

# Advanced Learning Algorithms - Week 4

## Decision Tree Model

Every 'object' have some values. Face: Round, Square - Ear: Pointy, Floppy...

Algorithm iterates over the tree to classify an object.

**Pick One Tree That Fits The Best**

**How to build:**

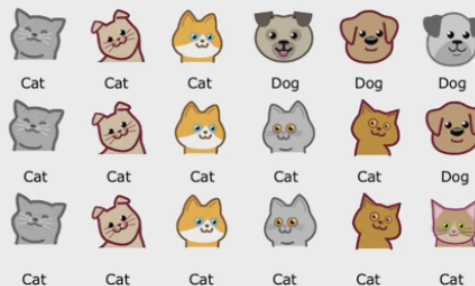
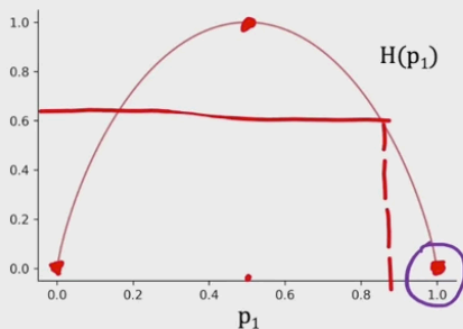
**Features** : Aim is to maximize purity.

**When to Stop Splitting?** : Max depth of the tree (overfitting), When a node becomes %100 one class, threshold of minimum # of examples in node.

**Entropy as a measure of impurity**

## Entropy as a measure of impurity

$p_1$  = fraction of examples that are cats



$$p_1 = 3/6 \quad H(p_1) = 1$$

$$p_1 = 5/6 \quad H(p_1) = 0.65$$

$$p_1 = 6/6 \quad H(p_1) = 0$$

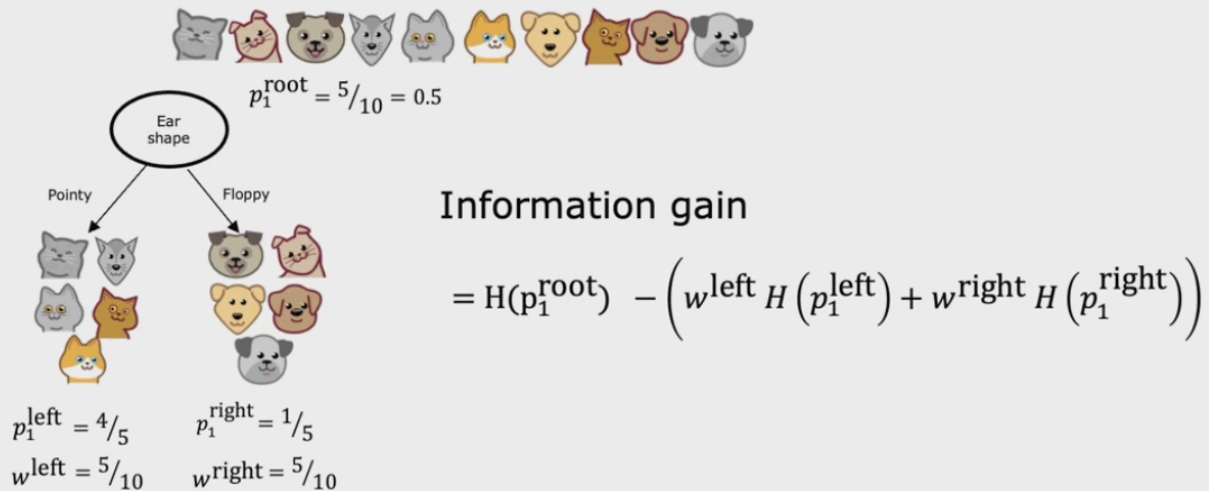
*Entropy Function*

**Reduction of Entropy / Increase in Purity : Choosing a split**

Reduction in Entropy = Information Gain

Should bother to increase the size of the tree?

# Information Gain



- At the first step: Everything is at root.
- Check the highest information gain among the possible splitting options (features).
- Repeat - until hitting one of the thresholds.

## One Hot Encoding

-> One feature that can take up 3 values - Ear: Pointy/Floppy/Round

-> Three feature that can take up only 2 values (0 or 1) - Pointy Ear: 0/1

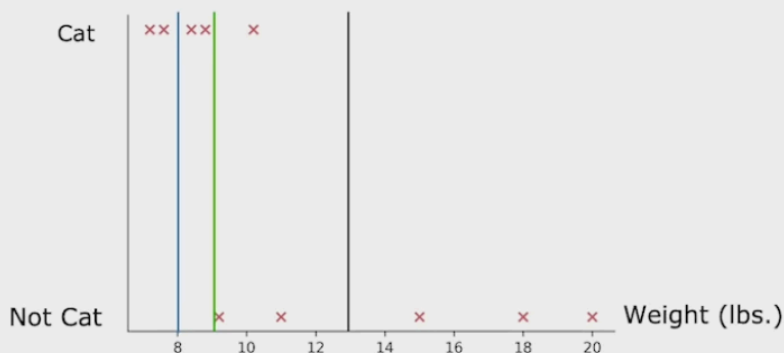
Also works with regression/neural

## Continuous Features

-> Weight (5.0 - 10.0)

Pick a value for threshold. Split the examples. (Weight < 8 lbs)

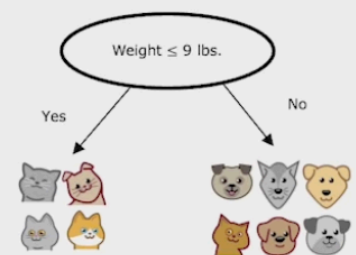
# Splitting on a continuous variable



$$H(0.5) - \left( \frac{2}{10} H\left(\frac{2}{2}\right) + \frac{8}{10} H\left(\frac{3}{8}\right) \right) = 0.24$$

$$H(0.5) - \left( \frac{4}{10} H\left(\frac{4}{4}\right) + \frac{6}{10} H\left(\frac{1}{6}\right) \right) = 0.61$$

$$H(0.5) - \left( \frac{7}{10} H\left(\frac{5}{7}\right) + \frac{3}{10} H\left(\frac{0}{3}\right) \right) = 0.40$$



One Decision Tree : Highly sensitive to small changes to the data.

