CSDN新首页上线啦,邀请你来立即体验! (http://blog.csdn.net/)

立即体验

CSDN

博客 (http://blog.csdn.net/?ref=toolbar)

学院 (http://edu.csdn.net?ref=toolbar)

下载 (http://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://dbtth://db





登录 (http:///passport.csdn.net/account/mobileregister?ref=toolbar&action=mobileRegister) /posteditieve/letothatar)

/activity?utm_source=csdnblog1) --完整总结(**ID3**,**C4.5**,**CART,**剪枝,替代)

2016年08月06日 14:24:47

7589

原文:

http://blog.csdn.net/zhaocj/article/details/50503450 (http://blog.csdn.net/zhaocj/article/details/50503450)

一、原理

决策树是一种非参数的监督学习方法,它主要用于分类和回归。决策树的目的是构造一种模型,使之能够 从样本数据的特征属性中,通过学习简单的决策规则——IF THEN规则,从而预测目标变量的值。

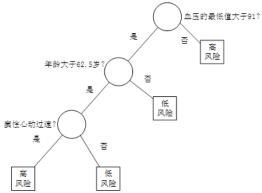


图1决策树

例如,在某医院内,对因心脏病发作而入院治疗的患者,在住院的前24小时内,观测记录下来他们的19个 特征属性——血压、年龄、以及其他17项可以综合判断病人状况的重要指标,用图1所示的决策树判断病人 是否属于高危患者。在图1中,圆形为中间节点,也就是树的分支,它代表IF THEN规则的条件;方形为终 端节点(叶节点),也就是树的叶,它代表IF THEN规则的结果。我们也把第一个节点称为根节点。 决策树往往采用的是自上而下的设计方法,每迭代循环一次,就会选择一个特征属性进行分叉,直到不能 再分叉为止。因此在构建决策树的过程中,选择最佳(既能够快速分类,又能使决策树的深度小)的分叉 特征属性是关键所在。这种"最佳性"可以用非纯度(impurity)进行衡量。如果一个数据集合中只有一种分 类结果,则该集合最纯,即一致性好;反之有许多分类,则不纯,即一致性不好。有许多指标可以定量的 度量这种非纯度,最常用的有熵,基尼指数(Gini Index)和分类误差,它们的公式分别为:

Entropy =
$$E(D) = -\sum_{j=1}^{J} p_j \log_2 p_j$$
 (1)

$$\text{Gini Index} = \text{Gini}(D) = \sum_{j=1}^{J} p_j (1 - p_j) = \sum_{j=1}^{J} p_j - \sum_{j=1}^{J} p_j^2 = 1 - \sum_{j=1}^{J} p_j^2 \qquad (2)$$

Classification Error = $1 - \max\{p_i\}$

上述所有公式中,值越大,表示越不纯,这三个度量之间并不存在显著的差别。式中D表示样本数据的分类 集合,并且该集合共有J种分类, p_i 表示第j种分类的样本率:

$$p_j = \frac{N_j}{N} \qquad (4)$$

式中N和 N_i 分别表示集合D中样本数据的总数和第i个分类的样本数量。把式4带入式2中,得到:



EmbeddedApp (http://bl...

(http://blog.csdn.net

/app_12062011) 原创 粉丝

码云 喜欢 未开通

219 1291

(https://git€

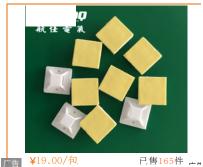
他的最新文章

更多文章 (http://blog.csdn.net/app_12062011)

人脸识别之人脸对齐(八)--LBF算法 (http://blog.csdn.net/App_12062011/art icle/details/78662404)

人脸识别之人脸对齐(七)--JDA算法 (http://blog.csdn.net/App_12062011/art icle/details/78662348)

人脸识别之人脸对齐(六)--ERT算法 (http://blog.csdn.net/App_12062011/art icle/details/78661992)



已售165件 二,4

在线课程



腾讯云容器服务架构实 现介绍() 讲师: 董晓杰



(http://edu.csdn.net

溶器技术。存58属域的实 践 (http://edu.csdn.net /73?utm_source=blog9) /hwiyicourse /series_detail /73?utm_source=blog9)

他的热门文章

$$\operatorname{Gini}(D) = 1 - \sum_{j=1}^{J} \left(\frac{N_j}{N}\right)^2 = 1 - \frac{\sum_{j=1}^{J} N_j^2}{N^2}$$
 (5)

目前常用的决策树的算法包括ID3(Iterative Dichotomiser 3, 第3代迭戈二叉树)、C4.5和 CART (ClassificationAnd Regression Tree, 分类和回归树)。前两种算法主要应用的是基于熵的方法, 而第三种性用的是基尼指数的方法。下面我们就逐一介绍这些方法。

ID3

ID3是由Ross Quinlan首先提出,它是基于所谓"Occam'srazor"(奥卡姆剃刀),即越简单越好,也就是越 是小型的决策树越优于大型的决策树。如前所述,我们已经有了熵作为衡量样本集合纯度的标准,熵越 大,越不纯,因此我们希望在分类以后能够降低熵的大小,使之变纯一些。这种分类后熵变小的判定标准 可以用信息增益(Information Gain)来衡量,它的定义为:

$$G(D \mathfrak{A}^{\mathfrak{D}}) = E(D) - \sum_{i=1}^{n} \frac{N_i}{N} E(D_i) \qquad (6)$$

该式表示在样本集合D下特征属性A的信息增益,n表示针对特征属性A,样本集合被划分为n个不同部分, 即A中包含着n个不同的值, N_i 表示第i个部分的样本数量, $E(D_i)$ 表示特征属性A下第i个部分的分类集合的 熵。信息增益越大,分类后熵下降得越快,则分类效果越好。因此我们在D内遍历所有属性,选择信息增益 最大的那个特征属性进行分类。在下次迭代循环中,我们只需对上次分类剩下的样本集合计算信息增益, 如此循环, 直至不能再分类为止。

C4.5

C4.5算法也是由Quinlan提出,它是ID3算法的扩展。ID3应用的是信息增益的方法,但这种方法存在一个问 题,那就是它会更愿意选择那些包括很多种类的特征属性,即哪个A中的n多,那么这个A的信息增益就可 能更大。为此, C4.5使用信息增益率这一准则来衡量非纯度, 即:

$$GR(D,A) = \frac{G(D,A)}{SI(D,A)}$$
 (7)

式中, SI(D, A)表示分裂信息值, 它的定义为(实际就分类熵):

$$SI(D,A) = -\sum_{i=1}^{n} \frac{N_i}{N} \log_2 \frac{N_i}{N}$$
 (8)

该式中的符号含义与式6相同。同样的,我们选择信息增益率最大的那个特征属性作为分类属性。

CART

CART算法是由Breiman等人首先提出,它包括分类树和回归树两种。我们先来讨论分类树,针对特征属性 A, 分类后的基尼指数为:

$$Gini_{sp}(D, A) = \sum_{i=1}^{n} \frac{N_i}{N} Gini(D_i)$$
 (9)

该式中的符号含义与式6相同。与ID3和C4.5不同,我们选择分类基尼指数最小的那个特征属性作为分类属 性。当我们每次只想把样本集合分为两类时,即每个中间节点只产生两个分支,但如果特征属性A中有多于 **2**个的值,即n>2,这时我们就需要一个阈值 β ,它把D分割成了 D_1 和 D_2 两个部分,不同的 β 得到不同的 D_1 和 D_2 ,我们重新设 D_1 的样本数为L, D_2 的样本数为R,因此有L+R=N,则式9可简写为:

$$\operatorname{Gini}_{\operatorname{sp}}(D, A^{(\beta)}) = \frac{L}{N} \operatorname{Gini}(D_1) + \frac{R}{N} \operatorname{Gini}(D_2) \quad (10)$$

我们把式5带入上式中,得到

CVBS视频信号解析 (http://blog.csdn.net/ app_12062011/article/details/19475741)

目标跟踪算法的分类(一) (http://blog.c sdn.net/app 12062011/article/details/484 36959)

33836

DSP之外部设备连接接口之EMIF (http:// blog.csdn.net/app_12062011/article/detai Is/7869589)

31400

50Hz工频干扰消除 (http://blog.csdn.net/ app_12062011/article/details/7770294) **21873**

DSP之中断总结篇 (http://blog.csdn.net/a pp 12062011/article/details/7864377) **19390**

$$\begin{aligned}
\operatorname{Gini}_{\mathrm{sp}}(D, A^{(\beta)}) &= \frac{L}{N} \left(1 - \frac{\sum_{j=1}^{J} L_{j}^{2}}{L^{2}} \right) + \frac{R}{N} \left(1 - \frac{\sum_{j=1}^{J} R_{j}^{2}}{R^{2}} \right) \\
&= \frac{1}{N} \left(L - \frac{\sum_{j=1}^{J} L_{j}^{2}}{L} + R - \frac{\sum_{j=1}^{J} R_{j}^{2}}{R} \right) \\
&\stackrel{1}{=} = 1 - \frac{1}{N} \left(\frac{\sum_{j=1}^{J} L_{j}^{2}}{L} + \frac{\sum_{j=1}^{J} R_{j}^{2}}{R} \right)
\end{aligned} \tag{11}$$

式中, $\sum_{i,j=L} \sum_{R_j=R} \infty$ 式11只是通过不同特征属性A的不同阈值 β 来得到样本集D的不纯度,由于D内的样本数量N是一定的,因此对式11求最小值问题就转换为求式12的最大值问题:

$$SG(D, A^{(\beta)}) = \frac{\sum_{j=1}^{J} L_j^2}{L} + \frac{\sum_{j=1}^{J} R_j^2}{R}$$
(12)

以上给出的是分类树的计算方法,下面介绍回归树。两者的不同之处是,分类树的样本输出(即响应值)是类的形式,如判断蘑菇是有毒还是无毒,周末去看电影还是不去。而回归树的样本输出是数值的形式,比如给某人发放房屋贷款的数额就是具体的数值,可以是0到120万元之间的任意值。为了得到回归树,我们就需要把适合分类的非纯度度量用适合回归的非纯度度量取代。因此我们将熵计算用均方误差替代:

$$LS(D) = \frac{1}{N} \sum_{i=1}^{N} (y_i - r_i)^2$$
 (13)

式中N表示D集合的样本数量, y_i 和 r_i 分别为第i个样本的输出值和预测值。如果我们把样本的预测值用样本输出值的平均来替代,则式**13**改写为:

$$LS(D) = \frac{1}{N} \sum_{i=1}^{N} \left(y_i - \frac{\sum_{i=1}^{N} y_i}{N} \right)^2$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left[y_i^2 - \frac{2y_i \sum_{i=1}^{N} y_i}{N} + \frac{\left(\sum_{i=1}^{N} y_i\right)^2}{N^2} \right]$$

$$= \frac{1}{N} \left[\sum_{i=1}^{N} y_i^2 - \frac{2\left(\sum_{i=1}^{N} y_i\right)^2}{N} + \frac{\left(\sum_{i=1}^{N} y_i\right)^2}{N} \right]$$

$$= \frac{1}{N} \left[\sum_{i=1}^{N} y_i^2 - \frac{\left(\sum_{i=1}^{N} y_i\right)^2}{N} \right]$$
(14)

上式表示了集合D的最小均方误差,如果针对于某种特征属性A,我们把集合D划分为s个部分,则划分后的均方误差为:

$$LS(D,A) = \sum_{i=1}^{3} \frac{N_i}{N} LS(D_i)$$
 (15)

式中 N_i 表示被划分的第i个集合 D_i 的样本数量。式15与式14的差值为划分为s个部分后的误差减小:

$$\Delta LS(D, A) = LS(D) - \sum_{i=1}^{s} \frac{N_i}{N} LS(D_i)$$
 (16)

与式6所表示的信息增益相似,我们寻求的是最大化的误差减小,此时就得到了最佳的s个部分的划分。同样的,当我们仅考虑二叉树的情况时,即每个中间节点只有两个分支,此时s=2,基于特征属性A的值,集合D被阈值 β 划分为 D_1 和 D_2 两个集合,每个集合的样本数分别为L和R,则:

$$\Delta LS(D, A^{(\beta)}) = LS(D) - \frac{L}{N}LS(D_1) - \frac{R}{N}LS(D_2)$$
 (17)

把式14带入上式, 得:

式中, y_i 是属于集合D的样本响应值, I_i 和 I_i 分别是属于集合 D_1 和 D_2 的样本响应值。对于某个节点来说,它的 样本数量以及样本响应值的和是一个定值,因此式18的结果完全取决于方括号内的部分,即:

$$\Delta LLS\left(\Sigma; A^{(\beta)}\right) = \frac{\left(\sum_{i=1}^{L} l_i\right)^2}{L} + \frac{\left(\sum_{i=1}^{R} r_i\right)^2}{R}$$
(19)

因此求式18的最大值问题就转变为求式19的最大值问题。

我们按照样本响应值是类的形式,还是数值的形式,把决策树分成了分类树和回归树,它们对应不同的计 算公式。那么表示特征属性的形式也会有这两种形式,即类的形式和数值的形式,比如决定是否出去踢 球,取决于两个条件:风力和气温。风力的表示形式是:无风、小风、中风、大风,气温的表示形式就是 具体的摄氏度,如-10 $^{\circ}$ $^{\circ}$ $^{\circ}$ $^{\circ}$ 之间。风力这个特征属性就是类的形式,而气温就是数值的形式。又比如决 定发放房屋贷款,其金额取决于两个条件:是否有车有房和年薪。有车有房的表示形式是:无车无房、有 车无房、无车有房、有车有房,而年薪的表示形式就是具体的钱数,如0~20万。有车有房这个特征属性就 是类的形式,年薪就是数值的形式。因此在分析样本的特征属性时,我们要把决策树分为四种情况:特征 为类的分类树(如决定是否踢球的风力)、特征为数值的分类树(如决定是否踢球的温度)、特征为类的 回归树(如发放贷款的有车有房)和特征为数值的回归树(如发放贷款的年薪)。由于特征形式不同,所 以计算方法上有所不同:

I、特征为类的分类树:对于两类问题,即样本的分类(响应值)只有两种情况:响应值为0和1,按照特 征属性的类别的样本响应值为1的数量的多少进行排序。例如我们采集20个样本来构建是否踢球分类树、设 出去踢球的响应值为1,不踢球的响应值为0,针对风力这个特征属性,响应值为1的样本有14个,无风有6 个样本,小风有5个,中风2个,大风1个,则排序的结果为:大风<中风<小风<无风。然后我们按照这个顺 序依次按照二叉树的分叉方式把样本分为左分支和右分支,并带入式12求使该式为最大值的那个分叉方 式,即先把是大风的样本放入左分支,其余的放入右分支,带入式12,得到A,再把大风和中风放入左分 支,其余的放入右分支,带入式12,得到B,再把大风、中风和小风放入左分支,无风的放入右分支,计算 得到C, 比较 $A \times B \times C$, 如果最大值为C, 则按照C的分叉方式划分左右分支, 其中阈值 β 可以设为 $3 \cdot$ 对于 非两类问题,采用的是聚类的方法。

 Π 、特征为数值的分类树:由于特征属性是用数值进行表示,我们就按照数值的大小顺序依次带入式12. 计算最大值。如一共有14个样本,按照由小至大的顺序为: abcdefghijklmn, 第一次分叉

为: a|bcdefghijklmn, 竖线"|"的左侧被划分到左分支,右侧被划分到右分支,带入式12计算其值,然后第 二次分叉:ablcdefghijklmn,同理带入式12计算其值,以此类推,得到这13次分叉的最大值,该种分叉方 式即为最佳的分叉方式,其中阈值 β 为分叉的次数。

III、特征为类的回归树: 计算每种特征属性各个种类的平均样本响应值, 按照该值的大小进行排序, 然后 依次带入式19, 计算其最大值。

Ⅳ、特征为数值的回归树:该种情况与特征为数值的分类树相同,就按照数值的大小顺序依次带入式19, 计算最大值。

在训练决策树时,还有三个技术问题需要解决。第一个问题是,对于分类树,我们还需要考虑一种情况, 当用户想要检测一些非常罕见的异常现象的时候,这是非常难办到的,这是因为训练可能包含了比异常多 得多的正常情况,那么很可能分类结果就是认为每一个情况都是正常的。为了避免这种情况的出现,我们 需要设置先验概率,这样异常情况发生的概率就被人为的增加(可以增加到0.5甚至更高),这样被误分类 的异常情况的权重就会变大、决策树也能够得到适当的调整。先验概率需要根据各自情况人为设置、但还 需要考虑各个分类的样本率,因此这个先验值还不能直接应用,还需要处理。设Qi为我们设置的第i个分类 的先验概率, N_i 为该分类的样本数,则考虑了样本率并进行归一化处理的先验概率 q_i 为:

$$q_{j} = \frac{Q_{j}/N_{j}}{\sum_{j=1}^{J}(Q_{j}/N_{J})}$$
 (20)

把先验概率带入式12中,则得到:

$$SG(D, A^{(\beta)}, q) = \frac{\sum_{j=1}^{J} (q_j L_j)^2}{\sum_{j=1}^{J} (q_j L_j)} + \frac{\sum_{j=1}^{J} (q_j R_j)^2}{\sum_{j=1}^{J} (q_j R_j)}$$
(21)

第二个需要解决的问题是,某些样本缺失了某个特征属性,但该特征属性又是最佳分叉属性,那么如何对 该样本进行分叉呢?目前有几种方法可以解决该问题,一种是直接把该样本删除掉;另一种方法是用各种 算法估计该样本的缺失属性值。还有一种方法就是用另一个特征属性来替代最佳分叉属性,该特征属性被 称为替代分叉属性。因此在计算最佳分叉属性的同时,还要计算该特征属性的替代分叉属性,以防止最佳 分叉属性缺失的情况。CART算法就是采用的该方法,下面我们就来介绍该方法。

寻找替代分叉属性总的原则就是使其分叉的效果与最佳分叉属性相似,即分叉的误差最小。我们仍然根据 特征属性是类还是数值的形式,也把替代分叉属性的计算分为两种情况。

当特征属性是类的形式的时候,当最佳分叉属性不是该特征属性时,会把该特征属性的每个种类分叉为不 同的分支,例如当最佳分叉属性不是风力时,极有可能把5个无风的样本分叉为不同的分支(如3个属于左 分支, 2个属于右分支), 但当最佳分叉属性是风力时, 这种情况就不会发生, 也就是5个无风的样本要么 属于左分支。要么属于右分支。因此我们把被最佳分叉属性分叉的特征属性种类的分支最大样本数量作为 该种类的分叉值,计算该特征属性所有种类的这些分叉值,最终这些分叉值之和就作为该替代分叉属性的 分叉值。我们还看前面的例子, 无风的分叉值为3, 再计算小风、中风、大风的分叉值, 假如它们的值分别 为4、4、39则把风力作为替代分叉属性的分叉值为14。按照该方法再计算其他特征属性是类形式的替代分 叉值,则替代性由替代分叉值按从大到小进行排序。在用替代分叉属性分叉时那些左分支大于右分支样本 数的种类被分叉为左分支,反之为右分支,如上面的例子,无风的样本被分叉为左分支。

当特征属性是数值的形式的时候, 样本被分割成了四个部分: LL、LR、RL和RR, 前一个字母表示被最佳 分叉属性分叉为左右分支,后一个字母表示被替代分叉属性分叉为左右分支,如LR表示被最佳分叉属性分 叉为左分支,但被替代分叉属性分叉为右分支的样本,因此LL和RR表示的是被替代分叉属性分叉正确的样 本,而LR和RL是被替代分叉属性分叉错误的样本,在该特征属性下,选取阈值对样本进行分割,使LL+RR 或LR+RL达到最大值,则最终max{LL+RR, LR+RL}作为该特征属性的替代分叉属性的分叉值。按照该方 法再计算其他特征属性是数值形式的替代分叉值,则替代性也由替代分叉值按从大到小进行排序。最终我 们选取替代分叉值最大的那个特征属性作为该最佳分叉属性的替代分叉属性。

为了让替代分叉属性与最佳分叉属性相比较,我们还需要对替代分叉值进行规范化处理,如果替代分叉属 性是类的形式,则替代分叉值需要乘以式12再除以最佳分叉属性中的种类数量,如果替代分叉属性是数值 的形式,则替代分叉值需要乘以式19再除以所有样本的数量。规范化后的替代分叉属性如果就是最佳分叉 属性时,两者的值是相等的。

第三个问题就是过拟合。由于决策树的建立完全是依赖于训练样本,因此该决策树对该样本能够产生完全 一致的拟合效果。但这样的决策树对于预测样本来说过于复杂,对预测样本的分类效果也不够精确。这种 现象就称为过拟合。

将复杂的决策树进行简化的过程称为剪枝、它的目的是去掉一些节点、包括叶节点和中间节点。剪枝常用 的方法有预剪枝和后剪枝两种。预剪枝是在构建决策树的过程中,提前终止决策树的生长,从而避免过多 的节点的产生。该方法虽然简单但实用性不强,因为我们很难精确的判断何时终止树的生长。后剪枝就是 在决策树构建完后再去掉一些节点。常见后剪枝方法有四种: 悲观错误剪枝(PEP)、最小错误剪枝

(MEP) 、代价复杂度剪枝 (CCP) 和基于错误的剪枝 (EBP) 。CCP算法能够应用于CART算法中,它 的本质是度量每减少一个叶节点所得到的平均错误,在这里我们重点介绍CCP算法。

CCP算法会产生一系列树的序列 $\{T_0,T_1,...,T_m\}$,其中 T_0 是由训练得到的最初的决策树,而 T_m 只含有一个根 节点。序列中的树是嵌套的,也就是序列中的 T_{tr1} 是由 T_t 通过剪枝得到的,即实现用 T_{tr1} 中一个叶节点来替代 T_i 中以该节点为根的子树。这种被替代的原则就是使误差的增加率 α 最小,即

$$\alpha = \frac{R(n) - R(n_t)}{|n_t| - 1}$$
 (22)

