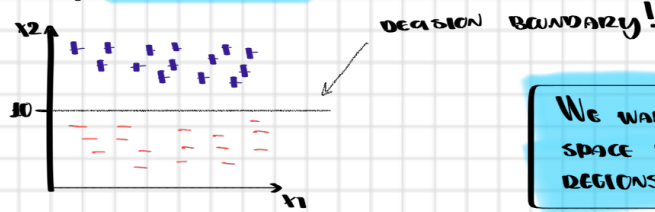
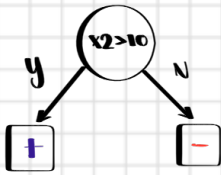


Lecture 11 - 22/10/24

Decision Trees

- A SUPERVISED LEARNING METHOD
- Use AN IF-ELSE STRUCTURE TO DEFINE A DECISION BOUNDARY
 - ↳ Each FOCUSES ON A SINGLE FEATURE

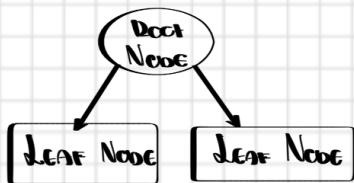


WE WANNA DIVIDE THE FEATURE SPACE INTO DISTINCT, NON-OVERLAPPING REGIONS OR "STRATA"

- Ask Question RECURSIVELY TO CREATE A TREE
- SEVERAL DECISION REGIONS CAN PREDICT THE SAME VALUE
- IN REGRESSION TREES THE OUTPUT LABEL / LABEL IS THE MEAN

GOOD CHARACTERISTICS:

- INTERPRETABILITY → Easy EXPLAINABLE
- EFFICIENCY → Save ONLY THE TREE NOT THE DATA, $O(\log(m))$ FOR TREE OPERATIONS.
- Every decision LEADS TO A BINARY SPLIT



↳ THIS IS CALLED **STUMP**. → IT IS A ONE-LEVEL TREE

SPLITS CAN BE BOTH:

- BINARY split ↘
- Multi-way split ↙

THERE ARE THREE TYPES OF INPUT FEATURES:

- o NUMERICAL
- o CATEGORICAL
- o MIXED

WHICH ALL LEAD TO

CLASSIFICATION TREES

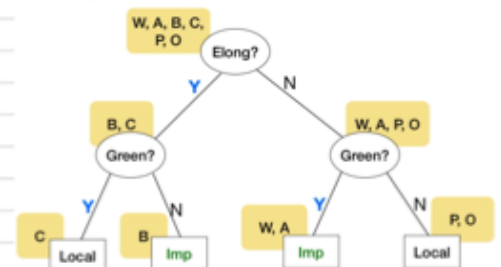
REGRESSION TREES

	Elongate	Big	Green	Class
Watermelon	N	Y	Y	Imported
Apple	N	N	Y	Imported
Banana	Y	N	N	Imported
Cucumber	Y	N	Y	Local
Pumpkin	N	Y	N	Local
Orange	N	N	N	Local



THIS TREE WORKS PERFECTLY, BUT, THERE ARE MULTIPLE TREES THAT ACHIEVE THE GOAL...

SO, WHICH ONE TO CHOOSE?



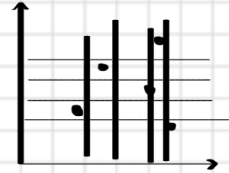
THIS LOOKS BETTER!

How to train a decision tree

How do we build a good decision tree?

- A **BALANCED** tree is good!
 - We **USE HEURISTIC** aiming for good trees
 - USE GREEDY algo** ~ At **EACH NODE**, make a **SPLIT** that **BEST DIVIDES** on **SPECIFIC CRITERION**
- {

 - Try** **all possible splits** for the current node
 - Use HEURISTIC** to measure how good each split is
 - Pick BEST SPLIT** for that node
 - REPEAT** **three steps above** **RECURSIVELY** for the **children**



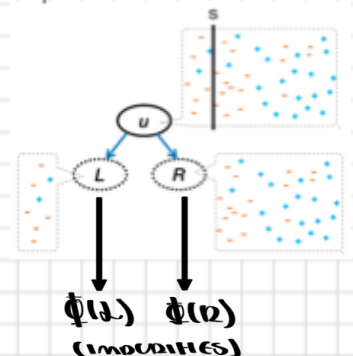
Criterion for evaluating quality of a split (in regression trees)

