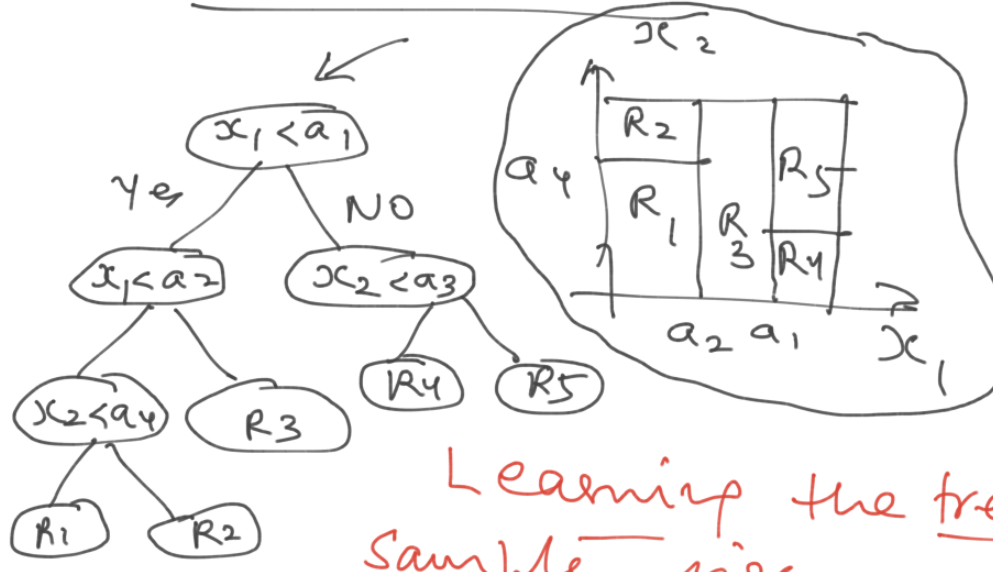


Lecture 13 FML April 6, 21
multiclass classification

Decision Tree

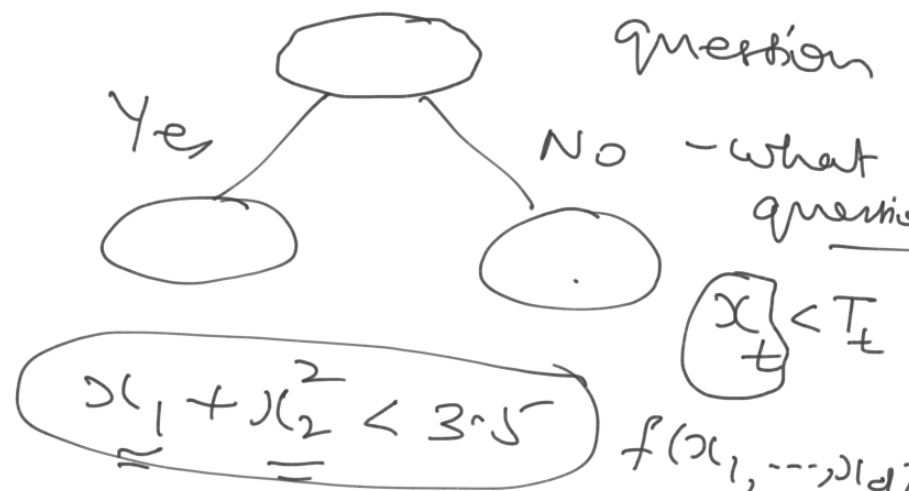


Learning the tree
sample size

$(x_1, y_1), \dots, (x_m, y_m)$ \parallel
 \uparrow \uparrow
classes. m classes

Getting optimal
 decision tree is NP
 complete

Heuristic algorithm
Greedy algorithm.



