

# Lecture 13 - Feb 21

## Decision Tree Classifiers

### Reading

[\*Data Mining and Machine Learning\*](#)

19 Decision Tree Classifier

### Upcoming Deadlines

Project 1 Proposal (coming soon)

## Background

Tree-based methods partition feature space into a set of rectangles and fit a simple model in each

Major methods: CART + C4.5 CS

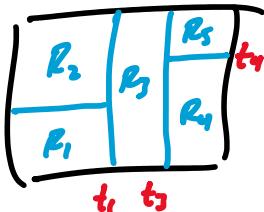
1D3

Ross Quinlan

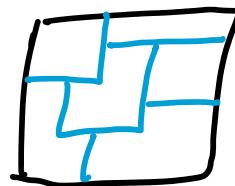
J48

We restrict to binary partitions  $X_2 \leq t_2$

This

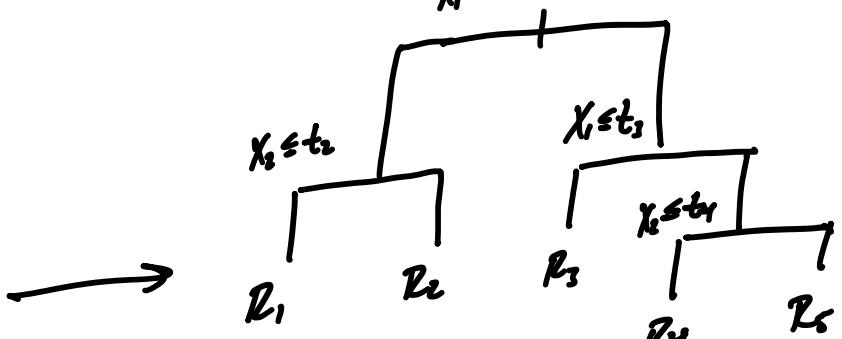


Not this



$X_1$

$X_1 \leq t_1$



Recursive binary tree  
is very interpretable  
and works best to  
represent the partitions  
in higher dimensions

(popular in medicine)

$$X \in \mathbb{R}^{n \times d} \quad x_i \in \mathbb{R}^d \quad y_i \in \{c_1, \dots, c_k\}$$

Decision tree predicts  $\hat{y}_i$  for each  $x_i$

$\mathcal{R}$  - data space where points lie

Decision trees split  $\mathcal{R}$  with axis-parallel hyperplanes until partitions are "pure" w.r.t. classes. The tree nodes represent splits + leaves are resulting regions

19.1

Axis-parallel hyperplanes  $h(x) : w^T x + b = 0$

↑ weight vector in  $\mathbb{R}^d$  normal to the hyperplane  
← offset of hyperplane to origin

$w$  restricted to  $\{e_1, \dots, e_d\}$

$$\Rightarrow h(x) : e_j^T x + b = 0 \\ x_j + b = 0$$

✓ 1.1 Points

## Split Points

A hyperplane specifies a decision or split point since it splits  $\mathcal{R}$  into 2 half-spaces:  $h(x) \leq 0$  vs  $h(x) > 0$

$$\begin{array}{ll} x_j + b \leq 0 & x_j + b > 0 \\ x_j \leq -b & x_j > -b \end{array}$$

Generic split:  $x_j = v^{\frac{-b}{\parallel}}$  → splits  $\mathcal{R} = \mathcal{R}_Y + \mathcal{R}_N$

