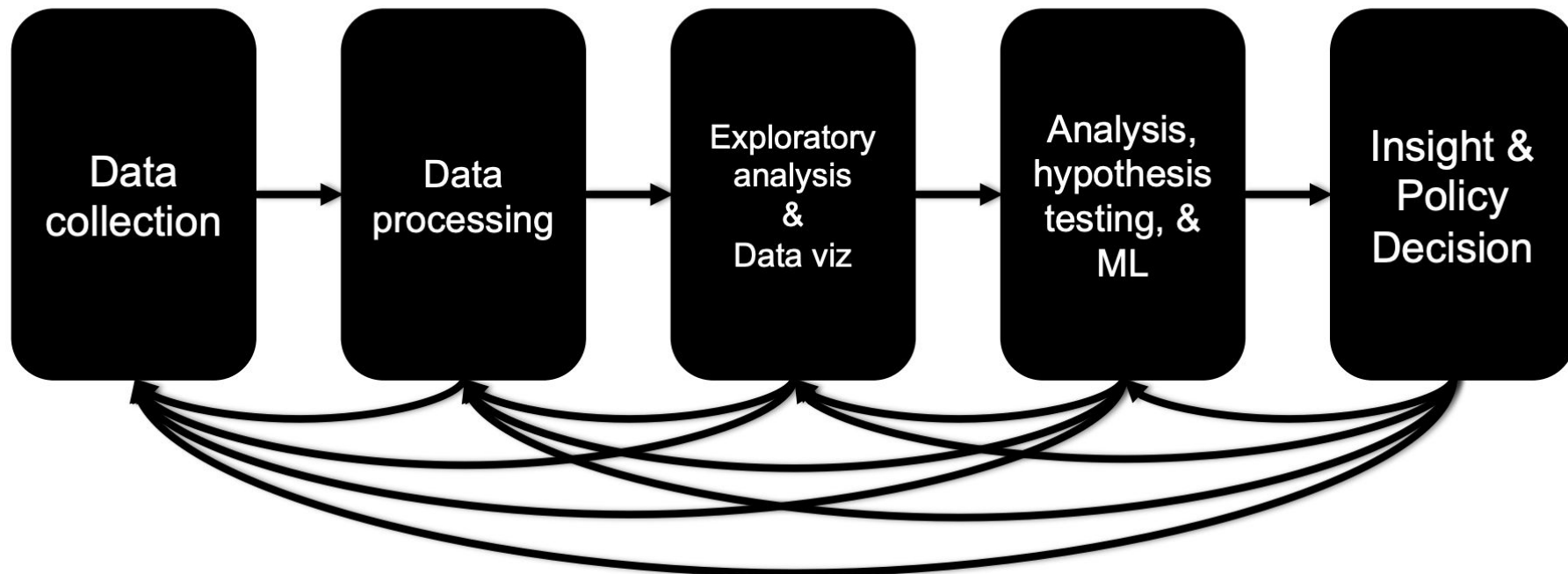


Cont ML Algorithms

CSC 380 - Principles of Data Science

Lecture 7.2

So far in the course



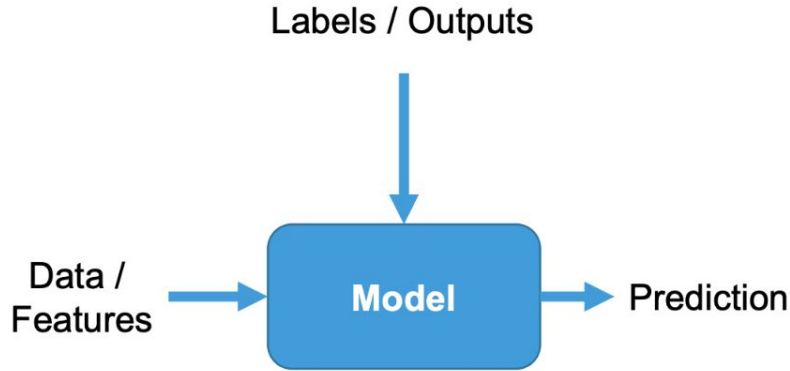
In our last lecture

- **Data Preprocessing:**
 - Cleaning
 - Integration
 - Transformation
 - Reduction
 - Discretization
 - Normalization
- **Machine Learning:**
 - **Supervised Learning :**
 - Linear Regression
 - Bayesian Classification
 - **Unsupervised Learning :**
 - Clustering - K-means