- Binary decision tree
- Questions will be $x_i < T_i$

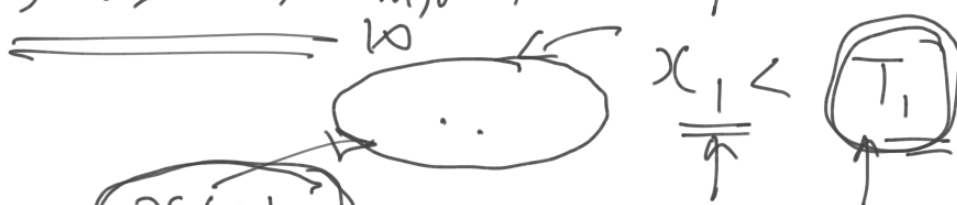
- Binary tree

- $x_i < T_i$

sample $(x_1, y_1), \dots$

$(x_{1(1)}, \dots, x_{1(d)})$

$(x_1, y_1), \dots, (x_m, y_m)$



$x(1) \dots x(d)$

Yes

No

class 3

class 1

...

$(x_4, y_4) \dots (x_{10}, y_{10})$

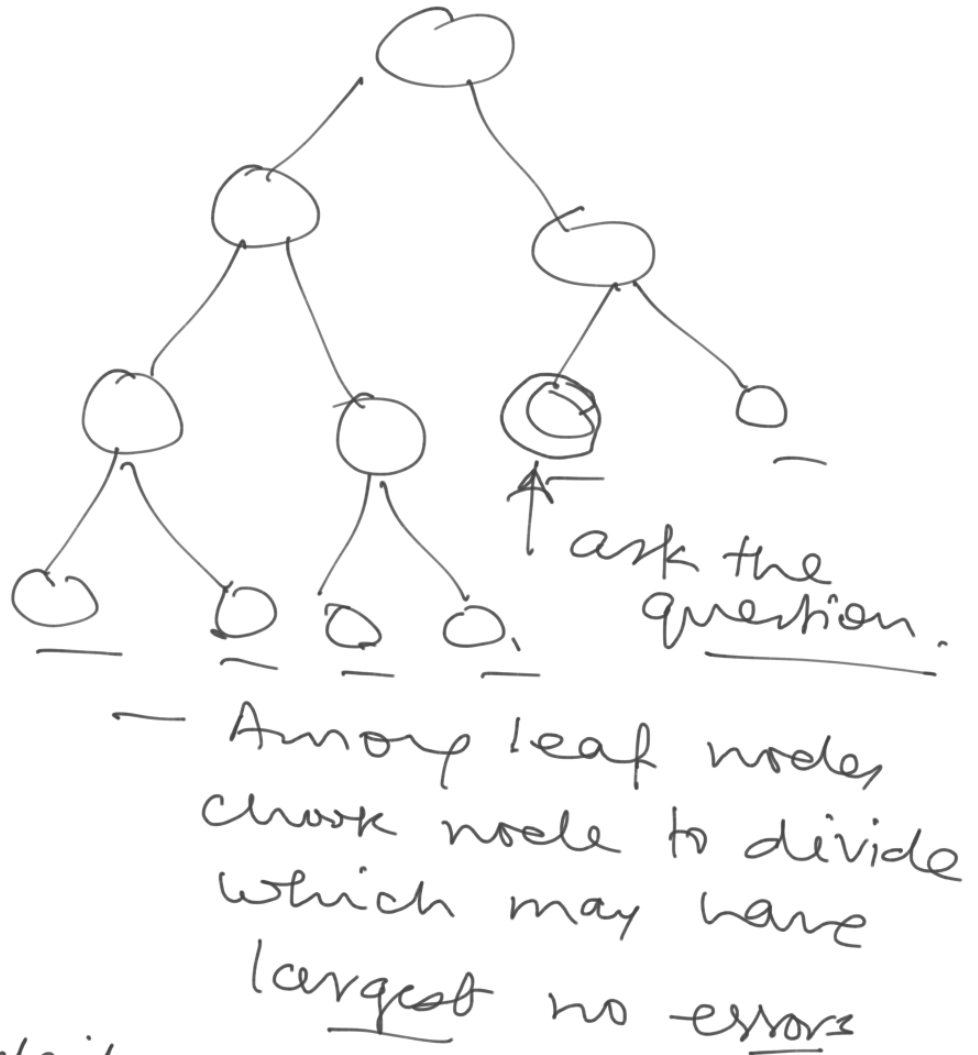
miss classified
train points -

= 1, 3



Greedy Algorithm

- ① Any leaf node will be labelled by the majority class.
- ② At any given node n_t we have to decide whether to divide it or not: if training sample is having large error then we may want to divide it



- After choosing next node to divide need to choose the question n_t

$$\underbrace{x_t}_{\text{what feature?}} < \underbrace{T_t}_{\text{what threshold?}}$$

- pick the feature and a good threshold
- Then change the feature & try with a good threshold---

- Pick the best feature & threshold.

