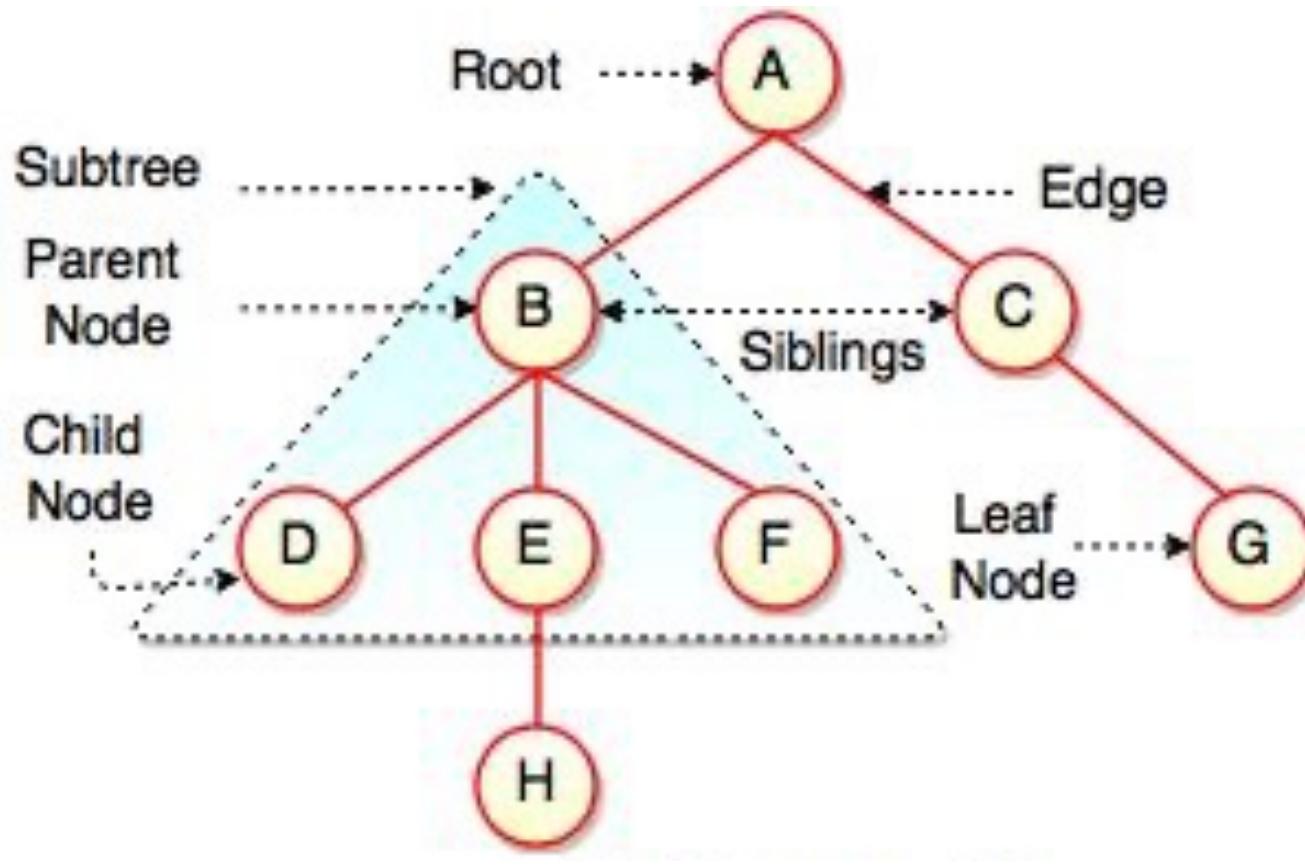


Topic 6

Decision Tree

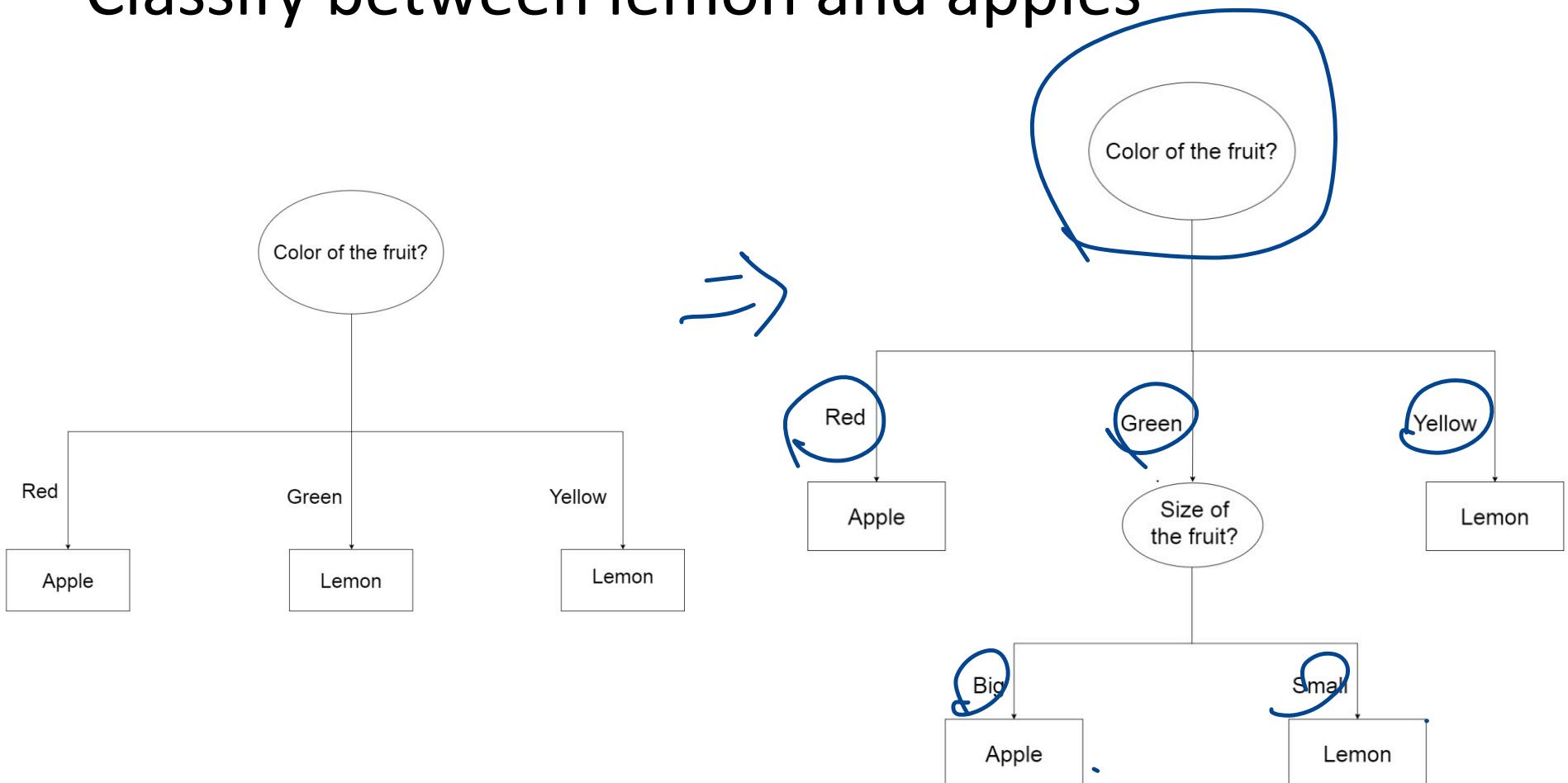
What are trees?

