

一个对不带类别标记文本进行分类的方法

蒋志方¹, 祝翠玲², 吴强¹

(1. 山东大学计算机科学与技术学院, 济南 250061; 2. 山东经济学院信息管理学院, 济南 250014)

摘要: 利用无监督聚类方法和朴素贝叶斯分类的特点, 把 UC 获得的预分类结果作为朴素贝叶斯分类器的训练样本, 将处在聚类结果中类属模糊区域的文本交给训练好的朴素贝叶斯分类器再行分类, 实现了对不带任何类别标记文本的准确分类, 可得到较准确的分类结果。
关键词: 文本分类; 无监督文本聚类; 朴素贝叶斯分类; 欧氏距离

Method of Unlabeled Texts Classification

JIANG Zhifang¹, ZHU Cuiling², WU Qiang¹

(1. School of Computer Science and Technology, Shandong University, Jinan 250061;

2. College of Information Management, Shandong Economic University, Jinan 250014)

【Abstract】 Using the specialty of the unsupervised clustering and the naïve Bayes classification, the paper gives a method that gains results of the text clusters and takes some of results as the training samples of the naïve Bayes classifier and let the trained naïve Bayes classifier reclassify those texts in illegible area of the clustering results. Consequently the method can classify the unlabeled text accurately and also can gain a better result of classification.

【Key words】 Text classification; Unsupervised text clustering; Naïve Bayes classification; Euclid distance

文本挖掘是对海量、异质、非结构化的文本数据源进行的涉及数据挖掘、自然语言处理、计算语言学、信息检索及分类、知识管理等综合研究领域。文本分类常由训练过程和测试过程构成。在训练过程中, 首先用训练文本训练分类器, 得到最优的文本特征集合。在测试过程中, 首先将测试文本用最优特征子集表示, 再经分类器分类, 得到测试文本所属的类别。文本聚类把一组文档按照相似性划分为若干类, 属于同类的文档间的距离尽可能小而不同类文档间的距离则尽可能大^[1]。

带标记训练样本的获得常由人工来完成的, 随着文本数据的快速增长, 这种方式的代价越来越高。故对文本的自动分类已成为组织和管理文本数据的关键技术, 如何实现在非人工获得带标记训练样本和非人工预分类的情况下的准确文本分类, 是分类方法研究的一个重要研究课题。

为实现对文本的自动分类, 可采用聚类方法来产生训练样本和得到对文本分类的类别标记。本文提出了一个对不带任何类别标记的文本进行准确分类的方法。采用向量空间模型(vector space model, VSM)来表示文本, 即将文本表示成在 n 维向量空间中的一个点。对在向量空间中的文本, 指定聚类半径 R , 利用无监督聚类(unsupervised clustering, UC)进行聚类, 获得文本类别标记集合和聚类的正例中心和反例中心, 然后把聚类结果中处在包含正例中心区域内的文本作为训练样本对朴素贝叶斯(naïve Bayes, NB)分类器进行训练, 最后将在聚类结果中处于类属模糊区的文本让训练好的 NB 分类器再分类, 可得到较准确的分类结果, 提高分类精度。

1 无监督文本聚类算法

无监督文本聚类算法(unsupervised text clustering, UTC)

是一种把 UC 算法^[2-4]用于文本聚类的方法。在此算法中, 通过指定聚类半径 R , 分别对每类文本进行聚类并获得文本聚类中心; 然后, 把文本聚类中心作为对文本的预分类的根据, 即对任意给定的文本, 计算其与各聚类中心的距离; 离其最近的聚类中心所对应的类就是该文本的所属类。

1.1 算法描述

设有 m 个参加聚类的文本, 一个文本由 n 个特征词来表示, 记为 $X_i(A_{i1}, A_{i2}, \dots, A_{in})$, 形成 n 维向量空间中的一个点; 记聚类结果为 $C_j, j=1, 2, \dots, L$, L 为到目前为止所形成的类的个数, C_j 是一个属于同一个类的文本集合; O_j 是类 C_j 的中心, 也是 n 维向量空间中的一个点, 记为 $O_j(B_{j1}, B_{j2}, \dots, B_{jn})$, 但 O_j 可能不是 m 个文本之一; 记 k_j 为属于类 C_j 的文本的个数; 记