式中,R(n)表示 T_i 中节点n的预测误差, $R(n_i)$ 表示 T_i 中以节点n为根节点的子树的所有叶节点的预测误差之 和, $|n_1|$ 为该子树叶节点的数量, $|n_1|$ 也被称为复杂度,因为叶节点越多,复杂性当然就越强。因此 α 的含义 就是用一个节点n来替代以n为根节点的所有 $|n_i|$ 个节点的误差增加的规范化程度。在 T_i 中,我们选择最小的 α 值的节点进行替代,最终得到 T_{H1} 。以此类推,每需要得到一棵决策树,都需要计算其前一棵决策树的 α 值,根据 α 值来对前一棵决策树进行剪枝,直到最终剪枝到只剩下含有一个根节点的 T_m 为止。 根据决策树是分类树还是回归树,节点的预测误差的计算也分为两种情况。在分类树下,我们可以应用上 面介绍过的式1~式3中的任意一个,如果我们应用式3来表示节点n的预测误差,则:

$$R(n) = \frac{\sum_{j=1}^{m} N_j - \max{\{N_j\}}}{\sum_{i=1}^{m} N_i} = \frac{N - \max{\{N_j\}}}{N} \tag{23}$$

式中, N_i 表示节点n下第j个分类的样本数,N为该节点的所有样本数, $\max\{N_i\}$ 表示在m个分类中,拥有样本

数最多的那个分类的样本数量。在回归树下,我们可以应用式14来表示节点n的预测误差:

$$R(n) = \frac{\sum_{i=1}^{N} (y_i^2) - \frac{(\sum_{i=1}^{N} y_i)^2}{N}}{N} \tag{24}$$

式中,y表示第i个样本的响应值,N为该节点的样本数量。我们把式23和式24的分子部分称为节点的风险

我们用全部样本得到的决策树序列为 $\{T_0,T_1,...,T_m\}$,其所对应的 α 值为 $\alpha_0<\alpha_1<...<\alpha_m$ 。下一步就是如何从这 个序列中最优的选择一颗决策树 T_i 。而与其说找到最优的 T_i ,不如说找到其所对应的 α_i 。这一步骤通常采用 的方法是交叉验证(Cross-Validation)。

我们把 L^{∞} 本随机划分为数量相等的V个子集 L_{v} , v=1,...,V。第v个训练样本集为:

$$L^{(v)} = L - L_v$$
 $v = 1, ..., V$ (25)

则 L_v 被用来做 $L^{(v)}$ 的测试样本集。对每个训练样本集 $L^{(v)}$ 按照CCP算法得到决策树的序列 $\{T_0^{\omega}, T_1^{\omega}, ..., T_m^{\omega}\}$, 其对应的 α 值为 α_0 $^{(n)}<\alpha_1$ $^{(n)}<...<\alpha_m$ $^{(n)}$ 。 α 值的计算仍然采用式22。对于分类树来说,第 ν 个子集的节点n的预测

$$R^{(v)}(n) = \frac{\sum_{j=1}^{m} N_j^{(v)} - \max\{N_j^{(v)}\}}{\sum_{j=1}^{m} N_j^{(v)}} = \frac{N^{(v)} - \max\{N_j^{(v)}\}}{N^{(v)}}$$
(26)

式中, $N_i^{(\prime)}$ 表示训练样本集 $L^{(\prime)}$ 中节点n的第i个分类的样本数, $N^{(\prime)}$ 为 $L^{(\prime)}$ 中节点n的所有样本数, $\max\{N_i^{(\prime)}\}$ 表示在m个分类中, $L^{(v)}$ 中节点n拥有样本数最多的那个分类的样本数量。对于回归树来说,第v个子集的节 点**n**的预测误差为:

$$R^{(v)}(n) = \frac{\sum_{i=1}^{N^{(v)}} ((y_i^{(v)})^2) - \frac{(\sum_{i=1}^{N^{(v)}} y_i^{(v)})^2}{N^{(v)}}}{N^{(v)}}$$
(27)

式中, $y_i^{(\prime)}$ 表示训练样本集 $L^{(\prime)}$ 中节点n的第i个样本的响应值。我们仍然把式26和式27的分子部分称交叉验 证子集中的节点风险值。

我们由训练样本集得到了树序列,然后应用这些树对测试样本集进行测试,测试的结果用错误率来衡量, 即被错误分类的样本数量。对于分类树来说,节点n的错误率为:

$$e^{(v)}(n) = N_v - N_{v,i}$$
 (28)

式中, N_{v} 表示测试样本集 L_{v} 中节点n的所有样本数, $N_{v,i}$ 表示 L_{v} 中第j个分类的样本数,这个j是式26中 $\max\{|L_i^{(v)}|\}$ 所对应的j。对于回归树来说,节点n的错误率为:

$$e^{(v)}(n) = \sum_{i=1}^{N_v} (y_{v,i}^2) - 2 \left(\frac{\sum_{i=1}^{N^{(v)}} y_i^{(v)}}{N^{(v)}} \right) \left(\sum_{i=1}^{N_v} y_{v,i} \right) + N_v \left(\frac{\sum_{i=1}^{N^{(v)}} y_i^{(v)}}{N^{(v)}} \right)^2$$
(29)

式中, $y_{\nu i}$ 表示 L_{ν} 的第i个样本响应值。决策树的总错误率 $E^{(\nu)}$ 等于该树所有叶节点的错误率之和。 虽然交叉验证子集决策树序列T(V)的数量要与由全体样本得到的决策树序列T的数量相同,但两者构建的形 式不同,它需要比较两者的 α 值后再来构建。而为了比较由全体样本训练得到 α 值与交叉验证子集的 α $^{(\prime)}$ 值之 间的大小, 我们还需要对 α 值进行处理, 即

$$\acute{\alpha}_i = \sqrt{\alpha_i \alpha_{i+1}} \qquad (30)$$

其中 $\alpha'_0 = 0$,而 α'_m 为无穷大。

我们设按照式22得到的初始训练子集 $L^{(\prime\prime)}$ 决策树序列为 $\{T_0^{(\prime\prime)},T_1^{(\prime\prime)},...,T_m^{(\prime\prime)}\}$, 其所对应的 $\alpha^{(\prime\prime)}$ 值为 $\{\alpha_0^{(\prime\prime)},\alpha_1^{(\prime\prime)},...,$ $\alpha_m^{(v)}$ 。最终的树序列也是由这些 $T^{(v)}$ 组成,并且也是嵌套的形式,但可以重复,而且必须满足:

$$\dot{\alpha}_k < \max \left\{ \alpha_j^{(v)} \right\}$$
(31)

该式的含义是 $T^{(v)}$ 中第k个子树的 $\alpha^{(v)}$ 值要小于 α^{\prime}_k 的最大的 $\alpha^{(v)}$ 所对应的子树,因此最终的树序列有可能是这 种形式: $T_0(0,T_1(0,T_1(0,T_2(0,T_2(0,T_2(0,T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0,T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T_2(0),T$

子集的决策树序列构建好了,下面我们就可以计算V个子集中树序列相同的总错误率之和,即

$$E_j = \sum_{v=1}^{V} E_j^{(v)}$$
 (32)

则最佳的子树索引**J**为:

$$J = \arg\min_{j} (E_j) \quad (33)$$

最终我们选择决策树序列 $\{T_0,T_1,...,T_m\}$ 中第J棵树为最佳决策树,该树的总错误率最小。

如果我们在选择决策树时使用1-SE(1 Standard Error of Minimum Error)规则的话,那么有可能最终的决 策树不是错误率最小的子树,而是错误率略大,但树的结构更简单的那颗决策树。我们首先计算误差范围 SF:

$$SE = \sqrt{E_J(N - E_J)} \tag{34}$$

式中, $E_{\mathbf{p}}$ 模示最佳子树的总错误率,N为总的样本数。则最终被选中的决策树的总错误率 $E_{\mathbf{K}}$ 要满足:

$$E_K < E_I + SE \tag{35}$$

并且决策村的结构最简单。

以上我们完整并详细的介绍了构建决策树,也就是训练决策树的过程,在这个过程中我们没有放过每个技 术细节。」而预测样本很简单,只要把样本按照最佳分叉属性(或替代分叉属性)一次一次分叉,直到到达 决策树的个节点,该叶节点的值就是该预测样本的响应值。

OpenCV中,决策树应用的是CART算法,并且分类树和回归树都实现了。而剪枝是应用的CCP算法。

我们先给出用于构建决策树的参数:

```
[cpp]
 1.
      CvDTreeParams::CvDTreeParams(): max_categories(10), max_depth(INT_MAX), min_sample_count(10),
2.
          cv_folds(10), use_surrogates(true), use_1se_rule(true),
3.
          truncate\_pruned\_tree(\textbf{true}), \ regression\_accuracy(0.01f), \ priors(0)
 4.
      {}
5.
 6.
     CvDTreeParams::CvDTreeParams( int _max_depth, int _min_sample_count,
                                     float _regression_accuracy, bool _use_surrogates,
 7.
8.
                                     int _max_categories, int _cv_folds,
 9.
                                     bool _use_1se_rule, bool _truncate_pruned_tree,
                                     const float* _priors ) :
10.
11.
          max_categories(_max_categories), max_depth(_max_depth),
12.
         min sample count( min sample count), cv folds ( cv folds).
13.
          use_surrogates(_use_surrogates), use_1se_rule(_use_1se_rule),
14
          truncate_pruned_tree(_truncate_pruned_tree),
          regression_accuracy(_regression_accuracy),
15.
16
          priors(_priors)
17. {}
```

其中参数的含义为:

max depth: 该参数指定决策树的最大可能深度。但由于其他终止条件或者是被剪枝的缘故, 最终的决策 树的深度可能要比max depth小。该值不能小于0,否则报错。如果用户把该值设置为大于25的数值,则系 统会重新把该值设为25, 即该值不能超过25

min sample count: 如果某个节点的样本数量小于该值,则该节点将不再被分叉

regression_accuracy: 该参数是终止构建回归树的一个条件,表示回归树的响应值的精度如果达到了该 值,则无需再分叉。该值不能小于0,否则报错

use_surrogates:如果该参数为真,则替代分叉节点需要产生,替代分叉节点用在解决缺失特征属性和进 行变量重要性的估计

max_categories:表示特征属性为类的形式的最大的类的数量,如果在训练的过程中,类的数量大于该 值,则采用聚类的方法会更高效,该参数只应用于非两类的分类树问题。该值一定要大于或等于2,否则报

cv folds: 如果该参数大于1,则应该交叉验证来进行剪枝,该值表示交叉验证的子集数量。该值一定不能 小于0,否则报错。如果用户把该值设置为1,则系统会重新把该值设为0

use 1se rule: 如果为真,表示在剪枝的过程中应用1SE规则,该规则的应用使决策树更简单,对训练数 据的噪声更有抵抗力, 但是精度会下降了

truncate pruned tree: 如果为真,要被剪掉的节点将被从树上移除,否则它们被保留

priors:设置决策树的特征属性的先验概率,该变量为矩阵的形式,并且与特征属性的排序相同

CvDTree构造函数:

```
[cpp]
  1.
      CvDTree::CvDTree()
  2.
         data = 0;
                  //样本数据
  3.
  4.
         var_importance = 0;
                           //重要数据
         default_model_name = "my_tree";
                                     //表明该类型为决策树
  5.
  6.
  7.
     clear();
                  //清空一些变量
  8.
训练样本! 构建决策树函数CvDTree::train:
      [qqɔ]
      bol CvDTree::train( const CvMat* _train_data, int _tflag,
  1.
                       const CvMat* _responses, const CvMat* _var_idx,
                       const CvMat* _sample_idx, const CvMat* _var_type,
const CvMat* _missing_mask, CvDTreeParams _params )
      <u>...</u>
  3.
  4.
      //_train_data训练样本数据,必须为32FC1类型的矩阵形式
  5.
      flag训练数据的特征属性类型,如果为CV_ROW_SAMPLE,表示样本是以行的形式储存的,即_train_data矩阵的每一行
  6.
      为一个样本(或特征矢量);如果为CV_COL_SAMPLE,表示样本是以列的形式储存的
      //_responses决策树的结果,即响应值,该值必须是32SC1或32FC1类型的一维矩阵(即矢量)的形式,并且元素的数量必须
  7.
      与训练样本数据_train_data的样本数一致
      //_var_idx标识感兴趣的特征属性,即真正用于训练的那些特征属性,该值的形式与_sample_idx变量相似
  8.
      // sample idx标识感兴趣的样本,即真正用于训练的样本,该值必须是一维矩阵的形式,即矢量的形式,并且类型必须是
      8UC1、8SU1或者32SC1。如果为8UC1或8SU1,则该值的含义是用掩码的形式表示对应的样本,即0表示不感兴趣的样本,其他
      数为感兴趣的样本,因此矢量的元素数量必须与训练样本数据 train data的样本数一致;如果为32SC1,则该值的含义是那些
      感兴趣的样本的索引,而不感兴趣的样本的索引不在该矢量中出现,因此该矢量的元素数量可以小于或等于_train_data的样本
      //_var_type特征属性的形式,是类的形式还是数值的形式,用掩码的形式表现对应特征属性的形式,0表示为数值的形式,1
 10.
      表示为类的形式。该值必须是一维矩阵,并且元素的数量必须是真正用于训练的那些特征属性的数量加1,多余的一个元素表示的
      是响应值的形式,即是分类树还是回归树
 11.
      //_missing_mask缺失的特征属性,用掩码的形式表现对应的特征属性,0表示没有缺失,而且必须与_train_data的矩阵尺
      寸大小一致
 12.
      //_params为构建决策树的参数
 13.
         bool result = false:
                             //标识变量
 14.
 15
         CV_FUNCNAME( "CvDTree::train" );
 16.
 17.
 18
         __BEGIN__;
 19.
 20
         clear();
                  //清空一些变量
```

下面的train函数与上面的train函数不同的地方就是矩阵的类型不同,一个是Mat,另一个是CvMat。

还要随机的把样本均等分 cv_folds 个子集。最后一个参数 $_share$ 赋值为false,表示重新构建一个决策树

data = new CvDTreeTrainData(_train_data, _tflag, _responses,

//调用do_train函数,对全体样本数据进行训练

CV_CALL(result = do_train(0));

END :

return result;

21.

22.

23.

24

25. 26.

27 28.

29. 30.

31. }

```
[qqɔ]
1.
     bool CvDTree::train( const Mat& _train_data, int _tflag,
 2.
                         const Mat& _responses, const Mat& _var_idx,
                         const Mat& sample idx, const Mat& var type,
3.
 4.
                         const Mat& _missing_mask, CvDTreeParams _params )
 5.
     {
 6.
         CvMat tdata = _train_data, responses = _responses, vidx=_var_idx,
             sidx=_sample_idx, vtype=_var_type, mmask=_missing_mask;
8.
         return train(&tdata, _tflag, &responses, vidx.data.ptr ? &vidx : 0, sidx.data.ptr ? &sidx
9.
                      vtype.data.ptr ? &vtype : 0, mmask.data.ptr ? &mmask : 0, _params);
10. }
```

//实例化CvDTreeTrainData类,并通过set_data函数设置训练样本数据,如果需要交叉验证(cv_folds > 1),则

_var_idx, _sample_idx, _var_type,

_missing_mask, _params, false);

第8页 共46页

在该train函数中,用于决策树的样本数据的一些参数保存在了 data中:

```
[cpp]
1.
    bool CvDTree::train( CvMLData* _data, CvDTreeParams _params )
2.
3.
       bool result = false;
    4.
5.
6.
7.
         BEGIN ;
     8.
    const CvMat* values = _data->get_values();
9.
        const CvMat* response = _data->get_responses();
10.
    const CvMat* missing = _data->get_missing();
const CvMat* var_types = _data->get_var_types();
11.
12.
13.
       const CvMat* train_sidx = _data->get_train_sample_idx();
    const CvMat* var_idx = _data->get_var_idx();
14.
15.
    16.
17.
18.
19.
        __END__;
20.
21.
        return result;
22. }
```

该train函数多了一个_subsample_idx参数

```
[cpp]
      \textbf{bool} \ \texttt{CvDTree}{::} \texttt{train(} \ \texttt{CvDTreeTrainData*} \ \_\texttt{data,} \ \textbf{const} \ \texttt{CvMat*} \ \_\texttt{subsample\_idx} \ )
1.
           bool result = false:
3.
 4.
           CV_FUNCNAME( "CvDTree::train" );
5.
 6.
 7.
           __BEGIN__;
8.
9.
           clear();
10.
           data = _data;
           data->shared = true;
                                      //表示共享决策树
11.
           //调用do_train函数,只对_subsample_idx指定的由_sample_idx参数确定的训练样本集的索引进行训练
12.
13.
           CV_CALL( result = do_train(_subsample_idx));
14.
           __END__;
15.
16.
           return result;
17.
18. }
```

do train函数,参数 subsample idx表示训练样本集合的索引

```
[cpp]
     bool CvDTree::do_train( const CvMat* _subsample_idx )
1.
2.
         bool result = false;
3.
4.
5.
        CV_FUNCNAME( "CvDTree::do_train" );
6.
         __BEGIN__;
        //提取出_subsample_idx指定的训练样本集合,如果_subsample_idx=0,则取所有样本数据
8.
9.
         //root表示根节点
10.
         root = data->subsample_data( _subsample_idx );
11.
12
        CV_CALL( try_split_node(root)); //递归调用try_split_node函数,构造决策树
13.
        if(root->split) //如果得到了决策树
14.
15.
            //分别确保左分支和右分支的正确性
16.
17.
            CV_Assert( root->left );
```

第9页 共46页

```
18.
           CV_Assert( root->right );
           //如果用户所设置的交叉验证的子集大于1,则表示需要对决策树进行剪枝。在前面已经提到,在实例化
19.
     CvDTreeTrainData时,会把cv_folds = 1重新设置为cv_folds = 0
20.
           if( data->params.cv_folds > 0 )
21.
               CV_CALL( prune_cv() ); //剪枝
            //如果决策树的数据不需要共享,则清空训练所需的数据
22.
23.
            if( !data->shared )
24.
               data->free_train_data();
25.
26.
            result = true; //正确返回标识
27.
28.
29.
        __END__;
30.
31.
        return result;
                        //返回
```

 try_split_node 函数的作用是在真正分叉之前,先预处理一下节点,试探一下该节点是否能够被分叉,能则 调用split wode data函数开始分叉,不能则退出,然后对得到的左分支和右分支再分别递归。



()

```
[cpp]
1.
     void CvDTree::try_split_node( CvDTreeNode* node )
2.
        CvDTreeSplit* best_split = 0;
3.
        //n表示该节点的样本数, vi表示特征属性
4.
        int i, n = node->sample_count, vi;
5.
6.
        bool can_split = true;
                               //标识该节点是否能够被分叉
        double quality_scale;
        // calc_node_value函数见后面的分析
8.
9
        calc_node_value( node );
        //如果该节点的样本数量小于所设定的阈值,或者目前决策树的深度超过了所设定的阈值,则该节点不能再被分叉
10.
11.
        if( node->sample_count <= data->params.min_sample_count ||
12.
            node->depth >= data->params.max_depth )
            can split = false:
                               //设置标识变量
13.
14
15.
        if( can split && data->is classifier )
                                             //分类树
16.
17.
            // check if we have a "pure" node,
            // we assume that cls_count is filled by calc_node_value()
18.
19.
            // cls_count数组存储着每个分类的样本数量,它是由calc_node_value()函数得到的
            int* cls count = data->counts->data.i;
20.
                                                  //m表示有多少个分类
21.
            int nz = 0, m = data->get_num_classes();
22.
            //遍历所有的分类,计算样本数量不为零的分类的个数
23.
            for( i = 0; i < m; i++ )
24.
               nz += cls_count[i] != 0;
25.
            //如果nz等于1,说明只有一个分类,因此该节点无需再分
26
            if( nz == 1 ) // there is only one class
               can_split = false; //设置标识变量
27.
28.
29
        else if( can_split ) //回归树
30.
            // node_risk,即节点的风险值由calc_node_value()函数得到的,对node_risk开根号并取平均,表示平均误
31.
     差量,即精度,其结果与所设定的阈值比较
            if( sqrt(node->node_risk)/n < data->params.regression_accuracy )
32.
33.
               can_split = false; //设置标识变量
34.
35.
        if(can_split) //如果该节点可以被分叉,则找到最佳的分类属性
37.
38.
            // find_best_split函数见后面的分析
39.
            best_split = find_best_split(node);
                                              //得到最佳分叉属性
40.
            // TODO: check the split quality ...
            //通过上面的注释可以看出,OpenCV的后续版本可能还会对这个最佳分类属性进行进一步的检测
41.
            node->split = best split;
                                     //赋值
42.
43.
        //该节点不能被分叉,或没有找到分类属性,也就是达到了也节点,则退出该函数,结束该分支的递归
44.
45.
        if( !can_split || !best_split )
46.
47.
            data->free node data(node);
                                     //清空该节点数据
48.
            return;
```

2017/11/30 上午10:51 第10页 共46页

```
49.
        // 用最佳分叉属性完成对节点的分叉,即标注各个样本的方向信息,calc node dir函数在后面给出了讲解,其中该函
50.
     数的返回值quality_scale为替代分叉属性值的规范化处理常数
        quality_scale = calc_node_dir( node );
51.
52.
        if( data->params.use_surrogates )
                                          //表示使用替代分叉属性
53
54.
            // find all the surrogate splits
55
            \ensuremath{//} and sort them by their similarity to the primary one
56.
            //遍历所有的特征属性,计算替代分叉属性
     \int_{\Gamma}
57
            for( vi = 0; vi < data->var_count; vi++ )
58
     1
                CvDTreeSplit* split;
59.
60
                //ci大于0表示特征属性是类的形式,ci小于0表示特征属性是数值的形式
                int ci = data->get_var_type(vi);
61.
62.
                //如果当前的特征属性为最佳分叉属性,则进行下一个vi的循环
63
                if( vi == best_split->var_idx )
                   continue:
64.
65
     [•••
66.
                if( ci >= 0 )
                              //特征属性是类的形式
67.
                   // find_surrogate_split_cat函数在后面给出详细的介绍
                   split = find_surrogate_split_cat( node, vi );
     &
69.
                       //特征属性是数值的形式
70.
                   // find_surrogate_split_ord函数在后面给出详细的介绍
                   split = find_surrogate_split_ord( node, vi );
71.
72.
                if( split )
73.
                //得到了替代分叉属性,对现有的替代分叉属性按从大到小的顺序进行排序
74.
75.
76.
                   // insert the split
                   //把替代分叉属性的分叉值保存在node->split->next内
77.
                   CvDTreeSplit* prev_split = node->split;
78
79.
                   //对替代分叉属性的分叉值进行规范化处理
80
                   split->quality = (float)(split->quality*quality_scale);
                   //排序
81.
82.
                   while( prev_split->next &&
83.
                         prev_split->next->quality > split->quality )
                       prev split = prev split->next;
84.
85.
                   split->next = prev_split->next;
86.
                   prev_split->next = split;
87.
               }
88
            }
89.
90
        //完成节点的分叉,得到左分支和右分支,并为下次分叉赋值给各个变量
91.
        split_node_data( node );
        //分别对左分支和右分支递归调用try_split_node函数
92.
93.
        try_split_node( node->left );
        try_split_node( node->right );
94.
95. }
```