Root
Stem
Branches
Twigs
Leaves

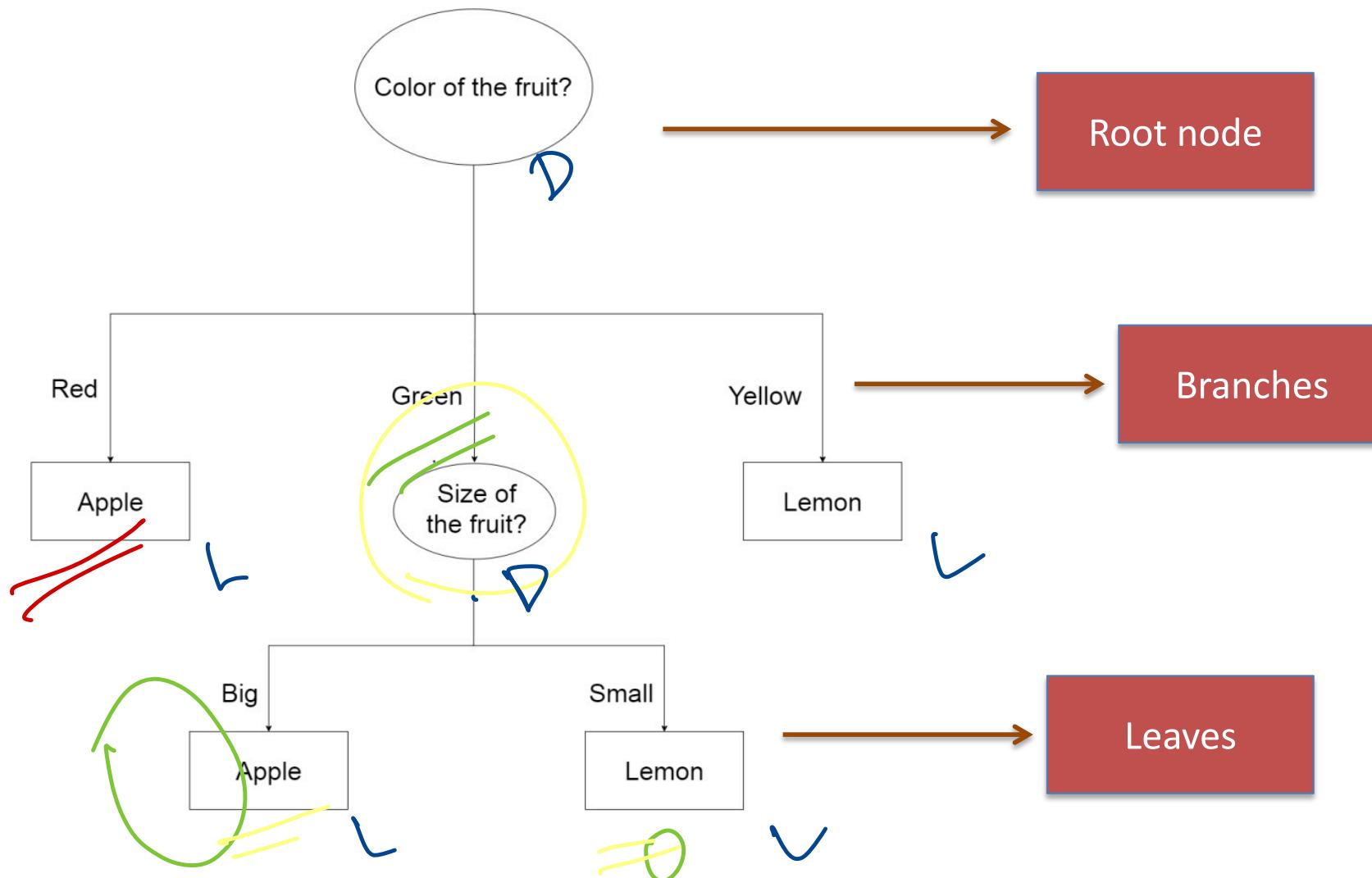


Decision Trees

- Classify between lemon and apples



Decision Trees

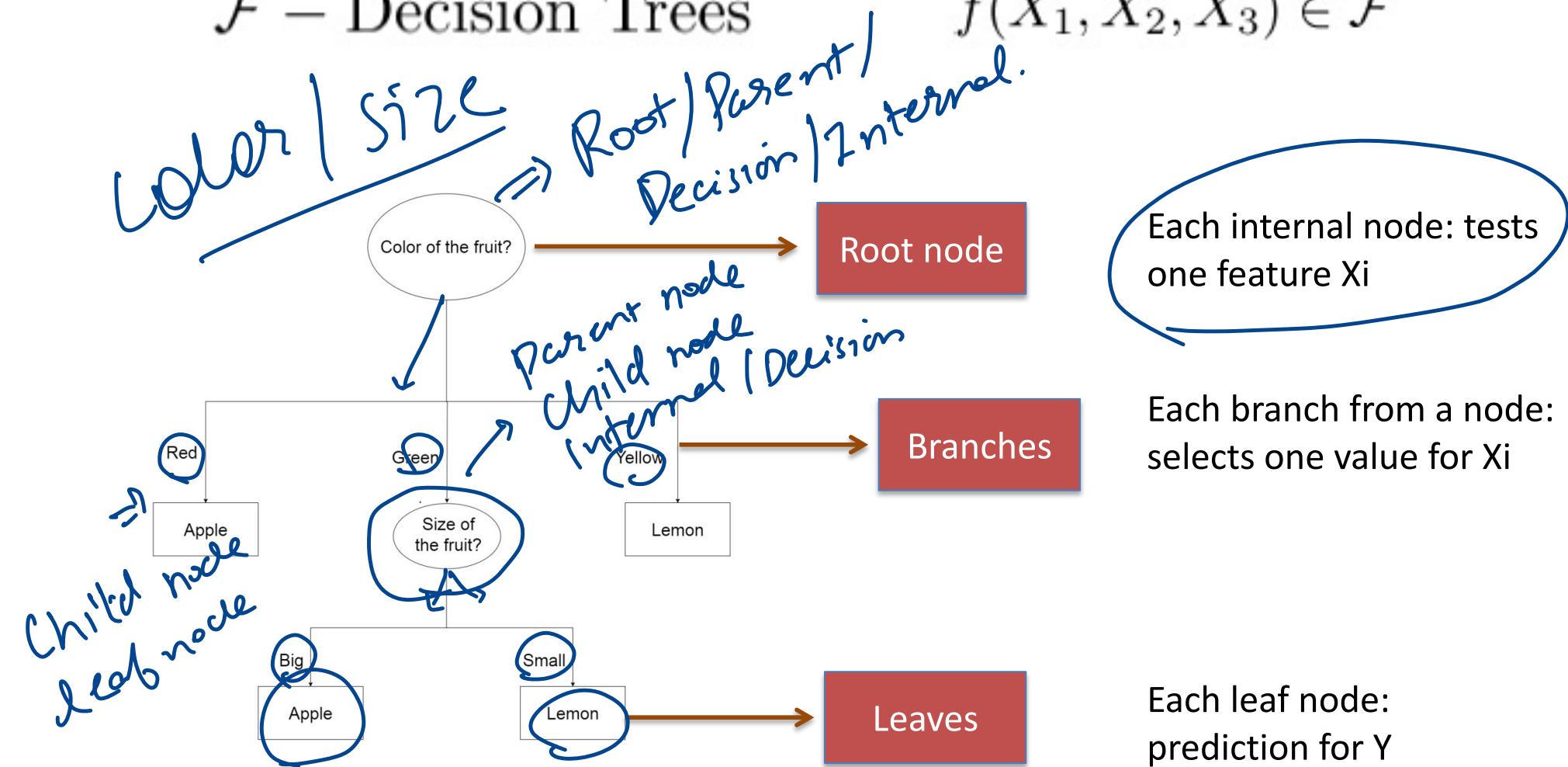


Rules for classifying data using attributes

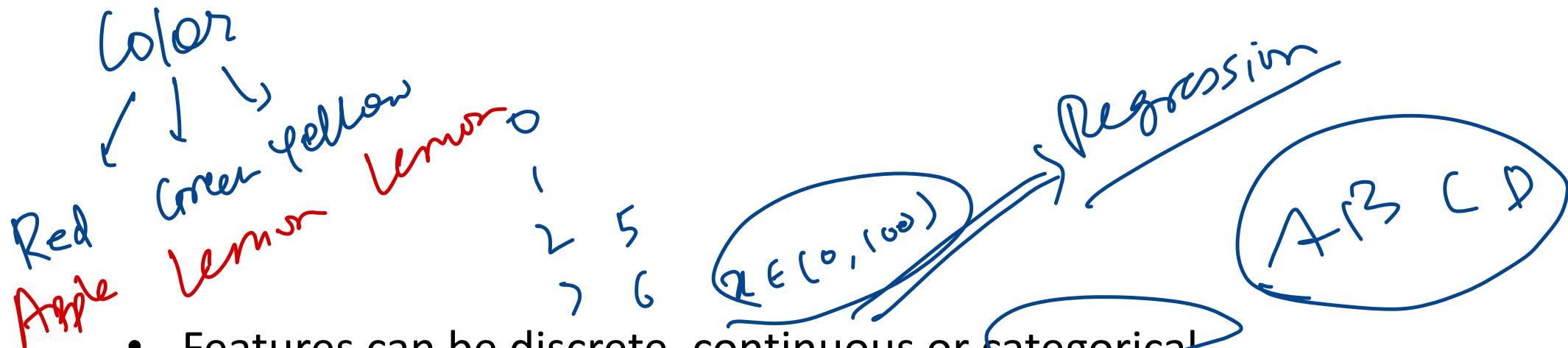
- The tree consists of decision nodes and leaf nodes.
- A decision node has two or more branches, each representing values for the attribute tested.
- A leaf node attribute produces a homogeneous result (all in one class), which does not require additional classification testing

\mathcal{F} – Decision Trees

$$f(X_1, X_2, X_3) \in \mathcal{F}$$



Features can be discrete, continuous or categorical

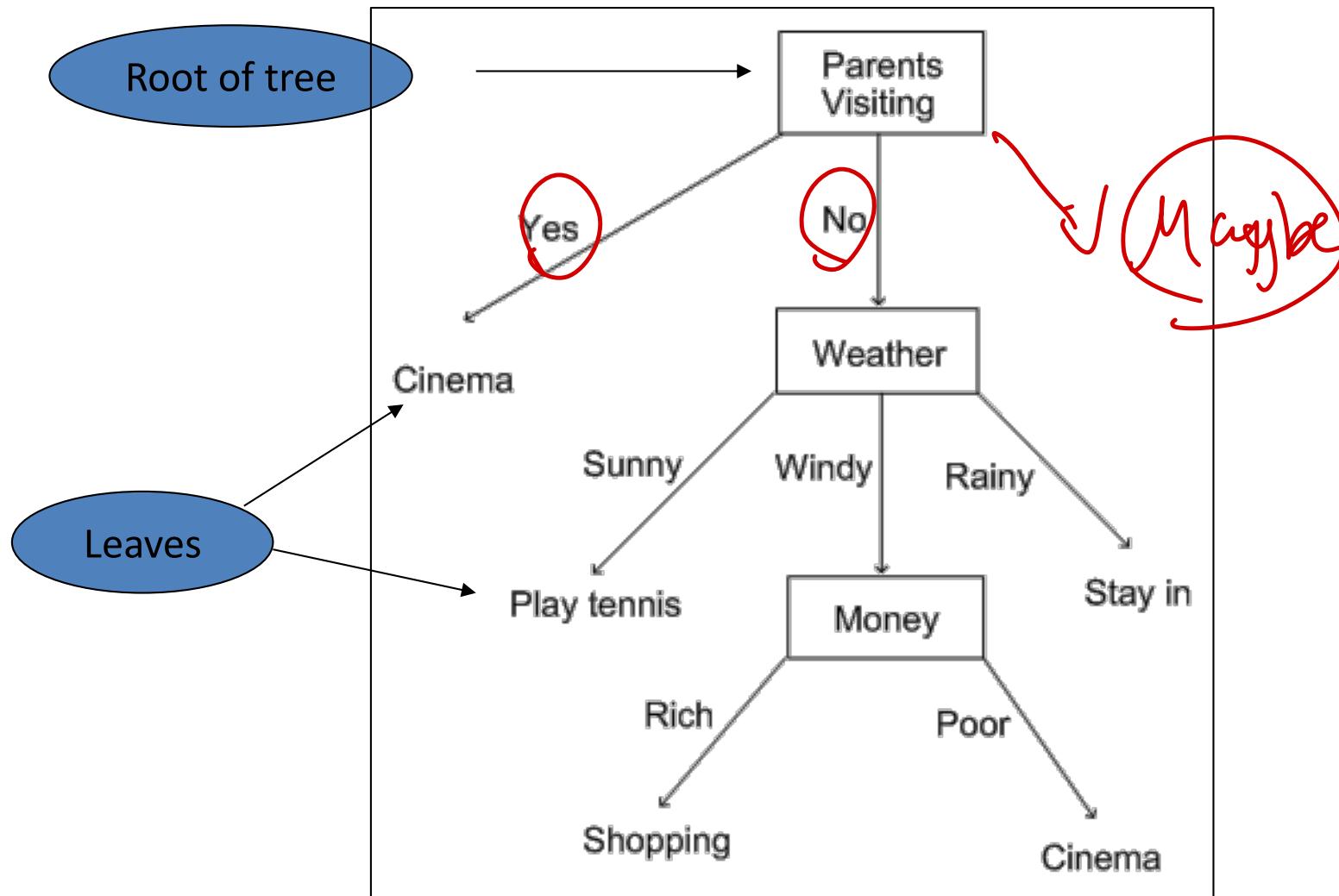


- Features can be discrete, continuous or categorical
- Each internal node: test some set of features $\{X_i\}$
- Each branch from a node: selects a set of value for $\{X_i\}$
- Each leaf node: prediction for Y

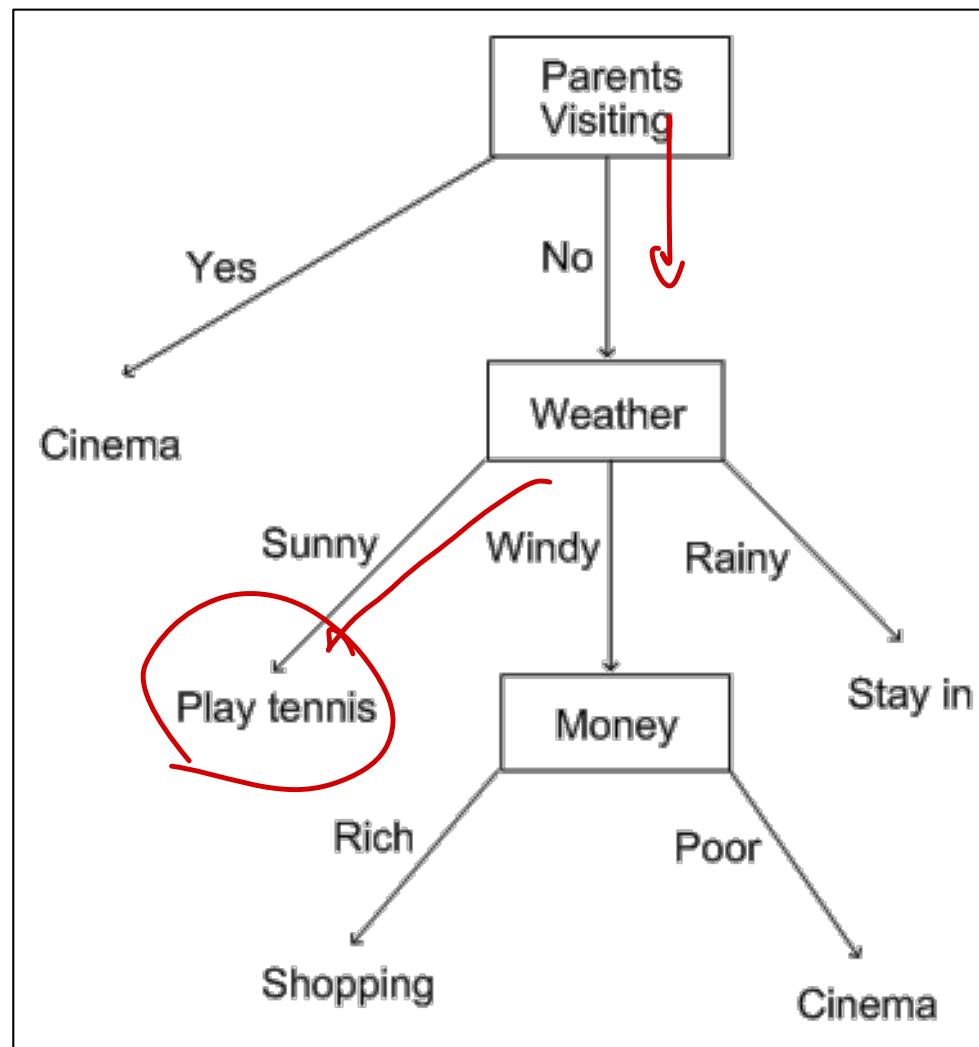
Example: What to do this Weekend?

