

tradaboost算法原理

原创 2017年05月14日 19:03:44

402

定义迁移学习的模型如下：设 $X_b$ 为源样例空间， $X_a$ 为辅助样例空间，源样例空间也就是我们的目标空间，就是想要去分类的样例空间。设 $Y=\{0, 1\}$ 为类空间，这里简化了多分类问题为二分类问题讨论，这样我们的训练数据也就是

$$T \subseteq \{(X = X_b \cup X_a) \times Y\}$$

测试数据：

$$S = \{(x_i^t)\}, \text{其中 } x_i^t \in X_b, \text{当 } i = 1, 2, \dots, k,$$

其中测试数据是未标注的，我么可以将训练数据划分为两个数据集：

$$T_a = \{(x_i^a, c(x_i^a))\}, \text{其中 } x_i^a \in X_a, \text{当 } i = 1, 2, \dots, n;$$

$$T_b = \{(x_j^b, c(x_j^b))\}, \text{其中 } x_j^b \in X_b, \text{当 } j = 1, 2, \dots, m.$$

其中 $c(x)$ 代表样本数据 $x$ 的真实所属的类别， $T_a$ 和 $T_b$ 的区别在于 $T_b$ 和测试数据 $S$ 是同分布的， $T_a$ 和测试数据是不同分布的，现在的任务就是给定很少的源数据 $T_b$ 和大量的辅助数据 $T_a$ 训练出一个分类器在测试数据 $S$ 上的分类误差最小。这里假设利用已有的数据 $T_b$ 不足以训练出一个泛化能力很强的分类器。

TrAdaBoost算法

我们利用AdaBoost算法的思想原理来解决这个问题，起初给训练数据 $T$ 中的每一个样例都赋予一个权重，当一个源域 $T_b$ 中的样本被错误的分类之后，我们认为这个样本是很难分类的，于是乎可以加大这个样本的权重，这样在下一次的训练中这个样本所占的比重就更大，这一点和基本的AdaBoost算法的思想是一样的。如果辅助数据集中的一个样本被错误的分类了，我们认为这个样本对于目标数据是很不同的，我们就降低这个数据在样本中所占的权重，降低这个样本在分类器中所占的比重，下面给出TrAdaBoost算法的具体流程：

算法 1 TrAdaBoost 算法描述

输入 两个训练数据集 $T_a$ 和 $T_b$ （根据公式(3.1)，合并的训练数据集 $T = T_a \cup T_b$ ），一个未标注的测试数据集 $S$ ，一个基本分类算法 $Learner$ ，和迭代次数 $N$ 。

初始化

- 1. 初始权重向量 $\mathbf{w}^1 = (w_1^1, \dots, w_{n+m}^1)^T$ ，其中，
$$w_i^1 = \begin{cases} 1/n & \text{当 } i = 1, \dots, n \\ 1/m & \text{当 } i = n + 1, \dots, n + m \end{cases}$$

- 2. 设置 $\beta = 1/(1 + \sqrt{2 \ln n/N})$ .

For  $t = 1, \dots, N$

- 1. 设置 $\mathbf{p}^t$ 满足

$$\mathbf{p}^t = \frac{\mathbf{w}^t}{\sum_{i=1}^{n+m} w_i^t}.$$

- 2. 调用 $Learner$ ，根据合并后的训练数据 $T$ 以及 $T$ 上的权重分布 $\mathbf{p}^t$ 和未标注数据集 $S$ ，得到一个在 $S$ 的分类器 $h_t : X \mapsto Y$ .



xiaocong1990 (http://b...

+ 关注

(http://blog.csdn.net/xiaocong1990) 码云

原创 345 粉丝 6 喜欢 0 (https://c... utm\_sour

他的最新文章

更多文章

(http://blog.csdn.net/xiaocong1990)

[leetcode]591. Tag Validator (/xiaocong1990/article/details/77917601)

[leetcode]664. Strange Printer (/xiaocong1990/article/details/77916536)

[leetcode]233. Number of Digit One (/xiaocong1990/article/details/77822731)

[leetcode]381. Insert Delete GetRandom O(1) - Duplicates allowed (/xiaocong1990/article/details/77822065)

在线课程



用户画像系统应用与技术分析 (http://edu.csdn.net/course/series\_detail/?utm\_source=wx2) (http://edu.csdn.net/hu... iyicourse/series\_detail/?utm\_source=wx2)