Better Option: Train whole bunch of different decision trees. - **Tree Ensemble**

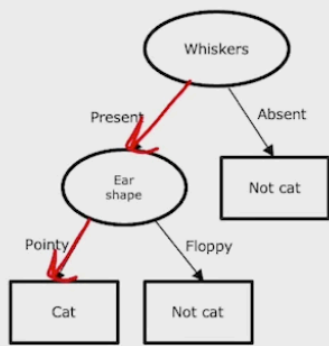
Make them vote for the output.

# Tree ensemble

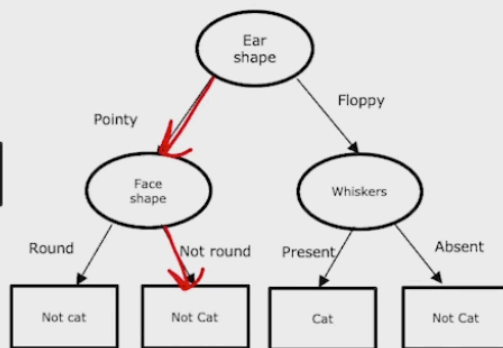
New test example



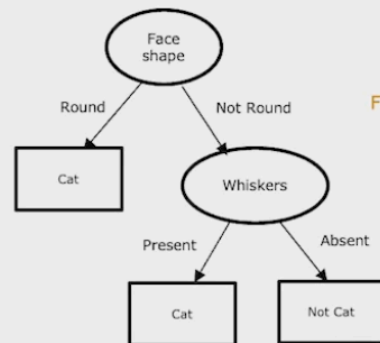
Ear shape: Pointy  
Face shape: Not Round  
Whiskers: Present



Prediction: Cat



Prediction: Not cat



**Sampling with replacement** -> replacing makes it possible to get the same output again. Tokens.

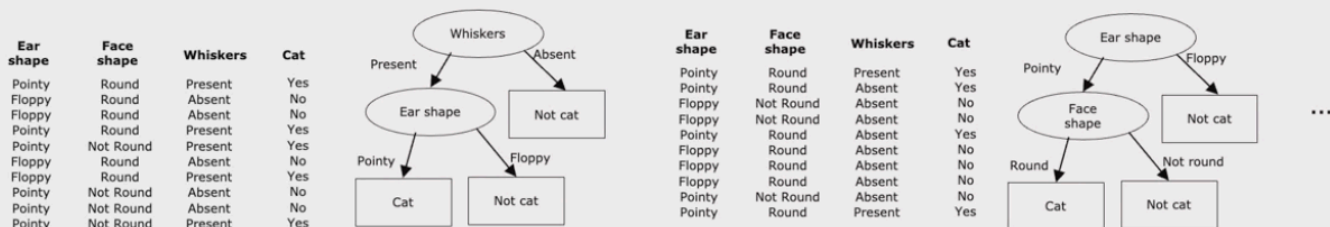
## Random Forest Algorithm - Tree Ensembles

### Generating a tree sample

Given training set of size  $m$

For  $b = 1$  to  $B$

Use sampling with replacement to create a new training set of size  $m$   
Train a decision tree on the new dataset



Bagged decision tree

after some # of decision trees, computation is high, more trees return into diminishing returns.

Randomizing the feature choice.

There are  $n$  features available at some split, create a subset from the  $n$ , say  $k$

make the algorithm choose from that subset of  $k$

$k = \text{root of } n$

## Boosted Decision Trees - ( XGB eXtreme Gradient Boosting )

Work more on the parts which you perform (tree performs) not well - idea of boosting.

# Boosted trees intuition

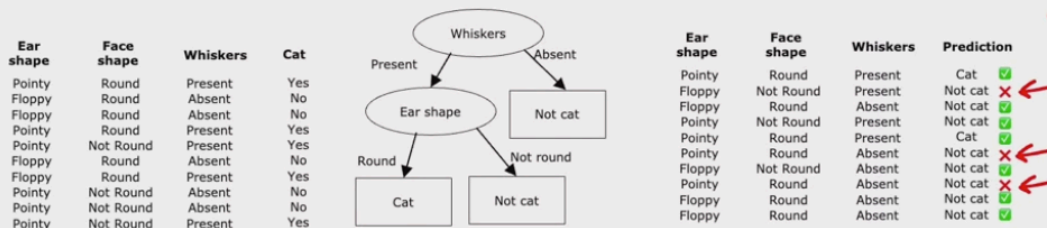
Given training set of size  $m$

For  $b = 1$  to  $B$ :

Use sampling with replacement to create a new training set of size  $m$

But instead of picking from all examples with equal  $(1/m)$  probability, make it more likely to pick misclassified examples from previously trained trees

Train a decision tree on the new dataset



Can use XGB both for classification and regression.

## Conclusions

Decision trees and tree ensembles work well on **Tabular (Structured) Data**

Performs poorly for **unstructured data (images, audio, text)**

## Neural Networks

Works well on **all types of data, including structured and unstructured data.**

May be slower.

Works with transfer learning, eliminate pre-training.

System of multiple models working together is better.