- If my parents are visiting ✓
 - We'll go to the cinema
- If not
 - Then, if it's sunny I'll play tennis
 - But if it's windy and I'm rich, I'll go shopping
 - If it's windy and I'm poor, I'll go to the cinema
 - If it's rainy, I'll stay in

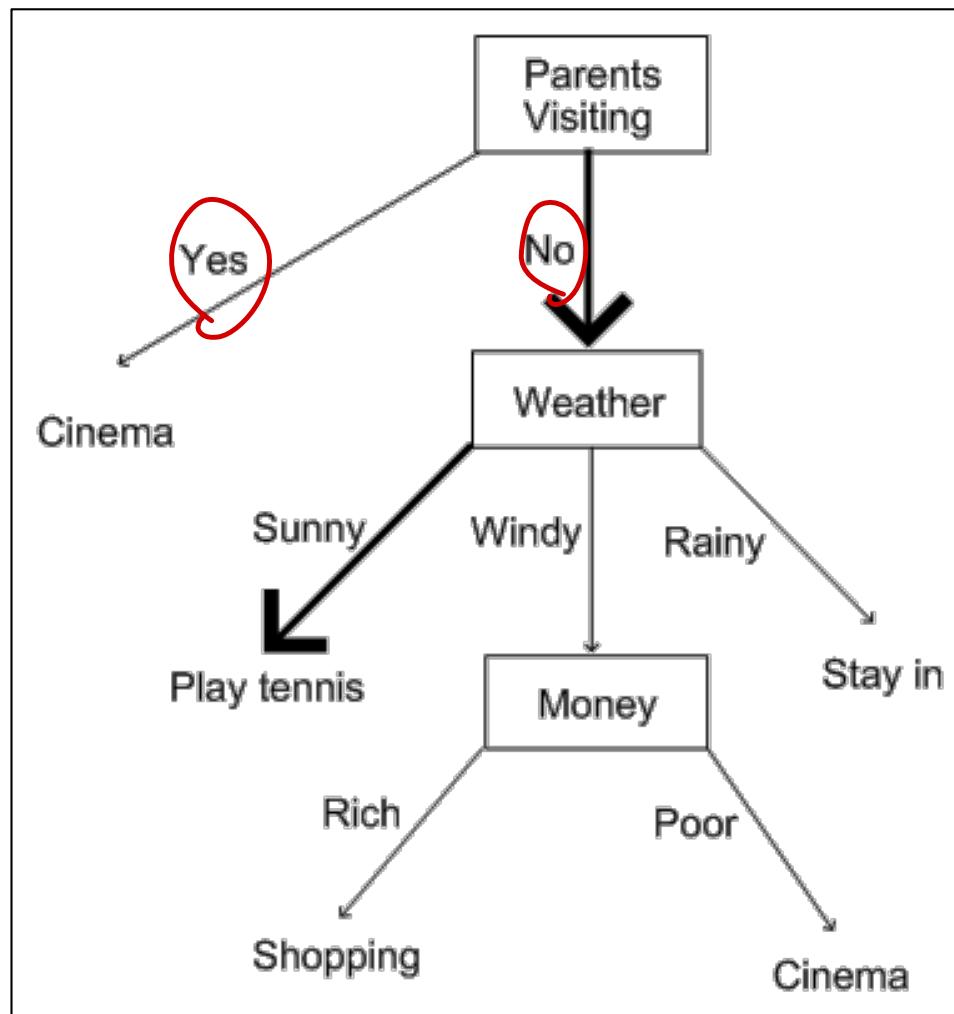
Written as a Decision Tree



Using the Decision Tree (No parents on a Sunny Day)



Using the Decision Tree (No parents on a Sunny Day)



From Decision Trees to Logic

- Read from the root to every tip
 - If this and this and this ... and this, then do this
- In our example:
 - If no_parents and sunny_day, then play_tennis
 - $\text{no_parents} \wedge \text{sunny_day} \rightarrow \text{play_tennis}$

How to design a decision tree

- Decision tree can be seen as rules for performing a categorisation
 - E.g., “what kind of weekend will this be?”
- Remember that we’re learning from examples
 - Not turning thought processes into decision trees
- The major question in decision tree learning is
 - Which nodes to put in which positions
 - Including the root node and the leaf nodes

Training and Visualizing a Decision Tree

```
from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier

iris = load_iris(as_frame=True)
X_iris = iris.data[["petal length (cm)", "petal width (cm)"]].values
y_iris = iris.target

tree_clf = DecisionTreeClassifier(max_depth=2, random_state=42)
tree_clf.fit(X_iris, y_iris)
```

```
{ from sklearn.tree import export_graphviz

    export_graphviz(
        tree_clf,
        out_file="iris_tree.dot",
        feature_names=["petal length (cm)", "petal width (cm)"],
        class_names=iris.target_names,
        rounded=True,
        filled=True
    )
}
```

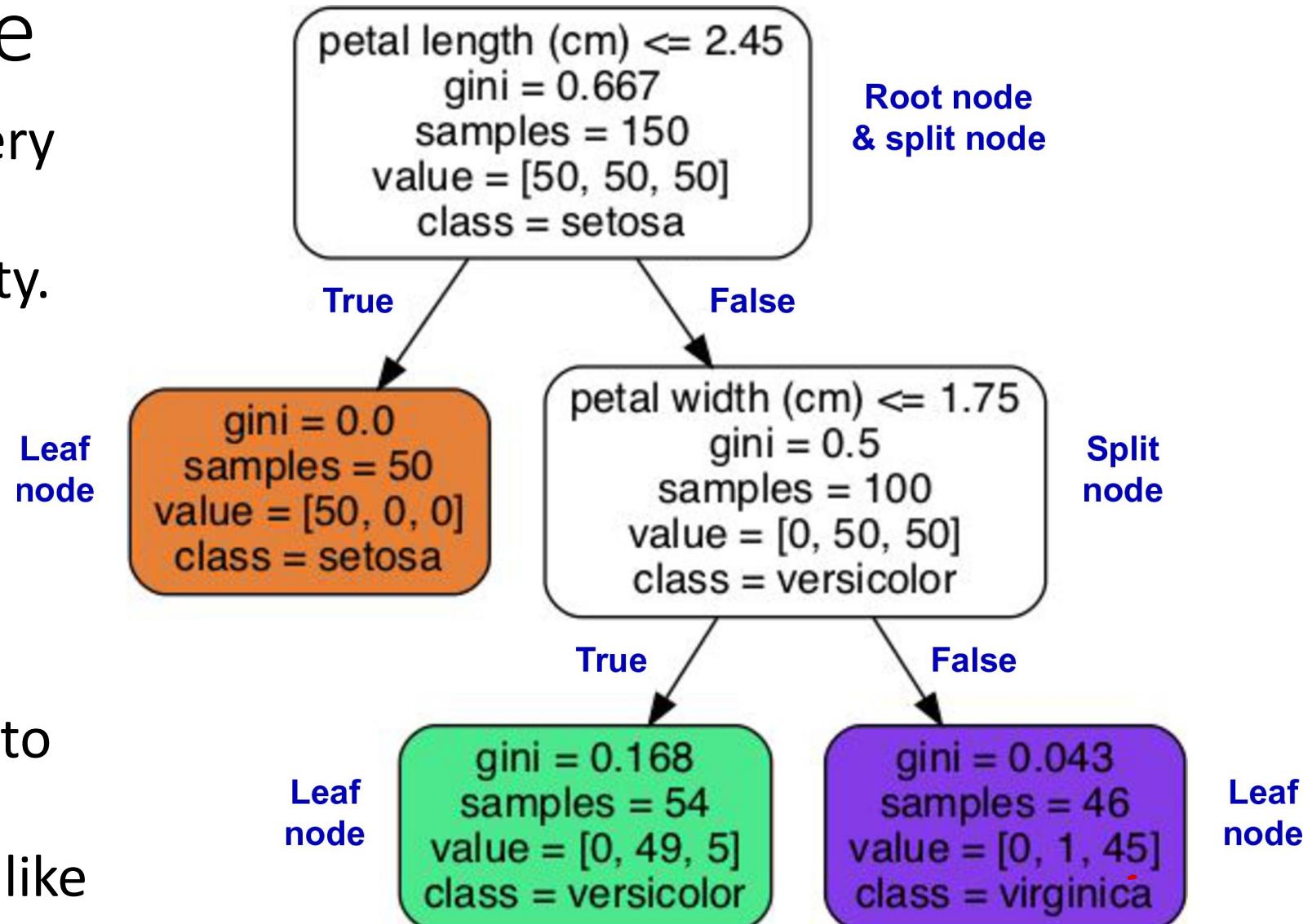
```
from graphviz import Source
Source.from_file("iris_tree.dot")
```

Iris Decision Tree

- Decision trees require very little data preparation.
- Measure of node impurity.
Pure (gini=0)
- Gini Impurity:

$$G_i = 1 - \sum_{k=1}^n p_{i,k}^2$$

- max_depth=2
- White box model – easy to interpret
- Not a black box model – like neural networks



$$1 - (0/54)^2 - (49/54)^2 - (5/54)^2 \approx 0.168.$$

Iris Decision Tree

- Decision trees require very little data preparation.
- Measure of node impurity. Pure (gini=0)
- Gini Impurity:
$$G_i = 1 - \sum_{k=1}^n p_{i,k}^2$$
- max_depth=2
- White box model – easy to interpret
- Not a black box model – like neural networks

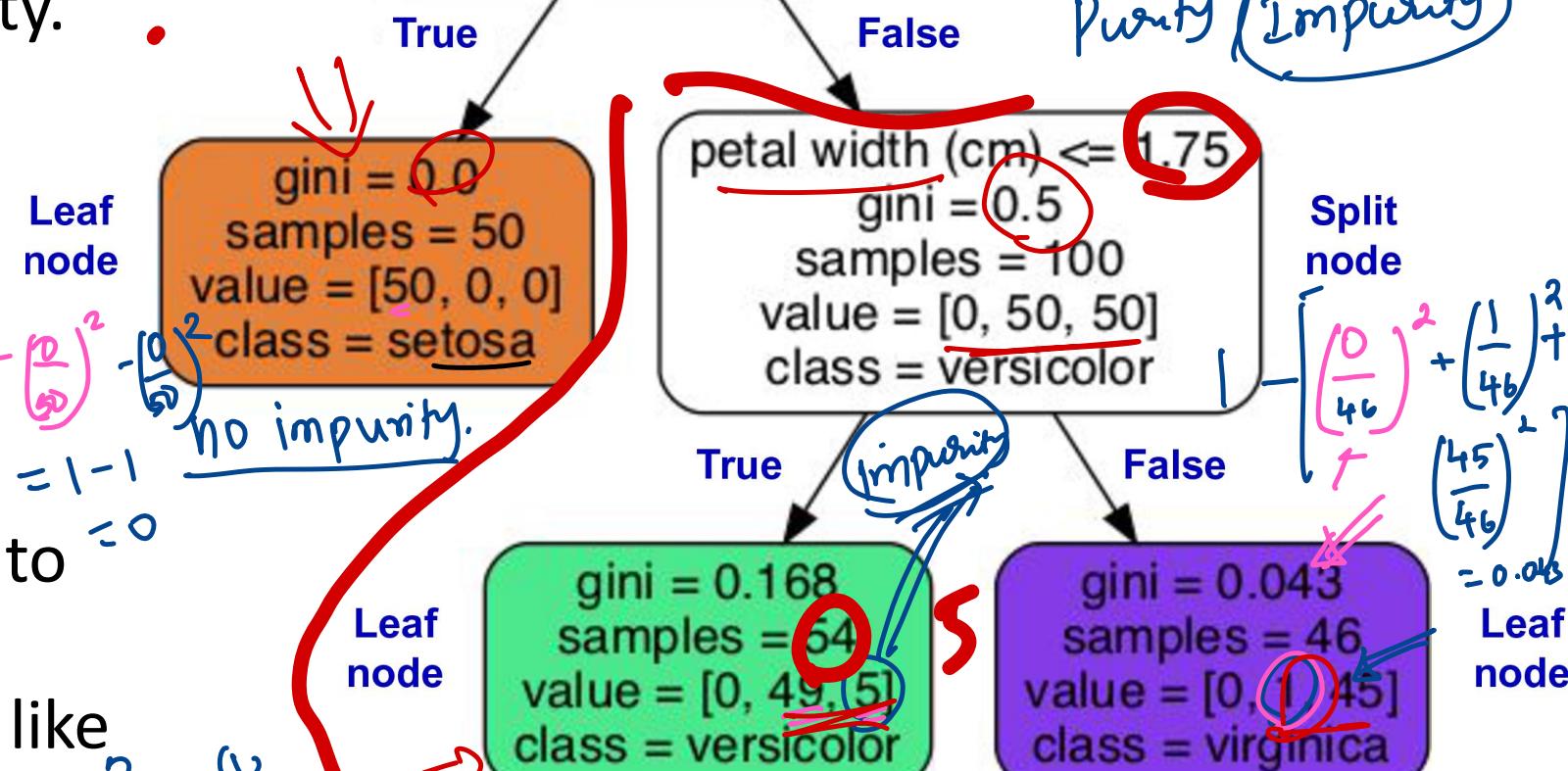
$$G_{\text{node}} = 1 - \sum_{k=1}^n p_k^{(1)} p_k^{(2)}$$

$$(1 - (0/54)^2 - (49/54)^2 - (5/54)^2) \approx 0.168$$

$$1 - \left\{ \left(\frac{50}{150}\right)^2 + \left(\frac{50}{150}\right)^2 + \left(\frac{50}{150}\right)^2 \right\}$$

petal length (cm) ≤ 2.45 cm
 \Rightarrow gini = 0.667
 samples = 150
 value = [50, 50, 50]
 class = setosa

