n, x_j) \right)^T \quad (12)$$

$$\eta_j = \left( k(x_1, x_j), k(x_2, x_j), \dots, k(x_n, x_j) \right)^T \quad (13)$$

$$k(x_i, x_j) = [\varphi(x_i), \varphi(x_j)] \quad (14)$$

因此,在  $K_W$  非奇异的情况下,式(8)的解向量  $\alpha$  是由  $K_W^{-1} K_B$  的前  $l(l = \min(c-1, n))$  个最大特征向量组成的。

于是,对于一个新的测试样本点  $x \in X$ ,它在特征空间  $F$  中特征向量  $w$  上的投影可表示为:

$$w \cdot \varphi(x) = \sum_{i=1}^n \alpha_i k(x_i, x) \quad (15)$$

对于文档分类问题,由于训练样本的维数远远大于样本数,因此矩阵  $K_W$  是奇异的,为了克服这一不足,我们用  $(K_W + \varepsilon E)$  来代替  $K_W$  来解决这一问题,其中  $\varepsilon$  是一个非常小的正常数,本文中我们取  $\varepsilon = 10^{-5}$ ,  $E$  是单位矩阵。

### 1.2 SVM 分类器

在利用 KDA 完成高维文档的降维处理后,为了高效地对降维后的低维特征空间进行分类,我们采用了具有良好泛化能力的支持向量机(SVM)<sup>[2]</sup>作为分类器。其实现方法如下:

设训练样本集  $\{(x_i, y_i)\}_{i=1}^l$ , 样本  $x_i \in R^n$ ,  $y_i \in \{-1, +1\}$  为类标签。SVM 通过求解式(16)找到一个具有最大间隔的超平面:

$$\min_{w, b, \xi} \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i \quad (16)$$

约束条件:

$$y_i[(w \cdot x + b) + \xi_i - 1] \geq 0, \xi_i \geq 0 \quad (17)$$

其中,  $C$  为一个用于控制误差的惩罚常数,  $\xi_i$  为非负松弛变量。

利用 Lagrange 乘子法,可以把式(16)转化为其对偶形式:

$$\max W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \quad (18)$$

约束条件为:

$$\sum_{i=1}^l \alpha_i y_i = 0, \alpha_i \in [0, C], i = 1, \dots, l \quad (19)$$

于是对于未知类别的数据点  $x$ , 可采用如下线形判决函数决定其所属类别:

$$f(x) = \text{sgn} \left( \sum_{i=1}^l \alpha_i y_i (x_i \cdot x) + b \right) \quad (20)$$

对于非线性 SVM 的情况,通过利用非线性映射  $\varphi$  把输入空间映射到高维特征空间,于是核函数  $K(x_i, x_j) = (\varphi(x_i) \cdot \varphi(x_j))$  可在特征空间计算,而无需知道映射  $\varphi$  的具体形式。

用核函数代替线形 SVM 中的点积形式,于是式(18)的对偶形式可变为:

$$\max W(\alpha) = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i,j=1}^l \alpha_i \alpha_j y_i y_j K(x_i, x_j) \quad (21)$$

约束条件为:

$$\sum_{i=1}^l \alpha_i y_i = 0, \alpha_i \in [0, C], i = 1, \dots, l \quad (22)$$

于是,非线性 SVM 的判决函数变为:

$$f(x) = \text{sgn} \left( \sum_{i=1}^l \alpha_i y_i K(x_i, x) + b \right) \quad (23)$$

由于 SVM 训练问题可归结为求解非负二次规划问题,常用的 SVM 训练方法存在如下不足:训练代价太大,在处理大规模高维数据时存在着收敛速度慢、分类准确率低。为了有效地提高 SVM 训练算法的性能,本文采用乘性更新规则方法<sup>[3]</sup>来优化求解 SVM 中的二次规划问题。该方法的优点是:无需参数设置、不需要选择工作子集,且具有并行更新变量的能力。实验表明该方法具有快速求解二次规划的能力,并具有较高的分类准确率。下面给出基于乘性更新方法的 SVM 训练方法。

由于式(18)的非负二次规划问题可归纳为如下标准形式:

$$\min F(v) = \frac{1}{2} v^T A v + b^T v \quad (24)$$

约束条件为:  $\forall i, v_i \geq 0$ 。矩阵  $A$  是正定和半对称的。由于对  $F(v)$  的优化为凸二次优化,加上非负  $v_i$  的约束,因此只能采用迭代的方法而不能采用分析的方法来获得全局最小值。

该迭代方法是通过矩阵  $A$  的正负元素来实现,首先在式(24)中,定义如下两个非负矩阵  $A^+$  和  $A^-$ :

$$A_{ij}^+ = \begin{cases} A_{ij}, & A_{ij} > 0 \\ 0, & \text{其他} \end{cases}, A_{ij}^- = \begin{cases} |A_{ij}|, & A_{ij} < 0 \\ 0, & \text{其他} \end{cases} \quad (25)$$

