

# Recursive Partitioning

**Practical Machine Learning (with R)**

UC Berkeley

# LOGISTIC METHODS

## Advantages

⇒ ...

⇒ ...

## Disadvantages

⇒ ...

⇒ ...

⇒ ...



# LINEAR METHODS: LIMITATIONS

## Advantages

- Interpretable
- Easy to train

## Disadvantages

- Logistic regression: multiclass problems
- Highly sensitive to inputs
- Linear functions →  
inflexible: do not model real data well





Partition Goal:

**PARTITION INPUT SO THAT THE  
RESULTING SMALLER GROUPS ARE  
MORE HOMOGENEOUS THAN THE  
PARENT.**

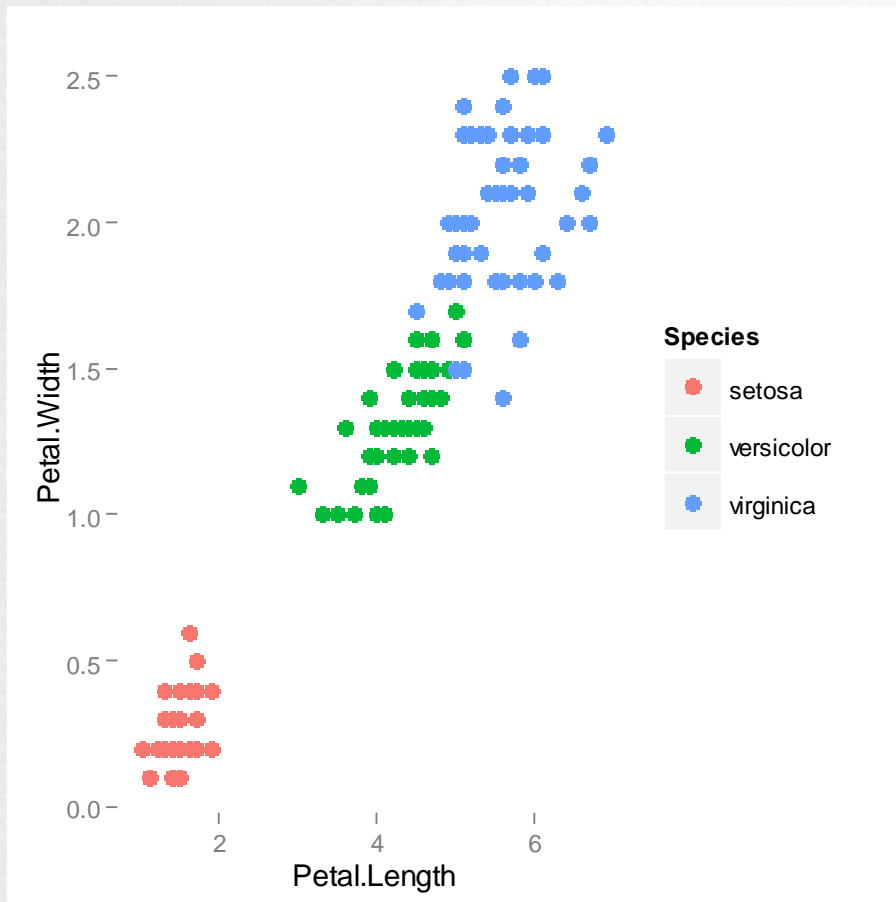


# PROCEDURE

- Find best **univariate plane (aka axis parallel)** to split the data into two subsets, such that the subsets are more alike than their parents
- 1. In each resulting subset ("node", "leaf") find the best univariate split, but only split the best one of these.
- 2. Repeat until stopping condition is met.