## Data Partition

$$D_Y = \{x^T \mid x \in D, x_j \leq v\}$$

$$D_N = \{x^T \mid x \in D, x_j > v\}$$

## Purity

$$\text{purity}(D_j) = \max_i \left( \frac{n_{ji}}{n_j} \right)$$

↑  
fraction of  
points in  $D_j$   
with the  
label  $i$

where  $|D_j| = n_j$  and  $n_{ji}$  is  
the number of points in  $D_j$  in  
class  $C_i$

with the  
majority  
label

## Categorical Attributes

If  $X_j$  is categorical,  $V \subseteq \text{range}(X_j)$

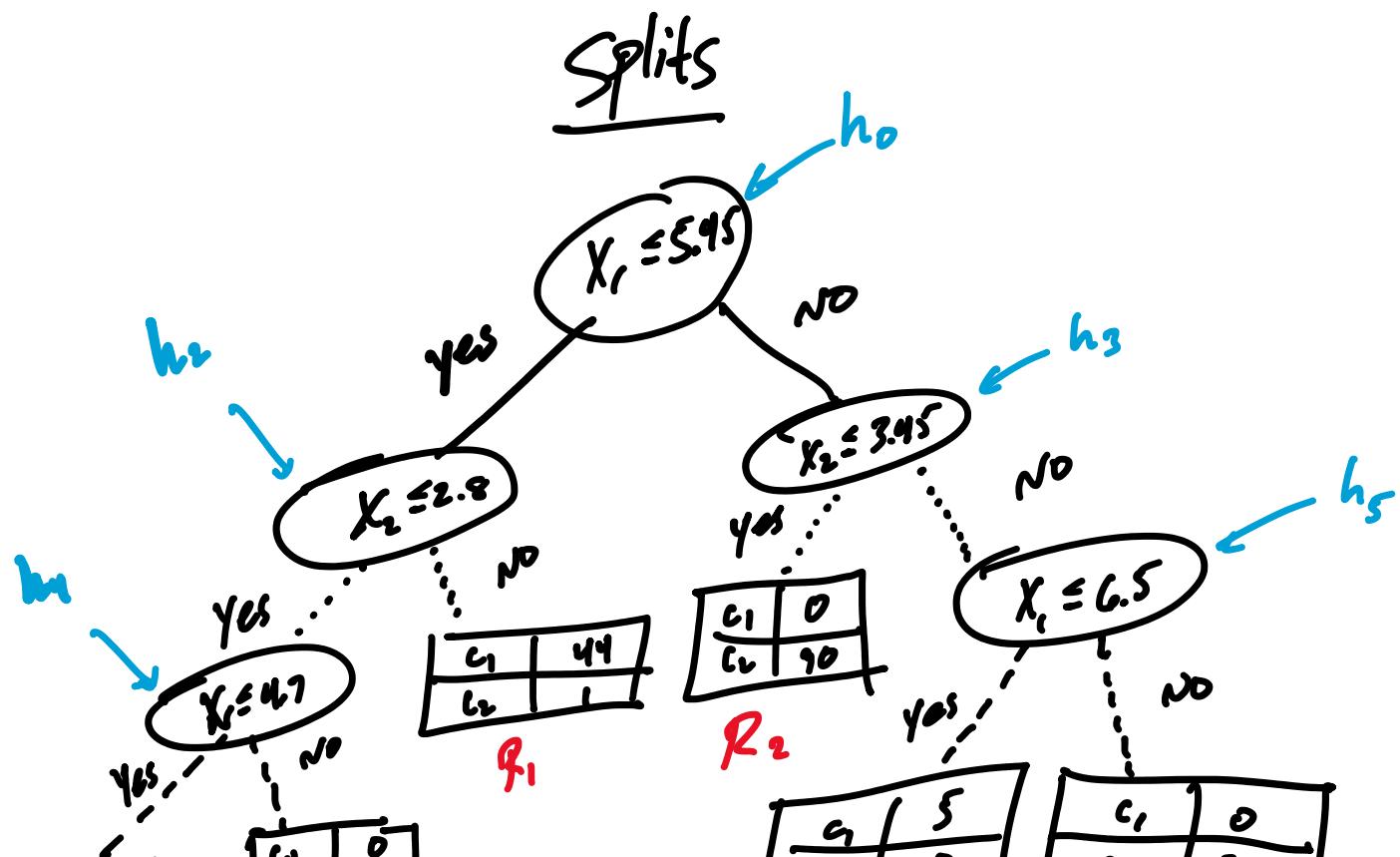
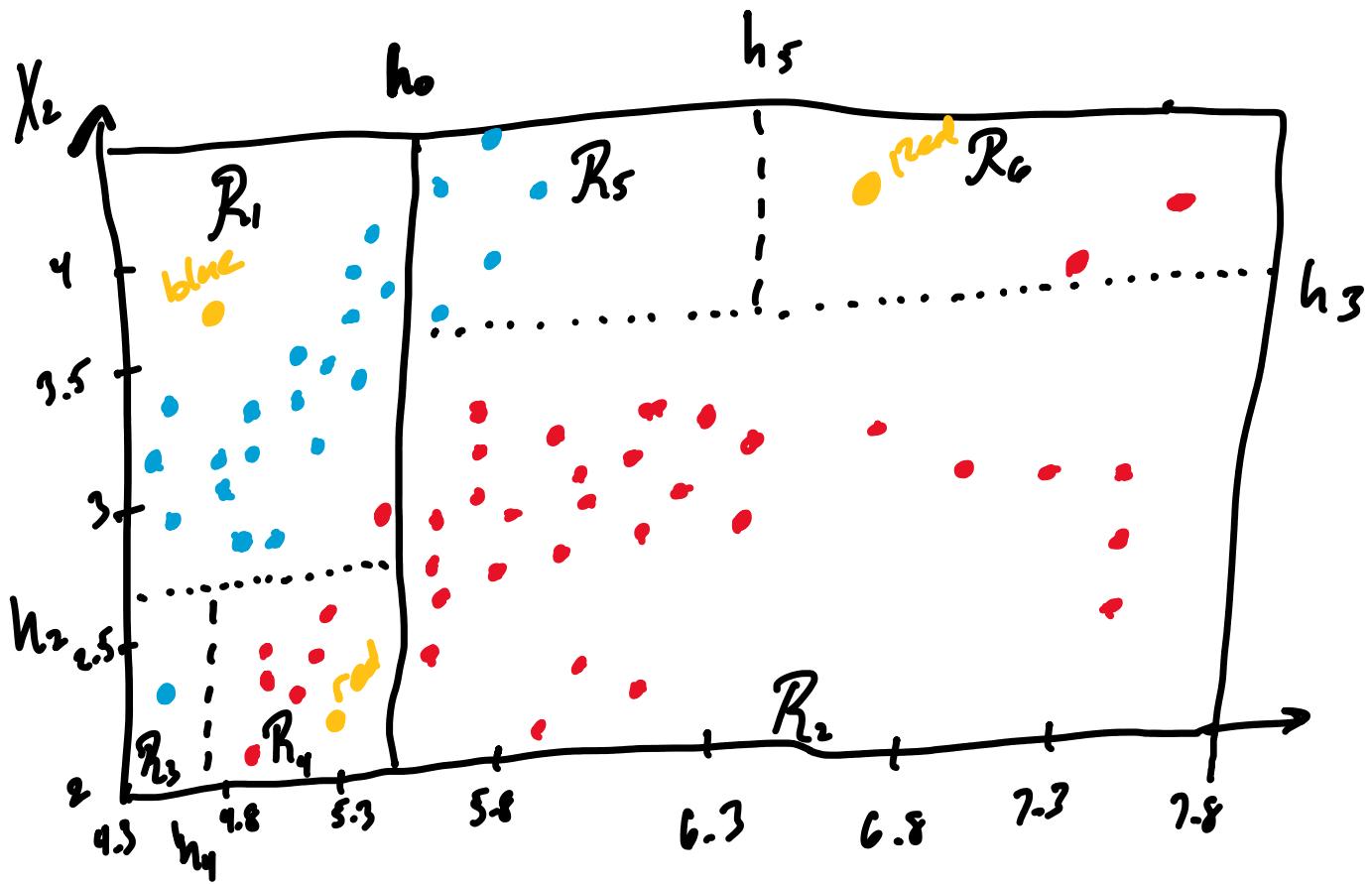
$R_Y \rightarrow$  points  $x$  with  $x_j \in V$

