

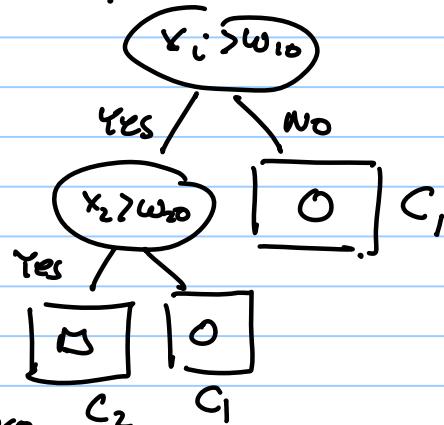
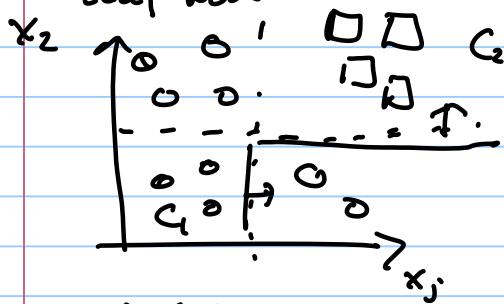
Decision Trees

3/30/2008

Divide and Conquer strategy for classifiers.

Tree. Decision Node in test function $f_m(x)$

Leaf Node

Each $f_m(x)$ defines a
disconnected in d-dim input space.Hierarchical structure \rightarrow Best case: can find
one of n regions with
 $\log n$ questionsDecision Trees are also
interpretable, because they can be converted to
a set of IF-THEN rules.Univariate Tree

test uses only one of the input dimension.

 \rightarrow e.g. attribute - color, size $f_m(x) = x_j > w_{m0}$ threshold value.binary split: $I_m = \{x | x_j > w_{m0}\}$ $R_m = \{x | x_j < w_{m0}\}$ Finding best tree is NP-complete, so use local
heuristic search procedures.

(2)

Classification Trees

Goodness of a split is quantified by an impurity measure.

Pure - if all instances after split belong to the same class.

Node m , N_m is no. instance reaching m (i.e. N for root node).

N_m^i of N_m belongs to class i , $\sum_i N_m^i = N_m$

The estimate for probability of class C_i is

$$\hat{p}(C_i | \Sigma, m) \equiv p_m^i = \frac{N_m^i}{N_m}.$$

Node m is pure if $p_m^i = 0$ or 1 , for all i .

Impurity Measure : $I_m = -\sum_{i=1}^K p_m^i \log p_m^i$.
entropy.

Other measures: Gini index $G(p, 1-p) = 2p(1-p)$

If node m is not pure, then the node should be split to decrease impurity

$$\hat{p}(C_i | \Sigma, m_j) \equiv p_{m_j}^i = \frac{N_{m_j}^i}{N_{m_j}} \quad \text{test result } j.$$

$$\text{Impurity after split } I'_m = -\sum_{j=1}^n \frac{N_{m_j}}{N_m} \sum_{i=1}^K p_{m_j}^i (\log p_{m_j}^i).$$

(3)

Classification & Regression Trees (CAKT)

CAKT: For all attributes and for all possible splits - calculate the impurity and choose the one with maximal purity.

Then repeat recursively and in parallel for all branches that are not pure. Continue till all branches are pure.

Danger - growing the tree until we have pure leafs risks overfitting the training data.

In practice, use a threshold θ_L

Don't split any node with impurity $< \theta_L$.

For leaf Nodes^m - output the posterior probability $P(C_i|m)$ - or output MAP $\hat{C}_i = \arg \max P(C_i|y)$

(4)

Regression Trees

Similar to classification tree, except we replace the impurity measure by one more suitable for regression.

node m , X_m subset of X reaching node m .

depth $b_m(x) = \begin{cases} 1, & \text{if } x \in X_m: x \text{ reaches node } m \\ 0, & \text{otherwise.} \end{cases}$

g_m is the estimated value in node m .

$$E_m = \frac{1}{N_m} \sum_t (r^t - g_m)^2 b_m(x^t)$$

$$N_m = |X_m| = \sum_t b_m(x^t)$$

$$\underline{\text{Select}} \quad g_m = \frac{\sum_t b_m(x^t) r^t}{\sum_t b_m(x^t)}$$

If E_m is below threshold, then create a leaf node and store its g_m value.

(Creates a piecewise approximation with discontinuities at leaf boundaries -

If E_m is above threshold, split the node so as to decrease the errors of the child nodes.

