

* Classification : is taking value and map/put to a discrete label.

eg: True or False
: male or Female

* Classification Learning:

* Instances : INPUT

* Concept (Function) \rightarrow maps Input to outputs

* Target concept: Output, answer, function that

* Hypothesis Class: all the functions I am willing to entertain (Classification)

* Sample : Training set, {Example input, paired with label} $\{x_i, y_i\}_{i=1}^n$

* Candidate : Target concept

* Testing set : little sample, for prediction.

That defines it
(e.g.: male or female: what male or female is!)

* Decision Trees:

1) Pick the best attribute , the best attribute then eliminate a lot of possibilities . Best \rightarrow splits the data

* Best is something that splits in half eg: gender \Rightarrow living or non-living

* You asked the next question

3) Follow the answer path

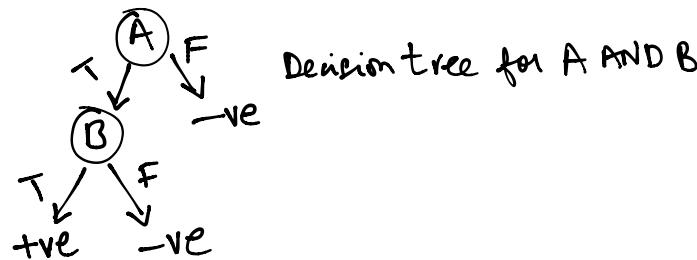
Go to 1 until we get an answer.

* you have to go through all possible questions and answers until the best solution is defined.

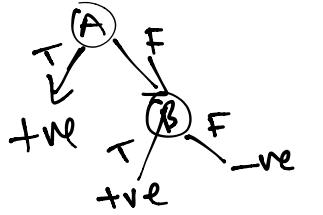
* Decision Trees: Expressiveness

* Boolean

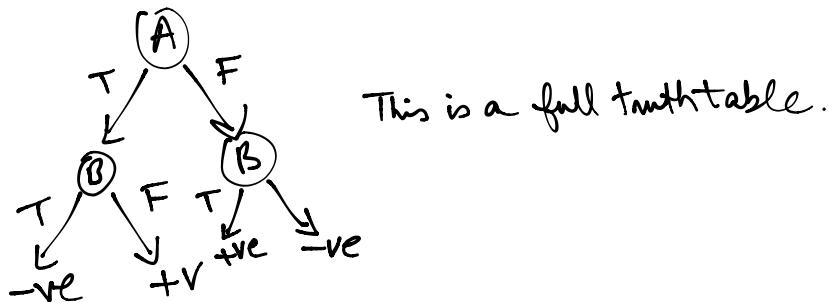
* AND & : A AND B



* A OR B decision
A' OR B all cumulative



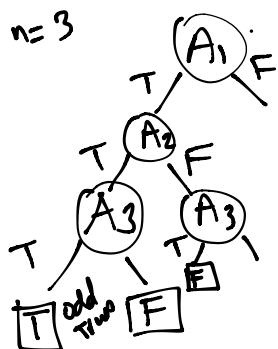
* A XOR B : exclusive OR, True when A True or B is True. when A True & B True, then False



* in case of OR we needed two nodes and in case of XOR we needed 3 it matter when the size of decision tree grows.

* n-OR : Any, the size of decision tree ^{nodes} is linear, n.

* n-XOR : odd Parity



$$O(2^N)$$

Exponential number of nodes, is also called hard.

* therefore in ML we look forward to problems like Any and less-like Parity.

* $\text{XOR} \rightarrow$ hard.

* exactly how expressive is decision tree.

* n attributes (boolean) $O(n!)$ nodes \Rightarrow exp

* how many trees?

* output is boolean

* how many trees

* Answer: # of nodes 2^n

* number of binary combination in the truth table : 2^{2^n}

* which grows very fast and is very expressive.

* by 6 there are billions of combinations, hence we need a intuitive combinations.

* ID3:

* we keep loops:

- 1) finding the best value attribute
- 2) Assign A as decision. attribute for node
- 3) For each value of A
create a descent of node
- 4) Sort Training Examples to leaves
- 5) If Examples perfectly classified. STOP

measure of information | Sample label
/ /

Information Gain (S, A)

Entropy: measure of randomness

$$\text{Gain}(S, A) =$$

$$\text{Entropy}(S) -$$

$$\leq \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Entropy}: \sum_{v=1}^n p(v) \log p(v)$$

* ID3 BIAs: restriction bias: A

* Inductive Bias: \hookrightarrow preference BIAs

- Good splits at the top
- Correct Over Inconet
- Shorter Trees.