$d(X_i, O_j) = \sqrt{\sum_{k=1}^n (A_{ik} - B_{jk})^2}$ 为第 i 个文本 X_i 和第 j 个聚类中心 O_j 的欧氏距离。

在对 m 个文本采用 UTC 算法进行聚类后, 得到了 L 个类别 $C = \{C_1, C_2, \dots, C_L\}$, 每个类 C_j 中的文本数 k_j 和所包含的文本 $\{X_{j,t} | t=1, 2, \dots, k_j\}$, 及其聚类中心 O_j 。把类 C_j 中包含的文本称为该类的正例集合, 记为 Ω_j^+ , m 个文本中不属于类 C_j 的其他文本称为类 C_j 的反例集合, 记为 Ω_j^- , 其中文本个数为 $m - k_j$ 。

算法 1 无监督文本聚类算法(UTC)

输入 聚类半径 R ; m 个 n 维文档集合 $Z = \{X_1, X_2, \dots, X_m\}, X_i \in R^n$ 为第 i 个文本;

作者简介: 蒋志方(1961-), 男, 副教授, 主研方向: 数据挖掘与信息可视化; 祝翠玲, 硕士、助教; 吴强, 博士生

收稿日期: 2006-08-24 **E-mail:** zfjiang@sdu.edu.cn

输出 聚类结果 C_j , 聚类中心 O_j , C_j 中所包含的文本集合。

步骤:

(1) $C_1 = \{X_1\}, L = 1, O_1 = X_1, Z \leftarrow Z - \{X_1\}$

(2) 若 $Z = \Phi$, 则转至(7)。

(3) 选择 $X_i \in Z$, 从已有聚类中心中寻找与 X_i 最接近的中心 O_j , 即

$$O_j = \arg \min_{k=1}^L d(X_i, O_k)$$

(4) 若 $d(X_i, O_j) < R$, 则将 X_i 加入类 C_j , 即

$$C_j \leftarrow C_j \cup \{X_i\}$$

调整 C_j 类的中心为

$$O_j \leftarrow \frac{k_j \times O_j + X_i}{k_j + 1}, k_j \leftarrow k_j + 1$$

然后转到(6)。

(5) 否则 $L \leftarrow L + 1, C_L \leftarrow \{X_i\}, O_L \leftarrow X_i$ 。

(6) $Z \leftarrow Z - \{X_i\}$, 然后转到(2)。

(7) 以类 C_j 的聚类中心 O_j 作为该类的正例中心 O_j^+ , 计算出每个类 C_j 的反例中心 O_j^- 。

(8) 对 m 个文本中的任意文本 X , 分别计算出其与类 C_j 的正例中心 O_j^+ 的欧氏距离和与反例中心 O_j^- 的欧氏距离:

$$d_{X,j}^+ = d(X, O_j^+)$$

$$d_{X,j}^- = d(X, O_j^-)$$

(9) 求 X 与所有正例中心之间欧氏距离的最小值和 X 与所有反例中心之间欧氏距离的最小值。

$$d_X^+ = \min_{j=1}^L d_{X,j}^+$$

$$d_X^- = \min_{j=1}^L d_{X,j}^-$$

(10) 决定文本 X 的类属:

若 $d_X^+ < d_X^-$, 则 $X \in C_j$; 若 $d_X^+ > d_X^-$, 则 $X \notin C_j$; 若 $d_X^+ = d_X^-$, 则 X 类属不确定。