于是有:  $A = A^+ - A^-$ 。根据以上定义的非负矩阵,迭代乘性更新规则定义如下:

$$v_i \leftarrow v_i \left[ \frac{-b_i + \sqrt{b_i^2 + 4(A^+ v)_i (A^- v)_i}}{2(A^+ v)_i} \right] \quad (26)$$

上述乘性更新规则的优点是实现简单,避免了常规 SVM 训练中的大量矩阵运算,且整个迭代更新过程不会违反非负性条件约束。而文献[3]中的定理1已证明:利用式(26)中的乘性更新规则,式(24)中的目标函数可以单调递减到全局最小值。因而保证了式(26)的更新迭代规则可以收敛到全局最小值。

对于 SVM 训练中的式(21)优化问题的具体实现,令  $A_{ij} = y_i y_j K(x_i, x_j)$ ,  $b_i = -1$ 。于是,优化求解式(20)的乘性更新规则如下:

$$\alpha_i \leftarrow \alpha_i \left[ \frac{1 + \sqrt{1 + 4(A^+ \alpha)_i (A^- \alpha)_i}}{2(A^+ \alpha)_i} \right] \quad (27)$$

由此可见,利用式(27)来求解 SVM 中的二次规划问题,避免了常规 SVM 训练中非常耗时的大量矩阵运算,使得训练速度得到很大提高,后面的实验结果有力地证明了这一点。

### 1.3 基于 KDA 和 SVM 的文档分类算法

设给定  $n$  个文档样本集  $X = \{x_1, x_2, \dots, x_n\}$  属于  $c$  个不同的类别。由于向量空间模型(Vector Space Model, VSM)是文档表示的主要形式,因此本文采用 VSM 中常用的 TF-IDF 项权重向量来表示每个文档  $x_i$ 。基于 KDA 和 SVM 的文档分类算法的主要思想为:首先利用核鉴别分析 KDA 把高维文档投影到低维特征空间,然后利用乘性更新规则优化的 SVM 分类器对降维的文档数据进行分类。算法的具体实现过程如下:

步骤1 利用向量空间模型 VSM 中的 TF-IDF 项权重向量来表示每个文档  $x_i$ 。

步骤2 通过求解式(8)获得最佳投影方向。

步骤3 对于新的测试文档样本,利用式(15)计算出它在特征空间  $F$  中特征向量上的投影。

步骤4 利用投影后的低维文档特征向量按照式(27)的乘性更新规则来优化训练式(21)中的 SVM 目标函数。

步骤 5 利用式(23)中的 SVM 判决函数来确定待测试文档所属的类别。

## 2 实验结果

为了测试本文提出的基于 KDA 和 SVM 的文档分类算法(简称 KDA-SVM)分类性能,我们将 KDA-SVM 与其他 4 种常用的文档分类算法:最近邻居法(K-Nearest Neighbor, KNN)、线性判别分析法(Linear Discriminant Analysis, LDA)、Bayes 方法和经典 SVM 算法进行了比较。算法中的核函数采用常用的 RBF 核函数,模型参数采用 5 次交叉验证法来确定。在利用 SVM 进行分类时,对于超过两个类别的情况,采用一对剩余法(One-vs-Rest)<sup>[7]</sup>进行多类分类。

实验中的测试数据采用了文档分类领域两个著名的测试集 Reuters-21578 和 20-Newsgroup。对于 Reuters-21578 数据集<sup>[8]</sup>,我们遵循常用的“ModApte”切分方式,选择文档最多的 10 个类别进行实验,其中包括了 7 194 个训练文档和 2 788 个测试文档。20-Newsgroup 数据集<sup>[9]</sup>包含大致均匀地分布在 20 个类别上的 18 846 个文档,我们采用每个类别中的一半数据作为训练集,另一半作为测试集。分类器性能评价指标使用常用的查准率(Precision)、查全率(Recall)及微平均  $F_1$  值,这三个评价指标的值越大,说明分类器的性能越好。它们的定义如下:

$$Precision = \frac{\text{正确分为某类的文档数}}{\text{测试集中分为该类别的文档总数}} \times 100\%$$

$$Recall = \frac{\text{正确分为某类的文档数}}{\text{测试集中属于该类别的文档总数}} \times 100\%$$

$$F_1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \times 100\%$$

表 1 和表 2 分别给出了 5 种分类方法在测试集 Reuters-21578 和 20-Newsgroup 上的实验结果,从中可以看出:本文提出的 KDA-SVM 分类算法的查准率、查全率和  $F_1$  评价指标普遍高于 KNN、LDA、Bayes 和 SVM 方法。原因在于 KDA-SVM 方法首先使用核鉴别分析法 KDA 对训练集中的高维文档空间进行非线性降维,然后在降维后的低维特征空间中利用优化的 SVM 进行分类预测。由于核鉴别分析 KDA 有效地减少了文档特征空间的维数,从而大大地提高了文档分类器 SVM 的计算效率。而利用乘性更新原则训练的 SVM 具有很快的训练速度和较好的分类效果,从而保证了 KDA-SVM 具有较高的分类准确率和较快的运行速度。