* Decision Trees: other consideration:

- Continuous Attributed?
- e.g. Age, weight, Distance



* what if we come across a value not in the training data?

* how about ranges:



* in truly continuous data there could be infinite number of checks.

* But this is how we evaluate continuous variables/data

* when do we really STOP? when everything is classified correctly!

* if there is noise in the data then it may run into an infinite loop.

* we have to be careful about over fitting / variance.

* when it violates Occam's Razor.

* how about using cross validation for ideal curve performance.

* we can go through the entire tree, and collapse the tree (Pruning) to see how much error is being produced.

* AFTER PRUNING, we can do Voting.

* Regression?

* information Gain may not work well, lets use a different splitting criteria

* Variance?

* Output: Avg, Local linear fit

* Recap:

Decision Trees

* Representation

* ID3: A Top Down Learning Algorithm

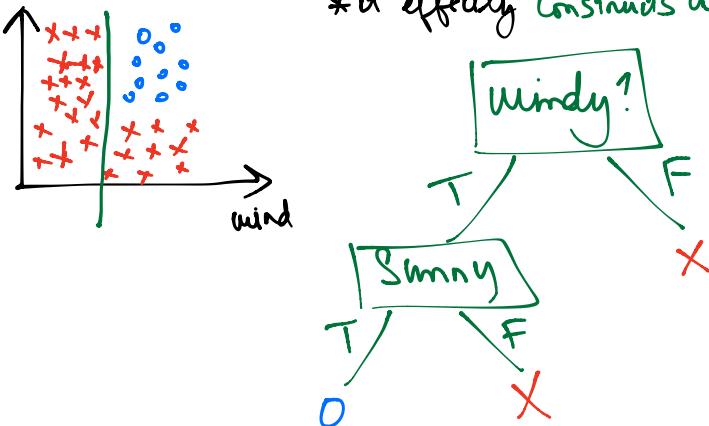
* Expressions of DTS

* BIAS of ID3

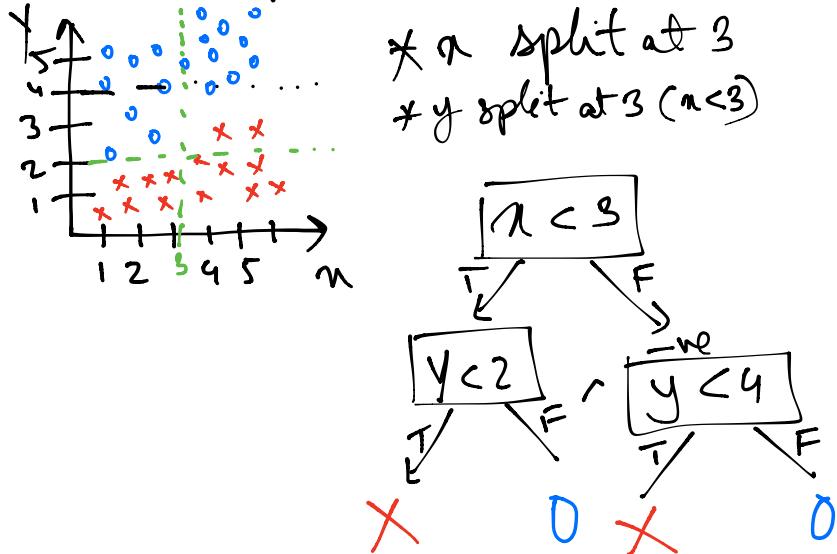
* "BEST ATTRIBUTES" $\xrightarrow{\text{like}}$ GAIN(S,A)

* Dealing with overfitting.

Decision Trees: we can ask multiple linear questions
 Sun ↑
 * it effectively constructs a linear barrier.



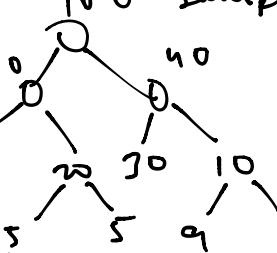
* what is the most useful split, with max Gain - attribute



* Decision Tree fn.: min_samples_split

100 - Samples

* Deep in the decision tree complications can occur when data is overfitted. Hence overfitting the data.