$R_N \rightarrow$  points  $x$  with  $x_j \in V$

## Decision Rules

DTs are easy to interpret  $\rightarrow$  splits = decision rules

# How It Works



$$\begin{array}{|c|c|} \hline c_1 & 1 \\ \hline c_2 & 0 \\ \hline \end{array}$$

465

$R_3$

$$\begin{array}{|c|c|} \hline i & \\ \hline c_1 & 0 \\ \hline c_2 & 0 \\ \hline \end{array}$$

$R_4$

71

$$\begin{array}{|c|c|} \hline c_1 & 5 \\ \hline c_2 & 0 \\ \hline \end{array}$$

$R_5$

$$\begin{array}{|c|c|} \hline c_1 & 0 \\ \hline c_2 & 3 \\ \hline \end{array}$$

$R_6$

## Decision Tree Algorithm Splits

Gini index is another measure of purity (CART)

$$G(D) = 1 - \sum_{i=1}^k P(c_i|D)^2$$

$P(c_j|D)=1$  for some  $c_j$

$P(c_i|D) = \frac{1}{k}$

high index = high disorder  $(0 \leq G(D) \leq \frac{k-1}{k})$

Weighted Gini index for split

$$G(D_p, D_n) = \frac{n_p}{n} G(D_p) + \frac{n_n}{n} G(D_n)$$

Low = less disorder  
= good

Classification And Regression Trees (CART) measure can also measure quality of splits

$$CART(D_p, D_n) = 2 \cdot \frac{n_p}{n} \cdot \frac{n_n}{n} \sum_{i=1}^k |P(c_i|D_p) - P(c_i|D_n)|$$

high = good

↳ try to maximize difference between class PMF for the partitions

Hyperparameters:  $\eta$  - leaf size  
 $\pi$  - leaf purity threshold

The "best split" points are chosen

### Split Point Evaluation

How do we objectively score split point candidates

Entropy measures the disorder in a system. A partition has lower entropy if it is purer

$$H(D) = - \sum_{i=1}^k P(c_i|D) \log_2 P(c_i|D) \quad (\text{entropy})$$

↑      ↑  
if near 1      prob. of  $c_i$   
in  $D$   
(pure), this  
is zero  
maximized  
at  $\log_2 k$   
if  $k$

Split entropy:  $H(D_Y, D_N) = \frac{n_Y}{n} H(D_Y) + \frac{n_N}{n} H(D_N)$

Information gain from a split:  $\leftarrow (ID3, C4.5, CS.0)$

$$Gain(D, D_r, D_n) = H(D) - H(D_r, D_n)$$

*high = good*

If this is high, entropy is reduced more

Gini index is another measure of purity  $\leftarrow (CART)$

$$G(D) = 1 - \sum_{i=1}^k P(c_i|D)^2$$

$P(c_j|D) = 1$  for some  $c_j$

$P(c_i|D) = \frac{1}{k}$

high index = high disorder  $\leftarrow \left( 0 \leq G(D) \leq \frac{k-1}{k} \right)$

Weighted Gini index for split

$$G(D_r, D_n) = \frac{n_r}{n} G(D_r) + \frac{n_n}{n} G(D_n)$$

*Low = less disorder  
= good*

Classification And Regression Trees (CART) measure can

also measure quality of splits

$\leftarrow$  "goodness"

$$CART(D_r, D_n) = 2 \cdot \frac{n_r}{n} \cdot \frac{n_n}{n} \sum_{i=1}^k |P(c_i|D_r) - P(c_i|D_n)|$$

