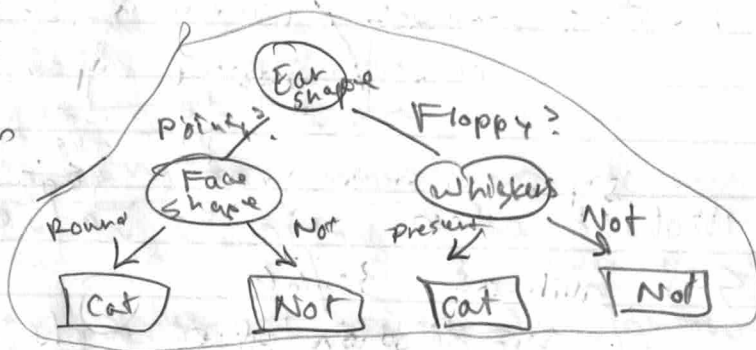


DECISION TREES

Super practical, great results, but not much academic interest!

Example:

- Build from data set
- Can build various trees from the data!



Learning a decision tree

- Pick root node split feature
- For each node in layer $N+1$, choose split feature
- Recur until reaching "pure" nodes → fully segmented classes

Key decisions:

- Choosing split features → want to ^(min)max purity
- When to stop splitting?

- node is 100% one class
- reached max depth
- Purity score improvements don't meet threshold
- # of examples in node is below threshold

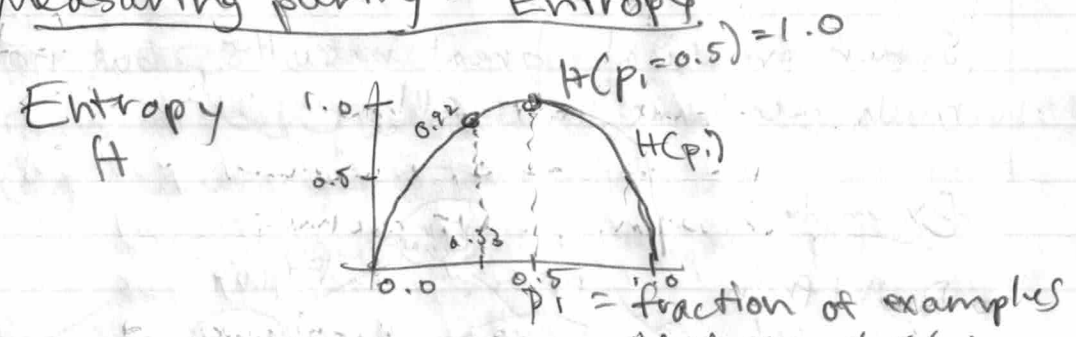
★ Still need to avoid overfitting!

$\log_2 x = y \Leftrightarrow e^y = x \Leftrightarrow 2^{-y} = x \Leftrightarrow 2^y = 1/x \Leftrightarrow \log_2 x = -\log_2 (1/x)$
 $2^a = e \Leftrightarrow \log_2 e = a$

So if $\log_2 x = y$
 then $\log_2 x = y \cdot \log_2 e$

$y = \frac{\log_2 x}{\log_2 e}$

Measuring purity - Entropy



Highest when ratio is 50/50!
 Still high at 33/66!

Define $p_1 =$ fraction of examples w/ label 1
 Define $p_0 = 1 - p_1$

Define:

$$H(p_1) = -p_1 \log_2(p_1) - p_0 \log_2(p_0)$$

$$= -p_1 \log_2(p_1) - (1-p_1) \log_2(1-p_1)$$

$H(0) = 0$

Choosing a Split - Information Gain = Reduction of Entropy

For each possible split feature, compute H with that feature - but weighted by #examples - for both branches it generates. Then compare to "no split" to get Information Gain i.e. reduction in entropy.

$$\text{Information Gain} = H(p_1^{\text{root}}) - (w^{\text{left}} H(p_1^{\text{left}}) + w^{\text{right}} H(p_1^{\text{right}}))$$
 where $w^{\text{left}}, w^{\text{right}}$ are the fraction of the examples falling in the left/right nodes.

Putting it together

Repeatedly until stopping criteria met ^(max depth, info gain below ϵ_1 , examples below ϵ_2)

- for each feature, compute info gain (IG)
- choose the one that maximizes IG
- recursion on L/R subtrees

Feature Encoding for non-binary features

- Categorical (non-binary) - use one-hot encoding
- Continuous values - pick a decision threshold by maximizing info gain

REGRESSION TREES

- Pick splits to reduce variance in subsets ^(weighted sum of)
- maximize (root node var) - ^(weighted sum)
- Final prediction = mean of leaf node examples

TREE ENSEMBLES

→ Decision trees are highly sensitive to changes in training data. So we train a tree ensemble and have them vote on a prediction. This increases robustness.

How? Random Forest algorithm.

- Use sampling with replacement to build data sets
- For each, build a decision tree (Bagged Decision Tree)
- For each node, only allow splits from $k < n$ random features (BGT → Random Forest). $\frac{k}{n}$

XGBoost / Boosted Decision Trees

Most common decision tree algorithm.

Modify Random Forest alg by, for each new random tree, increase prob. of selecting samples that are misclassified by previous trees.

→ Emphasizes samples we're getting wrong

★ Extreme Gradient Boost

- built in regularization
- good default hyperparams
- highly competitive
- fast to train!

(Also: might be more interpretable?)

When to use decision trees?

Tabular / structured data → XGBoost
(regression and classification)
(categorical and continuous)

Unstructured data (images, audio, text) → NNs
(works well on structured & semi-structured too)
Pros: transfer learning, can combine many models
Cons: slower