* hence we needed to keep the tuning limited by limiting min_sample_split to higher numbers.

Entropy: the measure which decides where to split the data.

* it is a measure of impurity in a bunch of examples

$$\text{Entropy} = \sum_i -P_i \log_2 (P_i) \quad \text{fraction of examples}$$

* it is opposite of purity

* Entropy : same class : entropy = 0.0

* evenly split between classes \rightarrow entropy = 1.0 (max value of entropy)

* Entropy : entropy = $\sum_i (P_i) \log_2 (P_i)$

* eg. slow / fast : ssff : slow nodes: 2 , total nodes : 4

entropy $\Rightarrow P_i$ is fraction of slow examples $\Rightarrow P_{\text{slow}} = \frac{2}{4} = 0.5$

$$P_{\text{fast}} = \frac{2}{4} = 0.5$$

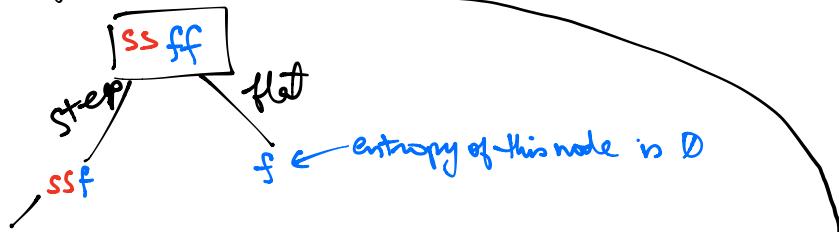
what is the entropy of this node : SSff

$$-0.5 \log_2(0.5) - 0.5 \log_2(0.5) = 1$$

Gain \Rightarrow information Gain =

* entropy of parent = 1

grade	bumpiness	speed limit	speed
steep	bumpy	yes	slow
steep	smooth	yes	slow
flat	bumpy	no	fast
steep	smooth	no	fast



$$P_{\text{slow}} = \frac{2}{3}$$

$$P_{\text{fast}} = \frac{1}{3}$$

$$\text{entropy} = \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$\text{entropy} = 0.91829 //$$

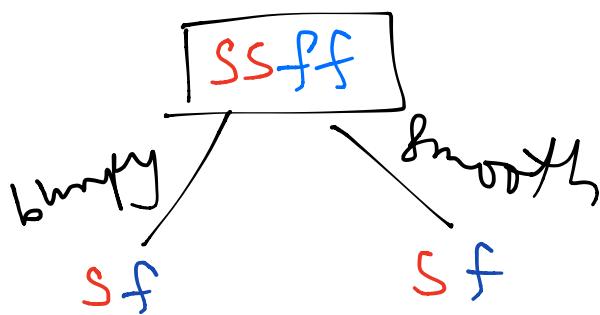
$$\text{entropy of children} = \frac{3}{4} (0.91829) + \frac{1}{4} (0) = 0.68872$$

$$* \text{Gain} = 1 - 0.68872 = 0.31127 \quad \text{if we split on grade}$$



entropy on bumpiness:

entropy (bumpiness):



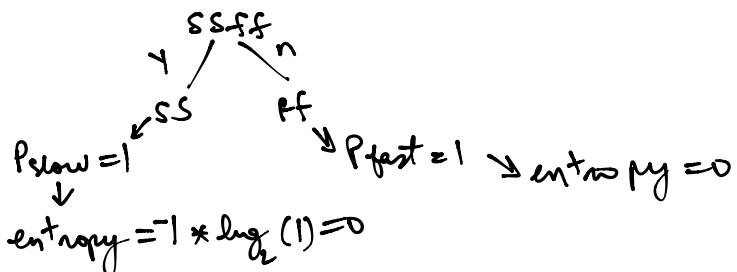
$$P_{slow} = \frac{1}{2}$$

$$P_{fast} = \frac{1}{2}$$

$$\text{entropy} = -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right)$$
$$= 1$$

$$\text{entropy of children} = \frac{1}{2} (1) + \frac{1}{2} (1) = 1$$

Gain on bumpiness = $1 - 1 = 0$, we do not want to start splitting starting here.



$$\text{Gain} = 1 - 0 = 1 // \Rightarrow \text{max gain}$$

* Bias - Variance trade off - somewhere in the middle, so it is able to generalize the new data and predict.

* DT - stop the growth of the decision trees before it goes too deep by limiting the depth of the tree and min-samples-split