

Reduction in training error:

$$f_{minority label}^S - \left(t_A^S \cdot f_{minority label}^A + t_B^S \cdot f_{minority label}^B \right)$$

error rate after the split

Gain as "Information Gain"

How much information gain (reduction in entropy) happens if the node is split?

Entropy: Measure of randomness (Uncertainty)

Highly uncertain events carry more information

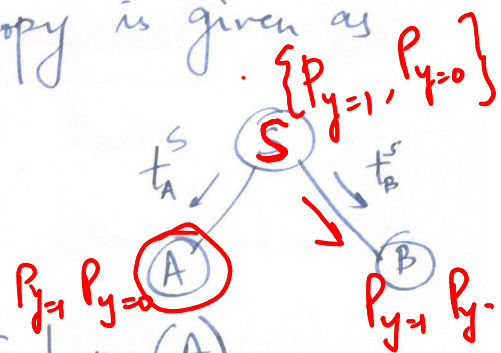
For a set of events $\{e_1, e_2, \dots, e_N\}$ which occur with probability $\{p_1, p_2, \dots, p_N\}$ the entropy is given as

$\{1, 0, 0, 0, 0\}$
No information

$$Entropy = \sum_i p_i \log \frac{1}{p_i}$$

Information Gain:

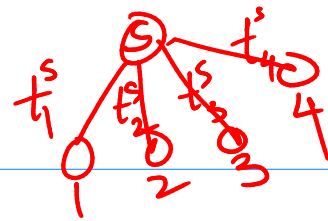
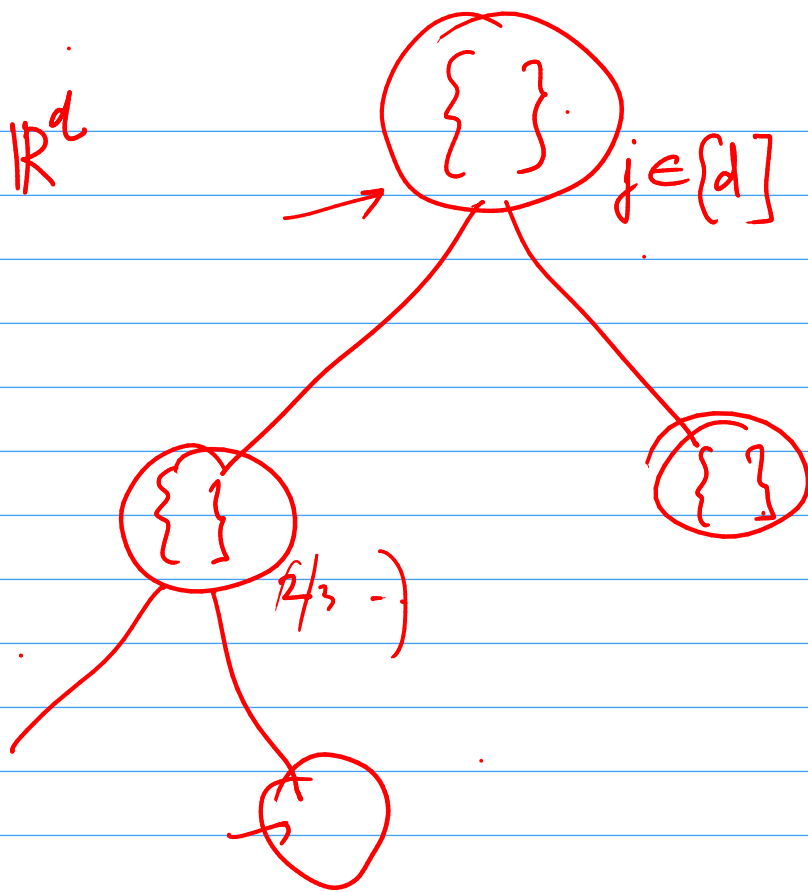
$$Entropy(S) - \left(t_A^S \cdot Entropy(A) + t_B^S \cdot Entropy(B) \right)$$



Probability is computed as the frequency of target labels for the examples in the set.

Ground truth

$$\underline{x} \in \mathbb{R}^d$$



Entropy(S)

