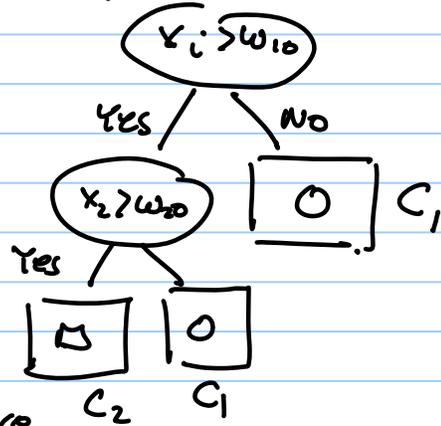
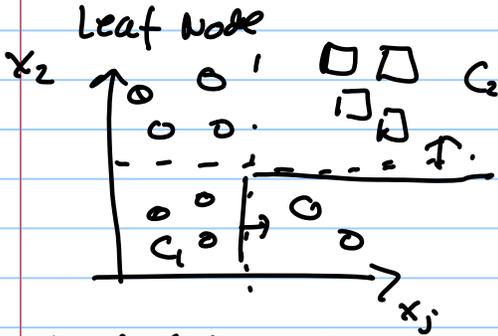


Decision Trees

Divide and Conquer strategy for classification.

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



Each $f_m(x)$ defines a disjunction in d -dim input space.

Hierarchical structure \rightarrow

Best case: can find one of N regions with $\log N$ questions

Decision 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 dimensions.

\rightarrow e.g. attributes - color, size

$$f_m(x) : x_j \geq w_{m0} \quad \text{threshold val.}$$

binary split:

$$L_m = \{x \mid x_j \geq w_{m0}\}$$

$$R_m = \{x \mid x_j < w_{m0}\}$$

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. instances 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

$$\hat{p}(C_i | \mathcal{X}, m) \equiv p_m^i = N_m^i / N_m.$$

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 $\phi(p, 1-p) = 2p(1-p)$

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

$$\hat{p}(C_i | \mathcal{X}, m, j) \equiv p_{m,j}^i = N_{m,j}^i / N_{m,j} \quad \text{test result } j.$$

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 (CART)

CART 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 leaves risks overfitting the training data.

In practice, use a threshold θ_e

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

For leaf nodes^m - output the posterior probability $P(C_i | m)$ - or output MAP $\hat{C}_i = \text{ARGMAX}_j P(C_j | m)$

(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 .

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

g_m is the estimated value in node m .

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

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

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

If E_m is below threshold, then create a leaf node and store the 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 ; $\cup_{j=1}^n X_{m,j} = X_m$
 $b_{m,j}(x) = \begin{cases} 1 & \text{if } x \in X_{m,j}: x \text{ reaches } m \text{ and take branch } j \\ 0 & \text{otherwise} \end{cases}$

$$g_{m,j} = \frac{\sum_t b_{m,j}(x^t) r^t}{\sum_t b_{m,j}(x^t)}$$

$$E_m' = \frac{1}{N_m} \sum_j \sum_t (r^t - g_{m,j})^2 b_{m,j}(x^t) \quad \text{error after split.}$$

(5)

Pruning

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

post-pruning - try to find and remove unnecessary sub-trees.

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 subtree on the pruning set, then prune the subtree and keep the leaf node.

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

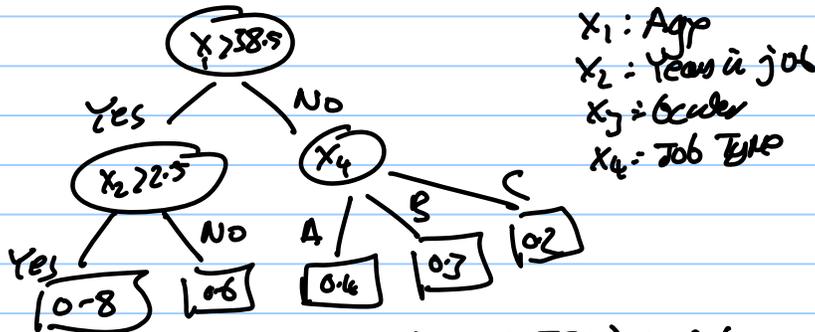
(6)

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-BASE



R1: IF (age > 38.5) AND (years-in-job > 22.5) THEN $y = 0.8$

R2: IF (age > 38.5) AND (years-in-job ≤ 22.5) THEN $y = 0.6$

R3: IF (age ≤ 38.5) AND (job-type = A) THEN $y = 0.4$

R4: AND (job-type = B) THEN $y = 0.3$

R5: AND (job-type = C) 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 gives 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 Alpaydm p188 for pseudo code)

Rules are added to explain positive examples
Rule R to rule R' .

Change in gain:

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

N instances cover by R N_+ no. true positives
 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 rule that maximizes the

rule value metric $\text{rvm} = \frac{p-n}{p+n}$

p - no. true |
 n - no. false |
positives on
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

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

Alternatives are the sphere case

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

Algorithm proceed as before.