Root node & split node



Gini

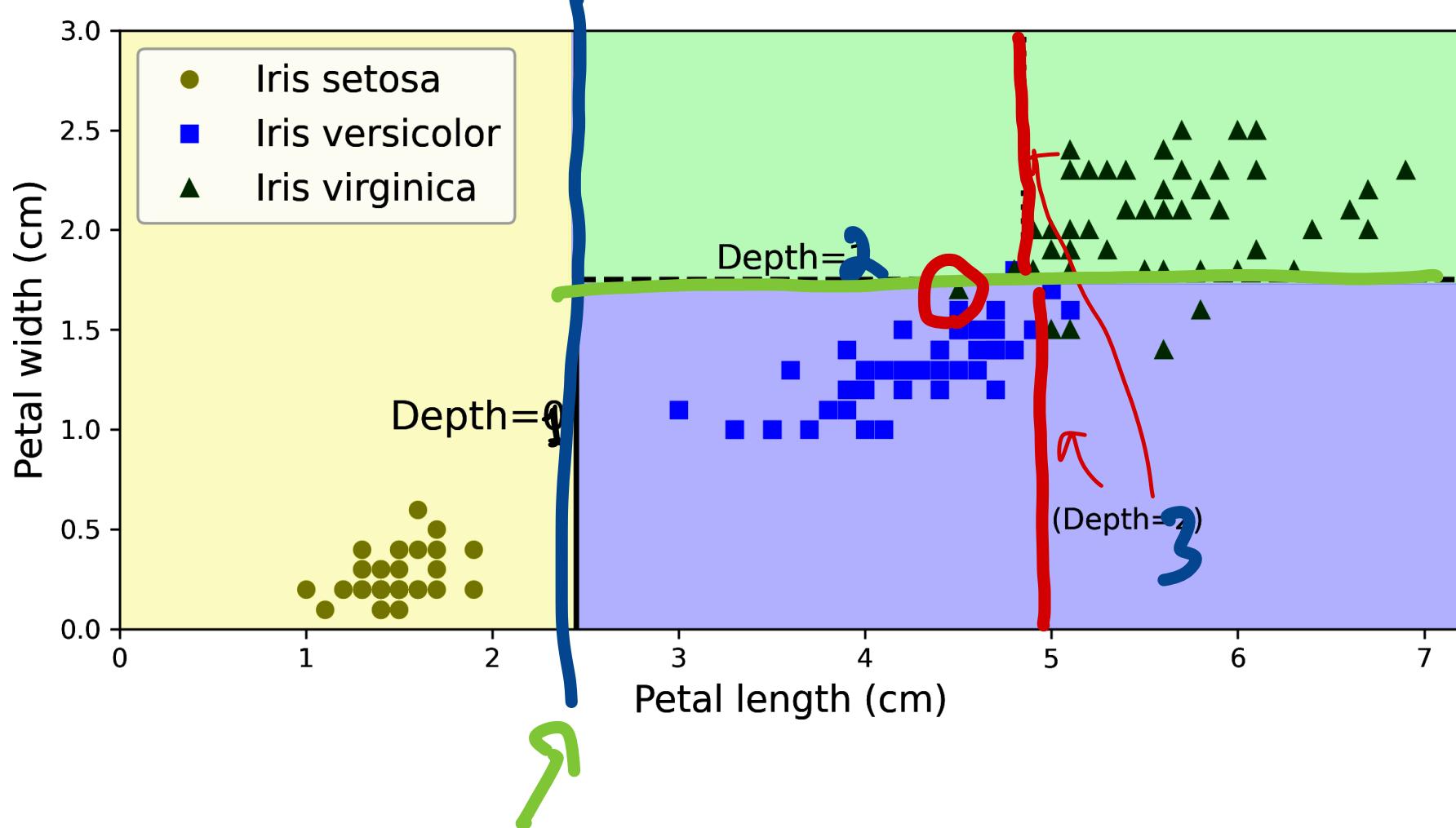
50	50	50
----	----	----

$$1 - \left\{ \left(\frac{50}{150} \right)^2 + \left(\frac{50}{150} \right)^2 + \left(\frac{50}{150} \right)^2 \right\}$$
$$1 - \left(\frac{1}{9} + \frac{1}{9} + \frac{1}{9} \right) = \frac{2}{3}$$

$$1 - \left\{ \left(\frac{50}{100} \right)^2 + \left(\frac{50}{100} \right)^2 \right\}$$
$$1 - \left(\frac{1}{2} + \frac{1}{2} \right) = \frac{1}{2}$$

Decision Tree Boundaries

- For `max_depth=2`
- `Max_depth=3` add the vertical dotted lines.



CART Training Algorithm

- Classification and regression Algorithm
- Splits training set into two subsets using a single feature k and a threshold t_k in this case petal length less than 2.45.
- Searches for (k, t_k) combination that produces the purest subsets, weighted by their size. The cost function is:

$$J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$

where $\begin{cases} G_{\text{left/right}} \text{ measures the impurity of the left/right subset} \\ m_{\text{left/right}} \text{ is the number of instances in the left/right subset} \end{cases}$

$k \Rightarrow$ feature
 $t_k \Rightarrow$ threshold. t_k

$$\begin{aligned} & \frac{54 \times 0.168}{100} + \frac{46 \times 0.043}{100} \\ &= (0.1105) \end{aligned}$$

Decision Tree Regularisation

Hyperparameters

- Recall in Data pre-processing lecture – decision tree gave a model with 0 error – they have a high tendency to overfit. Hence regularisation!!!

1 **max_depth**: Maximum depth of the tree in terms of layers.

3

- Max_features**: Maximum number of features that are evaluated for splitting at each node
- Max_leaf_nodes**: Maximum number of leaf nodes
- min_samples_split**: Minimum number of samples a node must have before it can be split
- min_samples_leaf**: Minimum number of samples a leaf node must have to be created
- min_weight_fraction_leaf**: Same as min_samples_leaf but expressed as a fraction of the total number of weighted instances

4

γ

2

Non-regularized vs. Regularized Tree

- Test accuracy:

No restrictions:

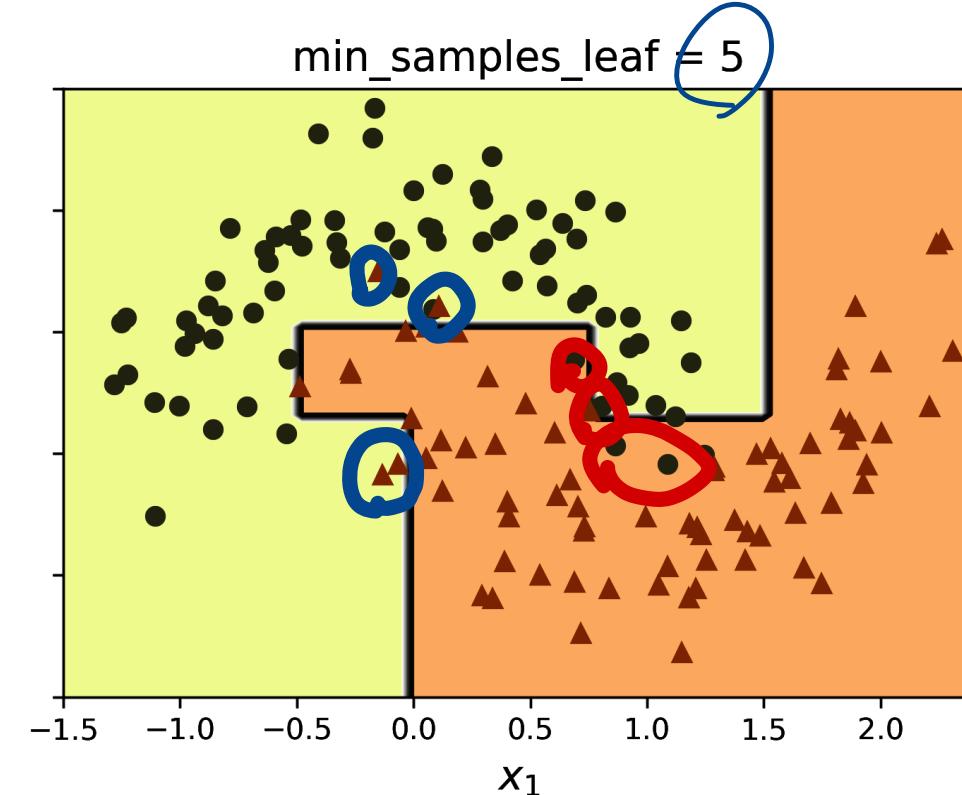
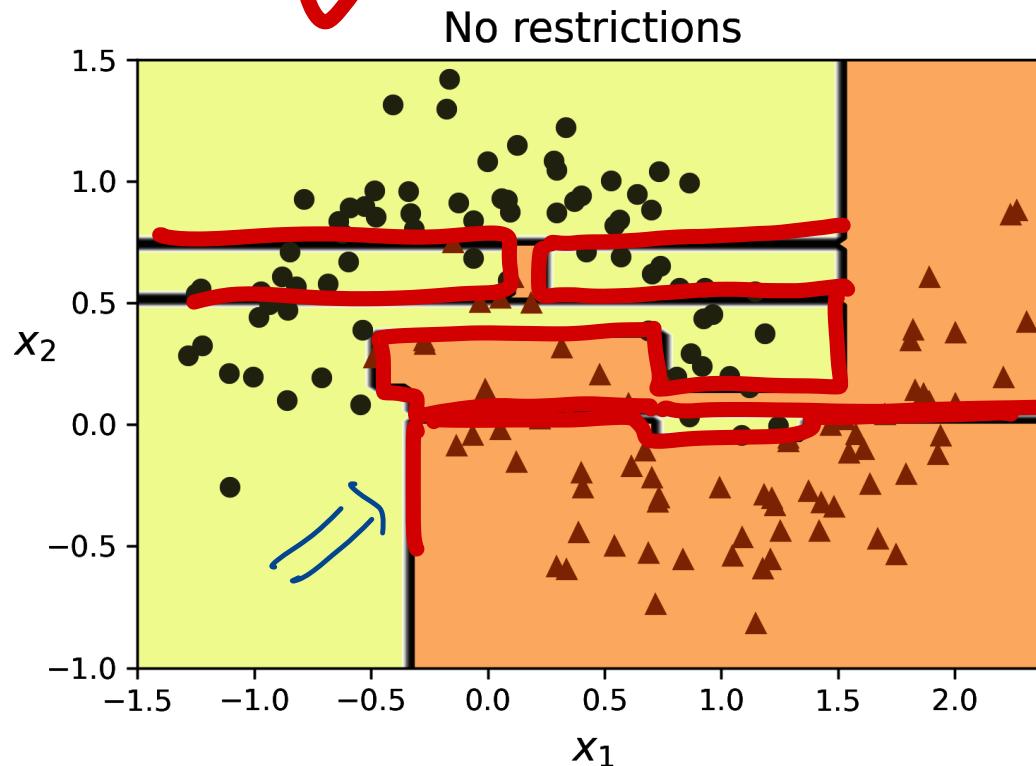
0.898

Restricted:

0.92

```
from sklearn.datasets import make_moons  
X_moons, y_moons = make_moons(n_samples=150, noise=0.2, random_state=42)
```

```
tree_clf1 = DecisionTreeClassifier(random_state=42)  
tree_clf2 = DecisionTreeClassifier(min_samples_leaf=5, random_state=42)  
tree_clf1.fit(X_moons, y_moons)  
tree_clf2.fit(X_moons, y_moons)
```