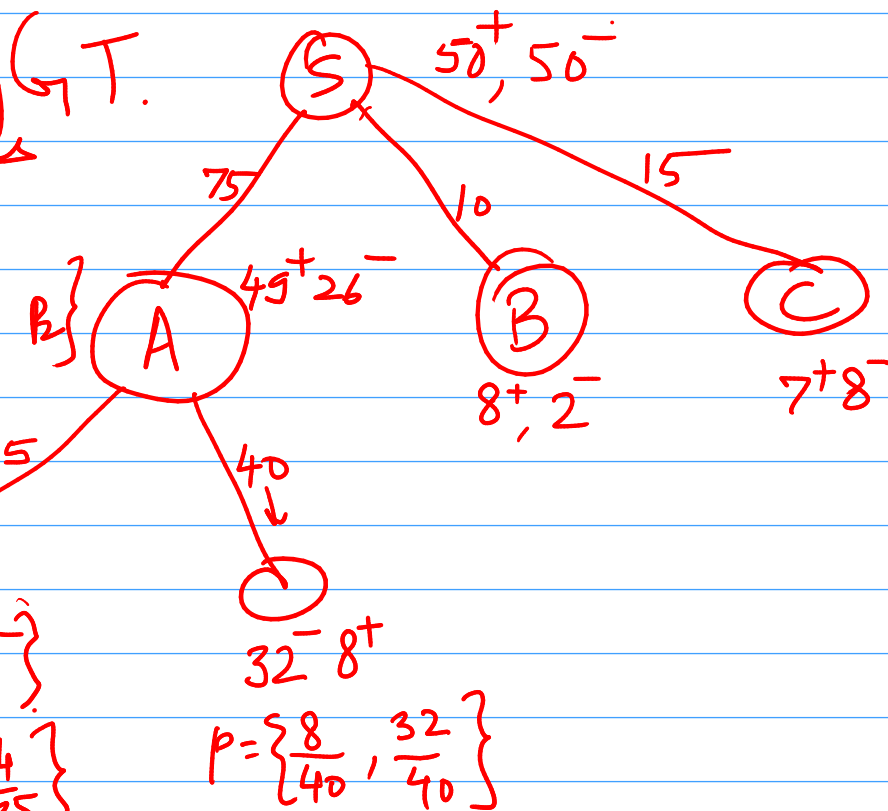
$$= \sum_{i=1}^c t_i^s \text{Entropy}(i)$$

Probability of labels
of examples
in a
node(set).

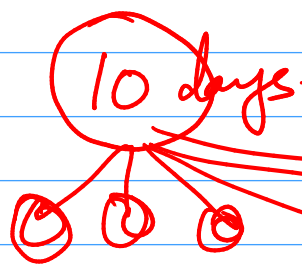
$$\underline{P} = \{P_1, P_2\}$$

$$\underline{P} = \left\{ \frac{31}{35}, \frac{4}{35} \right\}$$

$$\underline{P} = \left\{ \frac{8}{40}, \frac{32}{40} \right\}$$



Weekend (Yes, No)
B'day (Yes, No)
Exams over (Yes, No)
Date 14/Sep (Unique value)