-
- The criteria chosen to decide among the different nodes & features is no of misclassified sample points.
 - prob of error

Other criteria
for selection of nodes
and questions at a
time :-

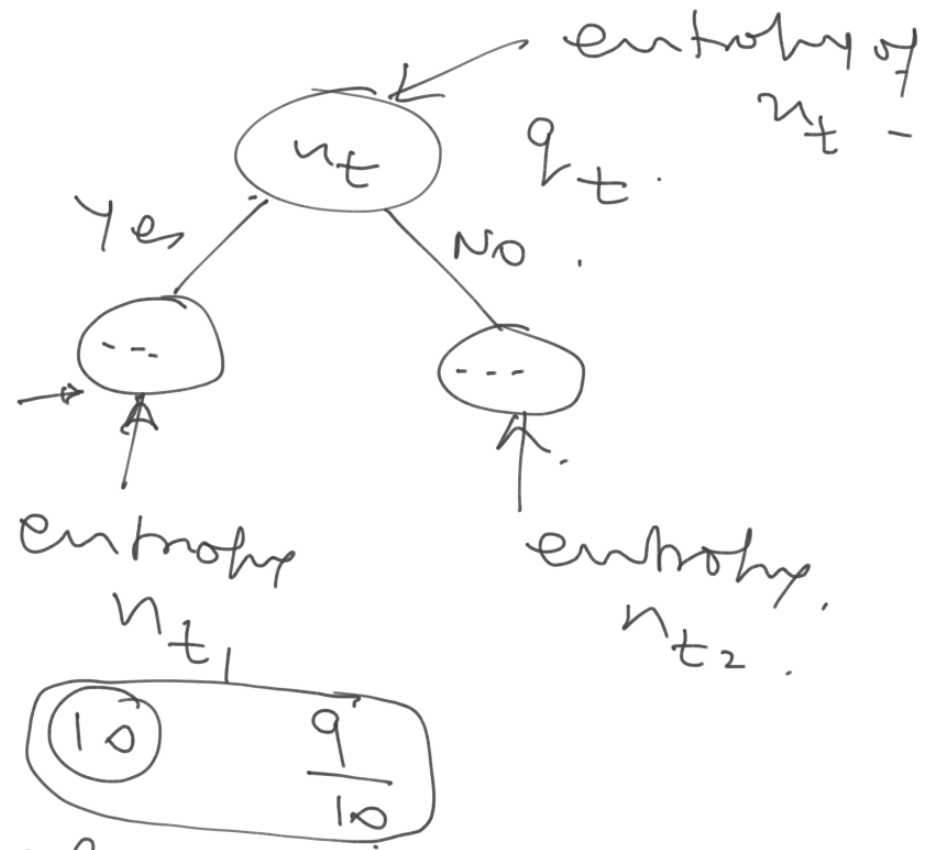
- Entropy
- Gini index.

Given a node n_t

This node has m_t
sample points and

$$p_l = \frac{m_t(l)}{m_t} \quad m_t(l) = \text{no of sample points of class } l$$

$$\text{entropy of } n_t = - \sum_{l=1}^k p_l \log p_l$$



Gini index 0_0

p_l = prob of class l

$$\sum_{l=1}^k p_l (1 - p_l).$$

entropy & Gini index

are large when
the node i $p_l = \frac{1}{k}$
 $\forall l$.

and then are 0 if

$$p_l = 1 \text{ for some } l.$$
$$= 0 \text{ others.}$$

- All these three criteria are convex fns.
- But misclassification prob is not diff but entropy & Gini index are diff and are also strictly convex
- entropy & Gini index are more commonly used methods/criteria.

