

基于 SVM 藏文文本分类的研究与实现

文/贾宏云 群诺 苏慧婧 次仁罗增 巴桑卓玛

摘要

利用 SVM(支持向量机)技术对复杂繁琐的汉文文本资源进行快速分类已经相当的成熟,但其在藏文文本分类中的应用还处于研究阶段,因此实验目的在于测试该方法是否在藏文文本分类中具有良好的性能。主要过程包括:文本向量空间模型化,获取 SVM 中核函数的参数并进行常用核函数分类性能对比,最后与 Logistic 回归分类器进行同等条件下的实验对比,验证了支持向量机模型在藏文文本分类中具有良好的分类效果。

【关键词】藏文文本分类 支持向量机 Logistic 回归

1 引言

由于文本记录着时代变迁的痕迹,文本的数量在历史的长河中不断增加,因人们在查

阅和学习的过程中对相关文本的寻找显得十分麻烦,所以文本分类的有效性对上述问题的解决具有重要作用。同时伴随着科学技术的发展,人们开始利用计算机的高效性进行文本自动分类,因此对分类数学模型的选择变得更为重要。

目前,机器学习算法已成为主流的方法,尤其在中文文本分类算法的研究上已经相当成熟,特别是 SVM 算法利用最小结构风险的原理使得分类性能更加优异,在武汉理工大学熊浩勇[1]的硕士毕业论文中已经详细描述,虽然其具有对模型参数的设置相当复杂并且耗时间长等不足之处,但所获取的模型参数十分精确。由于 SVM 的核函数很多,因此不同结构的文本使用的核函数具有差异性,藏文文本也存在这种情况。因此实验目的在于测试该方法是否在藏文文本分类中具有良好的性能。主要过程包括:文本向量空间模型化,获取 SVM 中核函数的参数并进行常用核函数分类性能对比,最后与 Logistic 回归分类器进行同等条件下的实验对比,验证了支持向量机模型在藏文文本分类中具有良好的分类效果。

2 SVM模型分类原理

SVM 是一种二分类模型,但可以在多分

类中进行多次二分类,它的基本模型定义在样本特征数据空间上的间隔最大的线性分类器,有效的解决样本特征数据在低维空间中非线性(线性也是非线性的特殊情况)的情况下,通过核函数把样本数据映射到高维空间中,利用经验风险和结构风险最小化原理找到线性超平面实现样本分类。由于藏文文本特征的高稀疏性和低维空间中样本的不可分性,所以本文选择一定的惩罚参数 C 和核函数来构造 SVM 分类器。

2.1 SVM分类器构建算法

2.1.1 确定目标函数

构建最优分类面来分割属于两类的训练集: $(x_i, y_i), i=1, 2, \dots, n, x_i \in R^d, y_i \in \{+1, -1\}$ 的问题,可转化为解下述二次规划问题:在约束条件 $y_i(w \cdot x_i + b) \geq 1, i=1, 2, \dots, n$ 下,为了使分类器具有更好的泛化能力和良好的分类效果,求 w 和 b 的优化条件是使两类样本到超平面的距离之和

$\frac{2}{\|w\|}$ 最大值,其等价于求目标函数 $\phi(w) = \frac{1}{2} \|w\|^2$ 的最小值。

2.1.2 引入松弛变量,构建惩罚参数

<< 上接 143 页

4 实验结果与分析

4.1 实验设计和结果

才旦夏茸大师文集第一册至十三册作为实验语料,其中第一册到第三册为训练语料,用于建设消歧知识库和训练属格的 5 个助词的检错正则表达式,第四册至第六册内容作了修改作为测试语料。语料的规模如表 3。

衡量指标选用了准确率 P:

$$P = \frac{\text{返回正确结果数}}{\text{返回结果总数}} \times 100\% \quad (8)$$

方法 1 代表基于正则表达式的属格助词自动检错算法 1,方法 2 代表基于正则表达式和知识库的属格助词自动检错算法,实验结果如表 4。

4.2 实验结果分析

从计算的结果可以看出,采用方法 2 检错率比方法 1 的检错率高,虽然只增加了消歧知识库,但准确率明显提高,但方法 2 比方法 1 增加了时间复杂度 $T(n)=O(\log(2n))$ 。