表 1 在 Reuters-21578 上的分类性能比较 %

分类方法	查准率	查全率	$F_1$
KNN	82.9	78.3	80.5
LDA	84.2	79.5	81.8
Bayes	85.6	81.4	83.4
SVM	87.3	85.7	86.5
KDA-SVM	90.5	87.2	88.8

表 2 在 20-Newsgroup 上的分类性能比较 %

分类方法	查准率	查全率	$F_1$
KNN	78.6	72.5	75.4
LDA	81.3	77.8	79.5
Bayes	89.5	80.3	84.7
SVM	93.2	85.6	89.2
KDA-SVM	96.4	87.8	91.9

另外,我们还对这 5 种分类算法在不同文档大小下的运行时间进行了比较。从图中的实验结果可以看出本文提出的分类算法 KDA-SVM 的运行时间都低于 KNN、LDA、Bayes 和 SVM 方法,并且随文档大小的增长幅度比较平缓,说明本文

提出的分类算法在处理大规模数据集时具有较好的扩展性。原因在于我们采用的核鉴别分析法 KDA 有效地减少了文档特征空间的维数,避免了高维文档引起的“维数灾难”问题。另外,采用乘性更新规则训练的优化 SVM 避免了常规 SVM 训练中非常耗时的大量矩阵运算,使得训练速度得到很大提高,从而保证了 KDA-SVM 具有很快的运行速度。

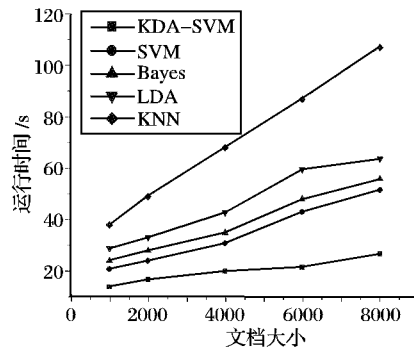


图 1 在 Reuters-21578 上的运行时间比较

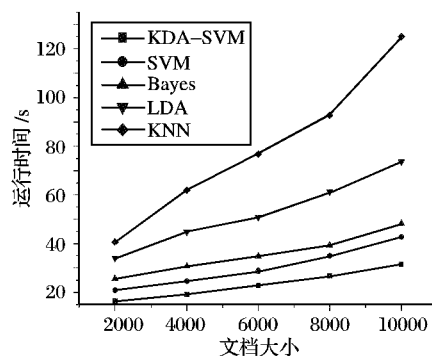


图 2 在 20-Newsgroup 上的运行时间比较

## 3 结语

为了有效地解决高维文档分类时的“维数灾难”问题,本文提出了基于核鉴别分析方法 KDA 和 SVM 的 Web 文档分类算法。实验结果表明,该算法有效地提高了文档分类的准确率,而且具有很快的运行速度。

### 参考文献:

- [1] SEBASTIANI F. Machine learning in automated text categorization [J]. ACM Computing Surveys, 2002, 34(1): 1-47.
- [2] VAPNIK V N. The nature of statistical learning theory [M]. New York: Springer, 1995.
- [3] SHA F, LIN Y Q, SAUL L K, LEE D D. Multiplicative updates for nonnegative quadratic programming [J]. Neural Computation, 2007, 19(8): 2004-2031.
- [4] BAUDAT G, ANOUAR F. Generalized discriminant analysis using a kernel approach [J]. Neural Computation, 2000, 12(10): 2385-2404.
- [5] MIKA S, RATSCH G, WESTON J, et al. Fisher discriminant analysis with kernels [C] // Proceedings of the 1999 IEEE Signal Processing Society Workshop on Neural Networks for Signal Processing. New York: IEEE Press, 1999: 41-48.
- [6] DUDA R O, HART P E, STORK D G. Pattern classification [M]. Second edition. Hoboken: Wiley-Interscience, 2000.
- [7] HSU C W, LIN C J. A comparison on methods for multi-class support vector machines [J]. IEEE Transactions on Neural Networks, 2002, 13(2): 415-425.
- [8] LEWIS D D. Reuters-21578 text categorization collection [EB/OL]. [2008-06-22]. <http://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html>.
- [9] 20 News Group [EB/OL]. [2008-06-22]. <http://kdd.ics.uci.edu/databases/20newsgroups/20newsgroups.htm>.