计算节点的先验概率、节点的值、节点的风险值和节点的错误率:

()

```
1.
                   void CvDTree::calc_node_value( CvDTreeNode* node )
  2.
   3.
                                //n为该节点的样本数, cv_n为交叉验证的子集数量
   4.
                                int i, j, k, n = node->sample_count, cv_n = data->params.cv_folds;
                                                                                                                                                    //分类树的分类数量,即它的响应值的数量
   5.
                                int m = data->get_num_classes();
                                //为缓存开辟空间,base_size为基本大小,ext_size为扩增大小
   6.
                                //根据是分类树还是回归树设置不同的大小
   7.
   8.
                                int base_size = data->is_classifier ? m*cv_n*sizeof(int) : 2*cv_n*sizeof(double)+cv_n*Siz
                                int \ ext\_size = n*(sizeof(int) + (data->is\_classifier ? sizeof(int) : sizeof(int) + sizeof(int) +
                                //开辟空间
10.
11.
                                cv::AutoBuffer<uchar> inn_buf(base_size + ext_size);
                                uchar* base_buf = (uchar*)inn_buf;
12.
                                                                                                                                                       //基本空间
13.
                                uchar* ext_buf = base_buf + base_size; //扩展空间
14.
                                int* cv_labels_buf = (int*)ext_buf;
15.
16.
                                //得到样本所对应的交叉验证的子集
                                const int* cv_labels = data->get_cv_labels(node, cv_labels_buf);
17.
18.
```

第11页 共46页 2017/11/30 上午10:51

```
19.
         if( data->is_classifier )
20.
21.
             // in case of classification tree:
             // * node value is the label of the class that has the largest weight in the node.
22.
23.
             // \ ^{\star} node risk is the weighted number of misclassified samples,
24
             //\ ^{*} j-th cross-validation fold value and risk are calculated as above,
25.
                  but using the samples with cv labels(*)!=j. 不包括该子集
             //\ ^{*} j-th cross-validation fold error is calculated as the weighted number of
26
27.
                  misclassified samples with cv_labels(*)==j.
                                                                仅是该子集
     \int_{\Gamma}
28
29.
             // compute the number of instances of each class
      1
30.
             //cls count数组用于表示每个分类中的样本数量
31.
             int* cls_count = data->counts->data.i;
32.
             // responses_buf指向的空间为样本数据的响应,即分类结果
33.
             int* responses_buf = cv_labels_buf + n;
34
             //得到样本的响应值,也就是样本的分类结果
     h
35.
             const int* responses = data->get class labels(node, responses buf);
36
             // cv_cls_count数组用于表示交叉验证时的每个子集的每类样本的数量
     <u>...</u>
37.
             int* cv_cls_count = (int*)base_buf;
38.
             double max_val = -1, total_weight = 0;
39
             int \max_k = -1;
      &
40.
             // priors指向存储着先验概率
41.
             double* priors = data->priors_mult->data.db;
             // cls_count数组清零
42.
43.
             for( k = 0; k < m; k++)
44.
                 cls_count[k] = 0;
45.
46.
             if( cv_n == 0 ) //不需要交叉验证
47.
                 //遍历所有训练样本,得到不同分类结果的样本数量
48.
                 for( i = 0; i < n; i++ )</pre>
49
50.
                    cls_count[responses[i]]++;
51.
                    //需要交叉验证
             else
52.
53.
54
                 // cv_cls_count数组清零
                 for( j = 0; j < cv_n; j++ )</pre>
55.
56.
                    for( k = 0; k < m; k++ )
                        cv_cls_count[j*m + k] = 0;
57.
58.
                 for( i = 0; i < n; i++ )</pre>
59.
60.
                 {
61.
                    //j表示第i个样本所在的交叉验证的子集,k表示该样本的分类结果
                    j = cv_labels[i]; k = responses[i];
62.
                    //对不同交叉验证子集的不同分类结果进行计数
63.
64
                    cv_cls_count[j*m + k]++;
65.
                }
66.
                 //统计所有样本的不同分类的数量
67.
                 for( j = 0; j < cv_n; j++ )</pre>
                    for( k = 0; k < m; k++ )
68.
69
                        cls_count[k] += cv_cls_count[j*m + k];
70.
             //如果有先验概率,并且该节点是根节点,计算先验概率
71.
72.
             if( data->have_priors && node->parent == 0 )
73.
74.
                 // compute priors_mult from priors, take the sample ratio into account.
                 //计算先验概率
75.
76
                 double sum = 0;
77.
                 for( k = 0; k < m; k++ )
78.
                 {
79
                    int n_k = cls_count[k];
                                             //得到第k个分类的样本数量,即式20的Nj
                    // data->priors->data.db[k]为式20的Qj,priors[k]为式20的Qj/Nj
80.
81.
                    priors[k] = data->priors->data.db[k]*(n_k ? 1./n_k : 0.);
                    sum += priors[k];
                                       //累加求和,表示式20中的∑j(Qj/Nj)
                }
83.
84
                 sum = 1./sum;
85.
                 //得到归一化的优先率
86.
                 for( k = 0; k < m; k++ )
87.
                    priors[k] *= sum; //得到式20中的qj
88.
89.
             //遍历所有分类,得到最多的样本数量max_val,以及对应的分类max_k
90.
             for( k = 0; k < m; k++ )
91.
92.
                 //得到加权后的第k个分类的样本数
93.
                 double val = cls count[k]*priors[k];
```

第12页 共46页 2017/11/30 上午10:51

```
total_weight += val;
95.
                 if( max val < val )</pre>
                                       //寻找最大值
96
                     max_val = val;
                                      //得到最大值
97.
98.
                     max_k = k; //最大值对应的分类索引
99
                 }
100.
             }
101.
102.
              node->class idx = max k;
                                       //赋值
      \int_{\Gamma}
103.
              //得到该节点的值,它等于拥有最大数量的那个分类的响应值
104.
              node->value = data->cat_map->data.i[
       1
105.
                 \label{lambdata-cat_ofs-data} \verb| data->cat_var_count | + max_k |;
106
              //该节点的风险值,即式23的分子部分
              node->node_risk = total_weight - max_val;
107.
108
              //遍历所有交叉验证的子集
109
              for( j = 0; j < cv_n; j++ )
      h
110.
111.
                 double sum_k = 0, sum = 0, max_val_k = 0;
      […
112.
                 max_val = -1; max_k = -1;
113.
                 //遍历该子集的各个分类
114.
                 for( k = 0; k < m; k++ )
      &
115.
116
                     double w = priors[k];
                                           //该类的先验概率
                     //得到加权后第j个子集的第k个分类的样本数,即式25中Lv中的第k个分类的样本数,v就是这里j,也就
117.
      是第i个测试样本集
118.
                     double val_k = cv_cls_count[j*m + k]*w;
                     //得到加权后的第j个子集中的不包括第k个分类的所有样本数,即式25中L(v)中第k个分类的样本数,v
119.
      就是这里j,也就是第j个训练样本集
120.
                     double val = cls_count[k]*w - val_k;
121.
                     sum_k += val_k; // val_k的累加,即Lv中的样本数
                                   // val的累加,即L(v)中的样本数
122.
                     sum += val;
                     //在当前子集下寻找具有最大分类样本数的val
123.
124
                     if( max_val < val )</pre>
                                          //寻找最大值
125.
126.
                         max val = val; //最大的val
                         max_val_k = val_k; //此时对应的最大val_k
127.
                         \max_k = k; //此时对应的分类索引
128.
129.
                     }
130.
                 }
                 //赋值
131.
                 node->cv_Tn[j] = INT_MAX; //最大常数
132.
                                                         //风险值,式26的分子部分
133.
                 node->cv node risk[i] = sum - max val:
134.
                 node->cv_node_error[j] = sum_k - max_val_k; //错误率,式28
135.
             }
136.
          else
                 //回归树
137.
138.
          {
139.
              // in case of regression tree:
140.
              // * node value is 1/n*sum_i(Y_i), where Y_i is i-th response,
141.
                  n is the number of samples in the node.
142.
              // * node risk is the sum of squared errors: sum_i((Y_i - node_value))^2)
              //\ ^{*} j-th cross-validation fold value and risk are calculated as above,
143.
144.
                   but using the samples with cv_labels(*)!=j.
              //\ ^{*} j-th cross-validation fold error is calculated
145.
146.
              //
                   using samples with cv_labels(*)==j as the test subset:
147
              //
                   error_j = sum_(i, cv_labels(i)==j)((Y_i - < node_value_j >)^2),
              //
                   where node value j is the node value calculated
148.
149.
              //
                   as described in the previous bullet, and summation is done
150.
                   over the samples with cv_labels(*)==j.
151.
              double sum = 0, sum2 = 0;
              float* values_buf = (float*)(cv_labels_buf + n);
153.
154
              int* sample_indices_buf = (int*)(values_buf + n);
              //得到所有样本的响应值, values指向该内存
156.
              const float* values = data->get_ord_responses(node, values_buf, sample_indices_buf);
157
              double *cv_sum = 0, *cv_sum2 = 0;
158.
              int* cv_count = 0;
159.
160.
              if( cv_n == 0 ) //不需要交叉验证
161.
162.
                 for( i = 0; i < n; i++ )</pre>
                                            //遍历所有样本
163.
164.
                     double t = values[i];
                                            //得到该样本的响应值
165.
                     sum += t; //响应值累加
                     sum2 += t*t;
                                   //响应值平方的累加
166.
```

```
167.
               }
168.
169
            else
                   //需要交叉验证
170.
171.
               cv_sum = (double*)base_buf;
                                         //表示各个子集的样本响应值的和
                                      //表示各个子集的样本响应值平方和
172.
               cv_sum2 = cv_sum + cv_n;
               cv_count = (int*)(cv_sum2 + cv_n); //表示各个子集的样本数量
173.
174.
               //以上三个数组清零
175.
               for( j = 0; j < cv_n; j++ )</pre>
176.
                   cv_sum[j] = cv_sum2[j] = 0.;
177.
      1
178.
                   cv_count[j] = 0;
179.
180.
181.
               for( i = 0; i < n; i++ ) //遍历所有样本
182.
                   j = cv_labels[i]; //得到当前样本所在的交叉验证的子集
183.
184.
                   double t = values[i]; //当前样本的响应值
185.
                   double s = cv_sum[j] + t; //计算各个子集的样本响应值的和
186.
                   double s2 = cv_sum2[j] + t*t; //计算各个子集的样本响应值平方和
187.
                   int nc = cv_count[j] + 1; //统计各个子集的样本数量
188.
                   //针对当前样本所在的交叉验证的子集
189.
                   cv_sum[j] = s;
                                 //响应值累加
                   cv_sum2[j] = s2; //响应值平方累加
190.
191.
                   cv_count[j] = nc;
                                    //样本计数
192.
               //遍历所有交叉验证的子集
193.
194.
               for( j = 0; j < cv_n; j++ )
195.
               {
                   sum += cv_sum[j]; //所有样本响应值的和
196.
                   sum2 += cv_sum2[j]; //所有样本响应值平方和
198.
               }
199
            //该节点的风险值,即式24的分子部分
200.
201.
            node->node risk = sum2 - (sum/n)*sum;
202.
            //该节点的值等于所有样本的平均响应值
203.
            node->value = sum/n:
204
            //遍历所有交叉验证的子集,计算每个子集的风险值和错误率
205
            for( j = 0; j < cv_n; j++ )</pre>
206.
               //s为该子集的样本响应值之和,即测试样本集的响应值之和;si为不包括该子集的所有样本的响应值之和,即
207.
     训练样本集的响应值之和
208
               double s = cv_sum[j], si = sum - s;
209.
               //s2为该子集的样本响应值平方之和,即测试样本集的响应值的平方之和;s2i为不包括该子集的所有样本的响
     应值平方之和,即训练样本集的响应值的平方之和
210.
               double s2 = cv_sum2[j], s2i = sum2 - s2;
               //c为该子集的样本数量,即测试样本集的数量;ci为不包括该子集的所有样本数量,即训练样本集的数量
211.
212.
               int c = cv_count[j], ci = n - c;
213.
               //r表示训练样本集的平均响应值
214.
               double r = si/MAX(ci,1);
215.
               //该子集的风险值,式27的分子部分
216.
               node->cv_node_risk[j] = s2i - r*r*ci;
217.
               //该子集的错误率,式29
               node->cv_node_error[j] = s2 - 2*r*s + c*r*r;
218.
219.
               node->cv_Tn[j] = INT_MAX; //最大常数
220.
221.
         }
222. }
```

找到最佳分叉属性find_best_split函数:

()

```
[cpp]
    CvDTreeSplit* CvDTree::find_best_split( CvDTreeNode* node )
1.
2.
        //实例化DTreeBestSplitFinder类,初始化一些变量
3.
        DTreeBestSplitFinder finder( this, node );
4.
        //并行处理,寻找最佳的分类特征属性,它调用DTreeBestSplitFinder类的重载运算符(),该函数的分析见下面
5.
6.
        cv::parallel_reduce(cv::BlockedRange(0, data->var_count), finder);
7.
        CvDTreeSplit *bestSplit = 0;
8.
        if( finder.bestSplit->quality > 0 ) //找到了最佳分类特征属性
9.
```

第14页 共46页

```
10.
 11.
               bestSplit = data->new split cat( 0, -1.0f );
 12.
               memcpy( bestSplit, finder.bestSplit, finder.splitSize );
                                                                       //赋值
 13.
 14.
                               //返回最佳分类特征属性
 15.
           return bestSplit;
 16. }
DTreeBestSplitFinder类的重载运算符():
        1
0
       √oid DTreeBestSplitFinder::operator()(const BlockedRange& range)
  1.
       ••• //vi表示样本数据的特征属性, vi1和vi2分别为特征属性队列中的首和末属性
  3.

✓ int vi, vi1 = range.begin(), vi2 = range.end();
  4.
       int n = node->sample_count; //得到该节点的样本数量CvDTreeTrainData* data = tree->get_data(); //得
  5.
  6.
                                                     //得到全部训练样本数据
  7.
           AutoBuffer<uchar> inn_buf(2*n*(sizeof(int) + sizeof(float))); //开辟一块内存
           //遍历所有特征属性
  8.
  9.
           for( vi = vi1; vi < vi2; vi++ )</pre>
 10.
 11.
               CvDTreeSplit *res;
               //得到该特征属性的类型,如果ci大于零,表示特征属性是类的形式,反之则是数值的形式
 12.
 13.
               int ci = data->get_var_type(vi);
 14
               //如果在该特征属性下只有一个样本,或者没有样本(即缺失),则退出该次vi循环
               if( node->get_num_valid(vi) <= 1 )</pre>
 15.
                  continue:
 16.
 17.
               if( data->is_classifier ) //分类树
 18.
 19.
 20.
                  if( ci >= 0 )
                                  //特征为类的分类树
                      {\tt res = tree->find\_split\_cat\_class(\ node,\ vi,\ bestSplit->quality,\ split,\ (uchar^*)}
 21.
                  else //特征为数值的分类树
 22
 23.
                      res = tree->find_split_ord_class( node, vi, bestSplit->quality, split, (uchar*)
 24
               else
                      //归纳树
 25.
 26.
 27
                  if( ci >= 0 )
                                  //特征为类的回归树
                      res = tree->find_split_cat_reg( node, vi, bestSplit->quality, split, (uchar*)in
 28.
 29
                  else //特征为数值的回归树
 30.
                      res = tree->find_split_ord_reg( node, vi, bestSplit->quality, split, (uchar*)in
 31.
               //如果计算得到的分类属性split的值大于bestSplit,则保存该split于bestSplit中,作为下次循环判断最佳
       分类属性的依据
 33.
               if( res && bestSplit->quality < split->quality )
 34.
                      memcpy( (CvDTreeSplit*)bestSplit, (CvDTreeSplit*)split, splitSize );
 35.
 36.
```

特征为类的分类树的分叉方法:

0

```
CvDTreeSplit* CvDTree::find_split_cat_class( CvDTreeNode* node, int vi, float init_quality,
1.
                                               CvDTreeSplit* _split, uchar* _ext_buf )
2.
     {
         //得到该节点的特征属性的形式
4.
5.
         int ci = data->get_var_type(vi);
         int n = node->sample_count;
                                     //得到该节点的样本数量
6.
7.
         //得到样本数据的分类数
8.
         int m = data->get_num_classes();
         //得到该特征属性vi的类别数量,也就是该特征有mi个种类
9.
10.
         int _mi = data->cat_count->data.i[ci], mi = _mi;
11.
         //为该特征属性vi开辟一块内存空间
         int base_size = m*(3 + mi)*sizeof(int) + (mi+1)*sizeof(double);
12.
13.
         if( m > 2 \&\& mi > data->params.max_categories )
```

第15页 共46页

```
base_size += (m*min(data->params.max_categories, n) + mi)*sizeof(int);
15.
         else
16.
             base_size += mi*sizeof(int*);
         cv::AutoBuffer<uchar> inn_buf(base_size);
17.
18.
         if( ! ext buf )
19.
            inn_buf.allocate(base_size + 2*n*sizeof(int));
20.
         uchar* base buf = (uchar*)inn buf;
         uchar* ext_buf = _ext_buf ? _ext_buf : base_buf + base_size;
21.
        //二叉树,1c和rc分别指向二叉树左、右分支,它们都是含有m个元素的一维数组,各存储着该节点各自分支的m个分类的
22.
     作本数量
23.
         int* lc = (int*)base_buf;
                                    //lc指向二叉树左分支
      1 int* rc = lc + m; //rc指向二叉树右分支
24.
25.
        ■// cjk表示该特征属性的mi个种类的m个分类的样本数量,是一个有mi个元素的m维数组(也可以看成是含有m个元素的mi
     维数组)
26.
         int* _cjk = rc + m*2, *cjk = _cjk;
27.
         // c_weights表示该特征属性中每个种类的权值,是含有为mi个元素的一维数组
     double* c_weights = (double*)alignPtr(cjk + m*mi, sizeof(double));
28.
29
     ••• int* labels_buf = (int*)ext_buf;
30.
       ✓//label指向该特征属性中各个样本所对应的种类
31.
     const int* labels = data->get_cat_var_data(node, vi, labels_buf);
int* responses_buf = labels_buf + n;
32.
33.
34
         //responses指向该节点样本的分类,即响应值
         const int* responses = data->get_class_labels(node, responses_buf);
35.
36.
37
         int* cluster_labels = 0;
         int** int_ptr = 0;
38.
39
         \quad \textbf{int} \ \textbf{i}, \ \textbf{j}, \ \textbf{k}, \ \textbf{idx}; \\
40.
         double L = 0, R = 0;
41.
         double best_val = init_quality;
         int prevcode = 0, best_subset = -1, subset_i, subset_n, subtract = 0;
         //得到先验概率,该值由calc_node_value函数计算得到
43.
44
         const double* priors = data->priors_mult->data.db;
                                                            //得到先验概率
45.
46.
         // init array of counters:
47.
         // c_{jk} - number of samples that have vi-th input variable = j and response = k.
         //初始化cjk数组,即清零。通过上面的英文注释可以看出,cjk中j代表着该特征属性vi的种类,k代表着响应值(即样本
48.
     分类)
49.
         for( j = -1; j < mi; j++ )</pre>
            for( k = 0; k < m; k++ )
50.
                cjk[j*m + k] = 0;
51.
         //遍历该节点的所有样本,为cjk赋值
52.
53
         for( i = 0; i < n; i++ )
54.
           //得到该特征属性的第i个样本的种类
55.
            j = ( labels[i] == 65535 \&\& data->is_buf_16u) ? -1 : labels[i];
56
            //得到该特征属性的第i个样本的分类,即响应值
57.
58.
            k = responses[i];
59.
           cjk[j*m + k]++;
                             //计数累加
60.
61.
         if( m > 2 )
                      //非两类问题,即多于2个的分类(响应值)
62.
63.
64.
             //如果该特征属性的种类数量mi大于所设定的阈值
65.
             if( mi > data->params.max_categories )
66
                //mi重新设置为阈值或样本数量的最小的那个值
67.
68.
                mi = MIN(data->params.max_categories, n);
69.
                cjk = (int*)(c_weights + _mi);
70.
                cluster_labels = cjk + m*mi;
                //聚类,这里就不再分析
                cluster_categories( _cjk, _mi, m, cjk, mi, cluster_labels );
72.
73.
             subset i = 1:
             subset_n = 1 << mi;
75.
76
77.
         else //两类问题
78.
79.
             assert( m == 2 ); //确定为两类问题
             int_ptr = (int**)(c_weights + _mi);
80.
81.
             //数组int_ptr[j]表示该特征属性的第j个种类的响应值为1的样本数量,因为是两类问题,另一个响应值是为0
82.
             for( j = 0; j < mi; j++ )</pre>
83.
                int_ptr[j] = cjk + j*2 + 1;
                                             //这里的2就是m
             //按样本数量由多到少进行排序,int_ptr指针的顺序改变了
85.
             icvSortIntPtr( int_ptr, mi, 0 );
```