$X_{m,j}$ subset of X_m taking branch j ; $V_j = X_{m,j} \subseteq X$,
 $b_{m,j}(x) = \begin{cases} 1 & \text{if } x \in X_{m,j}: x \text{ reaches } m \text{ and takes branch } j \\ 0 & \text{otherwise} \end{cases}$

$$g_{m,j} = \frac{\sum_t b_{m,j}(x^t) r^t}{\sum_t b_{m,j}(x^t)} \quad \text{error after split.}$$

$$E'_m = \frac{1}{N_m} \sum_j \sum_t (r^t - g_{m,j})^2 b_{m,j}(x^t)$$

(5)

Pruning

pre-pruning — stopping tree construction before we have reached pure leaf nodes.

post-pruning — try to find and remove unnecessary subtrees.

Grow tree until all nodes are pure.

Then find subtrees that cause overfitting.

To do so — obtaining pruning set of instances.

For each subtree,

replace it by a leaf node labeled with training instances covered by subtree.

If leaf node does not perform worse than the subset in the pruning set, then prune the subtree and keep the leaf node.

pre-pruning is faster — but post-pruning usually gives more accurate trees

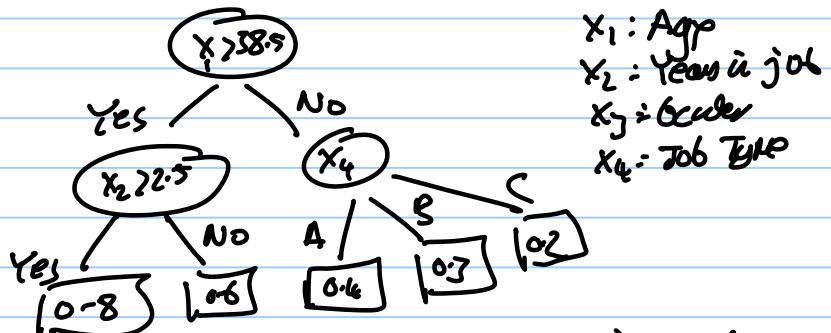
(6)

Rule Extraction from Trees

Interpretability

Each path from root to leaf corresponds to a conjunction of tests.

These paths can be written down as a set of IF-THEN rules — called a RULE-SET



x_1 : Age
 x_2 : Years in job
 x_3 : Gender
 x_4 : Job Type

$R_1 : \text{IF } (\text{age} > 38.5) \text{ AND } (\text{years-in-job} > 25) \text{ THEN } y = 0.8$

$R_2 : \text{IF } (\text{age} > 38.5) \text{ AND } (\text{years-in-job} \leq 25) \text{ THEN } y = 0.6$

$R_3 : \text{IF } (\text{age} \leq 38.5) \text{ AND } (\text{job-type} = A) \text{ THEN } y = 0.4$

$R_4 : \text{AND } (\text{job-type} = B) \text{ THEN } y = 0.3$

$R_5 : \text{AND } (\text{job-type} = C) \text{ THEN } y = 0.2$

Can validate the rules (by an expert)

— Can also calculate the percentage of training data covered by the rule (e.g. rule support).

Several different paths may have the same leaf values. This giving disjunction of conjunctions.

(7)

Learning Rules from Data.

Rule induction - works similarly to tree induction but does depth-first search and generates one path at a time.

Example: RIPPER (see Alpaydin p188
for pseudo code)

Rules are added to explain positive examples

Rule R to rule R'.

Change in gain:

$$\text{Gain}(R', R) = s \cdot \left(\log \frac{N'_+}{N'} - \log \frac{N_+}{N} \right)$$

N instances covered by R N_+ no. true positive
 $N' < N'$ for R.

Conditions are added to a rule until it covers no negative example.

Can prune back using a pruning set to find rules that maximizes the

rule value metric $\text{rum} = \frac{p-n}{n+n}$ p - no. true pos
 n - no. false pos
in pruning set

(8)

Propositional Rules & First Order Rules

IF Father(y,x) AND Female(y)
THEN Daughter(x,y)

Inductive logic Programming.
Learning first order rules is similar
to learning propositional rules

Multivariate Trees

All input variables can be used to split

$$f_m(\underline{x}) : \underline{w}_m^T \underline{x} + w_{m0} > 0$$

Alternates are the sphere areas

$$f_m(\underline{x}) : \|(\underline{x} - \underline{s}_m)\| \leq d_m$$

Algorithm proceed as before.