- How good is split S on these data points?
 - Pretend we've used S to **DIVIDE DATA POINTS** at the node
 - CALCULATE THE LABELS** \hat{y} of the two resulting regions (the mean) $\rightarrow y_{R1} \ y_{R2}$
 - Calculate **RSS** \forall data point at the node, **SUMMED OVER BOTH REGIONS**.

$$RSS = \sum_{i: x_i \in R_1} (y_i - y_{R1})^2 + \sum_{i: x_i \in R_2} (y_i - y_{R2})^2$$
 - RSS IS USED AS THE CRITERION FOR EVALUATING A SPLIT FOR A REGRESSION TREE!**
 - Pick the split with the minimum RSS** for that node

Criterion for evaluating quality of a split (in classification trees)

- How good is split S on these data points?
 - Pretend we've used S to **DIVIDE THE DATA POINTS** at **node u**
 - SPLIT CREATES** two child, **node L , node R**
 - Concept of **Impurity**:
 - A **MEASURE OF HOW MIXED THE DATA IS** w/ **class members**.



- What is then the **Quality of split S** ?
 - We **CANNOT** just **SUM UP THE IMPURITIES** of the two children!
 - USE WEIGHTS**:

$$\begin{matrix} \textcircled{L} & \textcircled{R} \\ P_L = N_L / N_u & P_R = N_R / N_u \end{matrix}$$

N_L = data points in left
 N_R = data points in right
 N_u = data points at u

$G = 0 \quad t = y \quad !!!$

- Combining the **impurities** using **WEIGHTED AVG'S**:

$$P_L \hat{\phi}(L) + P_R \hat{\phi}(R)$$

BEST SPLIT IS THE ONE THAT MINIMIZES THIS!

OR

STOP IF THERE IS NO SPLIT THAT GIVES US **MORE INFO.**

maximize $\hat{\phi}(t) - (P_R \hat{\phi}(t_R) + P_L \hat{\phi}(t_L))$

Impurity Function

- As said it is a **MEASURE OF HOW MIXED DATA IS**
- Φ CAN BE DEFINED AS A **FUNCTION OF POSTERIOR PROBABILITIES**:

$$\Phi(p_0, p_1, \dots, p_{K-1})$$



- Estimated posterior probability for **class k at node t**:

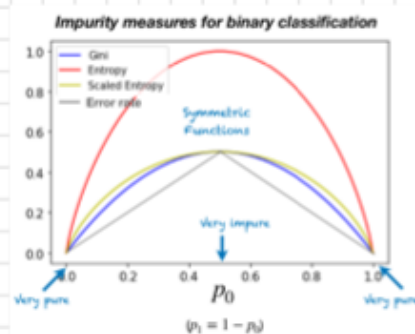
$$p_k = \hat{p}(k|t) = \frac{N_k}{N_t}$$

- Φ is **max** only when all p_k are equal.
- Φ is **minimum** if all points are in the same class.
- Φ is **symmetric**, i.e., it doesn't care about the order of input probs.

- Impurity Functions:**

- Gini impurity**: $\Phi(t) = 1 - \sum_{k=1}^K p_k^2$
- Entropy**: $\Phi(t) = - \sum_{k=1}^K p_k \log p_k$
-
- Classification Error Rate**

NOT A GREAT ONE

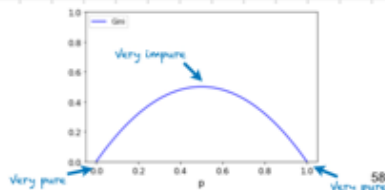


Gini Index

$$p_k = \hat{p}(k|t) = \frac{N_k}{N_t}$$

$$G(t) = \sum_{k=1}^K p_k(1-p_k) = 1 - \sum_{k=1}^K p_k^2$$

For a binary classification: $G(t) = 2p(1-p)$



Entropy

$$p_k = \hat{p}(k|t) = \frac{N_k}{N_t}$$

$$H(t) = - \sum_{k=1}^K p_k \log p_k$$

Entropy of random vars is the avg of "uncertainty" whereat to the vars possible outcomes

When everything is uniform = **Entropy maximized**

WE WANNA **minimize entropy** by every split!