1.2 算法复杂度分析

算法从步骤(2)到步骤(3)为一循环过程。算法要对每个文本 $X_i (i=1, 2, \dots, m)$ 计算其与 L 个类的聚类中心的距离, 所以该循环过程的时间复杂度是 $O(L^2mn)$ 。在步骤(7)中计算反例中心的时间复杂度为 $O(L^2)$, 而在步骤(8)至步骤(9)的时间复杂度为 $O(2mnL) + O(mn)$ 。所以整个算法的时间复杂度是 $O(L^2mn) + O(L^2) + O(2mnL) + O(mn) = O(L^2mn)$ 。因此该算法具有较高的效率。

1.3 分类准确性分析

m 个文本中的每一个 X 都可以用上述算法被归属于一个类 C_j 。事实上, 对于任意给定的文本集合, 在向量空间中, 存在一个 $\varepsilon (\varepsilon > 0)$ 区域, 其示意如图1所示。图中点划线所表示的区域就是 区域, 由欧氏距离公式和 $|d_X^+ - d_X^-| < \varepsilon$ 可知图中曲线为双曲线。

当 $|d_X^+ - d_X^-| > \varepsilon$ 时, 文本 X 在此区域之外。当 $d_X^+ - d_X^- > \varepsilon$ 时, 可以确定 $X \in C_j$, 而当 $d_X^- - d_X^+ > \varepsilon$ 时, 则可以确定 $X \notin C_j$ 。

当 $|d_X^+ - d_X^-| < \varepsilon$ 时, 文本 X 在此区域内时, X 与 C_j 的正例中心和反例中心距离相差不多, 对其的归类就不一定准确。

可把 ε 作为一个决策阈值, 来确定那些利用 UTC 方法可

进行准确分类的文本, 当 $|d_X^+ - d_X^-| > \varepsilon$ 时, 可保证分类的准确性。

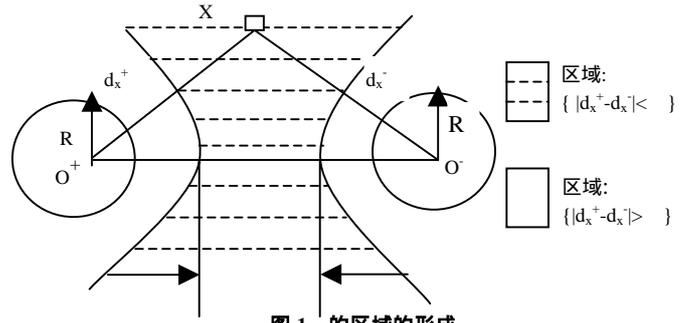


图1 区域的形成

2 朴素贝叶斯分类

2.1 分类方法

朴素贝叶斯分类(Naive Bayes classification, NBC)^[1,5,6]是一种基于概率的分类法, 它假设一个属性对给定类的影响独立于其他属性, 即特征独立性假设。当假设成立时, 与其他分类算法相比, NBC算法是最精确的^[9]。

假设有 L 个类 C_1, C_2, \dots, C_L 。给定一个没有类标记的未知数据样本 $X(a_1, a_2, \dots, a_n)$ (a_i 表示第 i 个属性的值), NBC将预测 X 属于具有最高后验概率 $P(C_j | X)$ (条件 X 下)的类 C_j , 即将未知的样本分配给类 C_j , 当且仅当

$$P(C_j | X) > P(C_i | X), 1 \leq i \leq L, i \neq j \quad (1)$$

根据贝叶斯定理, 有:

$$P(C_j | X) = \frac{P(X | C_j)P(C_j)}{P(X)} \quad (2)$$

为了降低计算 $P(X | C_j)$ 的开销, 由NBC的特征独立性假设, 有:

$$P(X | C_j) = \prod_{k=1}^n p(A_k = a_k | C_j) \quad (3)$$

其中, A_k 表示样本 X 的第 k 个属性; a_k 表示属性 A_k 的值, n 表示属性的个数。因此式(2)可以改写为

$$P(C_j | X) = \frac{\prod_{k=1}^n p(A_k = a_k | C_j)P(C_j)}{P(X)} \quad (4)$$

式(4)中, $P(X)$ 是一个定值, 由于NBC只需要比较 $P(C_j | X)$ 的大小, 而无需计算 $P(C_j | X)$ 的值, 因此对 $P(X)$ 可不予计算。