第16页 共46页 2017/11/30 上午10:51

```
86.
             subset_i = 0;
87.
             subset n = mi;
88
         //遍历样本的各个分类,初始化数组lc和rc,即式12的Lj和Rj
89
90
         for( k = 0; k < m; k++ )</pre>
91.
92.
             int sum = 0:
93.
             for( j = 0; j < mi; j++ )</pre>
                 sum += cjk[j*m + k];
94.
      ď
95
             rc[k] = sum;
                           //把所有样本数量先存储在rc中
                        //清零
96.
             lc[k] = 0;
      1 }
97.
98.
         ·//遍历特征属性的各个种类,得到右分支的样本数R
      : for( j = 0; j < mi; j++ )
99.
100.
             double sum = 0;
101.
      h
             //遍历第_{\rm j}个特征属性种类的所有分类,每个分类的样本个数乘以它所对应的先验概率,然后再累加求和,即为该特征
102.
      属性种类的权值
             for( k = 0; k < m; k++ )
103.
104
                 sum += cjk[j*m + k]*priors[k];
105
             c_weights[j] = sum; //第j个特征属性种类的权值
      ₡
                                //权值累加,得到右分支的样本数量R
106.
             R += c_weights[j];
107
         //遍历特征属性的各个种类
108.
109
         for( ; subset_i < subset_n; subset_i++ )</pre>
110.
111.
             double weight;
112.
             int* crow;
113.
             //分别表示式12中的∑Lj2和∑Rj2
114.
             double lsum2 = 0, rsum2 = 0;
115
116.
             if( m == 2 )
                            //两类问题
117.
                 //idx是相对地址指针,指向cjk数组的种类,这里的2就是m
                 //这里的int_ptr已经是排序后的指向
118.
119.
                 idx = (int)(int_ptr[subset_i] - cjk)/2;
                    //非两类问题
120.
121.
122.
                 int graycode = (subset_i>>1)^subset_i;
123.
                 int diff = graycode ^ prevcode;
124.
125.
                 // determine index of the changed bit.
126.
                 Cv32suf u:
127.
                 idx = diff >= (1 << 16) ? 16 : 0;
                 u.f = (float)(((diff >> 16) | diff) & 65535);
128.
129.
                 idx += (u.i >> 23) - 127;
130.
                 subtract = graycode < prevcode;</pre>
                 prevcode = graycode;
131.
132.
133.
             //crow为绝对地址指针,指向该节点的样本分类
134.
             crow = cjk + idx*m;
135.
             weight = c_weights[idx];
                                      //指向该种类的权值
             //该特征属性的种类权值太小,则该种类不予考虑
136.
137.
             if(weight < FLT_EPSILON)
                 continue;
138.
139.
140
             if( !subtract )
                              //subtract为零,即两类问题或聚类中graycode > prevcode
141.
142.
                 //遍历该种类的所有分类
143.
                 for( k = 0; k < m; k++ )
144.
145.
                    int t = crow[k];
                                      //第k类的样本数量
                    //lval和rval为式21的qjLj和qjRj,在初始化数组lc和rc的时候,定义lc为0,rc为全部样本,所以
146.
      这里左分支要是增加t个样本,那么右分支就要减少t个样本
147.
                    int lval = lc[k] + t;
148.
                    int rval = rc[k] - t;
149
                    //先验概率,以及它的平方
                    double p = priors[k], p2 = p*p;
150.
151.
                    // lsum2和rsum2分别为式21的两个分式的分子部分
152.
                    lsum2 += p2*lval*lval;
                    rsum2 += p2*rval*rval;
153.
154.
                    lc[k] = lval; rc[k] = rval;
                                                //更新qjLj和qjRj
155.
156.
                 //L和R分别为式21的两个分式的分母部分,两者是此消彼长的关系
157.
                 L += weight;
                 R -= weight;
158.
```

第17页 共46页

```
159.
160.
              else
                      //subtract不为零
161
                 for( k = 0; k < m; k++ )</pre>
162.
163.
                 {
164.
                     int t = crow[k];
                     int lval = lc[k] - t;
165.
166.
                     int rval = rc[k] + t;
                     double p = priors[k], p2 = p*p;
167.
168
                     lsum2 += p2*lval*lval;
                     rsum2 += p2*rval*rval;
169.
170.
                     lc[k] = lval; rc[k] = rval;
171.
                 L -= weight;
172.
173.
                 R += weight;
174.
              //如果L和R都大于某一常数
175.
176.
              if( L > FLT_EPSILON && R > FLT_EPSILON )
      <u>...</u>
177.
              {
178.
                 //式21的值
179.
                 double val = (lsum2*R + rsum2*L)/((double)L*R);
      &
180.
                 if( best_val < val )</pre>
                                       //替换成最大值
181.
                 {
                     best_val = val;
                                     //最佳值
182.
                                              //最佳特征属性的种类,起到式12中的β作用
183.
                     best_subset = subset_i;
184
185.
               //结束遍历特征属性的各个种类
186.
187.
188
          CvDTreeSplit* split = 0;
                                 //如果找到了最佳分类值
189
          if( best_subset >= 0 )
190.
191
              split = _split ? _split : data->new_split_cat( 0, -1.0f );
              split->var_idx = vi; //得到该特征属性
192.
193.
              split->quality = (float)best_val;
                                               //最佳分叉值,式12或21的值
194
              //初始化split->subset为0
              \label{eq:memset} memset(\ split->subset,\ 0,\ (data->max\_c\_count\ +\ 31)/32\ ^*\ \textbf{sizeof(int))};
195
196.
              if( m == 2 )
                            //两类问题
197.
              {
                  //遍历前best subset个特征属性的种类
198.
                 for( i = 0; i <= best_subset; i++ )</pre>
199
200.
201
                     //idx是相对地址指针,指向cjk数组的该种类,它是排序后的顺序,这里的>>1就是除以2,2也就是m
202.
                     idx = (int)(int_ptr[i] - cjk) >> 1;
                     //按int的大小(32位)进行存储,每个特征属性的种类占一位,前32位存储在split->subset[0]中,
203.
      第二个32为存储在split->subset[1],以此类推,每个数组又由低位到高位进行存储,例如某特征属性有19个种类,则这19
      个种类存放在split->subset[0]的前19位中。被分配到左分支上的特征属性种类,该位置1,右分支上的特征属性种类,该位
204.
                     split->subset[idx >> 5] |= 1 << (idx & 31);
205
                 }
206
              else
                      //非两类问题
207.
208
209.
                 for( i = 0; i < _mi; i++ )</pre>
210.
211.
                     idx = cluster_labels ? cluster_labels[i] : i;
                     if( best subset & (1 << idx) )
212.
213
                         split->subset[i >> 5] |= 1 << (i & 31);
214.
215.
             }
216.
217.
          return split;
218.
```

特征为数值的分类树的分叉方法:

()

```
[cpp]
    CvDTreeSplit* CvDTree::find_split_ord_class( CvDTreeNode* node, int vi,
1.
2.
                                                 float init_quality, CvDTreeSplit* _split, uchar* _
3.
4.
        const float epsilon = FLT_EPSILON*2;
                                                 //定义一个最小常数
```

```
5.
        int n = node->sample_count;
                                   //该节点的样本数量
        //该节点具有vi特征属性的样本数量,因为有可能有些样本不具备vi这个特征属性,所以n1 \leq n
6.
7.
        int n1 = node->get_num_valid(vi);
        int m = data->get_num_classes();
                                        //样本数据的分类数
8.
9.
        //为该特征属性vi开辟一块内存空间
10.
        int base_size = 2*m*sizeof(int);
11.
        cv::AutoBuffer<uchar> inn_buf(base_size);
12.
        if(! ext buf)
     13.
14.
        uchar* ext_buf = _ext_buf ? _ext_buf : base_buf + base_size;
15.
     float* values_buf = (float*)ext_buf;
16.
        int* sorted_indices_buf = (int*)(values_buf + n);
17.
     int* sample_indices_buf = sorted_indices_buf + n;
18.
19.
        const float* values = 0:
        const int* sorted_indices = 0;
20.
     //对该节点的所有样本,按照特征属性vi的值的从小到大的顺序进行排序,分别得到两个数组:values和
21.
     sorted_indices,这个数组的索引值表示排序后的索引值,而values值为特征属性的值,sorted_indices值为未排序前的
     举本索引值
22.
        data->get_ord_var_data( node, vi, values_buf, sorted_indices_buf, &values,
23.
                             &sorted_indices, sample_indices_buf );
     int* responses_buf = sample_indices_buf + n;
24.
25
        //responses指向该节点样本的分类,即响应值
        const int* responses = data->get_class_labels( node, responses_buf );
26.
27.
        const int* rc0 = data->counts->data.i;
                                              //指向样本分类
28
        int* lc = (int*)base_buf; //表示左分支
29.
30.
        int* rc = lc + m;
                           //表示右分支
31.
        int i, best_i = -1;
32.
        // lsum2和rsum2分别表示式12的第一项和第二项分式的分子部分,即\SigmaLj2和\SigmaRj2,best_val即为式12的值
        double lsum2 = 0, rsum2 = 0, best_val = init_quality;
34.
        //得到不同分类的先验概率
35
        const double* priors = data->have_priors ? data->priors_mult->data.db : 0;
36.
37.
        // init arrays of class instance counters on both sides of the split
38.
        for( i = 0; i < m; i++ )</pre>
39.
40
41.
            lc[i] = 0;
                       //左分支每个分类都是0
            rc[i] = rc0[i]; //右分支每个分类都是相应的样本数
42.
43.
44.
45.
        // compensate for missing values
46.
        //补偿缺失特征属性的样本,减去右分支的一些数量
47.
        for( i = n1; i < n; i++ )</pre>
48.
49.
            rc[responses[sorted_indices[i]]]--;
50.
51.
        if(!priors) //样本没有先验概率
52.
53.
            //分别表示左、右分支的样本数量,即是式12的L和R
54.
55.
            int L = 0, R = n1;
            //初始化∑Rj2,此时∑Lj2为0
56.
57.
            for( i = 0; i < m; i++ )
58
               rsum2 += (double)rc[i]*rc[i];
            //按照特征属性值从小到大的顺序遍历该节点的所有样本,每循环一次,左分支的样本数加1,右分支样本数减1
59.
60.
            for( i = 0; i < n1 - 1; i++ )</pre>
61.
               int idx = responses[sorted_indices[i]]; //得到该样本的响应值,即分类
62.
               int lv, rv;
L++; R--;
                          //左分支加1,右分支减1
64.
65.
               //lv和rv分别表示式12中的Lj和Rj
               lv = lc[idx]; rv = rc[idx];
               //更新5Lj2和5Rj2
67.
68.
               //在循环之前左分支第j个分类的样本数为Lj,循环之后样本数为Lj +1,那么经过平方以后∑Lj2增加了
     2Lj +1,即(Lj +1)2 - Lj2 = Lj2+2Lj +1- Lj2 =2Lj +1,因此这里lsum2要累加lv*2 + 1;同理右分支从Rj减少到
     Rj - 1,平方后减少了2Rj - 1,因此rsum2要累减rv*2 - 1
69.
               lsum2 += lv*2 + 1;
               rsum2 -= rv*2 - 1;
70.
71.
               //更新Lj和Rj
               lc[idx] = lv + 1; rc[idx] = rv - 1;
72.
73.
               //只有当该特征属性的值与其排序后相邻的值有一定差距时才计算式12
               if( values[i] + epsilon < values[i+1] )</pre>
75.
```

第19页 共46页 2017/11/30 上午10:51

```
76.
                     double val = (lsum2*R + rsum2*L)/((double)L*R);
77.
                     if( best_val < val ) //得到最大值
78.
79.
                         best_val = val;
                                           //式12的值
                         best_i = i; //特征属性排序后的索引值 , 起到式12阈值\beta的作用
80.
81.
82.
                 }
83
             }
      L^{\frac{1}{2}}else
84
85
                  //样本有先验概率
86.
      1
                                    //分别表示式21中两个分式的分母部分
87.
              double L = 0, R = 0;
88
              //初始化式13中的∑qjRj和∑(qjRj)2
              for( i = 0; i < m; i++ )</pre>
89.
90.
91.
                 double wv = rc[i]*priors[i];
      h
92.
                 R += wv:
93
                 rsum2 += wv*wv;
      <u>...</u>
94.
95
              //按照特征属性值从小到大的顺序遍历该节点的所有样本,每循环一次,左分支的样本数{
m Im},右分支样本数减{
m p}
              for( i = 0; i < n1 - 1; i++ )
      &
97.
98
                 int idx = responses[sorted_indices[i]]; //得到该样本的响应值,即分类
99.
                 int lv. rv:
100.
                 double p = priors[idx], p2 = p*p;
                                                    //先验概率,及它的平方
101.
                 L += p; R -= p; //更新式21中两个分式的分母部分
                 lv = lc[idx]; rv = rc[idx];
102.
                                             //得到式21中的qjLj和qjRj
103.
                 //更新式13中的∑(qjLj)2和∑(qjRj)2
104.
                 lsum2 += p2*(lv*2 + 1);
                 rsum2 -= p2*(rv*2 - 1);
105.
                 lc[idx] = lv + 1; rc[idx] = rv - 1;
                                                     //更新式21中的qjLj和qjRj
107.
108.
                 if( values[i] + epsilon < values[i+1] )</pre>
109.
110.
                     double val = (lsum2*R + rsum2*L)/((double)L*R);
111.
                     if( best_val < val )</pre>
112.
                     {
113.
                         best_val = val;
114.
                         best_i = i;
115.
                     }
116.
                 }
117.
             }
118.
119.
          CvDTreeSplit* split = 0;
120.
121.
          if( best_i >= 0 ) //如果有可分叉的阈值,则表明得到了最佳分叉属性
122.
123.
              split = _split ? _split : data->new_split_ord( 0, 0.0f, 0, 0, 0.0f );
                                                                                  //实例化split
      变量
              split->var_idx = vi;
                                    //特征属性
124.
125.
              //特征属性的值,这里取了平均
126.
              split->ord.c = (values[best i] + values[best i+1])*0.5f;
127.
              //分叉所对应的特征属性中的索引值,起到阈值\beta的作用
128.
              split->ord.split_point = best_i;
                                   //表示该分叉无需反转
129.
              split->inversed = 0;
130
              split->quality = (float)best_val;
131.
132.
          return split; //返回
133. }
```

特征为类的回归树的分叉方法:

()

```
[cpp]
1.
    CvDTreeSplit* CvDTree::find_split_cat_reg( CvDTreeNode* node, int vi, float init_quality, CvDTr
2.
3.
        //得到特征属性的形式,即类别形式
4.
        int ci = data->get_var_type(vi);
       int n = node->sample_count; //得到该节点的样本数量
5.
6.
        //得到该特征属性的种类数量,即有mi种不同的类别形式
        int mi = data->cat_count->data.i[ci];
7.
8.
        //为该特征属性开辟一块内存空间
```

```
int base_size = (mi+2)*sizeof(double) + (mi+1)*(sizeof(int) + sizeof(double*));
9.
         cv::AutoBuffer<uchar> inn_buf(base_size);
10.
11.
         if( !_ext_buf )
            inn_buf.allocate(base_size + n*(2*Sizeof(int) + Sizeof(float)));
12.
13.
         uchar* base_buf = (uchar*)inn_buf;
         uchar* ext_buf = _ext_buf ? _ext_buf : base_buf + base_size;
         int* labels_buf = (int*)ext_buf;
15.
         //label指向该特征属性中各个样本所对应的种类
16.
     Const int* labels = data->get_cat_var_data(node, vi, labels_buf);
float* responses_buf = (float*)(labels_buf + n);
17.
18.
19.
         int* sample_indices_buf = (int*)(responses_buf + n);
      1 //responses指向该节点样本的分类,即响应值,它已经是按照响应值的大小排序后的序列,对应的排序后的索引值存储
20.
     在eample_indices_buf数组内
21.
     const float* responses = data->get_ord_responses(node, responses_buf, sample_indices_buf);
22.
         //sum可以认为是拥有mi个元素的一维数组,它的值表示每个特征属性种类的样本响应值之和
23.
         double* sum = (double*)cv::alignPtr(base_buf, SizeOf(double)) + 1;
     //counts可以认为是拥有mi个元素的一维数组,它的值表示每个特征属性种类的样本数量
24.
25.
         int* counts = (int*)(sum + mi) + 1;
26.
     ••• double** sum_ptr = (double**)(counts + mi);
       ~ int i, L = 0, R = 0; //L和R分别表示式19中的L和R
27.
28.
         // best_val表示最佳的分类属性,lsum和rsum分别表示式19中的∑li和∑ri
     double best_val = init_quality, lsum = 0, rsum = 0;
29.
30
         int best_subset = -1, subset_i;
         //遍历所有该特征属性的种类,初始化数组sum和counts为0
31.
32.
         for( i = -1; i < mi; i++ )</pre>
            sum[i] = counts[i] = 0;
33.
34.
35.
         \ensuremath{/\!/} calculate sum response and weight of each category of the input var
36.
         //遍历该节点的所有样本,为数组sum和counts赋值
37.
         for( i = 0; i < n; i++ )
38
            //得到该特征属性的第i个样本的种类
39.
40
            int idx = ( (labels[i] == 65535) && data->is_buf_16u ) ? -1 : labels[i];
            //得到第idx个种类的响应值之和
41.
42.
            double s = sum[idx] + responses[i];
            //得到第idx个种类的样本数量
43.
            int nc = counts[idx] + 1;
44.
45.
            sum[idx] = s; //更新sum
46.
            counts[idx] = nc; //更新counts
47.
48.
49.
         // calculate average response in each category
50.
         //遍历该特征属性的所有种类,先把所有的样本赋予右分支,左分支为零
51.
         for( i = 0; i < mi; i++ )</pre>
52.
                            //计算式19中的R
//计算式19中的∑ri
53.
            R += counts[i];
            rsum += sum[i];
54.
55.
            //现在数组sum的含义为每个特征属性种类的平均样本响应值,为了后面的排序之用
            sum[i] /= MAX(counts[i],1);
56.
                                  //指向数组sum
            sum_ptr[i] = sum + i;
57.
58
         //根据特征属性种类的平均样本响应值的大小进行排序, sum_ptr为排序后的指向
59.
60.
         icvSortDblPtr( sum_ptr, mi, 0 );
61.
62.
         // revert back to unnormalized sums
63.
         // (there should be a very little loss of accuracy)
         //数组sum的含义再次恢复为每个特征属性种类的样本响应值之和,这么做虽然会带来一定的误差,但比重新计算效率要
64.
65.
         for( i = 0; i < mi; i++ )</pre>
            sum[i] *= counts[i];
66.
         //按照平均样本响应值大小的顺序,遍历该特征属性的所有种类
68.
         for( subset_i = 0; subset_i < mi-1; subset_i++ )</pre>
69.
            //idx为第subset_i个排序所对应的那个特征属性种类
            int idx = (int)(sum_ptr[subset_i] - sum);
71.
72
            int ni = counts[idx];
                                   该种类的样本数量
73.
            if( ni ) //如果该种类有样本
74.
76.
                double s = sum[idx]:
                                      //该种类的样本响应值之和
77.
                //更新式19中的各个参数
                lsum += s; L += ni;
78.
79.
                rsum -= s; R -= ni;
                if( L && R ) //如果L和R都有值
81.
```

第21页 共46页 2017/11/30 上午10:51

```
82.
                 {
83.
                    double val = (lsum*lsum*R + rsum*rsum*L)/((double)L*R);
                                                                          //式20
84
                    if( best_val < val )</pre>
                                         //找到式19的最大值
85.
86.
                        best val = val;
                                       //更新最佳分类属性值,即式19的值
                        //最佳分类属性值对应的属性索引值,它的起到式19阈值β的作用
87.
88.
                        best subset = subset i;
89
90
                }
      \mathbb{L}
91
             }
92.
93.
94
         CvDTreeSplit* split = 0;
      = if( best_subset >= 0 )
                                 //如果找到了最佳分类值
95.
96
97.
             split = _split ? _split : data->new_split_cat( 0, -1.0f);
             split->var idx = vi: //分叉的特征属性
98.
99
             split->quality = (float)best_val; //分叉的值,即式19
100.
      [•••
             //初始化split->subset为0
101.
             memset( split->subset, 0, (data->max_c_count + 31)/32 * sizeof(int));
             //遍历前best_subset个特征属性的种类
103.
             for( i = 0; i <= best_subset; i++ )</pre>
104
             {
                 //idx是相对地址指针,它表示特征属性排序后的顺序
105.
106
                 int idx = (int)(sum_ptr[i] - sum);
                 //按int的大小(32位)进行存储,每个特征属性的种类占一位,在左分支的种类,则该位置1,在右分支的种
107
      类,该位清零
108
                 split->subset[idx >> 5] |= 1 << (idx & 31);
109.
110.
          return split; //返回
111.
112. }
```

特征为数值的回归树的分叉方法:

()

```
1.
     CvDTreeSplit* CvDTree::find_split_ord_reg( CvDTreeNode* node, int vi, float init_quality, CvDTr
2.
        const float epsilon = FLT_EPSILON*2;
                                             //定义一个最小的常数
3.
        int n = node->sample_count; //该节点的样本数量
4.
5.
        int n1 = node->get_num_valid(vi); //该节点具有vi特征属性的样本数量
        //为该特征属性开辟一块内存空间
6.
7.
        cv::AutoBuffer<uchar> inn buf:
8.
        if( !_ext_buf )
            inn_buf.allocate(2*n*(sizeof(int) + sizeof(float)));
9.
10.
        uchar* ext_buf = _ext_buf ? _ext_buf : (uchar*)inn_buf;
11.
        float* values_buf = (float*)ext_buf;
12.
        int* sorted_indices_buf = (int*)(values_buf + n);
        int* sample_indices_buf = sorted_indices_buf + n;
14.
        const float* values = 0;
15.
        const int* sorted_indices = 0;
        //对该节点的所有样本,按照特征属性vi的值的从小到大的顺序进行排序,分别得到两个数组:values和
     sorted_indices,这个数组的索引值表示排序后的索引值,而values值为特征属性的值,sorted_indices值为未排序前的
     样本索引值
17.
        data->get_ord_var_data( node, vi, values_buf, sorted_indices_buf, &values, &sorted_indices,
18.
        float* responses_buf = (float*)(sample_indices_buf + n);
        //responses指向该节点样本的分类,即响应值,它已经是按照响应值的大小排序后的序列,对应的排序后的索引值存储
19
     在sample indices buf数组内
20
        const float* responses = data->get_ord_responses( node, responses_buf, sample_indices_buf
21.
22.
        int i, best_i = -1;
        // best_val表示最佳的分类属性,lsum和rsum分别表示式19中的\Sigmali和\Sigmari,rsum先被初始化为所有样本,对于回归
     树来说, node->value表示该节点样本平均响应值
24
        double best_val = init_quality, lsum = 0, rsum = node->value*n;
25.
        int L = 0, R = n1;
                           //L和R分别表示式19中的L和R,R先被初始化为所有样本
26.
27
        // compensate for missing values
        //补偿缺失的样本,减去右分支的一些数量
28.
29.
        for( i = n1; i < n; i++ )</pre>
30.
            rsum -= responses[sorted_indices[i]];
31.
```

