

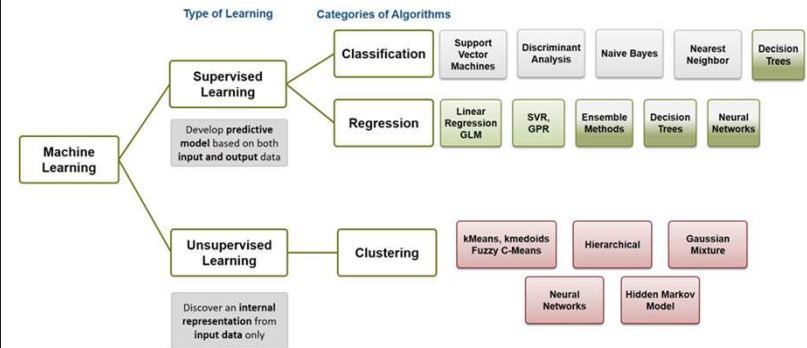
An Introduction to Decision Trees Classification Models with R

By Dario H. Romero, MS CS
Data Scientist / Production Optimization Engineer

git clone <https://github.com/darioromero/classifR.git>

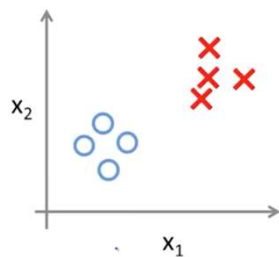
All Images may be subject to copyright

Machine Learning Types

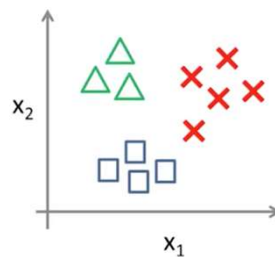


What is Classification?

Binary classification:



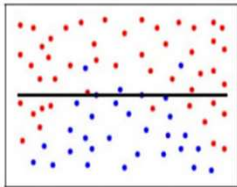
Multi-class classification:



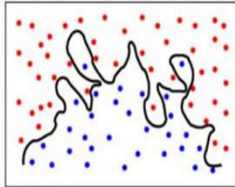
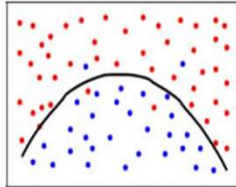
Matter of Interest

Generalization Problem

Underfitting



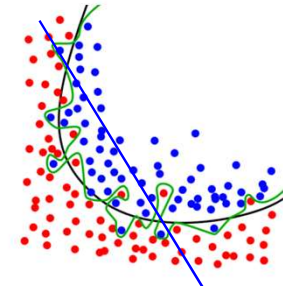
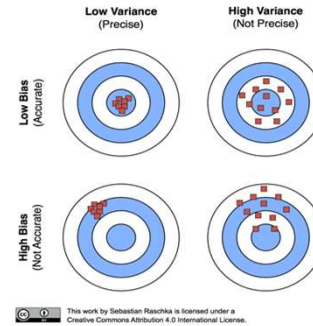
Overfitting



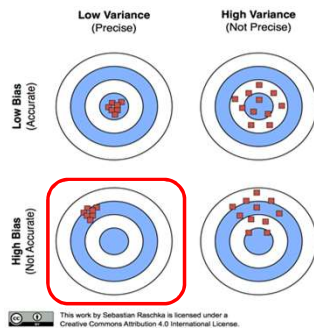
Bias & Variance

Bias occurs when an algo has *limited flexibility* to learn the true signal from a dataset.

Variance refers to an algo's *sensitivity* to specific sets of training data.

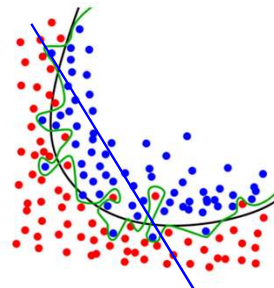


Bias & Variance



High bias, low variance algorithms train models that are consistent, but inaccurate on average.

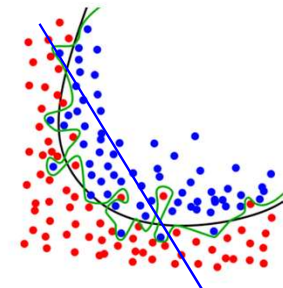
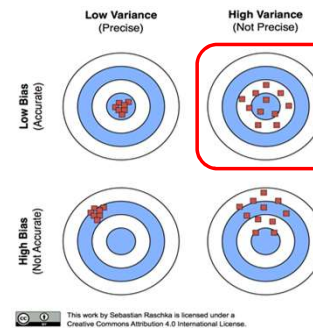
High variance, low bias algorithms train models that are accurate on average, but inconsistent.