贝叶斯理论把在训练数据中得到的概率作为条件概率, 因此式(4)中的一些概率可以直接从训练集获取。 $P(C_j)$ 为任意一个样本属于类 C_j 的概率, 若记类 C_j 中的样本个数为 $|C_j|$, 所有样本总数为 $|S|$, 则有:

$$P(C_j) = \frac{|C_j|}{|S|}$$

(1) 如果 A_k 是离散变量, 则有

$$P(A_k = a_k | C_j) = \frac{|C_j | A_k = a_k |}{|C_j|} \quad (5)$$

其中, $|C_j | A_k = a_k |$ 表示类 C_j 中属性 A_k 的值为 a_k 的样本数。

(2) 如果 A_k 是连续值属性, 则通常假定该属性服从高斯分布。因而,

$$P(A_k = a_k | C_j) = g(a_k, \mu_{c_j}, \sigma_{c_j}) = \frac{1}{\sqrt{2\pi}\sigma_{c_j}} e^{-\frac{(a_k - \mu_{c_j})^2}{2\sigma_{c_j}^2}} \quad (6)$$

其中, 给定类 C_j 的训练样本属性 A_k 的值; $g(a_k, \mu_{c_j}, \sigma_{c_j})$ 是属性

A_k 的高斯密度函数； μ_{c_j} 和 σ_{c_j} 分别为平均值和标准差。

那么，对未知样本 X 分类，对每个类 C_j ，计算 $P(X|C_j)P(C_j)$ ，样本被指派到 C_j ，当且仅当

$$P(X|C_j)P(C_j) > P(X|C_i)P(C_i), 1 \leq i \leq L, j \neq i \quad (7)$$

即 X 被指派到 $P(X|C_j)P(C_j)$ 最大的类 C_j 。

2.2 特点分析

- (1)在独立性假定成立时，可获得精度最优的分类效果；
- (2)具有更小的出错率、更高的健壮性和效率；
- (3)计算量较大。

3 无类别标记的文本的准确分类方法

3.1 方法的提出

无监督文本聚类的速度快但准确率较低，会出现不能准确判断其类别归属的模糊区域，而NBC在特征独立假设的前提下，具有分类准确度高，但计算量比较大的特点。可以只取出那些处于模糊区域的文本交给NBC来进行准确的分类。这样，将UTC方法和NBC方法结合起来，就有可能既保证有较快的分类速度，又有较高的分类准确率。由此本文提出一种基于无监督文本聚类 and NBC的分类方法，能够对不带任何类别标记的文本进行准确的分类。

3.2 方法的描述

对于任意给定的文本集合，可以给定一个决策阈值 $\varepsilon (\varepsilon > 0)$ ，首先对所有文本运用无监督文本聚类算法，产生一个聚类结果。在1.2节中，划分出了一个区域。处于区域之外，满足 $d_x^+ - d_x^- > \varepsilon$ 的文本由UTC就能得到比较准确的分类，可以把这些文本作为NB分类器的训练样本，来训练NB分类器。而处在区域之内的文本，用UTC方法对其进行分类出错的概率较高，所以可以把该文本交给经过训练的NB分类器进行重新分类。

算法2 对无类别标记的文本的准确分类算法(UNBTC)

输入 聚类半径 R ，决策阈值 ε ，待分类的文本集合 $Z = \{X_1, X_2, \dots, X_m\}, X_i \in R^n$ 。

输出 聚类结果 C_j, C_j 中所包含的文本。

步骤：

- (1)设定聚类半径 R ，决策阈值 ε 。
- (2)UTC算法对文本集合进行聚类，得到 $C_j(j=1,2,3,\dots,L)$ 。
- (3)对 $|d_x^+ - d_x^-| > \varepsilon$ 区域中的文本，得到具有明确类属的文本的集合