第22页 共46页

```
// find the optimal split
        //按照特征属性值从小到大的顺序遍历该节点的所有样本,每循环一次,左分支的样本数加1,右分支样本数减1
33.
34
        for( i = 0; i < n1 - 1; i++ )</pre>
35.
36.
            float t = responses[sorted_indices[i]]; //得到该样本的响应值,即分类
37
            L++; R--; //左分支加1,右分支减1
            //更新式20中的∑li和∑ri
38.
39
            lsum += t;
40.
            rsum -= t;
     \int_{\Gamma}
41.
            //只有当该特征属性的值与其排序后相邻的值有一定差距时才计算式19
42.
            if( values[i] + epsilon < values[i+1] )</pre>
     1
43.
                                                                    //式19
44.
               double val = (lsum*lsum*R + rsum*rsum*L)/((double)L*R);
45.
               if( best_val < val ) //得到最大值
46
47.
                   best_val = val; //式19的值
     h
                   best_i = i; //特征属性排序后的索引值,起到式19阈值\beta的作用
48.
49
50.
            }
51.
     CVDTreeSplit* split = 0;
53.
54
        if( best_i >= 0 ) //如果有可分叉的阈值,则表明得到了最佳分叉属性
55.
56.
            split = \_split ? \_split : data->new\_split\_ord( 0, 0.0f, 0, 0, 0.0f );
            split->var_idx = vi; //特征属性
57
            //特征属性的值,这里取了平均
58.
59.
            split->ord.c = (values[best_i] + values[best_i+1])*0.5f;
60.
            //分叉所对应的特征属性中的索引值,起到阈值β的作用
61.
            split->ord.split_point = best_i;
            split->inversed = 0; //表示该分叉无需反转
            split->quality = (float)best_val; //式19的值
63.
64.
        return split; //返回
65.
66. }
```

按最佳分叉属性标注该节点的所有样本是被分配到左分支还是右分支:

0

```
1.
     // calculate direction (left(-1), right(1), missing(0))
     // for each sample using the best split
     // the function returns scale coefficients for surrogate split quality factors.
    // the scale is applied to normalize surrogate split quality relatively to the
     // best (primary) split quality. That is, if a surrogate split is absolutely
5.
     // identical to the primary split, its quality will be set to the maximum value =
     // quality of the primary split; otherwise, it will be lower.
8.
     // besides, the function compute node->maxlr,
     // minimum possible quality (w/o considering the above mentioned scale)
10.
     // for a surrogate split. Surrogate splits with quality less than node->maxlr
11.
     // are not discarded.
12.
     double CvDTree::calc_node_dir( CvDTreeNode* node )
13.
14.
         //表示特征属性的种类是属于左分支还是右分支,-1为左分支,1为右分支,如果该特征属性缺失,则为0
15.
         char* dir = (char*)data->direction->data.ptr;
16.
         //n表示该节点的样本数量, vi表示分类的最佳特征属性
         int i, n = node->sample_count, vi = node->split->var_idx;
         double L. R:
18.
19.
         assert(!node->split->inversed); //确保分叉不反转
20.
21.
22.
         if( data->get_var_type(vi) >= 0 ) // split on categorical var
         //表示该特征属性是种类的形式
23.
24.
25.
             //开辟一块内存空间
26.
             cv::AutoBuffer<int> inn_buf(n*(!data->have_priors ? 1 : 2));
             int* labels_buf = (int*)inn_buf;
             //label指向该特征属性中各个样本所对应的种类
28.
29.
             const int* labels = data->get_cat_var_data( node, vi, labels_buf );
             // subset数组的每一位表示特征属性的种类,左分支的种类位是1,右分支的是0
30.
31.
             const int* subset = node->split->subset;
```

第23页 共46页 2017/11/30 上午10:51

```
if( !data->have_priors )
                                       //无先验概率
33.
34
                 int sum = 0, sum_abs = 0;
                 //遍历该节点的所有样本
35.
36.
                 for( i = 0; i < n; i++ )</pre>
37.
                 {
38.
                    int idx = labels[i]; //表示该样本的特征属性的种类
39
      #define CV_DTREE_CAT_DIR(idx, subset) \
40.
41.
         (2*((subset[(idx)>>5]&(1 << ((idx) & 31)))==0)-1)
42.
      1
                    //d为-1表示idx(特征属性的种类)属于左分支,为1表示属于右分支,如果没有该特征属性,ydh0
43.
44.
                    int d = ( ((idx >= 0)&&(!data->is_buf_16u)) || ((idx != 65535)&&
      (data->is_buf_16u)) ) ?
45
                        CV_DTREE_CAT_DIR(idx, subset) : 0;
                    //sum表示d累加求和,因为d也可能为负值,所以sum的含义为右分支比左分支多出的特征属性种
46
      sum_abs表示d的绝对值之和,表示的含义为被分叉的特征属性种类
47
                    sum += d; sum_abs += d & 1;
48.
                    dir[i] = (char)d; //赋值
49
50
                 //L和R分别表示左右分支的特征属性的种类数量
      <del>«</del>
                 R = (sum\_abs + sum) >> 1;
51.
52
                 L = (sum\_abs - sum) >> 1;
53.
54
             else
                     //有先验概率
55
             {
                 const double* priors = data->priors mult->data.db;
                                                                  //得到先验概率
56.
57.
                 double sum = 0, sum_abs = 0;
58.
                 int* responses_buf = labels_buf + n;
                 //responses指向该节点样本的分类,即响应值
59.
                 const int* responses = data->get_class_labels(node, responses_buf);
                 //遍历该节点的所有样本
61.
62.
                 for( i = 0; i < n; i++ )</pre>
63.
64.
                    int idx = labels[i];
                                         //表示该样本的特征属性的种类
65.
                    double w = priors[responses[i]];
                                                     //得到该响应值的先验概率
                    //d为-1表示idx(特征属性的种类)属于左分支,为1表示属于右分支,如果没有该特征属性,则d为0
66.
67.
                    int d = idx >= 0 ? CV_DTREE_CAT_DIR(idx, subset) : 0;
68.
                    sum += d*w; sum_abs += (d & 1)*w; //增加了先验概率
                    dir[i] = (char)d;
69.
 70.
                 //L和R分别表示左右分支的特征属性的种类数量
71.
72.
                 R = (sum\_abs + sum) * 0.5;
                 L = (sum\_abs - sum) * 0.5;
73.
74.
             }
75.
         else // split on ordered var
76.
77.
         //表示该特征属性是数值的形式
 78.
             // split_point表示式12或是19的阈值β
79.
80
             int split_point = node->split->ord.split_point;
             int n1 = node->get_num_valid(vi); //该节点具有vi特征属性的样本数量
81.
82
             //为该特征属性开辟一块内存空间
83.
             cv::AutoBuffer<uchar> inn_buf(n*(Sizeof(int)*
      (data->have_priors ? 3 : 2) + sizeof(float)));
84.
             float* val_buf = (float*)(uchar*)inn_buf;
             int* sorted_buf = (int*)(val_buf + n);
85.
86
             int* sample_idx_buf = sorted_buf + n;
87.
             const float* val = 0;
             const int* sorted = 0;
88.
             //对该节点的所有样本,按照特征属性vi的值的从小到大的顺序进行排序,分别得到两个数组:val和sorted,这个
      数组的索引值表示排序后的索引值,而val值为特征属性的值,sorted值为未排序前的样本索引值
90
             data->get_ord_var_data( node, vi, val_buf, sorted_buf, &val, &sorted, sample_idx_buf);
91.
             assert( 0 <= split_point && split_point < n1-1 );</pre>
92.
93
             if(!data->have_priors) //无先验概率
94
95
96.
                 for( i = 0; i <= split_point; i++ )</pre>
                    dir[sorted[i]] = (char)-1; //样本为左分支,赋值为-1
97.
98
                 for( ; i < n1; i++ )</pre>
                    dir[sorted[i]] = (char)1;
99.
                                              //样本为右分支,赋值为1
100.
                 for( ; i < n; i++ )</pre>
101.
                    dir[sorted[i]] = (char)0;
                                               //缺失的样本,赋值为0
                 //L和R分别为左右分支的样本数量
102.
```

第24页 共46页 2017/11/30 上午10:51

```
103.
                 L = split_point-1;
104.
                 R = n1 - split point + 1;
105
             else
                    //有先验概率
106.
107.
                                                                     //得到先验概率
108.
                 const double* priors = data->priors_mult->data.db;
                 int* responses_buf = sample_idx_buf + n;
109.
                 //responses指向该节点样本的分类,即响应值
110
111.
                 const int* responses = data->qet class labels(node, responses buf);
112.
                 L = R = 0:
113.
                 //遍历左分支的所有样本
114.
                 for( i = 0; i <= split_point; i++ )</pre>
115.
                     int idx = sorted[i];
116.
117.
                     double w = priors[responses[idx]];
                                                       //得到该样本所在的分类的先验概率
                     dir[idx] = (char)-1;
118.
                                           //左分支赋值为-1
                     L += w: //左分支的加权样本数量
119.
120.
121.
      <u>...</u>
                 //遍历右分支的所有样本
122.
                 for( ; i < n1; i++ )</pre>
123.
                     int idx = sorted[i];
124.
125.
                     double w = priors[responses[idx]];
                                                        //得到该样本所在的分类的先验概率
                     dir[idx] = (char)1;
126.
                                         //右分支赋值为1
127.
                     R += w;
                             //右分支的加权样本数量
128.
                 //遍历缺失的样本
129.
130.
                 for( ; i < n; i++ )</pre>
131.
                     dir[sorted[i]] = (char)0;
132.
             }
          node->maxlr = MAX(L, R);
                                    //表示左右分支最大值
134.
135
          //返回的值为式12或式19的值除以左右分支数量之和,即规范化常数
          return node->split->quality/(L + R);
136.
137. }
```

特征属性为数值形式的寻找替代分叉属性的函数:

0

```
1.
     CvDTreeSplit* CvDTree::find_surrogate_split_ord( CvDTreeNode* node, int vi, uchar* _ext_buf )
 2.
         const float epsilon = FLT EPSILON*2;
                                                //定义一个较小的常数
 3.
         //得到特征属性的样本是属于左分支(-1),右分支(1),还是缺失(0)
         const char* dir = (char*)data->direction->data.ptr;
 5.
 6
         //n表示该节点的样本数量,ni表示具有vi特征属性的样本数量
         int n = node->sample_count, n1 = node->get_num_valid(vi);
 7.
         //开辟一块内存空间
 8.
 9.
         cv::AutoBuffer<uchar> inn_buf;
         if( !\_ext\_buf )
10.
11.
             inn\_buf.allocate(\ n^*(\textbf{Sizeof}(\textbf{int})^*(\textbf{data->have\_priors}\ ?\ 3\ :\ 2)\ +\ \textbf{Sizeof}(\textbf{float}))\ );
         uchar* ext_buf = _ext_buf ? _ext_buf : (uchar*)inn_buf;
         float* values_buf = (float*)ext_buf;
13.
         int* sorted_indices_buf = (int*)(values_buf + n);
15.
         int* sample indices buf = sorted indices buf + n;
16.
         const float* values = 0;
         const int* sorted_indices = 0;
         //对该节点的所有样本,按照特征属性vi的值的从小到大的顺序进行排序,分别得到两个数组:values和
18.
     sorted_indices,这个数组的索引值表示排序后的索引值,而values值为特征属性的值,sorted_indices值为未排序前的
19.
         data->get_ord_var_data( node, vi, values_buf, sorted_indices_buf, &values, &sorted_indices,
20.
         {\prime\prime}{\prime} LL - number of samples that both the primary and the surrogate splits send to the left
         // LR - ... primary split sends to the left and the surrogate split sends to the right
21.
22.
         /\!/ RL - ... primary split sends to the right and the surrogate split sends to the left
23.
         // RR - ... both send to the right
24.
         // best i表示替代分叉点,即阈值B
25.
         int i, best_i = -1, best_inversed = 0;
         double best_val; //表示替代分叉值
26.
27.
         if(!data->have priors) //无先验概率
28.
29.
```

第25页 共46页 2017/11/30 上午10:51

```
//四个变量LL, RL, LR, RR由两个字母组成,前一个字母表示被最佳分叉属性分叉为左右分支,后一个字母表示被
30.
     替代分叉属性分叉为左右分支,即LL表示既属于最佳分叉属性分叉为左分支,又属于替代分叉属性分叉为左分支的样本数量,LR
     表示属于最佳分叉属性分叉为左分支,属于替代分叉属性分叉为右分支的样本数量,RL表示属于最佳分叉属性分叉为右分支,属
     于替代分叉属性分叉为左分支的样本数量,RR表示既属于最佳分叉属性分叉为右分支,又属于替代分叉属性分叉为右分支的样本
31.
            //也就是样本被分割成了LL, RL, LR, RR这四个部分
32.
            int LL = 0, RL = 0, LR, RR;
33
            // worst_val为二叉树左右分支中样本数量最多的值,它是一个阈值,替代分叉值只有大于该值才会被考虑
34.
            int worst_val = cvFloor(node->maxlr), _best_val = worst_val;
    \int_{\Gamma}
35
            int sum = 0, sum_abs = 0;
            //遍历所有样本
36
     1
37.
            for( i = 0; i < n1; i++ )</pre>
38
               //d表示方向,-1是左分支,1是右分支,0是缺失
39.
40
               int d = dir[sorted_indices[i]];
               //sum表示d累加求和,因为d也可能为负值,所以sum的含义为右分支比左分支多出的样本数;sum_abs表示d
41
     的绝对值之和,表示的含义为被分叉的样本数
42
               sum += d; sum_abs += d & 1;
43.
     [•••
           }
44
45
            // sum abs = R + L; sum = R - L
     &
            //LR和RR分别初始化为被最佳分叉属性分叉为左右分支的样本数,但对于替代分叉属性来说,所有样本都属于了右分
46.
47.
            RR = (sum\_abs + sum) >> 1;
48
            LR = (sum\_abs - sum) >> 1;
49
            // initially all the samples are sent to the right by the surrogate split,
50.
51.
            // LR of them are sent to the left by primary split, and RR - to the right.
52.
            // now iteratively compute LL, LR, RL and RR for every possible surrogate split value.
53.
            //按照该特征属性值从小到大的顺序遍历所有样本
            for( i = 0; i < n1 - 1; i++ )
55.
            {
56
               int d = dir[sorted_indices[i]];
                                             //取得该样本的方向
57.
58.
               if( d < 0 )
               //该样本被最佳分叉属性分叉为左分支,而且由于该样本位于best_i之前,所以也被替代分叉属性分叉为左分
59
     支,因此该样本属于LL
60.
61.
                   //把原本初始化为LR的样本重新划归到LL内
62.
                  LL++: LR--:
                   //比较并更新替代分叉点best_i和替代分叉值best_val
63
                   if( LL + RR > _best_val && values[i] + epsilon < values[i+1] )</pre>
64.
65
                      best_val = LL + RR;
66.
67.
                      best_i = i; best_inversed = 0;
68
                   }
69.
70.
               else if( d > 0 )
71.
               //该样本被最佳分叉属性分叉为右分支,而且由于该样本位于best_i之前,所以是被替代分叉属性分叉为左分
     支 . 因此该样本属于RL
72
               {
                   //把原本初始化为RR的样本重新划归到RL内
73.
74
                   RL++; RR--;
75.
                   //比较并更新替代分叉点best_i, 计算替代分叉值best_val
76.
                   if( RL + LR > _best_val && values[i] + epsilon < values[i+1] )</pre>
77
                   {
                      best val = RL + LR;
78.
79
                      best_i = i; best_inversed = 1;
80.
                  }
81.
               }
            //此处是否有误?应该为_best_val = best_val
83.
84
           best_val = _best_val;
                                //更新
               //有先验概率
        else
86.
87
88.
            double LL = 0, RL = 0, LR, RR;
89.
            // worst_val为二叉树左右分支中样本数量最多的值
90.
            double worst_val = node->maxlr;
91.
            double sum = 0, sum abs = 0;
92.
            const double* priors = data->priors_mult->data.db;
                                                            //获得先验概率
93.
            int* responses buf = sample indices buf + n;
94.
            //得到样本响应值
95.
            const int* responses = data->get_class_labels(node, responses_buf);
96.
            best val = worst val;
```

第26页 共46页

```
//遍历所有样本,得到带有先验概率的sum和sum_abs
97.
98.
              for( i = 0; i < n1; i++ )</pre>
99
                 int idx = sorted_indices[i];
100
101.
                 double w = priors[responses[idx]];
102.
                 int d = dir[idx];
                 sum += d*w; sum_abs += (d \& 1)*w;
103.
104.
105
      凸
106
              // sum abs = R + L; sum = R - L
              //LR和RR分别初始化为被最佳分叉属性分叉为左右分支的数量,但对于替代分叉属性来说,所有样本都属于了右分
107.
      ₹
108
              RR = (sum\_abs + sum)*0.5;
              LR = (sum\_abs - sum)*0.5;
109.
110.
111.
              // initially all the samples are sent to the right by the surrogate split,
              // LR of them are sent to the left by primary split, and RR - to the right.
112.
113.
              // now iteratively compute LL, LR, RL and RR for every possible surrogate split value.
      [•••
114.
              //按照该特征属性值从小到大的顺序遍历所有样本
115.
              for( i = 0; i < n1 - 1; i++ )</pre>
116.
      &
                 int idx = sorted_indices[i];
117.
118.
                 //得到该样本所属分类的先验概率
119.
                 double w = priors[responses[idx]];
120.
                 int d = dir[idx]; //取得该样本的方向
121.
                 if( d < 0 )
122.
123.
                 //该样本被最佳分叉属性分叉为左分支,而且由于该样本位于best_i之前,所以也被替代分叉属性分叉为左分
      支,因此该样本属于LL
124.
                 {
                     LL += w; LR -= w;
125.
126.
                     if( LL + RR > best_val && values[i] + epsilon < values[i+1] )</pre>
127
                         best_val = LL + RR;
128.
129.
                         best_i = i; best_inversed = 0;
130.
131.
132.
                 else if( d > 0)
                 //该样本被最佳分叉属性分叉为右分支,而且由于该样本位于best_i之前,所以是被替代分叉属性分叉为左分
133.
      支,因此该样本属干RL
134
135.
                     RL += w: RR -= w:
136
                     if( RL + LR > best_val && values[i] + epsilon < values[i+1] )</pre>
137.
138.
                         best val = RL + LR:
139
                         best_i = i; best_inversed = 1;
140.
141.
                 }
142.
             }
143.
144.
          //如果找到了替代分叉属性,则返回该替代分叉属性
145.
          return best i >= 0 && best val > node->maxlr ? data->new split ord( vi,
146.
              (values[best\_i] + values[best\_i+1])*0.5f, best\_i, best\_inversed, (float)best\_val ) : 0; \\
147. }
```

特征属性为类形式的寻找替代分叉属性的函数:

0

```
[cpp]
     CvDTreeSplit* CvDTree::find_surrogate_split_cat( CvDTreeNode* node, int vi, uchar* _ext_buf )
1.
2.
 3.
         const char* dir = (char*)data->direction->data.ptr;
         int n = node->sample count;
                                      //节点的样本数
 4.
 5
         //mi表示该特征属性的种类数量
         int i, mi = data->cat count->data.i[data->get var type(vi)], 1 win = 0;
 6.
 7.
         //为vi开辟一块内存空间
 8.
         int base_size = (2*(mi+1)+1)*sizeof(double) + (!data->have_priors ? 2*(mi+1)*sizeof(int)
9.
         cv::AutoBuffer<uchar> inn_buf(base_size);
10.
         if( ! ext buf )
             inn_buf.allocate(base_size + n*(sizeof(int) + (data->have_priors ? sizeof(int) : 0)))
11.
12.
         uchar* base_buf = (uchar*)inn_buf;
```