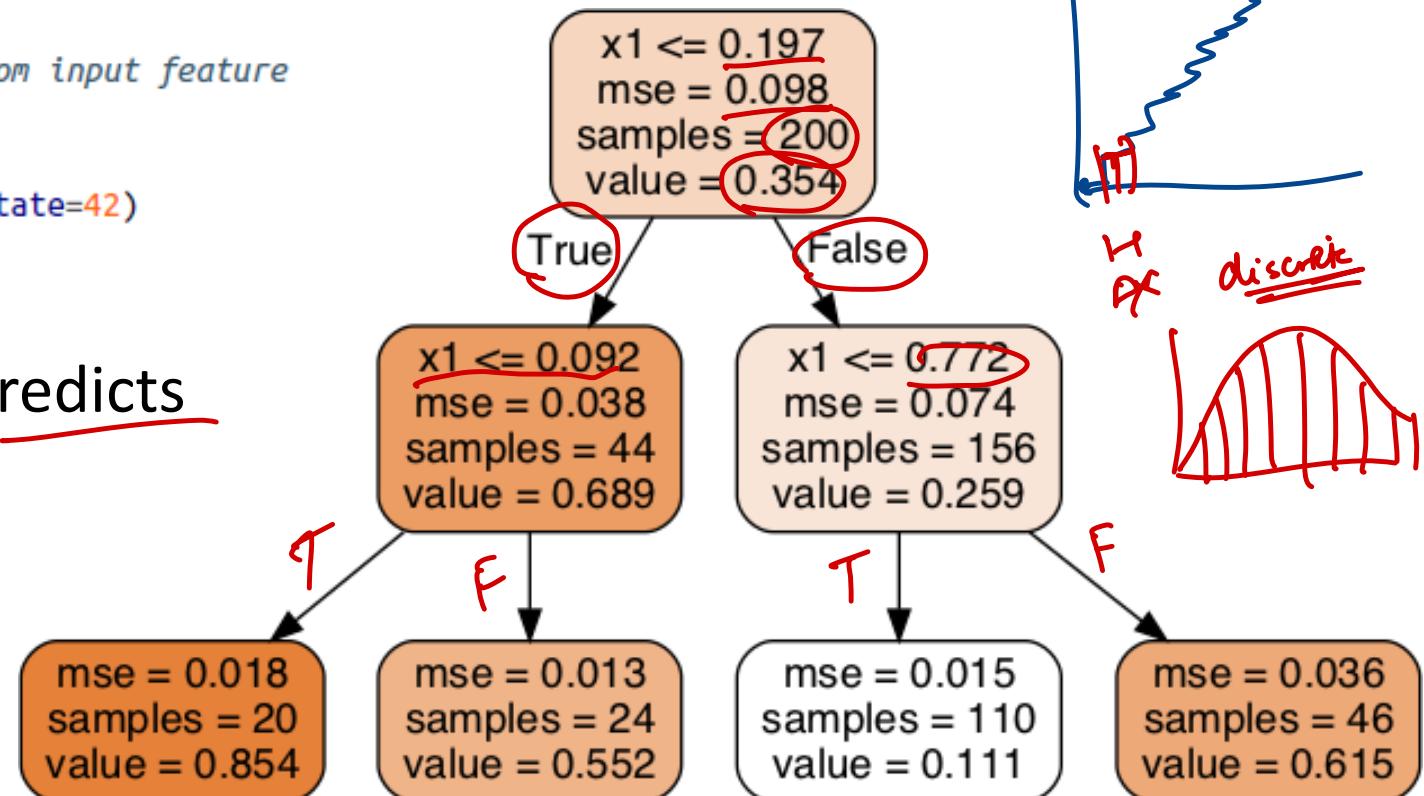
Regression with Decision Trees

```
import numpy as np
from sklearn.tree import DecisionTreeRegressor

np.random.seed(42)
X_quad = np.random.rand(200, 1) - 0.5 # a single random input feature
y_quad = X_quad ** 2 + 0.025 * np.random.randn(200, 1)

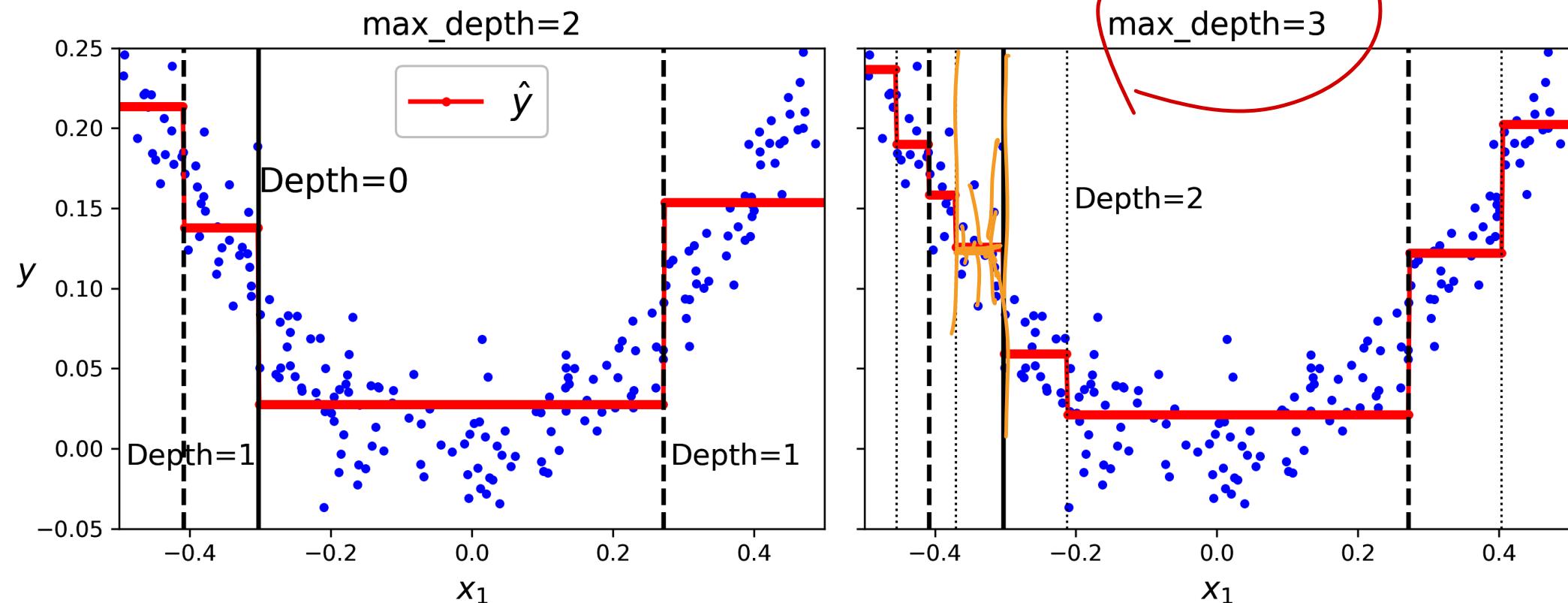
tree_reg = DecisionTreeRegressor(max_depth=2, random_state=42)
tree_reg.fit(X_quad, y_quad)
```

- Instead of predicting a class, it predicts a value.
- Example $x_{\text{new}} = 0.2$



Regression with Decision Trees

- Predicted values for each region is the average target value of the instances in that region.



CART for Regression

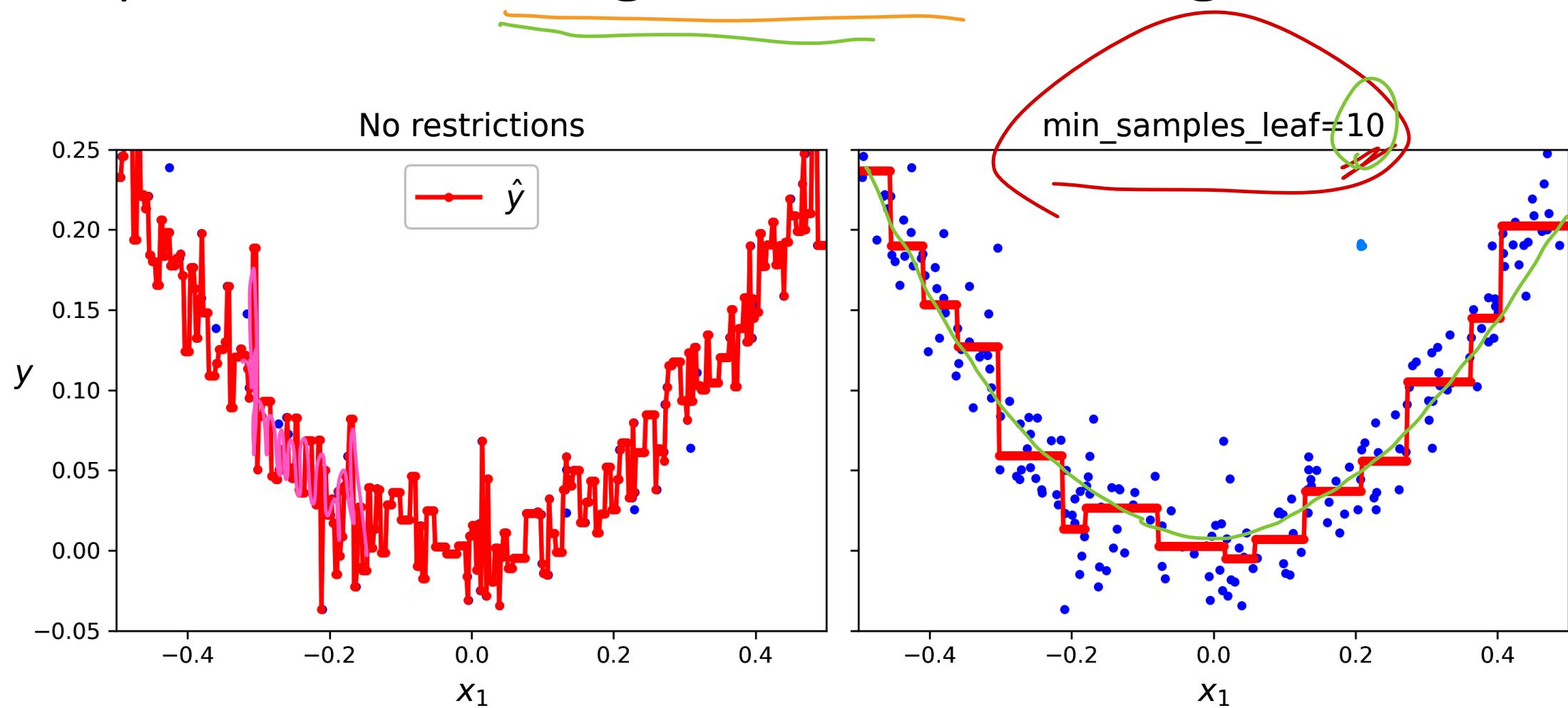
- Instead of trying to minimize impurity, tries to split the training data in order to minimize the MSE.

$$J(k, t_k) = \frac{m_{\text{left}}}{m} \text{MSE}_{\text{left}} + \frac{m_{\text{right}}}{m} \text{MSE}_{\text{right}}$$

where
$$\begin{cases} \text{MSE}_{\text{node}} = \frac{\sum_{i \in \text{node}} (\hat{y}_{\text{node}} - y^{(i)})^2}{m_{\text{node}}} \\ \hat{y}_{\text{node}} = \frac{\sum_{i \in \text{node}} y^{(i)}}{m_{\text{node}}} \end{cases}$$