*high = good*

$\hookrightarrow$  try to maximize difference between class PMF for the partitions

## Evaluating Split Points

need to compute PMFs  $\hat{P}(c_i | D)$  for each possible split points, which is expensive, so we can incrementally compute them instead

### Numerical Attributes

$\exists \text{ } m \text{ possible split points... one shortcut is to only take midpoints between successive } X_j \text{ coordinates of points}$

midpoints:  $\{v_1, \dots, v_m\}, m \leq n-1$

*Bayes & LTP*

$$\begin{aligned}\hat{P}(c_i | D_p) &= \hat{P}(c_i | X \leq v) \\ \hat{P}(c_i | D_n) &= \hat{P}(c_i | X > v) \\ &= \frac{\hat{P}(X \leq v | c_i) \hat{P}(c_i)}{\hat{P}(X \leq v)} = \frac{\hat{P}(X \leq v | c_i) \hat{P}(c_i)}{\sum_{j=1}^k \hat{P}(X \leq v | c_j) \hat{P}(c_j)}\end{aligned}$$

$\frac{n_i}{n}$

$$\hat{P}(X \leq v | c_i) = \frac{\hat{P}(X \leq v \text{ and } c_i)}{\hat{P}(c_i)} = \frac{\frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{x_j \leq v \text{ and } y_j = c_i\}}}{\frac{n_i}{n}} = \frac{N_{v,i}}{n_i}$$

$$\hat{P}(c_i | D_N) = P(c_i | X \leq v) = \frac{\frac{N_{v,i}}{n_i} \cdot \frac{n_i}{n}}{\sum_{j=1}^k \frac{N_{v,j}}{n_j} \cdot \frac{n_j}{n}} = \boxed{\frac{N_{v,i}}{\sum_{j=1}^k N_{v,j}}}$$

$$P(X > v | c_i) = 1 - \hat{P}(X \leq v | c_i) = 1 - \frac{N_{v,i}}{n_i} = \frac{n_i - N_{v,i}}{n_i}$$

$$\Rightarrow \hat{P}(c_i | D_N) = \frac{\frac{n_i - N_{v,i}}{n_i} \cdot \frac{n_i}{n}}{\sum_{j=1}^k \frac{n_j - N_{v,j}}{n_j} \cdot \frac{n_j}{n}} = \boxed{\frac{n_i - N_{v,i}}{\sum_{j=1}^k n_j - N_{v,j}}}$$

Pseudocode

DecisionTree( $X, y, n, \pi$ )

$$n = |X|$$

$$n_i = |X_i|$$

$$\text{purity}(X) = \max_i \frac{n_i}{n}$$

if  $n \leq n$  or  $\text{purity}(X) \geq \pi$

$$c^* = \arg \max_{c_i} \frac{n_i}{n}$$

create leaf node, label with class  $c^*$   
return

$$(\text{split point}^*, \text{score}) = (\emptyset, 0)$$

for  $i$  in range( $d$ )

if  $X_j$  is numeric:

$$(v, \text{score}) = \text{eval-numeric-attribute}(X, X_i)$$

find best  
cut in  
 $X_i$

$$\text{if } \text{score} > \text{score}^*: (\text{split point}^*, \text{score}^*) = (X_j \leq v, \text{score})$$

if  $X_j$  is categorical:

$$(V, \text{score}) = \text{eval-cat-attribute}(X, X_i)$$

$$\text{if } \text{score} > \text{score}^*: (\text{split point}^*, \text{score}^*) = (X_j \in V, \text{score})$$