第27页 共46页

```
13.
         uchar* ext_buf = _ext_buf ? _ext_buf : base_buf + base_size;
14.
15
         int* labels_buf = (int*)ext_buf;
         //label指向该特征属性中各个样本所对应的种类
16.
17.
         const int* labels = data->get_cat_var_data(node, vi, labels_buf);
18.
         {\prime\prime}{\prime} LL - number of samples that both the primary and the surrogate splits send to the left
19.
         // LR - ... primary split sends to the left and the surrogate split sends to the right
         // RL - ... primary split sends to the right and the surrogate split sends to the left
20.
     // RR - ... both send to the right //定义一个CvDTreeSplit变量
21.
22.
23.
         CvDTreeSplit* split = data->new_split_cat( vi, 0 );
      1 double best_val = 0;
24.
25.
         double* lc = (double*)cv::alignPtr(base_buf, sizeof(double)) + 1;
     double* rc = lc + mi + 1;
26.
         //初始化数组lc和rc,这两个数组的长度都为mi+1,作用是保存左右分支该特征属性每个种类的样本数量
27.
28.
         for( i = -1; i < mi; i++ )
     lc[i] = rc[i] = 0;
29.
30.
     \boxed{\bullet \bullet \bullet} // for each category calculate the weight of samples
31.
       \checkmark // sent to the left (lc) and to the right (rc) by the primary split
32.
         //第一步,先计算被最佳分叉属性分叉的该特征属性的左右分支各个种类的样本数量
      if(!data->have_priors)
34.
                                   //无先验概率
35
         {
             //此处是否有误?应该为int* _lc = (int*)rc + mi + 1;
36.
37.
             int* _lc = (int*)rc + 1;
             int* _rc = _lc + mi + 1;
38
             //初始化数组_lc和_rc,这两个数组的长度都为mi+1,作用是_lc保存着各个特征属性种类的方向信息的累加
39.
     和,_rc保存着各个特征属性种类的方向信息的绝对值累加和
40.
             for( i = -1; i < mi; i++ )</pre>
                 _lc[i] = _rc[i] = 0;
41.
             //遍历所有样本
42.
             for( i = 0; i < n; i++ )</pre>
43.
44.
                 //得到该样本的特征属性的种类
45.
46.
                 int idx = ( (labels[i] == 65535) && (data->is_buf_16u) ) ? -1 : labels[i];
                                  //得到该样本的方向信息
47.
                 int d = dir[i];
                 //计算方向信息的累加和sum,以及绝对值的累加和sum abs
48.
49.
                 int sum = lc[idx] + d;
50.
                 int sum_abs = _rc[idx] + (d & 1);
                                                      //更新数组 lc和 rc
51.
                 _lc[idx] = sum; _rc[idx] = sum_abs;
52.
             //遍历各个特征属性的种类,计算左右分支每个特征属性的种类的样本数量
53.
54
             for( i = 0; i < mi; i++ )</pre>
55.
             {
56.
                 int sum = _lc[i];
57.
                 int sum_abs = _rc[i];
                 lc[i] = (sum abs - sum) >> 1:
58.
59.
                 rc[i] = (sum_abs + sum) >> 1;
60.
             }
61.
62.
         else
                 //有先验概率
63.
         {
64.
             const double* priors = data->priors_mult->data.db;
                                                                  //获得先验概率
65.
             int* responses_buf = labels_buf + n;
             //得到样本的响应值
66.
67
             const int* responses = data->get_class_labels(node, responses_buf);
             //遍历所有样本
68.
69.
             for( i = 0; i < n; i++ )
70.
                 //得到该样本的特征属性的种类
71.
72.
                 int idx = ( (labels[i] == 65535) && (data->is_buf_16u) ) ? -1 : labels[i];
                 double w = priors[responses[i]];
                                                  //得到该样本响应值的先验概率
73.
74.
                 int d = dir[i];
                                  //方向信息
                 //计算方向信息的累加和sum以及绝对值累加和sum_abs
                 double sum = lc[idx] + d*w;
76.
77
                 double sum_abs = rc[idx] + (d & 1)*w;
78.
                 lc[idx] = sum; rc[idx] = sum_abs;
79.
80.
             //遍历各个特征属性的种类,计算左右分支每个特征属性的种类的数量
             for( i = 0; i < mi; i++ )</pre>
81.
82.
83.
                 double sum = lc[i]:
84.
                 double sum_abs = rc[i];
                 lc[i] = (sum\_abs - sum) * 0.5;
                 rc[i] = (sum\_abs + sum) * 0.5;
86.
```

第28页 共46页 2017/11/30 上午10:51

```
87.
88.
89
90.
         // 2. now form the split.
91.
         \ensuremath{/\!/} in each category send all the samples to the same direction as majority
92.
         //第二步,计算替代分叉值
93.
         //遍历特征属性的各个种类
94
         for( i = 0; i < mi; i++ )</pre>
95.
      Ľ.
96
             //得到该特征属性种类的左右分支的数量
97.
             double lval = lc[i], rval = rc[i];
      1
98.
             if(lval > rval) //左分支大于右分支
99
             {
                 //在subset数组内的相应位置1,而如果左分支小于右分支,则相应位为0;被置1的种类将被替代分叉属性分
100.
      叉为左分支,被清0的种类将被替代分叉属性分叉为右分支
                 split->subset[i >> 5] |= 1 << (i & 31);
101
                 best_val += lval; //累加替代分叉值
102.
103
                 l_win++;
104.
105
             e1se
                    //左分支小于右分支
                 best_val += rval;
                                   //累加替代分叉值
107.
108
         split->quality = (float)best_val; //替代分叉值
109.
110.
         if( split->quality <= node->maxlr || l_win == 0 || l_win == mi )
111.
             cvSetRemoveByPtr( data->split_heap, split ), split = 0;
112.
113.
         return split;
                        //返回
114. }
```

把节点分叉,并完成相关的运算:

()

```
void CvDTree::split_node_data( CvDTreeNode* node )
 1.
2.
3.
         //n为该节点的样本数,scount为训练的所有样本数,nl和nr分别表示左分支和右分支的样本数
         int vi, i, n = node->sample_count, nl, nr, scount = data->sample_count;
         //样本的方向信息,即属于哪个分支,-1为左分支,1为右分支,0为缺失
 5.
         char* dir = (char*)data->direction->data.ptr;
 6
         CvDTreeNode *left = 0, *right = 0;
 7.
 8.
         int* new_idx = data->split_buf->data.i;
         int new_buf_idx = data->get_child_buf_idx( node );
10.
         // work_var_count = 样本的特征属性的数量 + 1(如果是分类树的话)+ 1(如果应用交叉验证的话)
         int work_var_count = data->get_work_var_count();
11
         CvMat* buf = data->buf;
12.
13.
         size_t length_buf_row = data->get_length_subbuf();
14
         \label{eq:cv::AutoBuffer} \textit{cv::AutoBuffer} < \textit{uchar} > \textit{inn\_buf}(\textit{n}^*(3*\textit{sizeof}(\textit{int}) + \textit{sizeof}(\textit{float})));
         int* temp_buf = (int*)(uchar*)inn_buf;
15.
         //用替代分叉属性完成对该节点的分叉,即把该节点分叉为左分支或右分支,该函数在后面有讲解
16
17.
         complete_node_dir(node);
18.
         //遍历该节点的所有样本
         for( i = nl = nr = 0; i < n; i++ )
19
20.
21.
             //得到该样本的方向信息,这里的方向信息是左分支为0,右分支为1
22.
             int d = dir[i];
23.
             // initialize new indices for splitting ordered variables
             //数组new_idx保存着样本i所在分支的顺序索引,如i属于左分支,则new_idx[i]等于当前的nl,i属于右分支
24
     时 . new idx[i]等于当前的nr
25
             new_idx[i] = (nl & (d-1)) | (nr & -d); // d ? ri : li
             nr += d; //累加d,含义是计算右分支的样本数
26.
27.
             nl += d^1; //d与1异或,并累加,含义是计算左分支的样本数
28
29.
30
         bool split_input_data;
                                  //分叉以后的标识
31.
         //定义左右分支的节点变量
32.
         node->left = left = data->new_node( node, nl, new_buf_idx, node->offset );
33.
         node->right = right = data->new_node( node, nr, new_buf_idx, node->offset + nl );
         //判断目前决策树的深度和左右分支的样本数,如果满足要求,则该标识为1,否则为0
34.
35.
         split_input_data = node->depth + 1 < data->params.max_depth &&
             (node->left->sample_count > data->params.min_sample_count ||
36.
37.
             node->right->sample_count > data->params.min_sample_count);
```

2017/11/30 上午10:51 第29页 共46页

```
38.
 39.
                    // split ordered variables, keep both halves sorted.
 40
                    //遍历样本的所有形式是数值的特征属性,对特征属性是数值形式的进行分叉
                    for( vi = 0; vi < data->var_count; vi++ )
 41.
 42.
 43.
                           int ci = data->get_var_type(vi);
                           //如果是类形式的特征属性,或分叉标识为0,则继续下一个特征属性的循环
 44.
 45.
                           if( ci >= 0 || !split_input_data )
                                  continue;
 46.
            凸
 47.
                           //n1表示在特征属性vi下的样本数量
 48.
                           int n1 = node->get_num_valid(vi);
             1
 49.
                           float* src_val_buf = (float*)(uchar*)(temp_buf + n);
 50
                           int* src_sorted_idx_buf = (int*)(src_val_buf + n);
                           int* src_sample_idx_buf = src_sorted_idx_buf + n;
 51.
 52.
                           const float* src_val = 0;
                           const int* src_sorted_idx = 0;
 53
                          //对该节点的所有样本,按照特征属性vi的值的从小到大的顺序进行排序,分别得到两个数组:src_val和
 54.
            src_sorted_idx,这个数组的索引值表示排序后的索引值,而src_val值为特征属性的值,src_sorted_idx值为未排序前
            的样本索引值
 55
                           data->get_ord_var_data(node, vi, src_val_buf, src_sorted_idx_buf, &src_val, &src_sorted_idx_buf, &src_val, &src_sorted_idx_buf, &src_val_buf, 
                           //初始化temp_buf数组为src_sorted_idx数组
             &
                           for(i = 0; i < n; i++)</pre>
 57.
 58
                                   temp_buf[i] = src_sorted_idx[i];
 59.
 60.
                           if (data->is_buf_16u)
                                                                      //所有样本的数量小于65536
 61.
                                  unsigned short *ldst, *rdst, *ldst0, *rdst0;
 62.
 63.
                                  //unsigned short tl, tr;
 64.
                                  ldst0 = ldst = (unsigned short*)(buf->data.s + left->buf_idx*length_buf_row +
 65.
                                         vi*scount + left->offset);
                                  rdst0 = rdst = (unsigned short*)(ldst + nl);
 67.
 68
                                  // split sorted
                                  //按照顺序遍历该节点的样本
 69.
 70.
                                  for( i = 0; i < n1; i++ )</pre>
  71.
                                         int idx = temp_buf[i];
 72.
                                         int d = dir[idx]; //得到该样本的方向信息
 73.
 74.
                                         idx = new_idx[idx];
                                                                                //idx为当前样本所在分支的顺序索引
                                         if (d)
 75.
                                                          //右分支
 76.
                                                  *rdst = (unsigned short)idx;
                                                                                                           //赋值
 77.
 78
                                                 rdst++:
                                                                    //指向下一个地址
 79.
                                         else
 80.
                                                        //左分支
 81.
                                         {
                                                  *ldst = (unsigned short)idx;
                                                                                                            //赋值
 82.
 83.
                                                 ldst++;
                                                                  //指向下一个地址
 84.
                                         }
 85.
                                  //设置该特征属性vi下的左右分支的样本数
                                  left->set_num_valid(vi, (int)(ldst - ldst0));
 87.
 88
                                  right->set_num_valid(vi, (int)(rdst - rdst0));
 89.
                                  // split missing
 90.
 91.
                                  for( ; i < n; i++ )</pre>
                                                                          //分叉缺失该特征属性的样本
 92.
 93
                                         int idx = temp_buf[i];
 94.
                                         int d = dir[idx];
                                         idx = new_idx[idx];
 95.
                                         if (d)
 97.
                                         {
 98
                                                  *rdst = (unsigned short)idx;
                                                 rdst++;
100.
                                         }
101.
                                         else
102.
103.
                                                 *ldst = (unsigned short)idx;
104.
                                                 ldst++;
105.
                                         }
106.
                                  }
107.
108.
                           else
                                          //样本数大于65536
109.
                           //与if(data->is_buf_16u)的情况相似,唯一的不同就是分配的内存空间大小不同
110.
```

第30页 共46页

```
int *ldst0, *ldst, *rdst0, *rdst;
111.
112.
                  ldst0 = ldst = buf->data.i + left->buf_idx*length_buf_row +
113.
                      vi*scount + left->offset;
                  rdst0 = rdst = buf->data.i + right->buf_idx*length_buf_row +
114.
115.
                      vi*scount + right->offset;
116.
117.
                  // split sorted
118.
                  for( i = 0; i < n1; i++ )</pre>
119.
120
                      int idx = temp_buf[i];
                      int d = dir[idx];
121.
122.
                      idx = new_idx[idx];
123.
                      if (d)
124.
125.
                           *rdst = idx:
126.
                          rdst++;
127.
                      }
128.
                      else
129.
                      {
130.
                           *ldst = idx;
131.
                          ldst++;
132.
                      }
133.
                  }
134.
135.
                  left->set_num_valid(vi, (int)(ldst - ldst0));
136.
                  right->set_num_valid(vi, (int)(rdst - rdst0));
137.
138.
                  // split missing
139.
                  for( ; i < n; i++ )</pre>
140.
                      int idx = temp_buf[i];
141.
142.
                      int d = dir[idx];
143
                      idx = new_idx[idx];
                      if (d)
144.
145.
                           *rdst = idx;
146.
147.
                          rdst++:
148.
                      else
149.
150.
151.
                           *ldst = idx;
152.
                          ldst++:
153
154.
                  }
155.
              }
156.
157.
158
          // split categorical vars, responses and cv_labels using new_idx relocation table
159.
          //遍历样本的所有形式是类的特征属性,对特征属性是类形式的,以及响应值和交叉验证进行分叉
          for( vi = 0; vi < work_var_count; vi++ )</pre>
160.
161.
          {
162.
              int ci = data->get_var_type(vi);
163
              //n1表示在特征属性vi下的样本数量
164.
              int n1 = node->get_num_valid(vi), nr1 = 0;
              //如果是数值形式的特征属性,或分叉标识为0,则继续下一个特征属性的循环
165.
166
              if( ci < 0 || (vi < data->var_count && !split_input_data) )
                  continue;
167.
168.
169.
              int *src_lbls_buf = temp_buf + n;
              // src_lbls指向该特征属性中各个样本所对应的种类
170.
171.
              const int* src_lbls = data->get_cat_var_data(node, vi, src_lbls_buf);
              //把数组temp_buf初始化为src_lbls
172.
173.
              for(i = 0; i < n; i++)
                  temp_buf[i] = src_lbls[i];
175.
176
              if (data->is_buf_16u)
177.
178.
                  unsigned short *ldst = (unsigned short *)(buf->data.s + left->buf_idx*length_buf_ro
179.
                      vi*scount + left->offset);
180.
                  unsigned short *rdst = (unsigned short *)(buf->data.s + right->buf_idx*length_buf_r
181.
                      vi*scount + right->offset);
                  //遍历所有样本
182.
183.
                  for( i = 0; i < n; i++ )
184.
185.
                      int d = dir[i];
                                        //得到方向信息
```

第31页 共46页

```
//得到该样本的特征属性的种类
186.
                      int idx = temp_buf[i];
                      if (d)
187.
                                //右分支
188
189.
                           *rdst = (unsigned short)idx;
190.
                          rdst++;
                                    //指向下一个地址
                          nr1 += (idx != 65535 )&d;
                                                        //右分支计数
191.
192.
                      3
193.
                      else
                             //左分支
194.
                      {
195
                           *ldst = (unsigned short)idx;
                                    //指向下一个地址
196.
                          ldst++;
197.
                      }
198.
                  // data->var_count为样本的特征属性数量
199.
200.
                  if( vi < data->var_count )
201.
                  {
                      //设置该特征属性vi下的左右分支的样本数
202.
203
                      left->set_num_valid(vi, n1 - nr1);
204.
                      right->set_num_valid(vi, nr1);
205.
                  }
206.
       &
              else
207.
208
              //与if(data->is_buf_16u)的情况相似,唯一的不同就是分配的内存空间大小不同
209.
210.
                  int *ldst = buf->data.i + left->buf_idx*length_buf_row +
211.
                      vi*scount + left->offset;
                  int *rdst = buf->data.i + right->buf_idx*length_buf_row +
212.
213.
                      vi*scount + right->offset;
214.
215.
                  for( i = 0; i < n; i++ )</pre>
216.
217.
                      int d = dir[i];
218
                      int idx = temp_buf[i];
                      if (d)
219.
220.
                      {
221.
                           *rdst = idx;
222.
                          rdst++:
223.
                          nr1 += (idx >= 0)&d;
224.
                      }
                      else
225.
226.
                      {
                           *ldst = idx;
227.
228.
                          ldst++:
229.
                      }
230.
231.
                  }
232.
233.
                  if( vi < data->var_count )
234.
                  {
                      left->set_num_valid(vi, n1 - nr1);
235
236
                      right->set_num_valid(vi, nr1);
237.
                  }
238
              }
239.
          }
240.
241.
          // split sample indices
242.
243.
          int *sample_idx_src_buf = temp_buf + n;
244.
          //得到样本序列索引
          const int* sample_idx_src = data->get_sample_indices(node, sample_idx_src_buf);
245.
246.
          //初始化temp_buf数组为sample_idx_src数组
247.
          for(i = 0; i < n; i++)</pre>
248.
              temp_buf[i] = sample_idx_src[i];
249.
          int pos = data->get_work_var_count();
250.
251
          if (data->is_buf_16u)
252.
253.
              unsigned short* ldst = (unsigned short*)
      (buf->data.s + left->buf_idx*length_buf_row +
                  pos*scount + left->offset);
254.
255
              unsigned short* rdst = (unsigned short*)
      (buf->data.s + right->buf idx*length buf row +
256
                  pos*scount + right->offset);
257.
               //遍历所有样本
258.
              for (i = 0; i < n; i++)</pre>
```

第32页 共46页 2017/11/30 上午10:51

```
259.
                                    //得到该样本的方向信息
260.
                  int d = dir[i];
261
                  unsigned short idx = (unsigned short)temp_buf[i];
                  if (d)
                           //右分支
262.
263.
264
                       *rdst = idx;
                                      //样本索引
                      rdst++; //指向下一个地址
265.
266
267.
                         //左分支
                  else
268
269.
                      *ldst = idx;
270.
                      ldst++;
271.
272.
              }
273.
274
          else
       🖊 //与if(data->is_buf_16u)的情况相似,唯一的不同就是分配的内存空间大小不同
275.
276
277.
              int* ldst = buf->data.i + left->buf_idx*length_buf_row +
278.
                  pos*scount + left->offset;
279.
              int* rdst = buf->data.i + right->buf_idx*length_buf_row +
                  pos*scount + right->offset;
280.
281
              for (i = 0; i < n; i++)
282.
283.
                  int d = dir[i];
284
                  int idx = temp_buf[i];
                  if (d)
285.
286.
287.
                       *rdst = idx;
288.
                      rdst++;
290.
                  else
291
                       *ldst = idx;
292.
293.
                      ldst++;
294
295.
              }
296.
297.
          \ensuremath{/\!/} deallocate the parent node data that is not needed anymore
298.
299.
          data->free_node_data(node); //释放该节点
300. }
```

用替代分叉属性对节点进行分叉:

0

```
[cpp]
     \textbf{void} \ \texttt{CvDTree}{::} complete\_node\_dir(\ \texttt{CvDTreeNode*} \ node \ )
1.
2.
        //n表示该节点的样本数
3.
4
        int vi, i, n = node->sample\_count, nl, nr, d0 = 0, d1 = -1;
        // node->split->var_idx为该节点的最佳分叉属性,get_num_valid是取得最佳分叉属性下的样本数量
        //nz表示缺失该特征属性(即为最佳分叉属性)的样本的数量,nz等于0说明该节点的所有样本都拥有该特征属性,所以都
6.
     可以应用最佳分叉属性进行分叉,nz不等于0说明在该节点内有一些样本不具有最佳分叉属性,因此对这些样本就必须使用替代
     分叉属性进行分叉
7.
        int nz = n - node->get_num_valid(node->split->var_idx);
8.
        char* dir = (char*)data->direction->data.ptr; //取得样本的方向信息
9.
10
        // try to complete direction using surrogate splits
        //如果有缺失样本,则对这些样本使用替代分叉属性进行分叉
11.
12.
        if( nz && data->params.use_surrogates )
13.
            cv::AutoBuffer<uchar> inn_buf(n*(2*sizeof(int)+sizeof(float)));
14.
15.
            CvDTreeSplit* split = node->split->next;
                                                   //得到替代分叉属性
            //按照替代分叉属性值从大到小的顺序遍历所有的替代分叉属性
16.
17.
            for( ; split != 0 && nz; split = split->next )
18.
                int inversed_mask = split->inversed ? -1 : 0;
                                                           //表示是否反转方向
19.
20.
                vi = split->var_idx;
                                    //得到属于该替代分叉属性的特征属性
21.
22.
                if( data->get_var_type(vi) >= 0 ) // split on categorical var
```

第33页 共46页 2017/11/30 上午10:51

```
//该特征属性是类的形式
23.
24.
                {
25
                   int* labels_buf = (int*)(uchar*)inn_buf;
                   //labels指向该特征属性中各个样本所对应的种类
26
27.
                   const int* labels = data->get_cat_var_data(node, vi, labels_buf);
                   //subset数组的每个位表示一个特征属性的种类,该位置1,说明左分支样本数量大于右分支,反之则该
28
     位为0
29
                   const int* subset = split->subset;
30
                   //遍历所有样本
     \int_{\Gamma}
31
                   for( i = 0; i < n; i++ )</pre>
32.
     1
33.
                       int idx = labels[i];
                                            //特征属性的种类
34
                       if( !dir[i] && ( ((idx >= 0)&&(!data->is_buf_16u)) || ((idx != 65535)&&
     (data->is_buf_16u)) ))
35.
                       // dir[i]等于0,说明该样本还没有被分叉,即它缺失最佳分叉属性,因此开始应用替代分叉属性
     进行分叉
36
37
                           //由替代分叉属性确定该样本的分叉情况,左分支为-1,右分支为1
38.
     [•••
                          int d = CV_DTREE_CAT_DIR(idx, subset);
39
                          dir[i] = (char)((d \land inversed\_mask) - inversed\_mask);
40
                           //nz累减,当nz为0时,说明缺失的样本已全部分叉完,则退出遍历所有样本循环,并且也退出
     通替代分叉属性的循环
41
                           if( --nz )
42.
                              break;
43.
                       }
44
                   }
45.
46
                else // split on ordered var
47.
                //该特征属性是数值的形式
48.
                {
                   float* values_buf = (float*)(uchar*)inn_buf;
                   int* sorted indices buf = (int*)(values buf + n);
50.
51
                   int* sample_indices_buf = sorted_indices_buf + n;
                   const float* values = 0;
52.
                   const int* sorted_indices = 0;
53.
                   //对该节点的所有样本,按照特征属性vi的值的从小到大的顺序进行排序,分别得到两个数组:values和
     sorted_indices,这个数组的索引值表示排序后的索引值,而values值为特征属性的值,sorted_indices值为未排序前的
     样本索引值
55.
                   data->get_ord_var_data( node, vi, values_buf, sorted_indices_buf, &values, &sor
                   int split_point = split->ord.split_point;
56.
                   //n1表示在特征属性vi下的样本数量
57.
                   int n1 = node->get_num_valid(vi);
58.
59
                   //确保阈值β正确
                   assert( 0 <= split_point && split_point < n-1 );</pre>
60.
                   //按顺序遍历所有样本
61.
                   for( i = 0; i < n1; i++ )</pre>
62
63.
64
                       int idx = sorted_indices[i];
65.
                       // dir[i]等于0,说明该样本还没有被分叉,即它缺失最佳分叉属性,因此开始应用替代分叉属性
     讲行分叉
66
                       if( !dir[idx] )
67.
68
                           //左分支为-1,右分支为1
69.
                          int d = i <= split_point ? -1 : 1;</pre>
                          //方向赋值
70.
71.
                          dir[idx] = (char)((d ^ inversed_mask) - inversed_mask);
                           // nz累减,当nz为0时,说明缺失的样本已全部分叉完,则退出遍历所有样本循环,并且也退
72.
     出遍历替代分叉属性的循环
73.
                           if( --nz )
74.
                              break;
75
                       }
76.
                   }
77.
               }
78
            }
        }
79.
80
        // find the default direction for the rest
81.
82.
        //仍然有样本没有被分叉,则按照当前左、右分支的样本数进行分叉
83.
        if( nz )
84.
        {
85
            for( i = nr = 0; i < n; i++ )</pre>
               nr += dir[i] > 0; //右分支的数量
86.
87.
            nl = n - nr - nz; //左分支的数量
88.
            //nl大于nr,d0=-1;nl小于nr,d0=1;nl等于nr,d0=0
            d0 = nl > nr ? -1 : nr > nl;
89.
```