Bias & Variance

High bias, low variance algorithms train models that are consistent, but inaccurate on average.

High variance, low bias algorithms train models that are accurate on average, but inconsistent.



Bias / Variance tradeoff

But why is there a tradeoff?

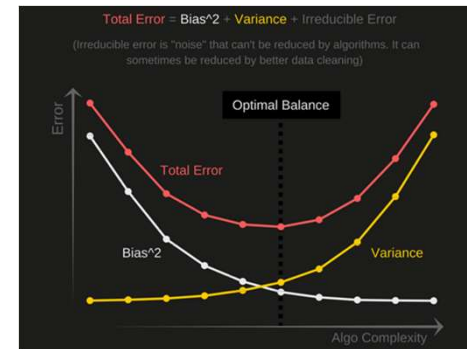
Low variance algos tend to be **less complex**, with simple or rigid underlying structure.

- e.g. Regression
- e.g. Naive Bayes
- *Linear algos*
- *Parametric algos*

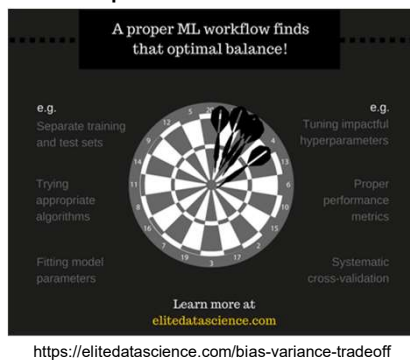
Low bias algos tend to be **more complex**, with flexible underlying structure.

- e.g. Decision trees
- e.g. Nearest neighbors
- *Non-linear algos*
- *Non-parametric algos*

Bias / Variance Optimal Balance

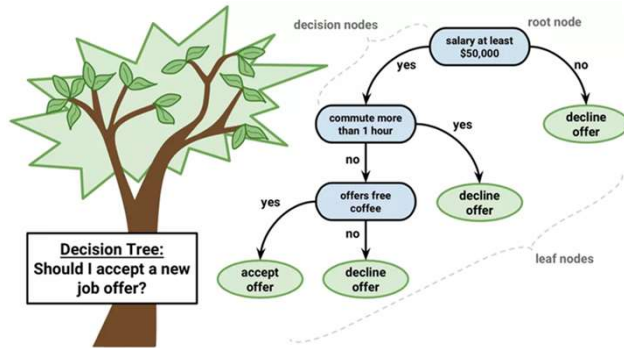


Bias / Variance Proper ML Workflow



Decision Trees for Classification

Decision Trees

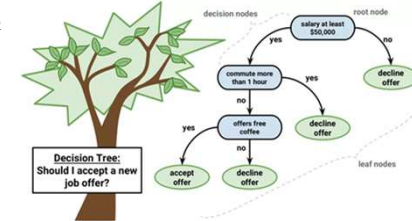


Decision Trees

The primary challenge in the decision tree implementation is to identify which attributes do we need to consider as the root node and each level. Handling this is known as the attributes selection.

- most popular attribute selection measures:
 - **Information Gain**
 - **Gini Index**

While using Information Gain as a criterion, we assume attributes to be categorical, and for Gini Index, attributes are assumed to be continuous.



Attribute Selection Methods

Entropy

The entropy $H(S)$ is a measure of the amount of uncertainty in the (data) set S .

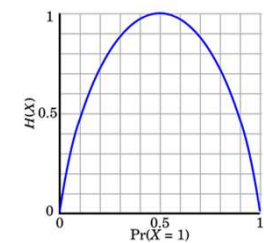
It is a measure of the disorder in a dataset.

Entropy for a partition $H(S, \vec{a})$

It is a measure of the disorder in a particular vector \vec{a} , within the dataset S .

Information gain

It is a measure of the decrease in disorder achieved by partitioning the original dataset.