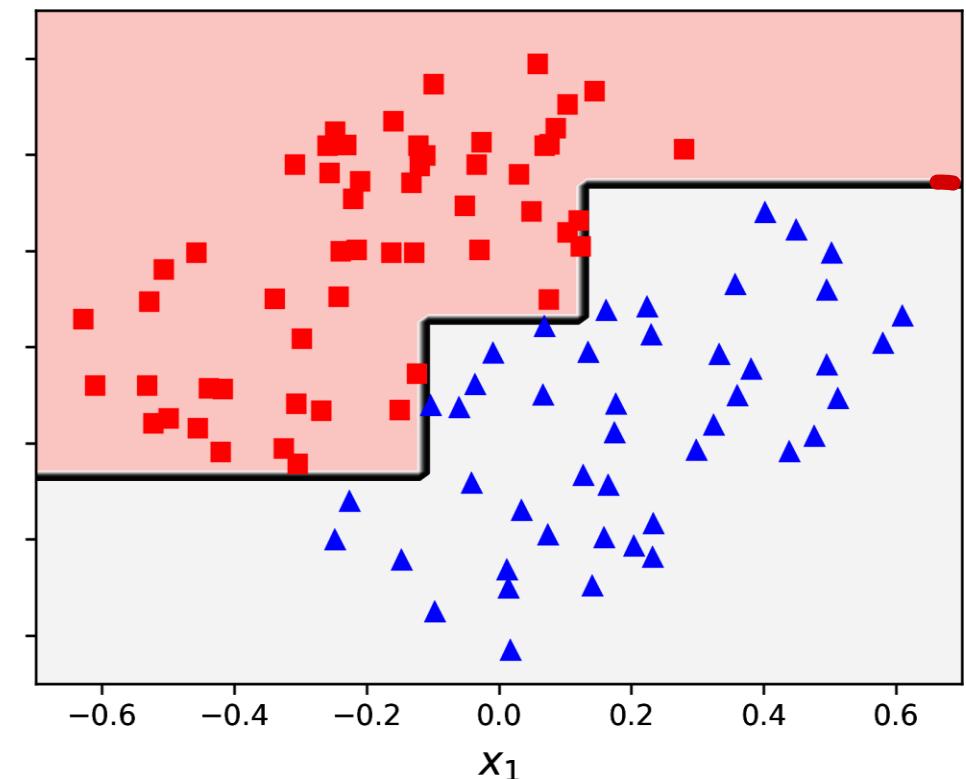
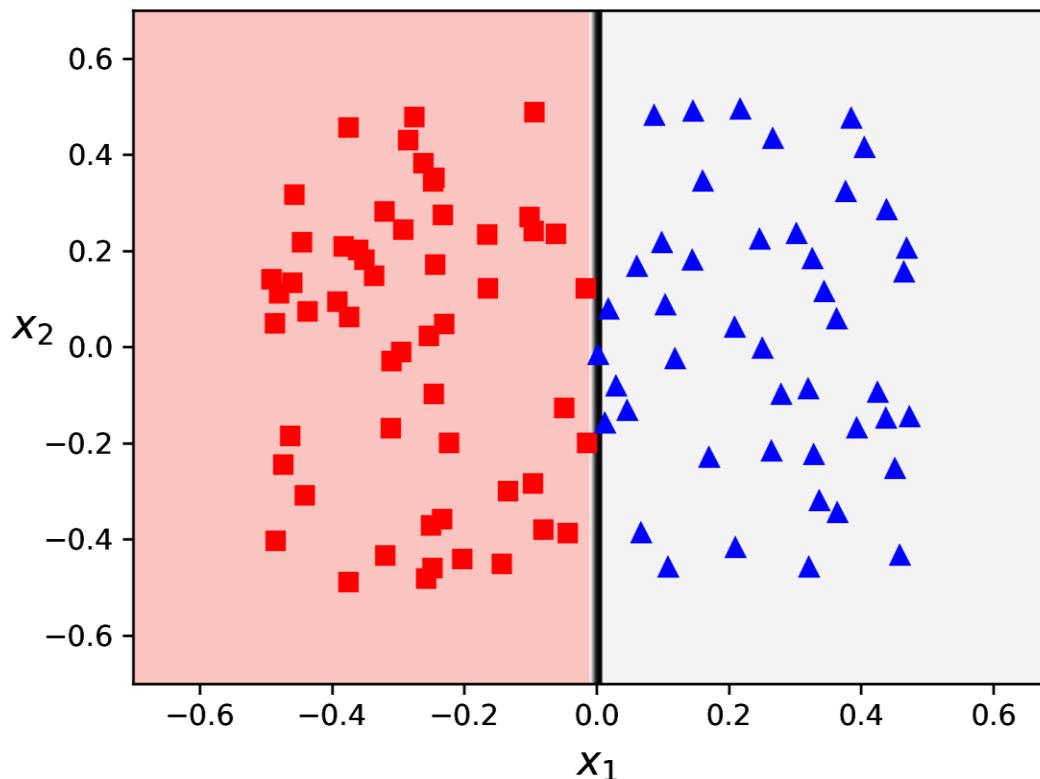


Importance of Regularization in Regression



Sensitivity to Axis Orientation

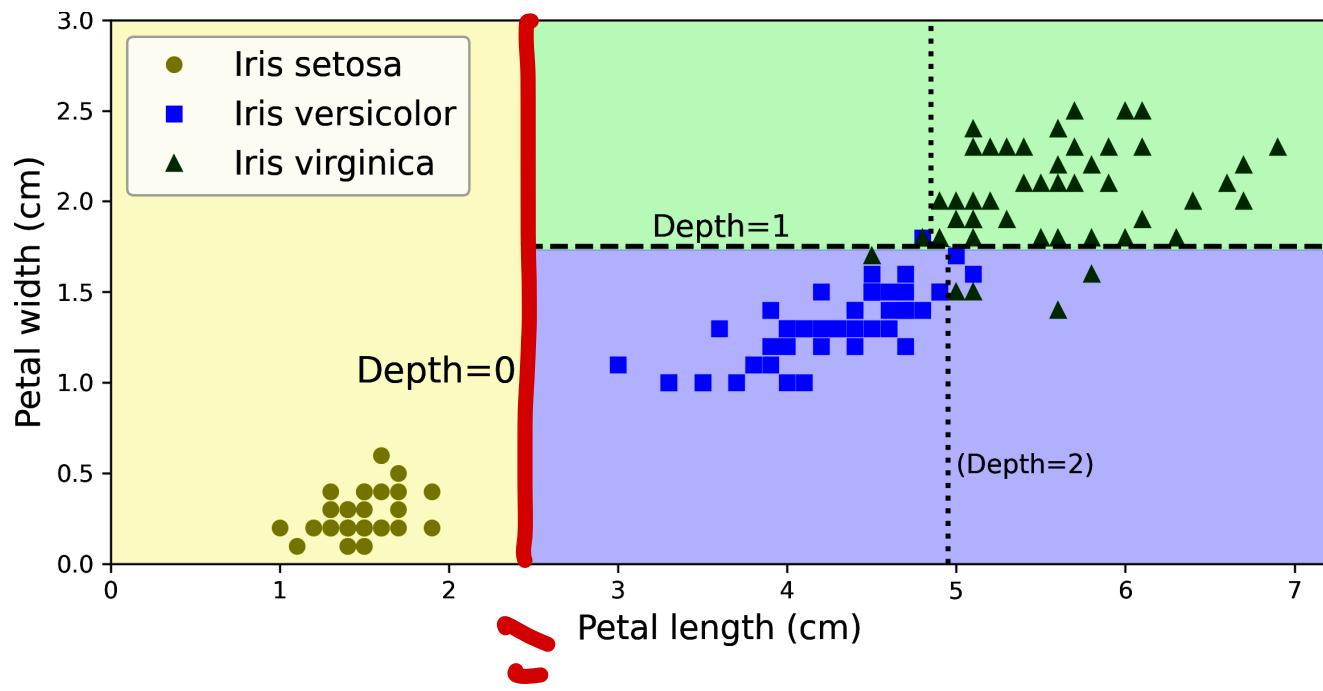
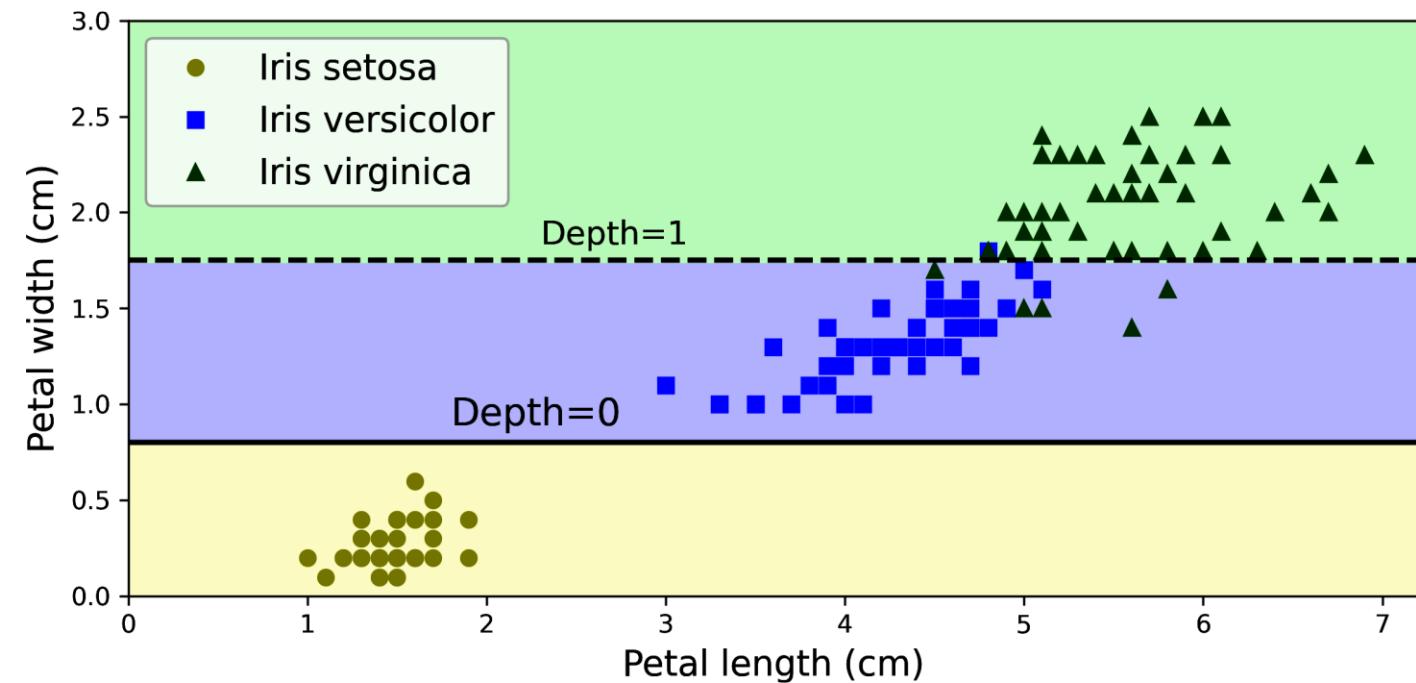
- Decision trees love orthogonal decision boundaries.
- Rotate the data by 45^0 and note the convoluted boundary.



High Variance in Decision Trees

- Decision Trees have high variance.
- Small changes in hyperparameters leads to very different models
- Since the training algorithm used by Sci-Kit learn is stochastic in nature, retraining the same model on the same data produces a very different model.
- By averaging predictions over many trees, it's possible to reduce the variance. Such an ensemble of trees is called a random forest.
- Next slide shows an example of the same dataset and two different tree configurations.

High Variance



The ID3 Algorithm

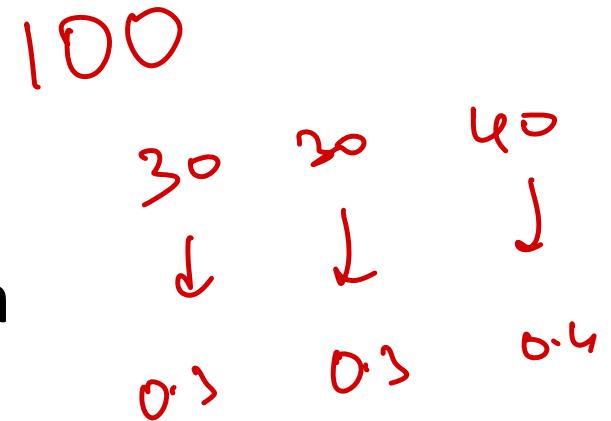
- Invented by J. Ross Quinlan in 1979
- ID3 uses a measure called Information Gain
 - Used to choose which node to put next
- Node with the highest information gain is chosen
 - When there are no choices, a leaf node is put on
- Builds the tree from the top down, with no backtracking
- Information Gain is used to select the most useful attribute for classification

Entropy – General Idea

- From Tom Mitchell's book:
 - “In order to define information gain precisely, we begin by defining a measure commonly used in information theory, called entropy that characterizes the (im)purity of an arbitrary collection of examples”
- A notion of impurity in data
- A formula to calculate the homogeneity of a sample
- A completely homogeneous sample has entropy of 0
- An equally divided sample has entropy of 1

Entropy - Formulae

- Given a set of examples, S
- For example, in a binary categorization
 - Where p_+ is the proportion of positives
 - And p_- is the proportion of negatives



$$Entropy(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_-$$

- For examples belonging to classes c_1 to c_n
 - Where p_n is the proportion of examples in c_n

$$Entropy(S) \equiv \sum_{i=1}^n -p_i \log_2 p_i$$

Entropy Example

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



Entropy Example

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

$$-(9/14) \log_2 \underline{(9/14)} - (5/14) \log_2 \underline{(5/14)}$$
$$= 0.940$$

Information Gain (IG)

- Information gain is based on the decrease in entropy after a dataset is split on an attribute.
- Which attribute creates the most homogeneous branches?
- ✓ First the entropy of the total dataset is calculated 0.94
- The dataset is then split on different attributes
- The entropy for each branch is calculated. Then it is added proportionally, to get total entropy for the split
- The resulting entropy is subtracted from the entropy before the split
- The result is the Information Gain, or decrease in entropy
- The attribute that yields the largest IG is chosen for the decision node

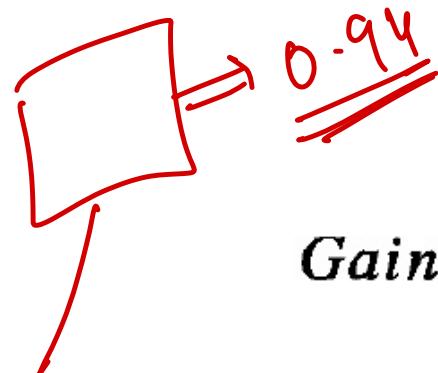
IG is the expected reduction in the entropy of target variable Y for data sample S due to sorting.