在实验过程中也发现一些有待解决的问题:消歧知识库的规模不大,还得扩充知识库。

5 结束语

纵观当前少数民族语言文字发展的现状,我们可以清楚地看到,要想句法分析向语义分析阶段顺利迈进,目前最重要的问题就是处理好藏文的虚词,藏文虚词的研究成果可以在各个层面上推广应用。下一步工作计划是,扩充藏文歧义虚词知识库,提高藏文虚词识别和检错的准确性。

总体而言藏文属格助词的识别及其自动检错算法的研究达到了可实用的水平。

参考文献

- [1] 卓玛吉,安见才让.藏文不自由虚词的自动识别研究[J].商业文化,2014(05).
- [2] 高定国,扎西加,赵栋材.计算机识别藏语虚词的方法研究[J].中文信息学报,2014,28(01):113-05.
- [3] 吴翔平.科技英语虚词分析法简介[J].

系统工程与电子技术学报,1986(07).

- [4] 杨慧玲.英语虚词在常规句和疑难句中的翻译比较分析[J].昆明师范高等专科学校学报,2006,28(01):86-88.
- [5] 多拉.藏语语义理解中功能性虚词研究[J].西藏大学学报(社会科学版),2011,26(04):106-107.
- [6] 胡书津.简明藏文文法:藏汉对照—2版[M].云南民族出版社,2000(10).
- [7] 格桑局冕,格桑央金.实用藏文文法教程[M].四川民族出版社,2004(11).
- [8] 才旦夏著.藏文文法详解[M].青海民族出版社,1954,5:18-45.

作者单位

拉毛措((1988-),女,藏族,西藏大学信息科学技术学院,硕士,主要从事自然语言信息处理研究。

作者单位

西藏大学信息科学技术学院 西藏自治区拉萨市 850000

表 4: 测试结果

C=0.1 时不同核函数下文本分类的准确率, 召回率, F1 值

| SVM | LINEAR | | | POLY | | | RBF | | | SIGMOID | | |
|------|--------|--------|------|--------|--------|------|--------|---------|------|---------|---------|------|
| 文本类别 | 精确率 | 召回率 | F1 值 | 精确率 | 召回率 | F1 值 | 精确率 | 召回率 | F1 值 | 精确率 | 召回率 | F1 值 |
| 人文类 | 97.58% | 94.53% | 0.96 | 99.17% | 93.75% | 0.96 | -- | -- | -- | -- | -- | -- |
| 教育类 | 98.54% | 94.84% | 0.97 | 96.62% | 93.90% | 0.95 | -- | -- | -- | -- | -- | -- |
| 政务类 | 97.68% | 90.36% | 0.94 | 95.88% | 91.43% | 0.94 | -- | -- | -- | -- | -- | -- |
| 时政类 | 98.56% | 83.66% | 0.91 | 93.06% | 81.71% | 0.87 | 26.83% | 100.00% | -- | -- | -- | -- |
| 经济类 | 99.29% | 92.11% | 0.96 | 97.96% | 94.74% | 0.96 | -- | -- | -- | -- | -- | -- |
| 法律类 | 99.52% | 93.72% | 0.97 | 99.53% | 94.17% | 0.97 | -- | -- | -- | -- | -- | -- |
| 民生类 | 98.99% | 80.33% | 0.89 | 97.96% | 78.69% | 0.87 | -- | -- | -- | 7.98% | 100.00% | -- |

表 5: 测试结果

C=10 时不同核函数下文本分类的准确率, 召回率, F1 值

| SVM | LINEAR | | | POLY | | | RBF | | | SIGMOID | | |
|------|--------|--------|------|--------|--------|------|--------|---------|------|---------|---------|------|
| 文本类别 | 精确率 | 召回率 | F1 值 | 精确率 | 召回率 | F1 值 | 精确率 | 召回率 | F1 值 | 精确率 | 召回率 | F1 值 |
| 人文类 | 97.58% | 94.53% | 0.96 | 99.17% | 93.75% | 0.96 | 0.78% | 0.78% | 0.01 | -- | -- | -- |
| 教育类 | 98.06% | 94.84% | 0.96 | 95.73% | 94.84% | 0.95 | 0.47% | 0.47% | 0.00 | -- | -- | -- |
| 政务类 | 97.32% | 90.71% | 0.94 | 98.03% | 88.93% | 0.93 | 3.57% | 3.57% | 0.04 | -- | -- | -- |
| 时政类 | 98.57% | 83.90% | 0.91 | 95.94% | 80.73% | 0.88 | 27.08% | 100.00% | 0.43 | -- | -- | -- |
| 经济类 | 98.60% | 92.76% | 0.96 | 99.30% | 93.42% | 0.96 | -- | -- | -- | -- | -- | -- |
| 法律类 | 98.57% | 92.83% | 0.96 | 99.06% | 94.17% | 0.97 | -- | -- | -- | -- | -- | -- |
| 民生类 | 98.06% | 82.79% | 0.90 | 97.92% | 77.05% | 0.86 | 1.64% | 1.64% | 0.02 | 7.98% | 100.00% | -- |