Entropy characterizes the (data) set S .

$$H(S) = \sum_{c \in C} -p(c) \log_2 p(c)$$

↑ Knowledge == Entropy ↓

Entropy

Compute the entropy for column [Class] $H(S)$: $H(\hat{S}) = - \sum_{i=1}^k p_i \log_2(p_i)$
 $H(S) = - (4/6 * \log_2(4/6) + 2/6 * \log_2(2/6))$
 $H(S) = (0.389975 + 0.528321) = 0.9183$

Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High

Predictors Target Variable

Entropy for Car Type

Compute the entropy for column [Car Type] $H(S, \text{Car Type})$:

$$Car_{Family} = - (2/3 * \log_2(2/3) + 1/3 * \log_2(1/3)) = 0.9183$$

$$Car_{Sports} = - 2/2 * \log_2(2/2) = 0.0$$

$$Car_{Truck} = - 1/1 * \log_2(1/1) = 0.0$$

Average Entropy for Car Type:

$$I(\text{Car Type}) = 3/6 * Car_{Family} +$$

$$2/6 * Car_{Sports} +$$

$$1/6 * Car_{Truck}$$

$$I(\text{Car Type}) = 0.5 * 0.9183 +$$

$$2/6 * 0 +$$

$$1/6 * 0 = 0.4591$$

Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High

Predictors Target Variable

Car Type

Family

High
High
Low

Sports

High
High

Truck

Low

Entropy for Age

Compute the entropy for column [Age] $H(S, \text{Age})$:

$$Age_{>20} = - (2/4 * \log_2(2/4) + 2/4 * \log_2(2/4)) = 1.0$$

$$Age_{\leq 20} = - 2/2 * \log_2(2/2) = 0.0$$

Average Entropy for Age:

$$I(\text{Age}) = 4/6 * Age_{>20} +$$

$$2/6 * Age_{\leq 20}$$

$$I(\text{Age}) = 0.6667 * 1.0 +$$

$$2/6 * 0 = 0.6667$$

Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High

Predictors Target Variable

Age

> 20

High
High
Low
Low

<= 20

High
High

Entropy for Color

Compute the entropy for column [Color] $H(S, \text{Color})$:

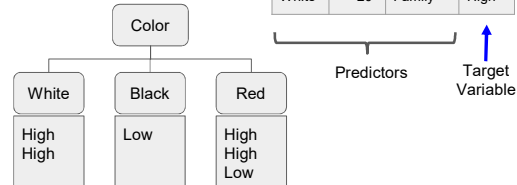
$$\text{Color}_{\text{White}} = -(2/2 * \log_2(2/2)) = 0.0$$

$$\text{Color}_{\text{Black}} = -(1/1 * \log_2(1/1)) = 0.0$$

$$\text{Color}_{\text{Red}} = -(2/3 * \log_2(2/3) + 1/3 * \log_2(1/3)) = 0.9183$$

Average Entropy for Age:

$$\begin{aligned} I(\text{Color}) &= 2/6 * \text{Color}_{\text{White}} + \\ &\quad 1/6 * \text{Color}_{\text{Black}} + \\ &\quad 3/6 * \text{Color}_{\text{Red}} \\ I(\text{Color}) &= 0.3333 * 0.0 + \\ &\quad 0.1667 * 0.0 + \\ &\quad 0.5 * 0.9183 = 0.4592 \end{aligned}$$



Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High

Information gain (IG)

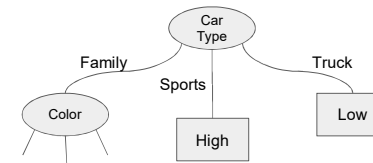
$$H(S) = 0.9183; I(\text{Car Type}) = 0.4591;$$

$$I(\text{Age}) = 0.6667; I(\text{Color}) = 0.4591$$

$$IG(\text{Car Type}) = H(S) - I(\text{Car Type}) = 0.9183 - 0.4591 = 0.4592$$

$$IG(\text{Age}) = H(S) - I(\text{Age}) = 0.9183 - 0.6667 = 0.2516$$

$$IG(\text{Color}) = H(S) - I(\text{Color}) = 0.9183 - 0.4591 = 0.4592$$



Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High

Gini Index or Gini Impurity

It is used for **impurity measure** instead of entropy.

Gini Index vs Gini Coefficient

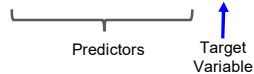
Despite their names they are not equivalent or even that similar.

Gini Index (impurity) is a measure of misclassification, which applies in a **multiclass classifier** context, while

Gini coefficient applies to **binary classification**.