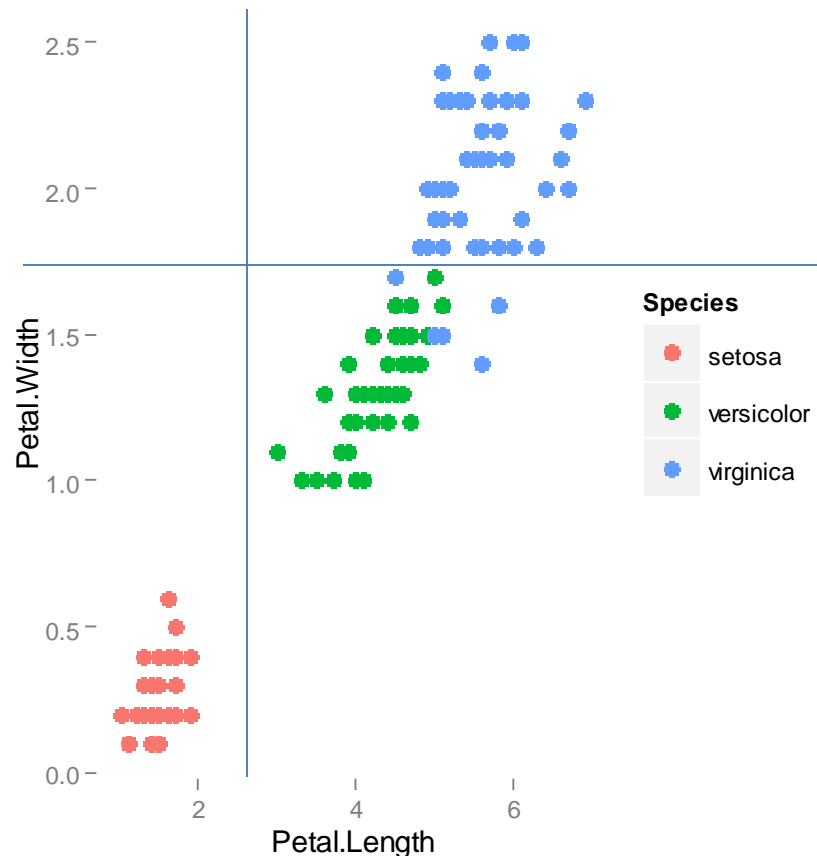
# A Simple Example





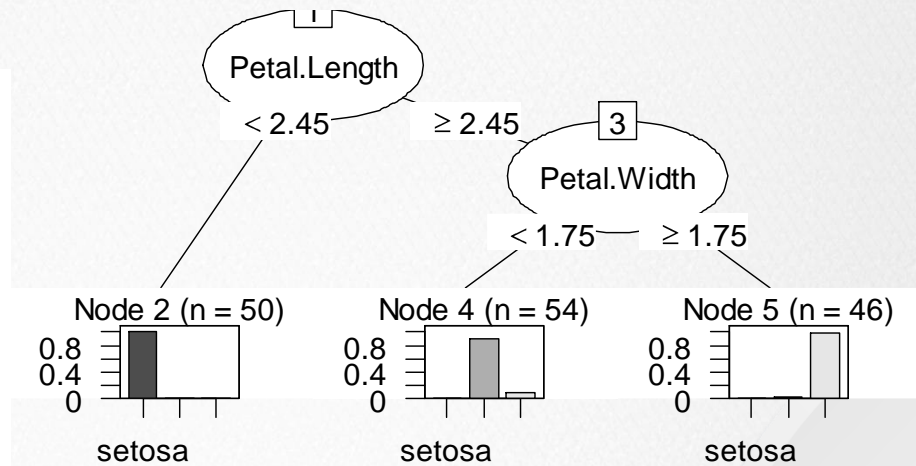
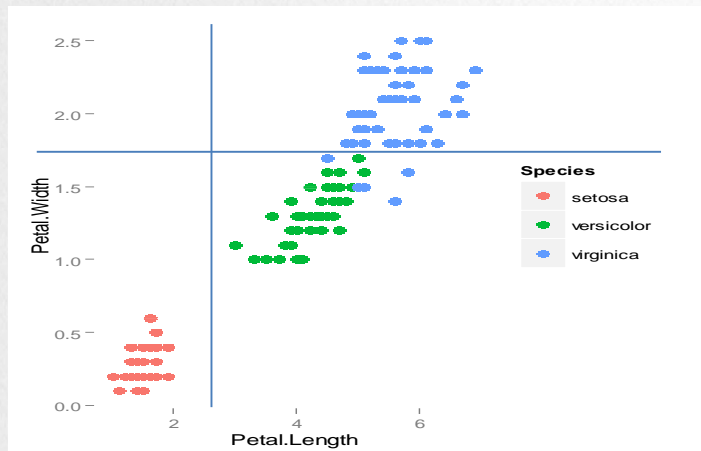
# A Simple Example

## Partitioning Requirements



# SOME NOTES

Splitting by planes is the same as a tree



Partitions define a rule\*


- Rules can be associated with outcomes → aggregation method

Trees always partition “all of of space”



# Splitting on Categorical Variable

- Select “metric”
- For each categorical variable
  - Find
$$\operatorname{argmin}_{S \in S}(\sum_{S_i} \operatorname{err}_i), i = 1..2$$
- Calculate:  $\sum_{S_i} \operatorname{err}_i$

- 
- Information Gain
  - *Gini* index (Class)