第34页 共46页 2017/11/30 上午10:51

46. 47.

ab->data.db[tree_count] = min_alpha;

```
90.
 91.
 92
          \ensuremath{//} make sure that every sample is directed either to the left or to the right
          //确保所有的样本都被分叉
 93.
 94.
          for( i = 0; i < n; i++ )</pre>
 95.
              int d = dir[i]; //得到方向,-1、1或0
 96.
 97
              if(!d) //d=0时,由d0值来确定它属于左分支还是右分支
 98.
       ď
 99
                 d = d0:
                  if(!d) //当d0=0时,交替为d赋值为-1和1
100.
       1
101.
                     d = d1, d1 = -d1;
 102.
              d = d > 0;
                         //方向信息重新定义:d=-1时,d为0;d=1时,d为1
103.
104
              //此时方向信息改为,左分支为0,右分支为1,这么做是为了后续程序的方便
105.
              dir[i] = (char)d; // remap (-1,1) to (0,1)
106.
107.
基于CCP算法的剪枝处理函数:
       &
()
       [cpp]
       void CvDTree::prune cv()
  1.
  2.
          CvMat* ab = 0;
                          //保存决策树序列的最小α值
  3.
          CvMat* temp = 0;
  4
          //保存交叉验证各个子集的决策树序列的错误率,即式28和式29
          CvMat* err_jk = 0;
  6.
  7.
  8.
          // 1. build tree sequence for each cv fold, calculate error_{Tj,beta_k}.
  9.
          // 2. choose the best tree index (if need, apply 1SE rule).
          // 3. store the best index and cut the branches.
 11.
 12
          CV_FUNCNAME( "CvDTree::prune_cv" );
 13.
 14.
          __BEGIN__;
 15.
          //tree_count表示构建的决策树序列的数量,cv_n表示交叉验证的子集数量,n表示根节点的样本数量
 16.
          int ti, j, tree_count = 0, cv_n = data->params.cv_folds, n = root->sample_count;
 17.
          // currently, 1SE for regression is not implemented
 18.
          //判断是否应用1SE规则,目前0penCV只实现了分类树的1SE算法
 19.
          bool use_1se = data->params.use_1se_rule != 0 && data->is_classifier;
          double* err;
          double min_err = 0, min_err_se = 0;
 21.
 22
          int min_idx = -1;
 23.
          CV_CALL(ab = cvCreateMat(1, 256, CV_64F));
 24.
                                                      //创建α的矩阵
 25.
 26.
          // build the main tree sequence, calculate alpha's
 27.
          //死循环,作用是构建CCP算法中的一系列树的序列{T0,T1,...,Tm},它是通过标注node->Tn的大小实现的,在该循环结
       束后形成的树序列中,构成Ti的所有节点的Tn ≤ i。在进入该死循环之前每个节点的Tn都为一个最大值。退出该死循环的条件
       是,得到了只有一个根节点的Tm
 28
          for(;;tree_count++)
                               //这里的tree_count表示树序列Ti的下标i
 29.
 30
              //计算Ti的最小α
 31.
              double min_alpha = update_tree_rnc(tree_count, -1);
 32.
              //标注Ti的各个节点的Tn值
 33.
              if( cut_tree(tree_count, -1, min_alpha) )
                  break:
                          //得到了Tm.则退出死循环
 34.
 35.
              //如果所构建的决策树序列的数量大于ab矩阵的行,则需要扩充ab矩阵,从而可以保存更多的最小\alpha值
 36.
              if( ab->cols <= tree_count )</pre>
 37.
 38.
                  //创建temp矩阵,比ab矩阵大
 39.
                  CV CALL( temp = cvCreateMat( 1, ab->cols*3/2, CV 64F ));
 40.
                  for( ti = 0; ti < ab->cols; ti++ )
 41.
                     temp->data.db[ti] = ab->data.db[ti];
                  cvReleaseMat( &ab );
 42.
                                      //释放ab矩阵
 43.
                  ab = temp;
                             //赋值
                  temp = 0;
                            //清零
 44.
 45.
```

第35页 共46页 2017/11/30 上午10:51

//保存最小α值

```
48.
         }
49.
50
         ab - data.db[0] = 0.; //\alpha'0 = 0
51.
52.
         if( tree count > 0 )
                              //如果得到了决策树序列
53.
             //遍历决策树序列, 计算式30
54.
55.
             for( ti = 1; ti < tree_count-1; ti++ )</pre>
56.
                 ab->data.db[ti] = sqrt(ab->data.db[ti]*ab->data.db[ti+1]);
      凸
57.
             //最后一个决策树序列的\alpha'm值赋于一个最大的值
58.
             ab->data.db[tree_count-1] = DBL_MAX*0.5;
      1
59.
                                                                      //创建矩阵
60
             CV_CALL( err_jk = cvCreateMat( cv_n, tree_count, CV_64F ));
             err = err_jk->data.db; //指向err_jk矩阵的首地址
61.
62.
63.
             //遍历交叉验证的所有子集,得到每个子集的决策树序列
      h
64.
             for( j = 0; j < cv_n; j++ )
65
             {
66.
                //tj表示当前子集的决策树序列的索引,即式31中下标j的含义;tk表示全体样本决策树序列的索引,即式31
      中下标k的含义
67
                 int tj = 0, tk = 0;
      &
                 //在当前子集内遍历决策树序列,按照式31找到最终的子集决策树序列
68.
69
                 for( ; tk < tree_count; tj++ )</pre>
70.
71.
                    double min_alpha = update_tree_rnc(tj, j); //得到\alpha值
                    if( cut_tree(tj, j, min_alpha) )
72.
                        min_alpha = DBL_MAX; //最大值,表示只有一个根节点的决策树
73.
74.
                    //按照tk的含义进行遍历
75.
                    for( ; tk < tree_count; tk++ )</pre>
76.
                    {
                        if( ab->data.db[tk] > min_alpha )
77.
78.
                           break: //退出for循环
79
                        //保存当前子树的总错误率,它等于子树中所有叶节点错误率之和
                        err[j*tree_count + tk] = root->tree_error;
80.
81.
                    }
82.
83.
84.
             //遍历树序列,得到最佳决策树
85.
             for( ti = 0; ti < tree_count; ti++ )</pre>
86.
                 double sum_err = 0;
                 //计算所有子集中树序列相同的子树的错误率
88.
89
                 for( j = 0; j < cv_n; j++ )
                    sum_err += err[j*tree_count + ti]; //式32
90.
                 if( ti == 0 || sum_err < min_err )
                                                  //得到最优的子树索引
91.
92.
                    min_err = sum_err; //更新最小错误率
93.
94
                    min_idx = ti; //得到子树的索引,式33
95.
                    if( use_1se )
                                   //应用1SE规则
                        min_err_se = sqrt( sum_err*(n - sum_err) );
96.
                                                                  //式34
97
                 else if( sum_err < min_err + min_err_se )</pre>
98.
                                                          //式35
                    //由于遍历树序列是按照从复杂到简单的顺序,因此越是后面的循环,所表示的决策树越简单,所以满足
99
      式35,而最终的决策树一定是最简单的
100.
                    min_idx = ti;
101.
         }
102.
103.
104.
         pruned_tree_idx = min_idx;
                                    //赋值
         //如果truncate_pruned_tree为true,则真正去掉那些被剪枝掉的节点
105.
         free_prune_data(data->params.truncate_pruned_tree != 0);
107.
108
          _END__;
         //释放三个矩阵
110.
         cvReleaseMat( &err_jk );
111.
         cvReleaseMat( &ab );
112.
         cvReleaseMat( &temp );
113. }
```

计算风险值的增加率 α 函数,其中参数fold为-1,表示计算右全部样本构建的决策树序列;fold大于等于0, 则计算交叉验证子集的决策树序列,而此时的fold表示子集的索引

()

```
[cpp]
1.
     double CvDTree::update_tree_rnc( int T, int fold )
2.
        CvDTreeNode* node = root; //根节点
3.
4.
        double min_alpha = DBL_MAX; //初始化为一个最大值
        for(;;)
                 //死循环
5.
6.
        //从根节点出发遍历决策树序列中第T个序列数的各个节点,搜索一遍后,退出该死循环
    \mathbb{L}_{\epsilon}
7.
                                //定义一个父节点
            CvDTreeNode* parent;
8.
9.
                     //死循环,从中间节点沿着左分支遍历决策树
            for(;;)
            //要想退出该死循环,必须满足下列条件之一:1、到达了叶节点;2、叶节点Tn或子集的树索引的节点cv_Tn小于形
10.
     参T_,即只遍历树序列中T的前一个树
11.
           {
               //当形参fold为负数时,t为node->Tn,否则t为node->cv_Tn[fold]
12.
13
               int t = fold >= 0 ? node->cv_Tn[fold] : node->Tn;
14.
     h
               if( t <= T | !node->left ) //退出死循环的两个条件
15.
                   //赋值叶节点的复杂度、风险值和错误率
16.
    <u>...</u>
                   node->complexity = 1; //节点的复杂度,即叶节点的数量
17.
18
                   node->tree_risk = node->node_risk;
                                                   //节点的风险值
     &
                   node->tree_error = 0.; //节点的错误率
19.
20.
                  if( fold >= 0 )
21.
                  {
                      //是子集的情况,则风险值和错误率被赋值为相应子集的相关值
22.
23.
                      node->tree_risk = node->cv_node_risk[fold];
24.
                      node->tree_error = node->cv_node_error[fold];
25.
                           //退出for死循环
26.
               3
27.
28.
               node = node->left;
                                 //得到该节点的左节点
29.
30.
            //从叶节点沿着右分支向中间节点遍历
31.
            for( parent = node->parent; parent && parent->right == node;
32.
               node = parent, parent = parent->parent )
33
               //赋值中间节点的复杂度、风险值和错误率,即该中间节点的所有子节点之和
34.
35.
               parent->complexity += node->complexity; //中间节点的复杂度累加
               parent->tree_risk += node->tree_risk; //中间节点的风险值累加
36.
               parent->tree_error += node->tree_error;
                                                   //中间节点的错误率累加
37.
38
               //计算α值,式22,这里需要注意的是,式22中的预测误差用风险值代替,即式23或式24中的分子部分代替整
     个分式, 我认为这么做的目的可能是为了减少误差, 提高效率
39.
               parent->alpha = ((fold >= 0 ? parent->cv\_node\_risk[fold] : parent->node\_risk)
40.
                   - parent->tree_risk)/(parent->complexity - 1);
               min\_alpha = MIN( min\_alpha, parent->alpha ); //得到最小的\alpha值
41.
42.
43.
44
            if( !parent )
                         //搜索到了根节点,则退出for死循环
            //此时,中间节点parent的左分支的所有子节点都已遍历过,下一步就是遍历parent的右分支的所有节点
46.
47.
            parent->complexity = node->complexity;
48.
            parent->tree_risk = node->tree_risk;
49.
            parent->tree_error = node->tree_error;
            node = parent->right;
50.
51.
52.
        return min_alpha; //返回最小的α值
53.
54. }
```

剪切节点,表现的形式是为Tn变量赋不同大小的值:

()

```
1.
    int CvDTree::cut_tree( int T, int fold, double min_alpha )
2.
        CvDTreeNode* node = root;
3.
        if(!node->left) //根节点没有左分支,则直接返回
4.
5.
           return 1;
6.
        for(;;) //死循环
7.
        //从根节点出发遍历决策树序列中第T个序列数的各个节点,搜索一遍后,退出该死循环
8.
9.
10.
           CvDTreeNode* parent;
                              //定义一个父节点
```

```
11.
           for(;;)
           //要想退出该死循环,必须满足下列条件之一:1、到达了叶节点;2、叶节点Tn或子集的树索引的节点cv Tn小于形
12.
     13.
14.
               int t = fold >= 0 ? node->cv Tn[fold] : node->Tn;
15.
               if( t <= T || !node->left )
                                        //退出死循环的前两个条件
                  break:
16.
               if( node->alpha <= min_alpha + FLT_EPSILON )</pre>
                                                       //退出死循环的第3个条件
17.
18
19
                  if( fold >= 0 )
20
                      node->cv_Tn[fold] = T;
                  else
21.
22
                      node -> Tn = T;
                  if( node == root )
23.
                                  //直接返回
24.
                      return 1;
25
                  break;
     h
               }
26.
27
               node = node->left;
                                 //下一个左节点
     <u>...</u>
28.
29.
           //从叶节点沿着右分支向中间节点遍历,不做其他任何操作
30.
           for( parent = node->parent; parent && parent->right == node;
     &
               node = parent, parent = parent->parent )
31.
32.
33.
34
           if( !parent )
                          //到达的了根节点,则退出for死循环
               break;
35
36.
37.
           node = parent->right;
                                //下一个右分支
38.
39.
        return 0;
                    //返回
41. }
```

下面我们梳理一下构建决策树的过程,也就是训练的过程:

do train函数实现构建决策树的作用,该函数主要就是递归调用try split node函数,完成决策树的构建,如 果需要剪枝,则再调用prune cv函数。

在try_split_node函数内,每次分叉之前,要用calc_node_value函数计算节点的样本分类的先验概率、节点 的值、节点的风险值和节点的错误率。然后判断退出递归的条件:节点的深度、节点的样本数,如果是分 类树,则是否为一个分类,如果是回归树,节点的风险值。如果上述条件不满足则退出递归。下一步就是 关键的调用计算最佳分叉属性的函数——find best split。如果还需要替代分叉属性,则根据特征属性是类 的形式还是数值的形式,分别调用find_surrogate_split_cat函数和find_surrogate_split_ord函数。由最佳分 叉属性得到的样本分叉方向信息是由calc_node_dir函数完成的,而由替代分叉属性得到的样本分叉方向信 息是由complete_node_dir函数实现的。分叉后的一些变量赋值是由split_node_data函数完成的。最后就是 分别对左分支和右分支递归调用try_split_node函数。

在计算最佳分叉属性时,是在find_best_split函数内,并行处理所有的特征属性。实现的程序是在 DTreeBestSplitFinder类的重载运算符()内,如果特征属性是类的形式的分类树,则调用 find split cat class函数;如果特征属性是数值的形式的分类树,则调用find split ord class函数,如果特 征属性是类的形式的回归树,则调用find_split_cat_reg函数,如果特征属性是数值的形式的回归树,则调用 find split ord reg函数,最终得到最佳分叉属性。

下面我们给出用决策树进行预测的函数。这里还需要说明一点的是,当特征属性是类的形式的时候,我们 还需要对这种形式进行量化处理,即按照构建决策树时训练样本的形式,把它们转换为: 0, 1, 2, ...。例如 在表示风力时,我们是把无风,小风,中风和大风量化为0, 1, 2, 3,则在预测样本时,如果当前风力是小 风,则它的输入值应写为1。当然不对其进行量化处理也是可以的,只要按照当初训练样本时,样本数据的 写法就可以。上述内容是通过参数preprocessed_input来识别的。

```
1.
    CvDTreeNode* CvDTree::predict( const CvMat* _sample,
       \textbf{const} \ \texttt{CvMat*} \ \underline{\texttt{missing}}, \ \textbf{bool} \ \texttt{preprocessed\_input} \ ) \ \textbf{const}
2.
    //_sample表示待预测样本,必须是32FC1类型的一维矢量形式,即一次只能预测一个样本,矢量中元素的数量必须与训练样本
    的特征属性的数量一致(具体还取决于参数preprocessed_input),并且要与训练样本的特征属性的排序相同
4.
    // missing表示缺失特征属性,为可选项,该参数必须是掩码的形式,1表示缺失,它的形式与参数 sample
    // preprocessed_input表示对特征属性是类的形式的数据是否进行预处理,即量化处理,如果为false,则表示正常的输
    入;如果为true,则表示进行了预处理,进行预处理的好处是可以加快预测速度
6.
    //该函数把最终到达的叶节点作为输出
7.
    {
8.
        cv::AutoBuffer<int> catbuf;
```

第38页 共46页 2017/11/30 上午10:51

```
9.
10.
        int i, mstep = 0;
11.
        const uchar* m = 0;
        CvDTreeNode* node = root;
                                 //得到决策树
12.
13.
14
        if( !node )
                      //尚没有构建好决策树,则报错
            CV_Error( CV_StsError, "The tree has not been trained yet" );
15.
    if( !CV_IS_MAT(_sample) || CV_MAT_TYPE(_sample->type) != CV_32FC1 || (_sample->cols != 1 && comple
16
17.
18.
19.
            (_sample->cols + _sample->rows - 1 != data->var_all && !preprocessed_input) ||
20.
            (\_sample->cols + \_sample->rows - 1 != data->var\_count \&\& preprocessed\_input) )
21.
        • // data->var_all表示训练时全部特征属性的数量,data->var_count表示训练时不包括缺失的特征属性的数量
22.
              CV_Error( CV_StsBadArg,
23.
            "the input sample must be 1d floating-point vector with the same " \,
            "number of elements as the total number of variables used for training" );
24.
25.
26
        const float* sample = _sample->data.fl; //指向预测样本数据的首地址
27.
     ••• //得到步长
       int step = CV_IS_MAT_CONT(_sample->type) ? 1 : _sample->step/sizeof(sample[0]);
28.
29.
         //为类的形式的特征属性开辟一块内存空间
     if (data->cat_count && !preprocessed_input ) // cache for categorical variables
30.
31
32.
            int n = data->cat count->cols:
33.
            catbuf.allocate(n);
34
            for( i = 0; i < n; i++ )
               catbuf[i] = -1; //初始化
35.
36.
37.
        if(_missing) //如果该样本缺失的某种特征属性
38.
39.
40.
            //检查输入参数 missing是否正确
41.
            if( !CV_IS_MAT(_missing) || !CV_IS_MASK_ARR(_missing) ||
                !CV_ARE_SIZES_EQ(_missing, _sample) )
42.
43.
                CV_Error( CV_StsBadArg,
44.
            "the missing data mask must be 8-bit vector of the same size as input sample" );
            m = _missing->data.ptr; //指向表示缺失特征属性的掩码首地址
45.
46.
            //得到步长
47.
            mstep = CV_IS_MAT_CONT(_missing->type) ? 1 : _missing->step/sizeof(m[0]);
48.
49.
        const int* vtvpe = data->var tvpe->data.i;
                                                  //训练样本特征属性的类型
50.
51.
        //预测样本要与训练样本的特征属性顺序一致,即第几个位置对应于哪个特征属性是固定的
        const int* vidx = data->var_idx && !preprocessed_input ? data->var_idx->data.i : 0;
52.
        //特征属性为类形式的各个属性的映射值,即当初表示该属性的数值
53.
54
        const int* cmap = data->cat_map ? data->cat_map->data.i : 0;
        //特征属性为类形式的各个属性的偏移量
55.
56.
        const int* cofs = data->cat_ofs ? data->cat_ofs->data.i : 0;
57.
        //遍历最佳决策树,进行样本的预测,达到了叶节点则预测结束
58.
        while( node->Tn > pruned_tree_idx && node->left )
59
60.
            CvDTreeSplit* split = node->split;
                                             //节点分叉
61.
            int dir = 0;
                         //表示分叉的方向信息,-1是左节点,1是右节点,
            //对预测样本进行分叉预测,第一次循环是应用最佳分叉属性进行分叉,以后都是应用替代分叉属性进行分叉
62.
            for( ; !dir && split != 0; split = split->next )
63.
64
                                        //得到的是该节点的分叉属性
                int vi = split->var_idx;
65.
66
                int ci = vtype[vi]; //分叉属性的类形索引,是数值形式还是类的形式
67.
                i = vidx ? vidx[vi] : vi; //得到分叉属性对应的特征属性索引
                float val = sample[(size_t)i*step]; //得到预测样本的分叉属性的值
68.
                if( m && m[(size_t)i*mstep] )
                               //该特征属性缺失,继续下一个循环,即应用替代分叉属性
70.
                   continue;
                if(ci < 0) // ordered 分叉属性为数值的形式
71.
                   // split->ord.c为分叉属性的值,如果预测样本的分叉属性的值val小于该值,则被分叉为左节点,否
     则为右节点
73.
                   dir = val <= split->ord.c ? -1 : 1;
74.
                else // categorical 分叉属性为类的形式
75.
76.
77.
                   if( preprocessed_input )
                                           //对预测样本数据已做预处理
78.
                      c = cvRound(val); //对预测样本响应值进行四舍五入取整
                   else //未对预测样本数据进行预处理,则需处理
79.
80.
81.
                       c = catbuf[ci];
                                      //c初始化为-1
                       if( c < 0 )
82.
```

```
83.
                       {
                                                //再次赋值,a为分叉属性的偏移量
84.
                           int a = c = cofs[cil:
85
                           //通常情况下,b为下一个特征属性的偏移量
                          int b = (ci+1 >= data->cat_ofs->cols) ? data->cat_map->cols : cofs[ci+1
86
87.
88
                          int ival = cvRound(val);
                           //输入的预测样本的特征属性为类的形式的值时必须是整数
89.
90
                          if( ival != val )
91.
                              CV_Error( CV_StsBadArg,
     \int_{\Gamma}
92
                              "one of input categorical variable is not an integer" );
93.
94.
                          int sh = 0;
95
                           //循环,一点点更新a和b值,直到找到分叉属性对应的量化值
96.
                           while( a < b )
97.
98
                              c = (a + b) >> 1; //取a和b的平均值
99.
100.
                              if( ival < cmap[c] )</pre>
                                                   //该属性小于c的映射值
101.
                                 b = c; //更新b值
102
                              else if( ival > cmap[c] )
                                                        //该属性大于c的映射值
                                 a = c+1; //更新a值
                              else //两值相等,则退出while循环
104.
105
                                 break;
106.
107
                          //如果没有得到合适的c值,则放弃该分叉属性,进行for循环,即应用下一个分叉属性
108.
                           if( c < 0 || ival != cmap[c] )
                              continue;
109.
110.
                           //得到真正的特征属性对应的值,即c要减去偏移量
111.
                          catbuf[ci] = c -= cofs[ci];
112.
                       }
                   //判断c的数据长度是否正确
114.
115
                   c = ((c == 65535) \&\& data->is_buf_16u)? -1 : c;
                   // split->subset是按位存储着不同特征属性的方向信息,通过宏定义CV_DTREE_CAT_DIR得到c对应
116.
      的方向信息,并赋值给dir
117.
                   dir = CV_DTREE_CAT_DIR(c, split->subset);
118.
                3
119.
120.
                if( split->inversed )
                                     //如果分叉方向需要反转
                   dir = -dir;
                                //反转方向
121.
122.
            //如果通过上面的for循环操作没有能够对预测样本进行分叉,那么就按照当初创建决策树时该节点的左、右分支的
123.
      样本数量进行分叉,哪个分支的样本多,就属于哪个分支
124.
            if( !dir )
125.
            {
126.
                //得到右分支与左分支的样本数量之差
                double diff = node->right->sample count - node->left->sample count;
127.
128.
                //左分支样本数量多,则dir为-1;反之dir为1
129.
                dir = diff < 0 ? -1 : 1;
130.
131.
            // dir为-1,指向左节点;dir为1,指向右节点,进入下一个while循环
            node = dir < 0 ? node->left : node->right;
132.
133.
134.
                       //返回叶节点
135.
         return node;
136.
```