Impurity is what is commonly used in decision trees.

Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High



Gini Index (Impurity)

Gini Index for **Binary Target variable** is:

$$\text{Gini Index} = 1 - \sum_{t=0}^{t=1} p_t^2$$

Similarly if Target Variable is categorical variable with multiple levels, the Gini Index will be still similar.

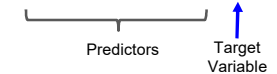
If Target variable takes **k** different values,

the Gini Index will be:

$$\text{Gini Index} = 1 - \sum_{t=0}^{t=k} p_t^2$$

Gini Index with Maximum value could be when all target values are equally distributed.

Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High



Gini Index: Attribute Selection Measure

- If a data set D contains examples from n classes, gini index, $gini(D)$ is defined as:

$$gini(D) = 1 - \sum_{j=1}^n p_j^2$$

where p_j is the relative frequency of class j in D

- If a data set D is split on A into two subsets D_1 and D_2 , the gini index $gini(D)$ is defined as

$$gini_A(D) = \frac{|D_1|}{|D|} gini(D_1) + \frac{|D_2|}{|D|} gini(D_2)$$

- Reduction in Impurity: $\Delta gini(A) = gini(D) - gini_A(D)$

Original by: Yingfan Liu, liuyf@se.cuhk.edu.hk

Gini Index (worked out example)

- Compute the Gini index for the overall collection of training examples on the target Variable (Class).

$$gini(D) = 1 - \sum_{j=1}^n p_j^2 = 1 - [(4/6)^2 + (2/6)^2] = 0.4444$$

- Compute the Gini index for the predictor Variable Car Type:

$$\text{Family} = 1 - [(2/3)^2 + (1/3)^2] = 0.4444$$

$$\text{Sports} = 1 - [(2/2)^2] = 0.0$$

$$\text{Truck} = 1 - [(1/1)^2] = 0.0$$

$$\text{Weighted Average} = (3/6) * 0.4444 + (2/6) * 0.0 + (1/6) * 0.0 = 0.2222$$

- Compute the Delta Gini index for the predictor Variable Car Type:

$$\text{Delta } gini(\text{Class}) = gini(D) - gini(\text{Car Type} | D) = 0.4444 - 0.2222 = 0.2222$$

Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High

Predictors Target Variable

Select attribute with lower Gini Index for splitting

<https://view.officeapps.live.com/op/view.aspx?src=http://dni-institute.in/blogs/wp-content/uploads/2014/11/Gini-Index-Calculation-Binary-Target-Variable.xlsx>

Gini Index (worked out example)

- Compute the Gini index for the overall collection of training examples on the target Variable (Class).

$$gini(D) = 1 - \sum_{j=1}^n p_j^2 = 1 - [(4/6)^2 + (2/6)^2] = 0.4444$$

- Compute the Gini index for the predictor Variable Age:

$$> 20: = 1 - [(2/4)^2 + (2/4)^2] = 0.5$$

$$<= 20: = 1 - [(2/2)^2] = 0.0$$

$$\text{Weighted Average} = (4/6) * 0.5 + (2/6) * 0.0 = 0.3333$$

- Compute the Delta Gini index for the predictor Variable Age:

$$\text{Delta } gini(\text{Class}) = gini(D) - gini(\text{Age} | D) = 0.4444 - 0.3333 = 0.1111$$

Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High

Predictors Target Variable

Select attribute with lower Gini Index for splitting

<https://view.officeapps.live.com/op/view.aspx?src=http://dni-institute.in/blogs/wp-content/uploads/2014/11/Gini-Index-Calculation-Binary-Target-Variable.xlsx>

Gini Index (worked out example)

- Compute the Gini index for the overall collection of training examples on the target Variable (Class).

$$gini(D) = 1 - \sum_{j=1}^n p_j^2 = 1 - [(4/6)^2 + (2/6)^2] = 0.4444$$

- Compute the Gini index for the predictor Variable Color:

$$\text{Red} = 1 - [(1/3)^2 + (2/3)^2] = 0.4444$$

$$\text{White} = 1 - [(2/2)^2] = 0.0$$

$$\text{Black} = 1 - [(1/1)^2] = 0.0$$

$$\text{Weighted Average} = (3/6) * 0.4444 + (2/6) * 0.0 + (1/6) * 0.0 = 0.2222$$