# TREATMENT OF CATEGORICAL VARIABLES

## ⇒ Grouped Categories

- Value treated as related

## ⇒ Independent Categories

- Values Treated as Independent



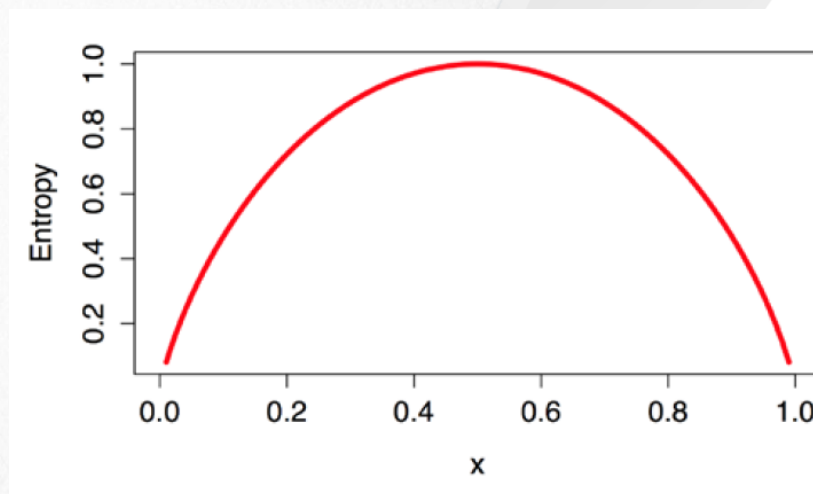
# Entropy (Two-Class Classification)

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

$$p_1 + p_2 = 1$$

*Entropy for node:*

$$-p_1 \log_2(p_1) - p_2 \log_2(p_2)$$





# Information Gain / Criteria

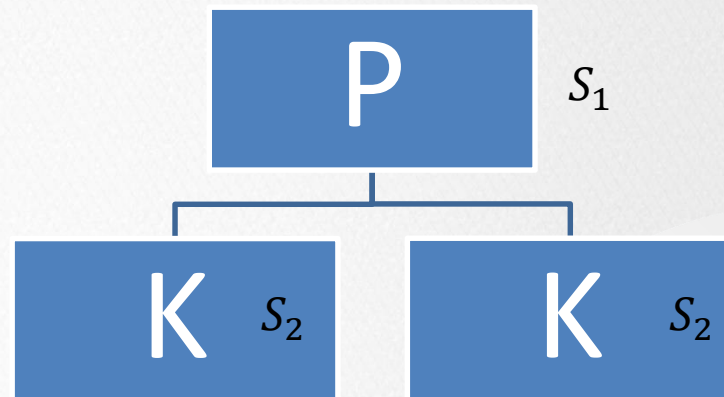
Change in Entropy of the System:

$$\text{InfoGain} = \text{Entropy}(S_1) - \text{Entropy}(S_2)$$

Change in system must consider all existing states. Splitting introduces an additional node.

$$\text{Entropy}(S) = \sum_{i=1}^n \omega_i \text{Entropy}(P_i)$$

$\omega_i := \text{node weight}$   
(proportion of observations in node)



# Gini Index (Two-Class Classification)

→ Measure node purity:

$$p_1(1 - p_1) + p_2(1 - p_2)$$

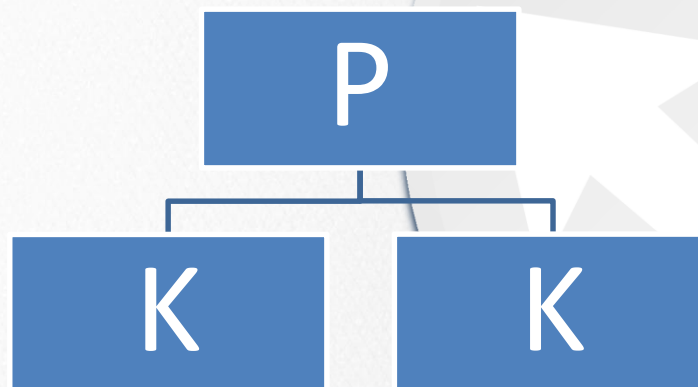
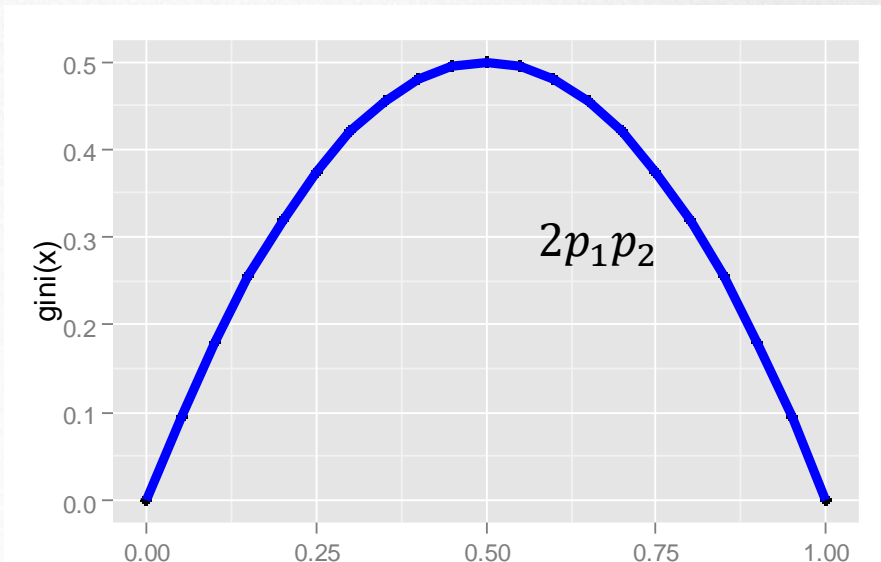
For two class:

$$p_1 + p_2 = 1$$

$$2p_1p_2$$

Minimize:

$$Gini(S) = \sum_{i=1}^n \omega_i 2p_1p_2$$



# SPLITTING ON CONTINUOUS VARIABLE

- ⇒ Determine Metric
- ⇒ Order data
  - If metric is a “cumulative” function calculate as cumulative function:

e.g.  $FPR = \text{cumsum}(FP) / \text{cumsum}(TN + FP)$

- Otherwise calculate at all possible split points or subset of split points

$$\text{argmin}_{x=n} \left( \sum_{i=1..2} \text{err}_i \right)$$

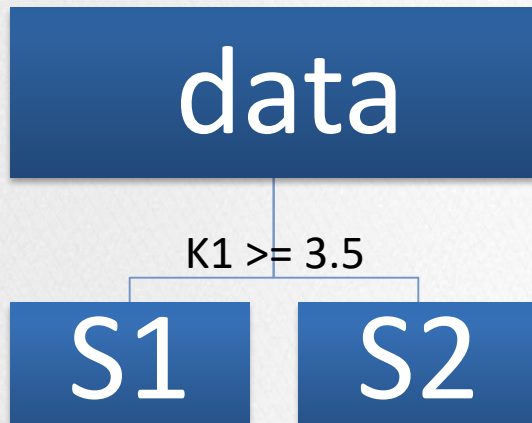




# data

Choose the split that  
minimizes the error  
 $\operatorname{argmin}_S(\text{Error})$

Ordinal							
Categorical			Continuous				
K1	K2	K3	V1	V2	V3	V4	V5
E <sub>K1</sub>	E <sub>K2</sub>	E <sub>K3</sub>	E <sub>V1</sub>	E <sub>V2</sub>	E <sub>V3</sub>	E <sub>V4</sub>	E <sub>V5</sub>



Choose the split that  
minimizes the error  
 $\text{argmin}_S(\text{Error})$

REPEAT WITH S1 AND S2

\* Very often predictor will be used again.

Ordinal							
Categorical			Continuous				
K1	K2	K3	V1	V2	V3	V4	V5
$E_{K1}$	$E_{K2}$	$E_{K3}$	$E_{V1}$	$E_{V2}$	$E_{V3}$	$E_{V4}$	$E_{V5}$

# MISSING DATA

- ➔ Missing values in predictors are common
- ➔ A split determines which observations go to the LHS and RHS. How to Handle `NA`s?
- ➔ `NA_Categorical`
  - Treat as separate category
- ➔ `NA` (in general)
  - Use **Surrogate Splits**





# SURROGATE SPLITS

- ⇒ Tree is built ignoring missing data
  - Any record with incomplete data (response or predictor) is rejected -or-
  - Missing data is rejected from determined the split
- ⇒ Variables are often collinear → splits are similar and send variables down the same path.
  - Choose a surrogate split that best approximates the chosen split (accuracy)
  - Very often this is also a good split.



# Tree Method Advantages I

- Highly interpretable
- Predict easy to implement (even in SQL)
- Handle many predictors (sparse, skewed, continuous, categorical) --> little need to pre-process them
- Non-parametric: do not require specification of predictor-response relationship



# Tree Method Advantages I

- Inherent method for handling missing data
- Trees insensitive to monotonic (order-preserving) transformation of inputs
  - $2^x$
  - No use in scaling and centering
- Intrinsic feature selection
- Computational simple and quick





# TREE DISADVANTAGES

- High Model Variance(sensitive to data)
  - Derives from each subsequent split is dependent on prior splits
- Less than optimal predictive performance
  - Rectangular regions!!!
- Limited number of outcome values
- Selection bias toward predictors with higher number of distinct values
- Tuning parameter,  $C_p$



# TREE VARIANTS

⇒ There are many tree variants

⇒ Tweaks

- change how splits are determined ? How many splits?
- when to stop growing the tree
- how the node value is determined



➤ Now Adapt to Regression





# APPENDIX



# Linear Models