在CvDTree类内, get var importance是一个很有用的函数,它可以判断在构建决策树的过程中,各个特征 属性的重要程度,它是通过计算分叉属性的分叉值得到的:

0

```
[qqa]
    const CvMat* CvDTree::get_var_importance()
1.
2.
        // var importance为CvDTree类的全局变量,用以表示各个特征属性的重要程度
3.
4.
        if( !var_importance )
                             //还没有对var_importance变量赋值
5.
           CvDTreeNode* node = root;
6.
                                    //根节点
           double* importance;
           if(!node) //还没有构建决策树,则返回该函数
8.
               return 0;
9.
```

第40页 共46页

```
//定义var_importance变量为一维数组的形式,数组内元素的个数就是特征属性的数量
10.
            var\_importance = cvCreateMat( 1, data->var\_count, CV\_64F );
11.
12.
            cvZero( var_importance );
                                    //初始化var_importance变量为0
            importance = var_importance->data.db;
                                               //指向var_importance数组的首地址
13.
            //遍历整个决策树,统计分叉属性的分叉值作为特征属性重要程度的衡量标准
14.
                     //死循环
15.
            for(;;)
16.
17.
               CvDTreeNode* parent;
                                    //父节点
               for(;; node = node->left ) //由父节点沿着左分支向叶节点遍历
18.
     凸
19
                   CvDTreeSplit* split = node->split;
20.
     1
                   //到达了叶节点,则退出for循环
21.
22.
                   if( !node->left || node->Tn <= pruned_tree_idx )</pre>
23.
                      break;
                   //统计节点的各个特征属性的所有分叉值(即式12或式19),该值作为特征属性的重要程度的度量,先计
24.
     算最佳分叉值,则计算其他的替代分叉值
25.
                   for( ; split != 0; split = split->next )
26
                      importance[split->var_idx] += split->quality;
     ···
27.
28.
               //由叶节点沿着右分支向父节点遍历
               for( parent = node->parent; parent && parent->right == node;
     &
30.
                   node = parent, parent = parent->parent )
31.
32.
33.
               if( !parent )
                              //到达了根节点,则退出for死循环
34.
                   break;
35.
36.
               node = parent->right;
                                    //下一个右分支
37.
38.
            //归一化var_importance变量
            cvNormalize( var_importance, var_importance, 1., 0, CV_L1 );
40.
41.
        return var_importance;
                                //返回var_importance变量
42.
43. }
```

三、应用实例

下面我们以贷款申请是否通过为例,它取决于5个条件:年龄(青年Y,中年M,老年O)、薪水(低L,中 M, 高H) 、有房子(有Y, 无N) 、有车(有Y, 无N) 和信贷情况(一般F, 好G, 优E) 。而分类结果N 表示贷款申请没通过,Y表示通过。具体的样本见下表(仅仅是虚构的例子,不具有任何实际作用):

序号	年龄	薪水	有房子	有车	信贷情况	发放贷款
1	Υ	L	N	N	F	Ν
2	Υ	L	Υ	N	G	Ν
3	Υ	М	Υ	Ν	G	Y
4	Υ	М	Υ	Υ	G	Υ
5	Υ	I	Υ	Υ	G	Υ
6	Υ	Μ	Ζ	Υ	G	Ζ
7	М	ш	Υ	Υ	E	Υ
8	М	I	Υ	Υ	G	Υ
9	М	L	N	Υ	G	Ν
10	М	М	Υ	Υ	F	N
11	М	Ι	Υ	Υ	E	Υ
12	М	Μ	Ν	Z	G	Ν
13	0	ш	Ν	Z	G	Ν
14	0	Ы	Υ	Υ	E	Υ
15	0	Ы	Υ	Z	E	Υ
16	0	М	N	Υ	G	Ν
17	0	L	N	Ν	E	Ν
18	0	Ι	N	Υ	F	Ν
19	0	Н	Υ	Υ	E	Y

2017/11/30 上午10:51 第41页 共46页

构建上表所示的决策树是一颗分类树,而且所有的特征属性都是类的形式,下面就是用以上样本构建决策 树并预测的程序为:

0

```
include "opencv2/core/core.hpp"
1.
     #include "opencv2/highgui/highgui.hpp"
2.
     #include "opencv2/imgproc/imgproc.hpp"
3.
     #Include "opencv2/ml/ml.hpp"
5.
6.
     #include <iostream>
     using namespace cv;
using namespace std;
7.
8.
9.
     //5个特征属性的描述
     static const char* var_desc[] =
10.
11.
     "Age (young=Y, middle=M, old=0)",
"Salary? (low=L, medium=M, high=H)",
12.
13.
14.
         "Own_House? (false=N, true=Y)",
         "Own_Car? (false=N, true=Y)",
15.
16
         "Credit_Rating (fair=F, good=G, excellent=E)",
17.
     };
18.
19.
     int main( int argc, char** argv )
20.
21.
22.
         float trainingData[19][5]={ {'Y', 'L', 'N', 'N', 'F'},
23.
24.
                               {'Y','L','Y','N','G'},
                               {'Y', 'M', 'Y', 'N', 'G'},
25.
                               {'Y','M','Y','Y','G'},
26
27.
                               {'Y', 'H', 'Y', 'Y', 'G'},
                               {'Y', 'M', 'N', 'Y', 'G'},
28.
29
                               {'M','L','Y','Y','E'},
                               {'M', 'H', 'Y', 'Y', 'G'},
30.
                               {'M','L','N','Y','G'},
31.
                               {'M', 'M', 'Y', 'Y', 'F'},
32.
                               {'M','H','Y','Y','E'},
33.
34.
                               {'M','M','N','N','G'},
35.
                               {'O', 'L', 'N', 'N', 'G'},
                               {'0','L','Y','Y','E'},
36.
                               {'0','L','Y','N','E'},
                               {'O', 'M', 'N', 'Y', 'G'},
38.
39
                               {'O', 'L', 'N', 'N', 'E'},
                               {'O', 'H', 'N', 'Y', 'F'},
40.
                               {'0','H','Y','Y','E'}
41.
         Mat trainingDataMat(19, 5, CV_32FC1, trainingData);
42.
                                                           //样本的矩阵形式
         //样本的分类结果,即响应值
43.
         44.
45.
         Mat responsesMat(19, 1, CV_32FC1, responses); //矩阵形式
46.
         float priors[5] = {1, 1, 1, 1, 1}; //先验概率,这里的每个特征属性的作用都是相同
47.
48.
         //定义决策树的参数
49.
         CvDTreeParams params( 15,
                                    // 决策树的最大深度
                             1, //决策树叶节点的最小样本数
50.
51.
                                 //回归精度,这里不需要
                                     //是否使用替代分叉属性,由于没有缺失的特征属性,所以这里不需要替代分
52.
     叉属性
53.
                             25.
                                  //最大的类数量
54.
                                  // 交叉验证的子集数,由于样本太少,这里不需要交叉验证
                             false.
55.
                                     //使用1SE规则,这里不需要
                             false,
                                      //是否真正的去掉被剪切的分支,这里不需要
56
                             priors
                                      //先验概率
57.
58
         //类形式的掩码,这里是分类树,而且5个特征属性都是类的形式,因此该变量都为1
59.
60.
        Mat varTypeMat(6, 1, CV_8U, Scalar::all(1));
61.
        CvDTree* dtree = new CvDTree(); //实例化CvDTree类
62.
63.
         //训练样本,构建决策树
64.
         dtree->train ( trainingDataMat,
                                           //训练样本
                                      //样本矩阵的行表示样本,列表示特征属性
65.
                       CV_ROW_SAMPLE,
```

2017/11/30 上午10:51 第42页 共46页

```
//样本的响应值矩阵
  66.
                          responsesMat,
                          Mat(), //应用所有的特征属性
  67.
  68.
                          Mat(),
                                   //应用所有的训练样本
  69.
                          varTypeMat, //类形式的掩码
                          Mat(), //没有缺失任何特征属性
  70.
  71.
                          params
                                  //决策树参数
  72.
                          );
  73.
           //调用get_var_importance函数
       Const CvMat* var_importance = dtree->get_var_importance(); //输出特征属性重要性程度
  74.
  75.
           for( int i = 0; i < var_importance->cols*var_importance->rows; i++ )
  76.
       1 {
  77.
  78.
               double val = var_importance->data.db[i];
  79.
               char buf[100];
  80.
               int len = (int)(strchr( var_desc[i], '(' ) - var_desc[i] - 1);
  81.
               strncpy( buf, var_desc[i], len );
               buf[len] = '\0';
  82.
  83.
               printf( "%s", buf );
               printf( ": %g%%\n", val*100. );
  84.
  85.
  86.
        float myData[5] = {'M', 'H', 'Y', 'N', 'F'}; //预测样本
  87.
  88.
           Mat myDataMat(5, 1, CV_32FC1, myData); //矩阵形式
           double r = dtree->predict( myDataMat, Mat(), false)->value;
                                                                       //得到预测结果
  89.
  90.
  91.
           cout<<endl<<"result: "<<(char)r<<endl;</pre>
                                                   //输出预测结果
  92.
  93.
           return 0;
  94.
  95. }
则最终的输出为:
Age: 5.32416%
Salary?: 5.18066%
Own House?: 40.8999%
Own_Car?: 23.6091%
Credit_Rating: 24.9861%
```

result: N

我们也可以用下列代码把决策树保存为xml文件:

dtree->save("dtree.xml");

如果要应用该决策树,只需要加载该文件即可,如:

CvDTree* dtree = new CvDTree();

dtree->load("dtree.xml");

Д

相关文章推荐

2017/11/30 上午10:51 第43页 共46页

决策树算法ID3、C4.5、CART (http://blog.csdn.net/taigw/article/details/44840771)

决策树是机器学习中非常经典的一类学习算法、它通过树的结构、利用树的分支来表示对样本特征的判断规则、从树的 叶子节点所包含的训练样本中得到预测值。决策树如何生成决定了所能处理的数据类型和预测性能。主要的决...



🧝 taigw (http://blog.csdn.net/taigw) 2015年04月03日 01:02 👊11116

\int_{Γ}

ID3、C4.5、CART三种决策树的区别 (http://blog.csdn.net/qq_27717921/article/details/7...

很早就想写写决策树,说起决策树做过数据挖掘的就不会感觉陌生,但是可能对ID3决策树算法、C4.5决策树算法以及C ART决策树之间的区别不太了解,下面就这三个比较著名的决策树算法分别写写决策树是如何...



🧿 qq_27717921 (http://blog.csdn.net/qq_27717921) 2017年07月09日 21:01 🖽562





不止20K, Python薪酬又飙升了??

2017年 Pytyhon薪酬曝光啦!看完后薪资报告后,同事说了一句:人生苦短,不学Python算白活....

(http://www.baidu.com/cb.php?c=IgF pyfqnHmknjnvPjc0IZ0qnfK9ujYzP1f4PjDs0Aw-

5Hc3rHnYnHb0TAq15HfLPWRznjb0T1YYmyD1nAFhnjmdm1nsnWu90AwY5HDdnHcsnH0knHR0IgF 5y9YIZ0IQzquZR8mLPbUB48ugfElAqspynETZ-YpAq8nWqdlAdxTvqdThP-

5yF_UvTkn0KzujY4rHb0mhYqn0KsTWYs0ZNGujYkPHTYn1mk0AqGujYknWb3rjDY0APGujYLnWm4n1c0ULl85H00TZbqnW0v0APzm1Ykn1T4P6)

决策树ID3、C4.5、CART算法: 信息熵, 区别, 剪枝理论总结 (http://blog.csdn.net/ljp8...

今天学习了决策树算法中的ID3、c4.5、CART算法,记录如下: 决策树算法:顾名思义,以二分类问题为例,即利用 自变量构造一颗二叉树,将目标变量区分出来,所有决策树算法的关键点如下:



决策树剪枝算法 (http://blog.csdn.net/yujianmin1990/article/details/49864813)

剪枝作为决策树后期处理的重要步骤,是必不可少的。没有剪枝,就是一个完全生长的决策树,是过拟合的。需要去掉 一些不必要的节点以使得决策树模型更具有泛化能力。...



🥻 yujianmin1990 (http://blog.csdn.net/yujianmin1990) 2015年11月16日 13:36 爲9936

决策树之剪枝原理与CART算法 (http://blog.csdn.net/u014688145/article/details/53326910)

决策树学习笔记(二)继续关于决策树的内容,本篇文章主要学习了决策树的剪枝理论和基于二叉树的CART算法。主 要内容: 1.理解决策树损失函数的定义以及物理含义 2.基尼指数的主要两个作用 3.理解CA...



🤦 u014688145 (http://blog.csdn.net/u014688145) 2016年11月24日 22:10 🕮 6456



程序员跨越式成长指南

完成第一次跨越, 你会成为具有一技之长的开发者, 月薪可能翻上几番; 完成第二次跨越, 你将 成为拥有局部优势或行业优势的专业人士,获得个人内在价值的有效提升和外在收入的大幅跃迁...

(http://www.baidu.com/cb.php?c=IgF pyfqnHmknjfzrjD0IZ0qnfK9ujYzP1f4PjnY0Aw-5Hc4nj6vPjm0TAq15Hf4rjn1n1b0T1Y4myw9m1uhn1l9P1u-

mWTd0AwY5HDdnHcsnH0knHR0lgF_5y9YIZ0lQzqMpgwBUvqoQhP8QvIGIAPCmgfEmvq_lyd8Q1R4uWc4uHf3uAckPHRkPWN9PhcsmW9huWqdIAdxTvqdThP-5HDknWFBmhkEusKzujY4rHb0mhYqn0KsTWYs0ZNGujYkPHTYn1mk0AqGujYkn10snjf10APGujYLnWm4n1c0ULl85H00TZbqnW0v0APzm1Yvnjbzr0)

第44页 共46页

机器学习笔记(九)——决策树的生成与剪枝 (http://blog.csdn.net/chunyun0716/articl...

基本的决策树生成算法主要有ID3和C4.5,它们生成树的过程大致相似,ID3是采用的信息增 一、决策树的牛成算法 益作为特征选择的度量,而C4.5采用信息增益比。构建过程如下: ...

🍆 chunyun0716 (http://blog.csdn.net/chunyun0716) 2016年05月08日 11:43 🛛 🖺 1796

决策树剪枝算法 (http://blog.csdn.net/ritchiewang/article/details/50254009)

算法目的: 决策树的剪枝是为了简化决策树模型,避免过拟合。算法基本思路: 减去决策树模型中的一些子树或者叶结 点,并将其根结点作为新的叶结点,从而实现模型的简化。模型损失函数 ...

🝙 ritchiewang (http://blog.csdn.net/ritchiewang) 2015年12月10日 21:19 🕮 1046

 $[\cdots]$

C4.5决策树 –剪枝的问题 (http://blog.csdn.net/PANHUBO/article/details/52014513)

C4.5决策树——剪枝的问题 一。先说点什么 今天是周日,作为实习生在公司不加班,只想多学点什么。同时 对未来也是有点迷茫, 就说这么多吧。 二。切...

🥁 PANHUBO (http://blog.csdn.net/PANHUBO) 2016年07月24日 16:15 🛛 1836

决策树剪枝的方法与必要性 (http://blog.csdn.net/yeting067/article/details/38614621)

zhuan 1 决策树剪枝的必要性本文讨论的决策树主要是基于ID3算法实现的离散决策树生成。ID3算法的基本思想是贪心 算法,采用自上而下的分而治之的方法构造决策树。首先检测训练数据集的所有特征,...

【机器学习】决策树(下)——CART算法及剪枝处理 (http://blog.csdn.net/HerosOfEar...

CART, 即分类与回归树(classification and regression tree), 也是一种应用很广泛的决策树学习方法...

💮 HerosOfEarth (http://blog.csdn.net/HerosOfEarth) 2016年09月03日 22:35 🖺1658



Delphi7高级应用开发随书源码 (http://download.csdn.net/detail/chenxh...

2003年04月30日 00:00 676KB

决策树归纳一般框架(ID3, C4.5, CART) (http://blog.csdn.net/u012314976/article/det...

构建决策树的目的是对已有的数据进行分类,得到一个树状的分类规则,然后就可以拿这个规则对未知的数据进行分类 预测。决策树归纳是从有类标号的训练元祖中学习决策树。决策树是一种类似于流程图的树结构,其中每...

🤛 u012314976 (http://blog.csdn.net/u012314976) 2014年12月31日 10:02 🕮 3584

机器学习自学之路-SVM 算法选择:三种算法优缺点比较(ID3、C4.5、CART) (http://...

ID3D3算法十分简单,核心是根据"最大信息熵增益"原则选择划分当前数据集的最好特征,信息熵是信息论里面的概念, 是信息的度量方式,不确定度越大或者说越混乱,熵就越大。在建立决策树的过程中,根据特征属性...

🍘 github_39261590 (http://blog.csdn.net/github_39261590) 2017年08月01日 16:55 🕮421

决策树(ID3、C4.5、CART、随机森林) (http://blog.csdn.net/gumpeng/article/details/...

本篇文章也是多个博客的综合,有待进一步整理。1、决策树-预测隐形眼镜类型(ID3算法,C4.5算法,CART算法,GI

第45页 共46页

NI指数,剪枝,随机森林) 1. 我们应该设计什么的算法,使得计算...

R语言-决策树算法(C4.5和CART)的实现 (http://blog.csdn.net/Hz_ZDeveloper/article/...

决策树算法的实现:一、C4.5算法的实现 a、需要的包: sampling、party library(sampling) library(party) sampling用于实

● Hz_ZDeveloper (http://blog.csdn.net/Hz_ZDeveloper) 2017年04月27日 09:53 □1907

09#R语事实现决策树分析 (http://blog.csdn.net/han2538740718/article/details/52300605)

决策树是附加概率结果的一个树状的决策图,是直观的运用统计概率分析的图法。机器学习中决策树是一个预测模型, 它表示对象属性和对象值之间的一种映射,树中的每一个节点表示对象属性的判断条件,其分支表示符合节点...

🚯 han25æ940718 (http://blog.csdn.net/han2538740718) 2016年08月24日 15:08 皿5208

后剪枝之悲观剪枝法 (http://blog.csdn.net/zhang4418876/article/details/47025029)

转自http://blog.csdn.net/woshizhouxiang/article/details/17679015 把一颗子树(具有多个叶子节点)的分类用一个叶子节 点来替代的话,在...

C zhang4418876 (http://blog.csdn.net/zhang4418876) 2015年07月23日 16:57 □2360

用 python实现 c4.5算法,并进行悲观剪枝 (http://blog.csdn.net/o1101574955/article/deta...

#coding=utf-8 import xlrd import xlwt import math import operator from datetime import date,datetime...

🚵 o1101574955 (http://blog.csdn.net/o1101574955) 2015年12月21日 13:43 🛛 🖺 2860

C4.5 树剪枝 (http://blog.csdn.net/fanbotao1209/article/details/44780597)

转自: http://www.cnblogs.com/superhuake/archive/2012/07/25/2609124.html 树剪枝 在决策树的创建时,由于数据中的 噪声和离群点...

fanbotao1209 (http://blog.csdn.net/fanbotao1209) 2015年03月31日 15:48 21782

C4.5算法详解(至今见过写的最好的算法详解) (http://blog.csdn.net/xuxurui007/articl...

C4.5是机器学习算法中的另一个分类决策树算法,它是基于ID3算法进行改进后的一种重要算法,相比于ID3算法,改进 有如下几个要点:用信息增益率来选择属性。ID3选择属性用的是子树的信息增益,这里...

xuxurui007 (http://blog.csdn.net/xuxurui007) 2014年01月09日 17:33
 □68292

2017/11/30 上午10:51 第46页 共46页