Information Gain (cont'd)

- A branch set with entropy of 0 is a leaf node.
- Otherwise, the branch needs further splitting to classify its dataset.
- The ID3 algorithm is run recursively on the non-leaf branches, until all the data is classified.

Information Gain (cont'd)

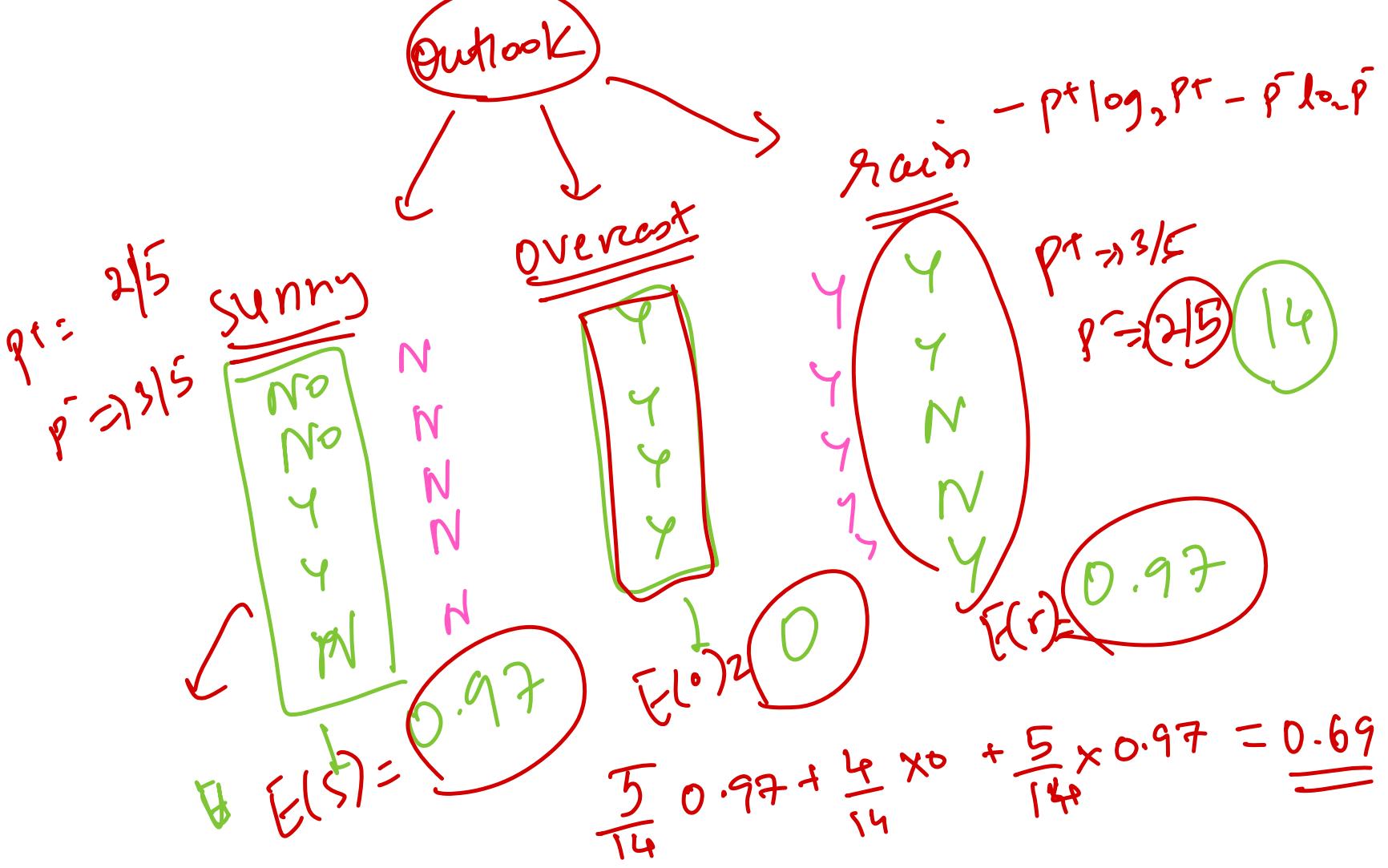
- Calculate $\text{Gain}(S, A)$
 - Estimate the reduction in entropy we obtain if we know the value of attribute A for the examples in S



$$\text{Gain}(S, A) \equiv \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

0.94

0.69

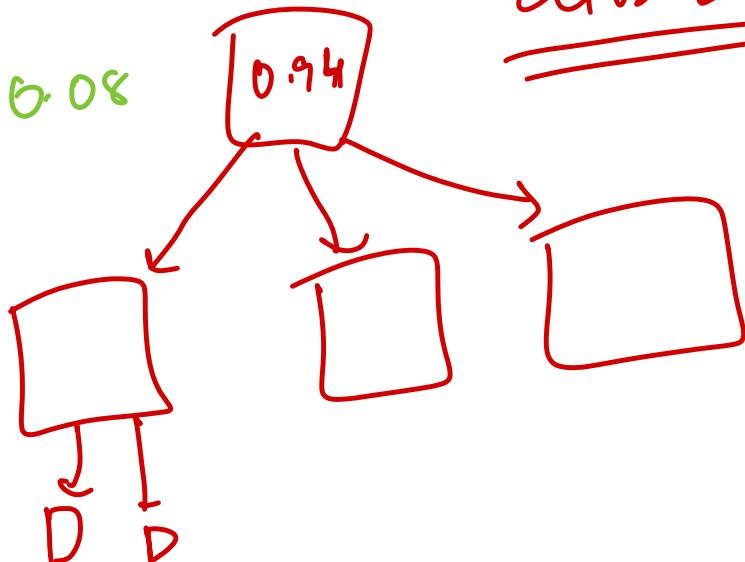


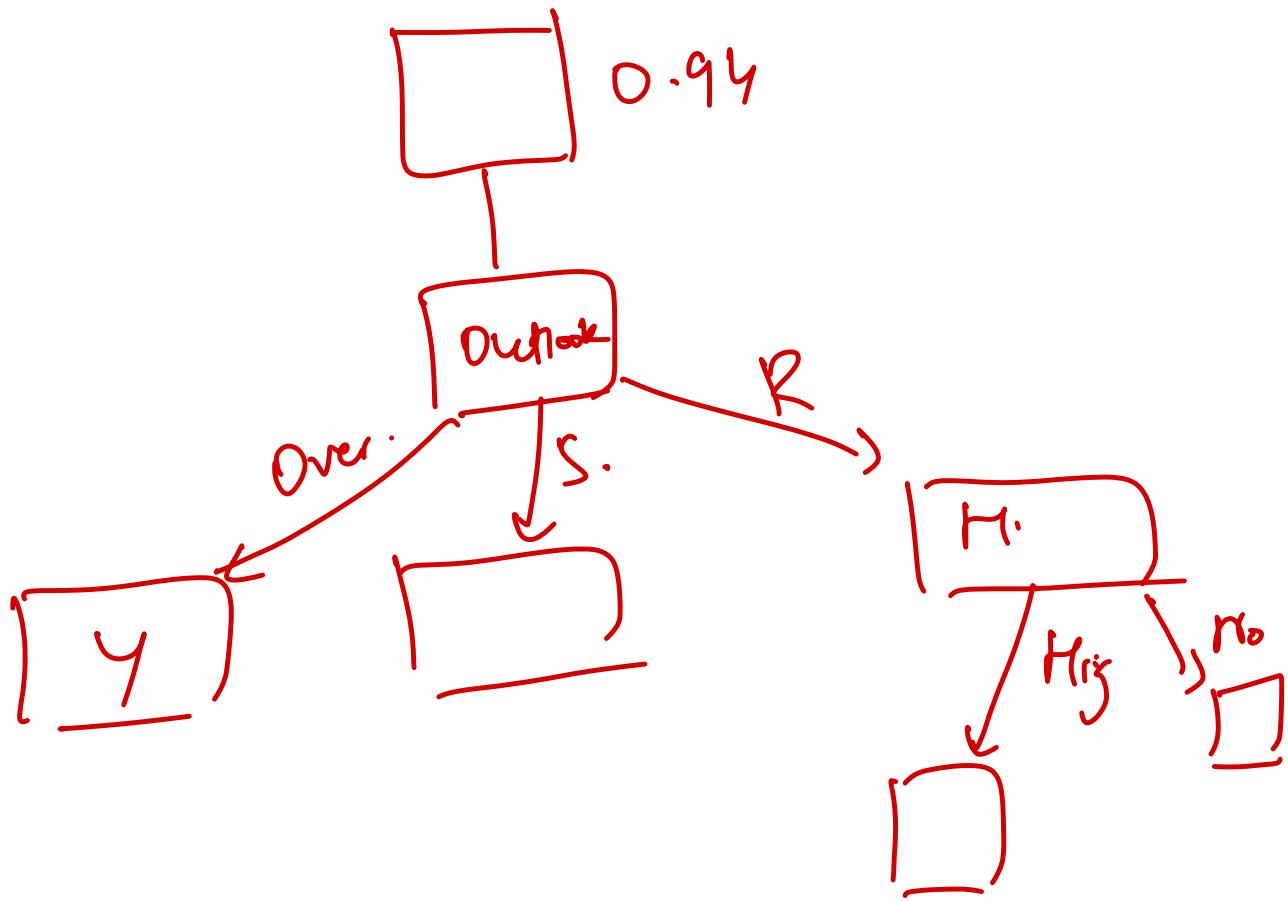
$$\underline{IG_0} = 0.94 - \underline{0.69} = \underline{0.25} \Rightarrow \underline{\text{outlook}}$$

$$\underline{IG_w} = 0.94 - \underline{0.93} = \underline{0.01} \quad () 0.2$$

Your task find IG for remaining attributes

$$\underline{IG_n} = () 0.08$$





An Example Calculation of Information Gain

- Suppose we have a set of examples
 - $S = \{s_1, s_2, s_3, s_4\}$
 - In a binary categorization
 - With one positive example and three negative examples
 - The positive example is s_1
- And Attribute A
 - Which takes values v_1, v_2, v_3
- s_1 takes value v_2 for A, s_2 takes value v_2 for A
 s_3 takes value v_3 for A, s_4 takes value v_1 for A

A		Final class
s_1	$v_2 \rightarrow$	+ve
s_2	v_2	-ve
s_3	v_3	-ve
s_4	v_1	- ve

First Calculate Entropy(S)

- Recall that

$$\text{Entropy}(S) = -p_+ \log_2(p_+) - p_- \log_2(p_-)$$

- From binary categorisation, we know that

$$p_+ = \frac{1}{4} \text{ and } p_- = \frac{3}{4}$$

- Hence, $\text{Entropy}(S) = -(1/4)\log_2(1/4) - (3/4)\log_2(3/4)$

$$= 0.811$$

Calculate Gain for each Value of A

- Remember that

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

- And that $S_v = \{\text{set of example with value } V \text{ for } A\}$

— So, $S_{v1} = \{s_4\}$, $S_{v2} = \{s_1, s_2\}$, $S_{v3} = \{s_3\}$

\hookrightarrow \hookrightarrow \hookrightarrow

$s_4 \approx N$
 $s_2 \approx N$
 $s_3 \approx N$
 $s_1 \approx Y$

- Now, $(|S_{v1}|/|S|) * Entropy(S_{v1})$

$$\begin{aligned} &= (1/4) * (-0/1) * \log_2(0/1) - (1/1) \log_2(1/1) \\ &= (1/4) * (0 - (1) \log_2(1)) = (1/4)(0-0) = 0 \end{aligned}$$

- Similarly, $(|S_{v2}|/|S|) = 0.5$ and $(|S_{v3}|/|S|) = 0$

$$E(S_{V2}) = -\frac{1}{2} \log_2(1/2) - \frac{1}{2} \log_2(1/2)$$

$$= \frac{1}{2} \log_2(2) + \frac{1}{2} \log_2(2)$$

$$= \frac{1}{2} + \frac{1}{2}$$

$$= 1$$

$$\cdot \frac{2}{4} \times 1 = \left(\frac{1}{2}\right)$$

Final Calculation

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

- So, we add up the three calculations and take them from the overall entropy of S:
- Final answer for information gain:
 - $Gain(S,A) = 0.811 - (0.25*0 + 1/2*1 + 0*0.25) = 0.311$

