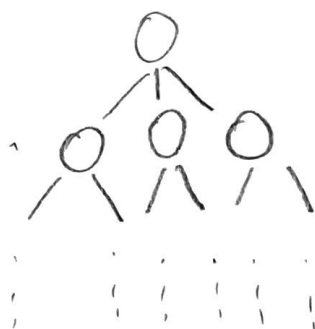


2022.02.28

# Decision Tree



从一开始的 input, 不断通过各种问题将 input 分类渐渐细化. 可 binary split / multi split.

## Discretization

将 continuous <sup>attribute</sup> 转化为 discrete 的  $n$  个 range

static: 在一开始一次性 discretize

dynamic: ranges can be found by equal interval bucketing

如何看哪种 split 更好?

最终各个 split 中 data point 的 portion 尽量不同.

例	$S_1$	$S_2$	$S_3$
	$\downarrow$	$\downarrow$	
包含	1个 $C_1$	0个 $C_1$	7个 $C_1$
	2个 $C_2$	8个 $C_2$	0个 $C_2$

Gini index: (1个 node 要 split 时怎样更好?)

$$GINI(t) = 1 - \sum_j [P(j|t)]^2$$

表示 data 原本的种类  $t$  表示最终 split.

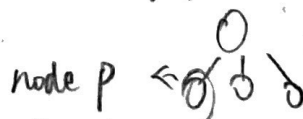
$\downarrow P(j|t)$  表示 node  $t$  中  $j$  的占比

例  $S_1$   
 $C_1$  为  $\frac{1}{3}$   
 $C_2$  为  $\frac{2}{3}$

Combine

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

$k$  表示 partition 个数  $n$  为 node  $p$  中总 records,  $n_i$  为 child  $i$  中的 records.



$GINI_{split}$  在算每一步 split 的好坏来决定此 node 如何 split.

例看split S. 则算[不split的 GINI index - split S 的 GINI index.]  
哪个split使  $\leftarrow$  最大, 哪个为最好.

Continuous attribute 中  
input 先要排序

Stop criteria

① node 中所有 records 属于同一 class

② node 中所有 records 有 similar attribute values.

Model evaluation

例 class 1 有 9990 samples, class 2 有 10 samples.

若 model 将所有都分成 class 1, accuracy 为 99%. Problem!

Cost Matrix

		predicted						
		Yes	No					
actual	Yes	$C(Y Y)$	$C(N Y)$	→ 给每个 cell 给个 cost. 例 <table border="1"> <tr> <td>10</td> <td>20</td> </tr> <tr> <td>1</td> <td>50</td> </tr> </table> ←	10	20	1	50
	10	20						
1	50							
No	$C(Y N)$	$C(N N)$						

$C(i|j)$  表示 actual 为 j, predict 成 i 的 cost.  
最终算总 cost.

Test Model

Holdout

$\frac{2}{3}$  for training,  $\frac{1}{3}$  for testing

Cross validation

将 data 分为 k 个 partitions (disjoint), k-1 会用来 training,  
求余 1 个用来 test.

例: 25个 classifier, 每个 error rate 为 0.35. 每个 unseen record 会根据 majority of  $\downarrow$  来分类.  
majority 错的 prob:  $\sum_{i=1}^{25} \xi^i (1-\xi)^{25-i} = 0.06$ .  
classifier 间需 independent.

· 用 GINI index 分某个 attribute (binary)

① sort the attribute values,

② 尝试所有 possible split 计算 GINI

③ 选最小值.

例 attribute values 3 7 11  
2 4 8 12  
分水岭

根据与“分水岭”大小关系来 split 成 2 组, 每组算 GINI