impurity index

of a node

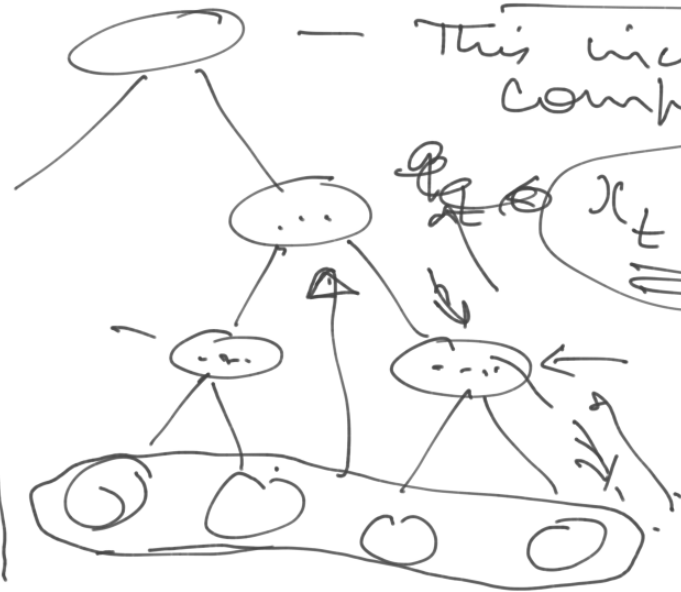
- The algorithm is stopped when ^{all} the nodes have really low impurity or no of points at each node are enough low

Problems with greedy

- It can give a bad split

This can be addressed by looking ahead.

- This increases computational complexity



- It can give a large tree
- This can do over fitting
- Grow and prune

Problem with Decision Tree

- These are unstable



Random Forest

- Bunch of tree
- These trees are formed in a random way.
- S sample $(x_1, y_1), \dots, (x_m, y_m)$
- Take a subsample S' randomly from S
- From the set of features $(x(1), \dots, x(d))$ randomly pick d' features, $d' \leq \sqrt{d}$.
- Form a decision tree S', d' features, d'

using greedy.

This way we get Tree T_1

- Repeat the process
- Take ^{different n times} subsample S' from S
- take d' features randomly
- Form Tree T_2 by greedy from (S', d')

T_1, T_2, \dots, T_L random tree

- For a new observation (OL) to classify :-

- Classify it using T_1, T_2, \dots, T_L trees
- overall classification via majority decision
- Random Forest can provide a good performance as avg.

- we have taken

$$d' \leq \sqrt{d}$$

This way the correlation between different trees are low.

- Theoretical analysis is hard. Very limited results available

- One reason is the use of greedy algorithm.

- Some studies on performance

are asking instead of using greedy algorithm for tree

randomly choose the node

- Using binary classification algorithms.

- OVA (One vs All.)
k classes $S = (x_1, y_1), \dots, (x_m, y_m)$

For each class $l \in Y$,
make classifier h_l taking two classes: $class_1 \rightarrow class_l$
 $class_2 \rightarrow$ The rest of classes.

— h_l is formed from any of the binary classification algorithms.

— We have k classifiers.

— For classifier l form score function f_l

— Suppose we used kernel method for binary classification

$$f_l(x) = \boxed{y (w^T \cdot \Phi(x))} \leftarrow$$

$$y \in \{1, -1\}$$

$$f_l(x) = p(x | \text{class } l) \leftarrow$$

— pick the class with largest $f_l(x)$

— calibration problem

$$\underline{\underline{f_l(x)}}$$

— This OVA works quite ~~at~~ well in practice

OVO method :

one vs one

- For each $l, l' \in \{1, \dots, K\}$ make a binary classifier $h_{ll'}$ using subsample of S which have labels l and l'
- Now finally when a new measurement x comes

$K(K-1)$ binary classifiers.

- classify x via all these classifiers and choose class with majority.
- Complexity of leap OVA is same as that of other multiclass classifiers
- But complexity of OVO can be very large

- But complexity of OVO still may not be very large because for $h_{l,l'}$ the sample size to train from can be quite small.

- One problem with OVO is that for some $h_{l,l'}$ the sample to train from can be quite small.

- But it does not have the calibration problem of OVA.

- If generalization error of $h_{l,l'} \leq r$
 $\forall l, l'$ then generalization error of overall classifier $\leq \underline{\underline{(k-1)r}}$.

- KNN and NN.
 \hookrightarrow k nearest nbr. \hookrightarrow Neural networks