$\cup$   
 if  $score > score^*$ :  $(split\ point, score) = (X_j \in v, score)$

# partition  $X$  into  $X_Y$  and  $X_N$  using  $split\ point^*$  and call recursively

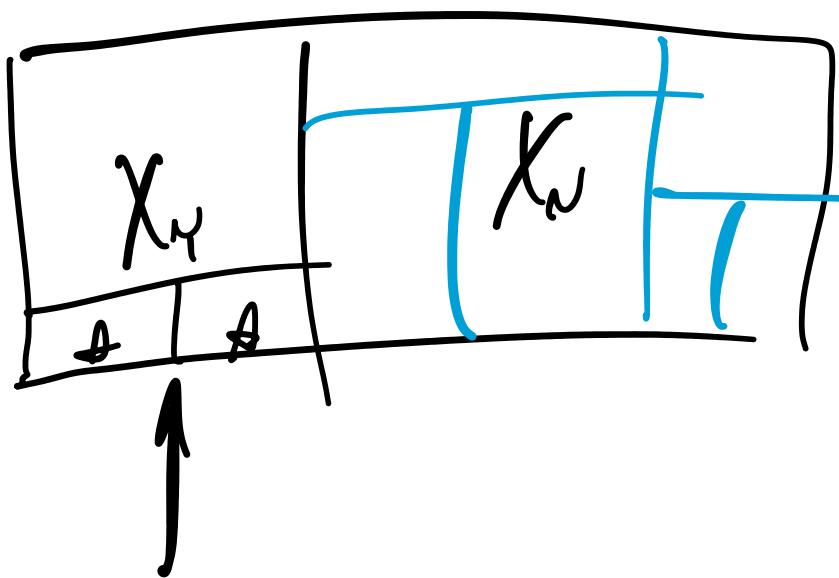
$$X_Y = \{x^T \mid x \in X \text{ satisfies } split\ point^*\}$$

$$X_N = \{x^T \mid x \in X \text{ does not } \dots\}$$

create internal node at  $split\ point^*$  with child nodes  $X_Y, X_N$

DecisionTree( $X_Y, Y_Y, n, \bar{n}$ )

DecisionTree( $X_N, Y_N, n, \bar{n}$ )



## Iris Example

```
import numpy as np
import matplotlib.pyplot as plt

from sklearn.datasets import load_iris
from sklearn.tree import DecisionTreeClassifier, plot_tree

# hyperparameters
n_classes = 3
plot_colors = "ryb"
plot_step = 0.02

# Load the iris dataset
iris = load_iris()
    # Load the iris dataset
iris = load_iris()

# run the model for each pair of features
for pairidx, pair in enumerate([[0, 1], [0, 2], [0, 3], [1, 2], [1, 3], [2, 3]]):
    # take the two corresponding features
    X = iris.data[:, pair]
    y = iris.target

    # fit the classifier
    clf = DecisionTreeClassifier().fit(X, y)

    # plot the decision boundary
    plt.subplot(2, 3, pairidx + 1)

    # predict the testing data on a tight mesh of points in space and color-code them
    x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
    y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
    xx, yy = np.meshgrid(np.arange(x_min, x_max, plot_step), np.arange(y_min, y_max, plot_step))

    plt.tight_layout(h_pad=0.5, w_pad=0.5, pad=2.5)
    Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
    Z = Z.reshape(xx.shape)
    cs = plt.contourf(xx, yy, Z, cmap=plt.cm.RdYlBu)

    plt.xlabel(iris.feature_names[pair[0]])
    plt.ylabel(iris.feature_names[pair[1]])

    # plot the training points
    for i, color in zip(range(n_classes), plot_colors):
        idx = np.where(y == i)
        plt.scatter(X[idx, 0], X[idx, 1], c=color, label=iris.target_names[i], cmap=plt.cm.RdYlBu, edgecolor='black', s=15)

    plt.suptitle("Decision surface of a decision tree using paired features")
    plt.legend(loc='lower right', borderpad=0, handletextpad=0)
    plt.axis("tight")

    plt.figure()
    clf = DecisionTreeClassifier().fit(iris.data, iris.target)
    plot_tree(clf, filled=True)
    plt.show()
```

Decision surface of a decision tree using paired features