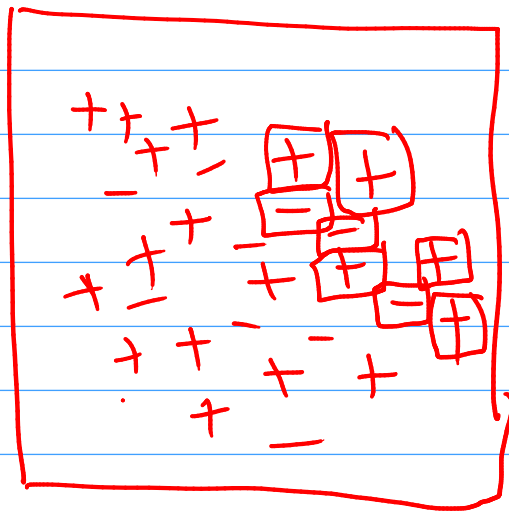


date

Memorized

dates
is a not
a good
attribute to
split

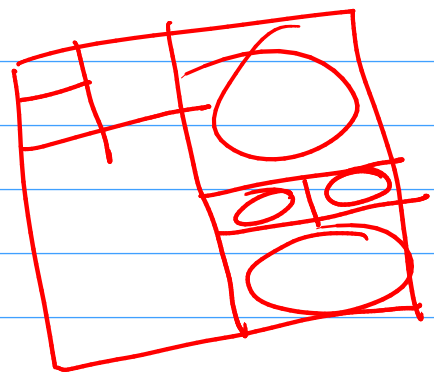
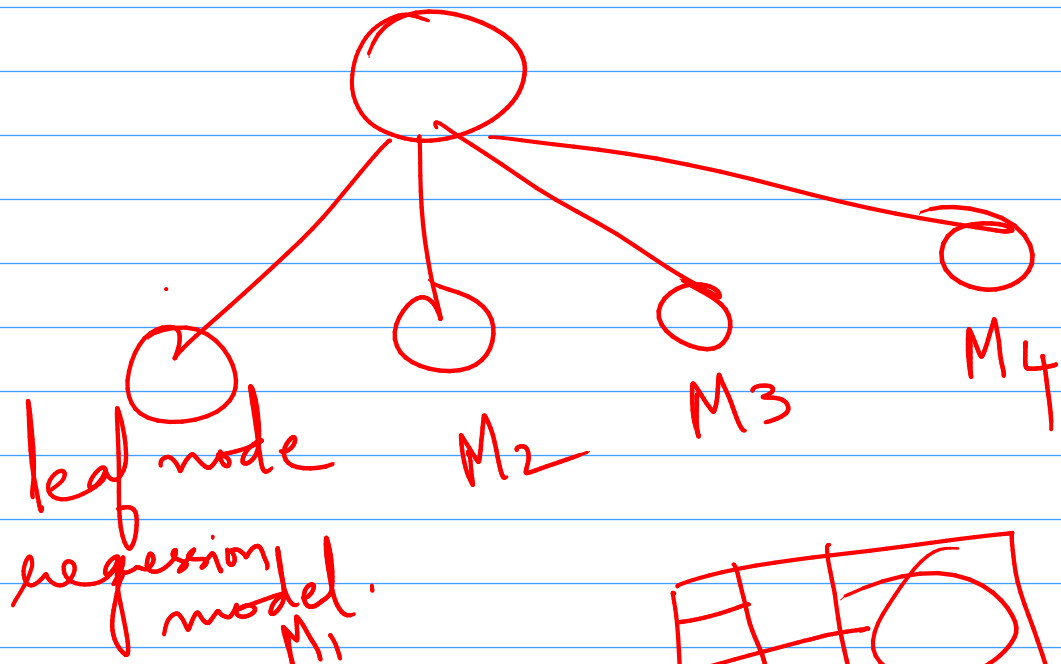
D Tree
creates a
partitioning
of the
input space



$VC_{dim} \infty$

leaf nodes with
single example

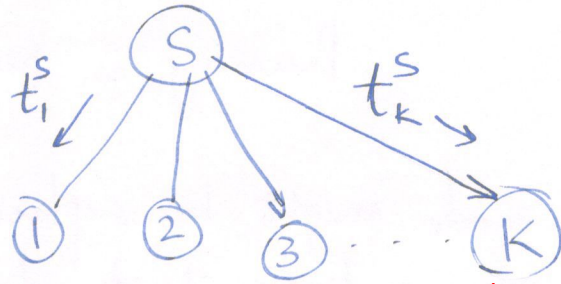
Useless machine



Information gain is more for attributes which have more number of values and therefore can lead to multiway splits.

- To avoid the bias in favour of attributes with a large number of values, we divide the information gain by a split gain.

$$\text{Split gain} = \sum_{i=1}^k t_i^s \log \frac{1}{t_i^s}$$



Information Gain Ratio:

Ratio

$$\text{split gain} = \frac{\left(\text{Entropy}(S) - \sum_i t_i^s \cdot \text{Entropy}(i) \right)}{\left(\sum_{i=1}^k t_i^s \log \frac{1}{t_i^s} \right)}$$

ID3

Iterative Dichotomizer (Version 3)

Works for binary valued attributes

Uses information gain for choosing the split attribute.

VC dim of a decision tree?

recursion

C4.5

CHAID. CART.

The learned trees are usually very large.

This leads to low error on the training data.

But such large trees do not perform well on the test data.
i.e. they show poor generalization performance.

Solution to the overfitting problem: Stop growing the tree

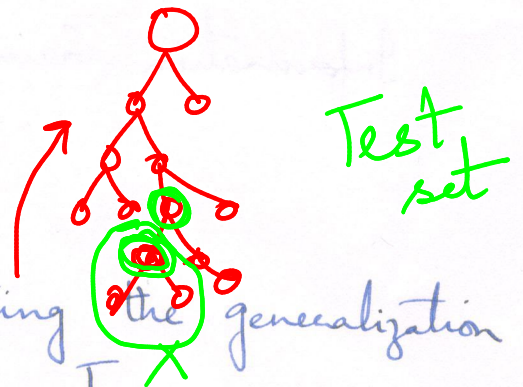
Early stopping. depth of the tree.
by limiting the number of iterations in the ID3 algorithm.
This results in a tree with bounded number of nodes.

Another way to address overfitting is to prune the tree after it is built.

Pruning: First construct a large decision tree.

Pruning is performed by a bottom-up walk on the large tree.

Each node might be replaced with one of its subtrees or with a leaf if it does not lead to much increase in the generalization error.



Pseudo Code for Pruning

Given: Some method $f(T)$ of estimating the generalization performance of the decision tree T .

For each node j in a bottom-up walk on T

find T' which minimizes $f(T')$ where T' is any generalization error.
one of the following:

- the current tree after replacing node j with a leaf 1.
- the current tree after replacing node j with a leaf 0.
- the current tree after replacing node j with its left subtree.
- the current tree after replacing node j with its right subtree.