In this lecture:

- Revise/new terms
- Evaluating models
- Decision Tree Classifier
- Knn (not to be confused with Kmeans from last lecture)
- Logistic regression

Supervised Learning



Unsupervised Learning



ML models distinguished by a number of factors

- Number of parameters needed (parametric / nonparametric)
- Whether they model uncertainty (probabilistic / nonprobabilistic)
- Do they model the data generation process? (generative / discriminative)

Revise : Linear v/s Non-Linear

- Linear models generate output as a linear combination of inputs.
 - PCA, linear regression, etc.
- Nonlinear models fit an arbitrary nonlinear function to map inputs to outputs:
 - Neural networks, support vector machine, nonlinear dimensionality reduction

Parametric vs Nonparametric

- Parametric : fixed number of parameters
- Nonparametric : either has an infinite number of parameters, or parameters grow with the amount of data

Parametric Methods	Non-Parametric Methods
Parametric Methods uses a fixed number of parameters to build the model.	Non-Parametric Methods use the flexible number of parameters to build the model.
Parametric analysis is to test group means.	A non-parametric analysis is to test medians.
It is applicable only for variables.	It is applicable for both – Variable and Attribute.
It always considers strong assumptions about data.	It generally fewer assumptions about data.
Parametric Methods require lesser data than Non-Parametric Methods.	Non-Parametric Methods requires much more data than Parametric Methods.
Parametric methods assumed to be a normal distribution.	There is no assumed distribution in non-parametric methods.

Parametric Methods	Non-Parametric Methods
Parametric data handles – Intervals data or ratio data.	But non-parametric methods handle original data.
Here when we use parametric methods then the result or outputs generated can be easily affected by outliers.	When we use non-parametric methods then the result or outputs generated cannot be seriously affected by outliers.
Parametric Methods can perform well in many situations but its performance is at peak (top) when the spread of each group is different.	Similarly, Non-Parametric Methods can perform well in many situations but its performance is at peak (top) when the spread of each group is the same.
Parametric methods have more statistical power than Non-Parametric methods.	Non-parametric methods have less statistical power than Parametric methods.
As far as the computation is considered these methods are computationally faster than the Non-Parametric methods.	As far as the computation is considered these methods are computationally slower than the Parametric methods.
Examples: Logistic Regression, Naïve Bayes Model, etc.	Examples: KNN, Decision Tree Model, etc.