- Compute the Delta Gini index for the predictor Variable Color:

$$\text{Delta } gini(\text{Class}) = gini(D) - gini(\text{Color} | D) = 0.4444 - 0.2222 = 0.2222$$

Color	Age	Car Type	Class
White	> 20	Family	High
Red	<= 20	Sports	High
Red	> 20	Sports	High
Black	> 20	Family	Low
Red	> 20	Truck	Low
White	<= 20	Family	High

Predictors Target Variable

Select attribute with lower Gini Index for splitting

<https://view.officeapps.live.com/op/view.aspx?src=http://dni-institute.in/blogs/wp-content/uploads/2014/11/Gini-Index-Calculation-Binary-Target-Variable.xlsx>

Errors and Model Accuracy

To be Or Not To be in Classification

Feature 1	Feature 2	... Feature n	Actual	Predicted	
...	1	1	TP
...	1	0	FN
...	0	1	FP
...	0	0	TN

Confusion Matrix: “the stem table” and Metrics

		Condition (as determined by “Gold standard”)		
		Condition Positive	Condition Negative	
Test Outcome	Test Outcome Positive	True Positive	False Positive (Type I error)	Positive predictive value = $\frac{\sum \text{True Positive}}{\sum \text{Test Outcome Positive}}$
	Test Outcome Negative	False Negative (Type II error)	True Negative	Negative predictive value = $\frac{\sum \text{True Negative}}{\sum \text{Test Outcome Negative}}$

$$\text{Sensitivity} = \frac{\sum \text{True Positive}}{\sum \text{Condition Positive}}$$

$$\text{Specificity} = \frac{\sum \text{True Negative}}{\sum \text{Condition Negative}}$$

$$\text{Accuracy} = \frac{tp + tn}{tp + tn + fp + fn}$$

$$\text{Precision} = \frac{tp}{tp + fp}$$

$$\text{Recall} = \frac{tp}{tp + fn}$$

$$F = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{True Negative Rate} = \frac{tn}{tn + fp}$$

Measures	Formulas
Accuracy	$\frac{\sum TP + \sum TN}{\sum \text{Total population}}$
Positive predicted value	$\frac{\sum TP}{\sum \text{Prediction positive}}$
Negative predicted value	$\frac{\sum TN}{\sum \text{Prediction negative}}$
False Omission rate	$\frac{\sum FN}{\sum \text{Prediction negative}}$
False discovery rate	$\frac{\sum FP}{\sum \text{Prediction positive}}$
Prevalence	$\frac{\sum \text{Condition positive}}{\sum \text{Total positive}}$
True positive rate	$\frac{\sum TP}{\sum \text{condition positive}}$
False positive rate	$\frac{\sum FP}{\sum \text{condition negative}}$
Positive likelihood ratio	$\frac{TPR}{FPR}$
Negative likelihood ratio	$\frac{FNR}{TNR}$

Gains and Losses

numerical form			table of costs			
predicted→ real ₄	Class_pos	Class_neg	predicted→ real ₄	Class_pos	Class_neg	
Class_pos	114	86	Class_pos	0	-3	=
Class_neg	7	93	Class_neg	-10	0	

We are interested only in potential losses

gains and losses form

predicted→ real ₄	Class_pos	Class_neg
Class_pos	0	-258
Class_neg	-70	0

$$\sum_{i,j=1}^2 n_{ij} \cdot w_{ij} = -328$$

End of Theory



Binary Classification using Decision Trees with R and CARET

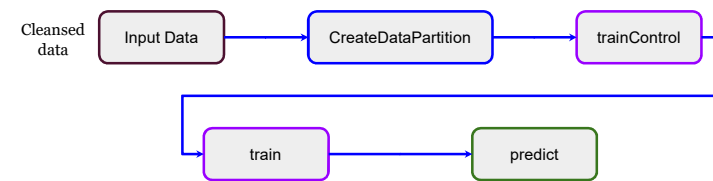
The CARET package in R

The caret package, short for classification and regression training, contains numerous tools for developing predictive models using the rich set of models available in R.

The package focuses on simplifying model training and tuning across a wide variety of modeling techniques.

It also includes methods for pre-processing training data, calculating variable importance, and model visualizations.

CARET General (Basic) Workflow



D:\Users\drome\gitrepos\classifR\jupyter lab {classif_R_decisiontree_gini_ig.ipynb}