(4)用 $d_x^+ - d_x^- > \varepsilon$ 区域中的已有类别归属的文本对NB分类器进行训练

(5)用经过训练的NBC分类器对满足 $|d_x^+ - d_x^-| < \varepsilon$ 的那些文本进行再分类。

(6)输出文本分类结果。

在本算法中，当待分类文本在文本向量空间中具有清晰类属时，应用了UTC算法的高效性与准确性，对处于向量空间中模糊区域的文本，则利用了NBC的分类错误率。

3.3 实验结果

在UCI知识发现数据库中的newsgroups数据集上对上述算法进行了试验，选取了5个分类中的123篇文档，选取了276个关键词对各个文档进行表示。表1是当聚类半径 R 取不同值时对guns类、mideast类、religion类的分类准确度(%)对比情况。

表1 NBC和UTC聚类准确率 (%)

accuracy	guns	mideast	religion
----------	------	---------	----------

method			
NBC	87.81	85.37	92.68
UTC(R=0.60)	77.24	65.04	82.93
UTC(R=0.50)	79.67	75.61	86.18
UTC(R=0.45)	83.74	82.11	88.61
UTC(R=0.40)	81.30	82.93	87.81

由表1可以看出，无论采用NBC还是采用UTC，对于religion类的数据得到的结果准确率都高于其余两类。这说明在文本向量空间中该类距离其它类比较远，界限较清晰。对于这3个类别而言，NBC的准确率均要高于UTC，故NBC确比较有精确的区分能力。在UTC中，随着聚类半径 R 的不同，聚类结果也有不同的变化，当 $R=0.45$ 时达到较好的分类效果。运用该方法，当 $R=0.45$ 时，取不同的值时最后得到的结果如表2所示。

表2 R=0.45时取不同值时的分类准确率 (%)

Class	guns	mideast	religion
0.10	84.55	82.11	88.61
0.15	85.37	82.93	89.43
0.18	86.99	83.74	91.06
0.20	85.37	83.74	91.87

由表2中可以看出，这3个类别的分类正确率总体呈增长趋势，但也出现一些异常的情况，这可能是由多种情况引起的，例如：在计算权重时忽略了关键词在文中位置对其在文中重要性的贡献度；在利用UTC算法得到的结果对NBC器进行训练时选取的训练文本的数量以及准确度方面的提高方面等，还有待进一步改进。

4 结语

将UTC法与NBC法相结合是一种有效的分类方法，它实现了对不带任何类别标记的文本的准确分类，并解决了UTC准确率低和NBC方法计算量大、需要预分类的问题，充分利用了UTC分类速度快和NBC精度高的优点，获得了较高的训练速度和较高的识别准确率。

参考文献

- 1 Han J, Kamber M. 数据挖掘概念与技术[M]. 范明, 孟小峰. 译. 北京: 机械工业出版社, 2001.
- 2 Eskin E, Arnold A, Prerau M, et al. A Geometric Framework for Unsupervised Anomaly Detection: Detecting Intrusions in Unlabeled Data[M]//Applications of Data Mining in Computer Security. Boston: Kluwer Academic Publisher, 2002.
- 3 罗敏, 王丽娜. 基于无监督聚类的入侵检测方法[J]. 电子学报, 2003, 31(11).
- 4 Dougherty J. Supervised and Unsupervised Discretization of Continuous Features[C]//Proceedings of the 12th International Conference on Machine Learning. 1995.
- 5 Plangley W, Thompson K. An Analysis of Bayesian Classifiers[C]//Proc. of the 10th National Conf. on Artificial Intelligence. 1992.
- 6 Nigam K, McCallum A, Thrun S. Learning to Classify Text from Labeled and Unlabeled Documents[C]//Proceedings of the 15th National Conference on Artificial Intelligence. 1998.
- 7 McCallum A, Nigam K. A Comparison of Event Models for Naive Bayes Text Classification[C]//Proc. of the Workshop on Learning for Text Categorization. 1998.
- 8 宫秀军, 刘少辉, 史忠植. 一种增量贝叶斯分类模型[J]. 计算机学报, 2002, 25(6).