2017 求职面试集锦之VIP服务版 (http://edu.csdn.net/course/series\_detail/?utm\_source=blog9) (http://edu.csdn.net/hu... iyicourse/series\_detail/70?utm\_source=blog9)

热门文章

主成分分析 (PCA) 原理详解 (/xiaocong1990/article/details/53584575) 1877

PyCharm配置Spark开发环境 (/xiaocong1990/article/details/53767722) 1875

联合分布 & 条件分布 & 边缘分布 (/xiaocong1990/article/details/72027006) 1327

奇异值分解(SVD)原理详解 (/xiaocong1990/article/details/54909126) 1319

3. 计算 $h_t$ 在 $T_b$ 上的错误率:

$$\epsilon_t = \frac{\sum_{i=n+1}^{n+m} w_i^t |h_t(x_i) - c(x_i)|}{\sum_{i=n+1}^{n+m} w_i^t}.$$

4. 设置 $\beta_t = \epsilon_t / (1 - \epsilon_t)$ .

5. 设置新的权重向量如下

$$w_i^{t+1} = \begin{cases} w_i^t \beta^{ |h_t(x_i) - c(x_i)| }, & \text{当 } i = 1, \dots, n \\ w_i^t \beta^{- |h_t(x_i) - c(x_i)| }, & \text{当 } i = n + 1, \dots, n + m \end{cases}$$

输出 最终分类器

$$h_f(x) = \begin{cases} 1, & \sum_{t=\lceil N/2 \rceil}^N \ln(1/\beta_t) h_t(x) \geq \frac{1}{2} \sum_{t=\lceil N/2 \rceil}^N \ln(1/\beta_t) \\ 0, & \text{其他} \end{cases}$$

可以看到, 在每一轮的迭代中, 如果一个辅助训练数据被误分类, 那么这个数据可能和源训练数据是矛盾的, 那么我们就可以降低这个数据的权重。具体来说, 就是给数据乘上一个 $\beta^{|h_t(x_i) - c(x_i)|}$ , 其中 $\beta$ 的值在0到1之间, 所以在下一轮的迭代中, 被误分类的样本就会比上一轮少影响分类模型一些, 在若干次以后, 辅助数据中符合源数据的那些数据会拥有更高的权重, 而那些不符合源数据的权重会降低。极端的一个情况就是, 辅助数据被全部忽略, 训练数据就是源数据 $T_b$ , 这样这时候的算法就成了AdaBoost算法了。在计算错误率的时候, 当计算得到的错误率大于0.5的话, 需要将其重置为0.5。

可以看到, TrAdaBoost算法在源数据和辅助数据具有很多的相似性的时候可以取得很好效果, 但是算法也有不足, 当开始的时候辅助数据中的样本如果噪声比较多, 迭代次数控制的不好, 这样都会加大训练分类器的难度。

版权声明: 本文为博主原创文章, 未经博主允许不得转载。



## 相关文章推荐

### 不懂爱情 (/starringnight/article/details/6915790)

今天上班没有见到他, 或许明天能见到吧。 我依然在等待着, 当我看到他, 内心里的忐忑让我恍惚回到了年轻时候的情感, 朦胧而美好, 可以只属于我一个人, 可能只属于我一个人。 所以...



starringnight (http://blog.csdn.net/starringnight) 2011-10-28 22:16 86

### 深度学习调参备忘 (一) (/yingyujianmo/article/details/45196333)

CNNs调参备忘: 在所有深度网络中, 卷积神经网络和图像处理最为密切相关, 卷积网络在很多图片分类竞赛中都取得了很好的效果, 但卷积网调参过程很不直观, 很多时候都是碰运气。为此, 卷积网络发明者Yann LeC...



yingyujianmo (http://blog.csdn.net/yingyujianmo) 2015-04-22 15:59 1470

### 准确率、召回率、F1 (/xiaocong1990/article/details/72584506)

准确率和召回率是广泛用于信息检索和统计学分类领域的两个度量值, 用来评价结果的质量。其中精度是检索出相关文档数与检索出的文档总数的比率, 衡量的是检索系统的查准率; 召回率是指检索出的相关文档数和文档库中所...



xiaocong1990 (http://blog.csdn.net/xiaocong1990) 2017-05-20 19:30 143