A Worked Example

Weekend	Weather	Parents	Money	Decision (Category)
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay in
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$FR \rightarrow 3C \\ 27 \\ \rightarrow 1S2 \\ 1Sh$$

$$-\frac{3}{7} \log_2 \frac{3}{7} - \frac{2}{7} \log_2 \left(\frac{2}{7}\right) - 2 \times \frac{1}{7} \log_2 \left(\frac{1}{7}\right)$$

$$0.7 \times () = 1.289$$

$$1.57 - 1.289 = 0.2816$$

A Worked Example

Weekend	Weather	Parents	Money	Decision (Category)
W1	Sunny	Yes	Rich	Cinema 1
W2	Sunny	No	Rich	Tennis 2
W3	Windy	Yes	Rich	Cinema 1
W4	Rainy	Yes	Poor	Cinema 1
W5	Rainy	No	Rich	Stay in 2
W6	Rainy	Yes	Poor	Cinema 1
W7	Windy	No	Poor	Cinema 1
W8	Windy	No	Rich	Shopping 4
W9	Windy	Yes	Rich	Cinema 1
W10	Sunny	No	Rich	Tennis 2

$$0.6 \Rightarrow C$$

$$T \Rightarrow 0.2$$

$$S2 \Rightarrow 0.1$$

$$Sh \Rightarrow 0.1$$

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$-0.6 \log_2 0.6 - 0.2 \log_2 0.2 - 2 \times (0.1 \log_2 0.1)$$

$$= \underline{\underline{1.571}}$$

Information Gain for All of S

- $S = \{W_1, W_2, \dots, W_{10}\}$
- Firstly, we need to calculate:
 - Entropy(S) = ... = 1.571
- Next, we need to calculate information gain
 - For all the attributes we currently have available
 - (which is all of them at the moment)
 - $\text{Gain}(S, \text{weather}) = 0.7$
 - $\text{Gain}(S, \text{parents}) = 0.61$
 - $\text{Gain}(S, \text{money}) = 0.2816$

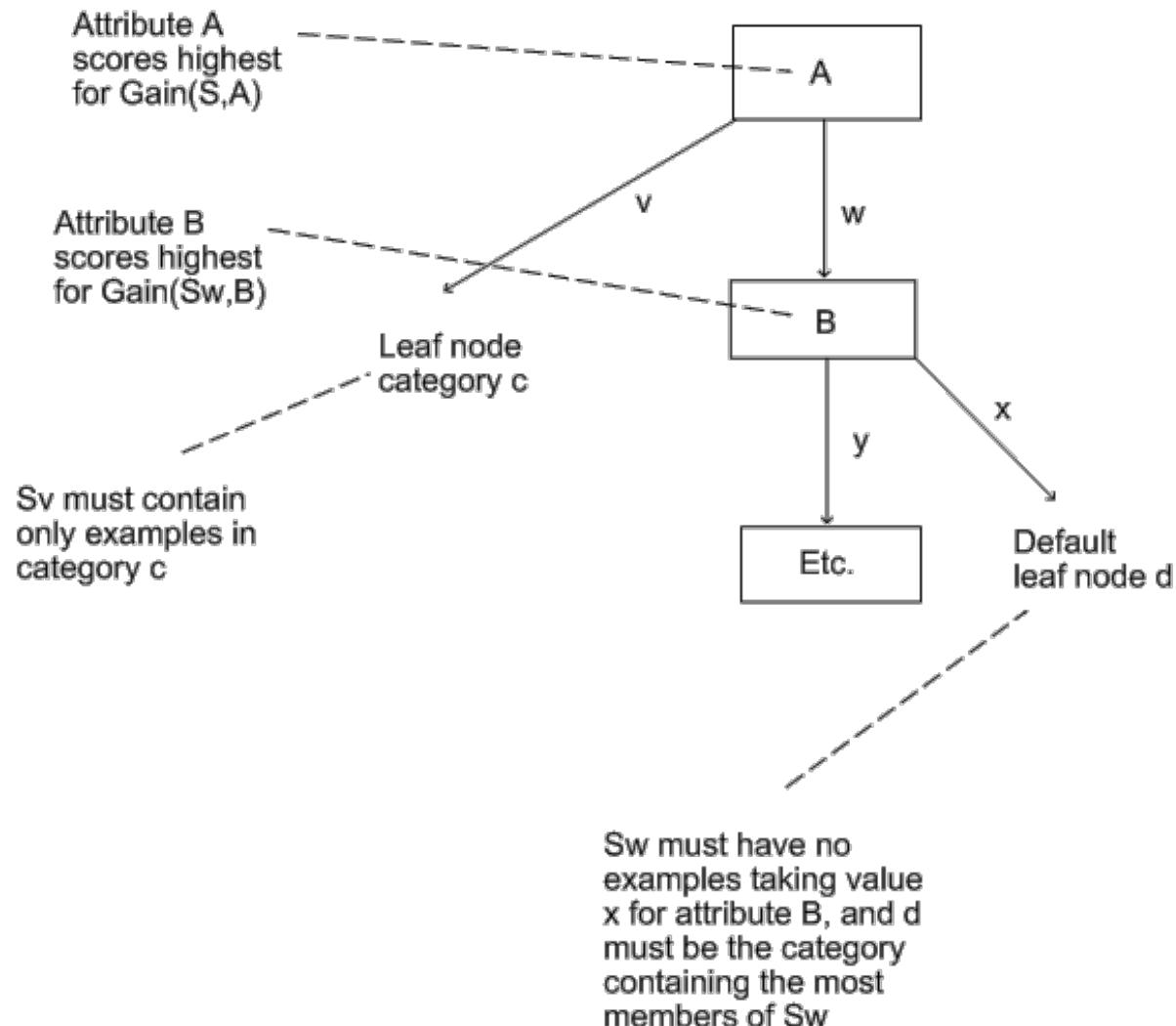
The ID3 Algorithm

- Given a set of examples, S
 - Described by a set of attributes A_i
 - Categorised into categories c_j
1. Choose the root node to be attribute A
 - Such that A scores highest for information gain
 - Relative to S , i.e., $\text{gain}(S,A)$ is the highest over all attributes
 2. For each value v that A can take
 - Draw a branch and label each with corresponding v

The ID3 Algorithm

- For each branch you've just drawn (for value v)
 - If S_v only contains examples in category c
 - Then put that category as a leaf node in the tree
 - If S_v is empty
 - Then find the default category (which contains the most examples from S)
 - Put this default category as a leaf node in the tree
 - Otherwise
 - Remove A from attributes which can be put into nodes
 - Replace S with S_v
 - Find new attribute A scoring best for $\text{Gain}(S, A)$
 - Start again at part 2
- Make sure you replace S with S_v

Explanatory Diagram



A Worked Example

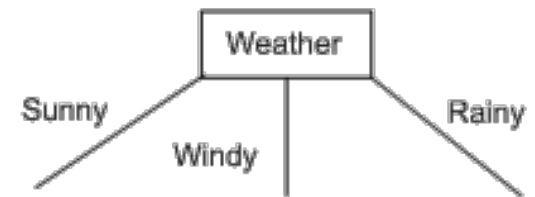
Weekend	Weather	Parents	Money	Decision (Category)
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay in
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

Information Gain for All of S

- $S = \{W_1, W_2, \dots, W_{10}\}$
- Firstly, we need to calculate:
 - $\text{Entropy}(S) = \dots = 1.571$
- Next, we need to calculate information gain
 - For all the attributes we currently have available
 - (which is all of them at the moment)
 - $\text{Gain}(S, \text{weather}) = \dots = 0.7$
 - $\text{Gain}(S, \text{parents}) = \dots = 0.61$
 - $\text{Gain}(S, \text{money}) = \dots = 0.2816$
- Hence, the weather is the first attribute to split on
 - Because this gives us the biggest information gain

Top of the Tree

- So, this is the top of our tree:
- Now, we look at each branch in turn
 - In particular, we look at the examples with the attribute prescribed by the branch
- $S_{\text{sunny}} = \{W1, W2, W10\}$
 - Categorisations are cinema, tennis and tennis for $W1, W2$ and $W10$
 - What does the algorithm say?
 - Set is neither empty, nor a single category
 - So we have to replace S by S_{sunny} and start again



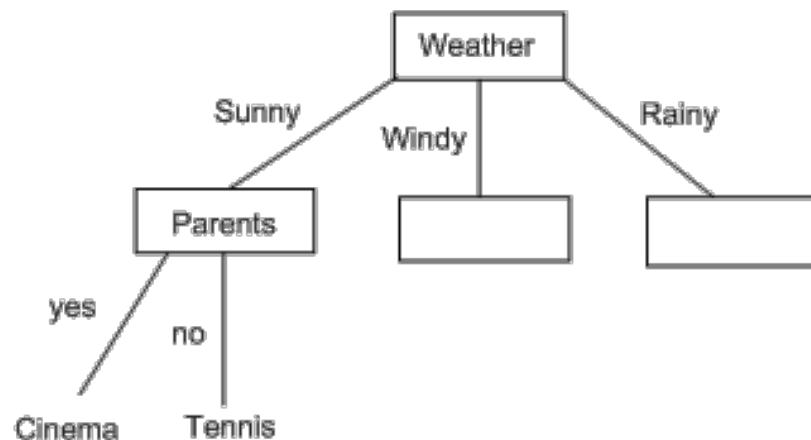
Working with S_{sunny}

Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W10	Sunny	No	Rich	Tennis

- Need to choose a new attribute to split on
 - Cannot be weather, of course – we've already had that
- So, calculate information gain again:
 - $\text{Gain}(S_{\text{sunny}}, \text{parents}) = \dots = 0.918$
 - $\text{Gain}(S_{\text{sunny}}, \text{money}) = \dots = 0$
- Hence we choose to split on parents

Getting to the leaf nodes

- If it's sunny and the parents have turned up
 - Then, looking at the table in previous slide
 - There's only one answer: go to cinema
- If it's sunny and the parents haven't turned up
 - Then, again, there's only one answer: play tennis
- Hence our decision tree looks like this:

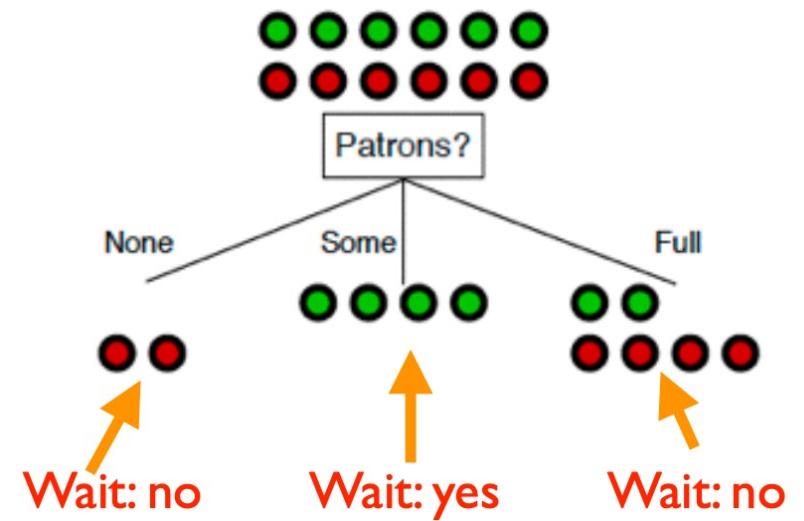


What is the optimal Tree Depth?

- We need to be careful to pick an appropriate tree depth.
- If the tree is too deep, we can overfit.
- If the tree is too shallow, we underfit
- Max depth is a hyper-parameter that should be tuned by the data. Alternative strategy is to create a very deep tree, and then to prune it.

Control the size of the tree

- If we stop early, not all training samples would be classified correctly.
- How do we classify a new instance:
 - We label the leaves of this smaller tree with the majority of training samples' labels



Summary of learning classification trees

- Advantages:
 - Easily interpretable by human (as long as the tree is not too big)
 - Computationally efficient
 - Handles both numerical and categorical data
 - It is parametric thus compact: unlike Nearest Neighborhood Classification, we do not have to carry our training instances around Building block for various ensemble methods (more on this later)
- Disadvantages
 - Heuristic training techniques
 - Finding partition of space that minimizes empirical error is NP-hard.
 - We resort to greedy approaches with limited theoretical underpinning.

Feature Space

- Suppose that we have p explanatory variables X_1, \dots, X_p and n observations.
 - a numeric variable: $n - 1$ possible splits
 - an ordered factor: $k - 1$ possible splits
 - an unordered factor: $\rightarrow 2(k-1) - 1$ possible splits.

Measures of Impurity

- At each node i of a classification tree, we have a probability distribution p_{ik} over k classes.

$$\hat{p}_{ik} = \frac{n_{ik}}{n_i}.$$

- Deviance: $D = \sum D_i$, where $D_i = -2 \sum_k n_{ik} \log(p_{ik})$.
- Entropy: $\sum p_{ik} \log(p_{ik})$.
- Gini index: $\sum_{j \neq k} p_{ij} p_{ik} = 1 - \sum_k p_{ik}^2$.
- Residual sum of squares

$$D = \sum_{\text{cases } j} (y_j - \mu_{[j]})^2$$

where $\mu_{[j]}$ is the mean of the values in the node that case j belongs to.