表 6: 测试结果对比

C=0.1 时不同核函数下 SVM 和 Logistic 回归文本分类性能对比

| SVM | LINEAR | | | POLY | | | Logistic | | |
|------|--------|--------|------|--------|--------|------|----------|--------|------|
| 文本类别 | 精确率 | 召回率 | F1 值 | 精确率 | 召回率 | F1 值 | 精确率 | 召回率 | F1 值 |
| 人文类 | 97.58% | 94.53% | 0.96 | 99.17% | 93.75% | 0.96 | 98.49% | 92.02% | 0.95 |
| 教育类 | 98.54% | 94.84% | 0.97 | 96.62% | 93.90% | 0.95 | 95.35% | 96.09% | 0.96 |
| 政务类 | 97.68% | 90.36% | 0.94 | 95.88% | 91.43% | 0.94 | 95.06% | 89.29% | 0.92 |
| 时政类 | 98.56% | 83.66% | 0.91 | 93.06% | 81.71% | 0.87 | 93.90% | 78.78% | 0.86 |
| 经济类 | 99.29% | 92.11% | 0.96 | 97.96% | 94.74% | 0.96 | 99.31% | 94.74% | 0.97 |
| 法律类 | 99.52% | 93.72% | 0.97 | 99.53% | 94.17% | 0.97 | 94.37% | 97.76% | 0.96 |
| 民生类 | 98.99% | 80.33% | 0.89 | 97.96% | 78.69% | 0.87 | 91.00% | 74.59% | 0.82 |

似。

(3) 由选择的特征向量中的值比较大, 使特征向量内积和差值相对很大, 因此 RBF 和 SIGMOID 的分类效果不好。

(4) 从表 6 测试结果可以看出, 当 SVM 核函数选择为 LINEAR 和 POLY 并且在上述参数下, 从整体参考值上看 SVM 的藏文本分类效果好于 Logistic 回归文本分类效果。

6 总结

本文采用基于 SVM 模型的藏文本分类实现过程中, 为了降低模型的复杂度, 对藏文本特征提取时, 忽略词与词之间联系, 因此假定词与词之间的互信息为零。实验验证了 SVM 模型对藏文本具有良好的效果, 因此后期会继续研究藏文本结构形式, 增大特征信息量, 提高分类的效果。

(通讯作者: 群诺)

参考文献

- [1] 熊浩勇. 基于 SVM 的中文文本分类算法研究与实现 [D]. 武汉理工大学, 2008.
- [2] 李航. 统计学习方法 [M]. 北京: 清华大学出版社, 2012.
- [3] 崔建明, 刘建明, 廖周宇. 基于 SVM 算法的文本分类技术研究 [J]. 计算机仿真, 2013.
- [4] 高定国, 珠杰. 藏文信息处理的原理与应用 [M]. 成都: 西南交通大学出版社, 2015.
- [5] 杨玉珍, 刘培玉, 朱振方, 邱烨. 应用特征项分布信息的信息增益改进方法研究 [J]. 山东大学学报 (理学版), 2009.
- [6] 杨杰明. 文本分类中文本表示模型和特征选择算法研究 [D]. 吉林大学, 2013.

作者简介

贾宏云 (1990-), 男, 四川省成都市人。硕士研究生在读, 西藏大学。主要研究方向为自然语言处理。

然语言处理。

苏慧婧 (1989-), 女, 四川省眉山市人。硕士研究生在读, 西藏大学。主要研究方向为自然语言处理。

次仁罗增 (1993-), 女, 藏族, 西藏自治区日喀则市人。硕士研究生在读, 西藏大学。主要研究方向为自然语言处理。

巴桑卓玛 (1991-), 女, 藏族, 西藏自治区山南市人。硕士研究生在读, 西藏大学。主要研究方向为自然语言处理。

通讯作者简介

群诺 (1972-), 男, 藏族, 西藏自治区拉萨市人。副教授, 西藏大学。主要研究方向为自然语言处理。

作者单位

西藏大学信息科学技术学院 西藏自治区拉萨市 850000