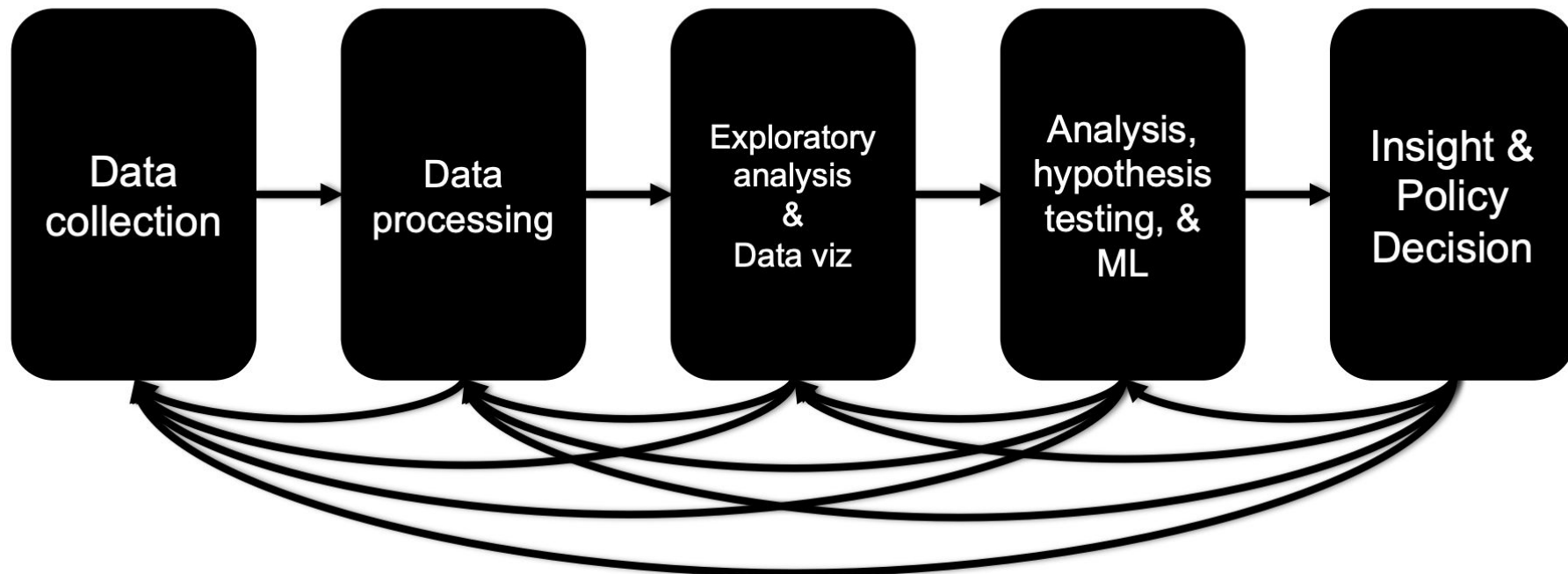
Probabilistic vs Non-Probabilistic

- non-probabilistic generates deterministic outputs / predictions from data ex: k-means
- A probabilistic model represents predictions as random variables, with a distribution that is fit to training data

Generative vs Discriminative

- **Discriminative models**
 - learn the (hard or soft) boundary between classes
 - Lesser assumptions
 - simply providing classification splits (and not necessarily in a probabilistic manner)
 - Assume some functional form for $P(Y|X)$
 - Estimate parameters of $P(Y|X)$ directly from training data
- **Generative models**
 - model the distribution of individual classes
 - providing a model of how the data is actually generated
 - make some kind of structure assumptions on your model
 - Assume some functional form for $P(Y)$, $P(X|Y)$
 - Estimate parameters of $P(X|Y)$, $P(Y)$ directly from training data
 - Use Bayes rule to calculate $P(Y|X)$.

So far in the course



Evaluation

1. Accuracy
2. Confusion matrix
3. Precision
4. Recall
5. F1 score
6. Precision-Recall or PR curve
7. ROC (Receiver Operating Characteristics) curve
8. PR vs ROC curve.

Accuracy

Pros

- Easy to interpret

Cons

- Paints an incorrect picture when classes are imbalanced
- there are different costs associated with the different mistakes.
- Threshold of accuracy has different costs for small differences ex: (0.51 and 0.99 is the same) versus (0.49 and 0.51 is not)

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

Confusion Matrix

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Confusion Matrix

		Actual Label			Total Predicted
		A	B	C	
Predicted Label	A	856 28.98%	58 1.96%	130 4.4%	1044 35.34%
	B	0	765 25.90%	136 4.6%	901 30.5%
	C	69 2.34%	33 1.12%	907 30.7%	1009 34.16%
Total Actual		925 31.31%	856 28.98%	1173 39.71%	2954 100%

Recall (or 1 - Sensitivity)

What proportion of actual positives was identified correctly?

- A measure of quality
- Higher recall means that an algorithm returns more of the correct results

$$\text{Recall} = \frac{\text{True Positive (TP)}}{\text{True Positive (TP)} + \text{False Negative (FN)}}$$

Precision (or Specificity)

What proportion of positive identifications was actually correct?

- A measure of quantity
- Higher precision means that an algorithm returns more correct results than incorrect ones

$$\text{Precision} = \frac{\text{True Positive (TP)}}{\text{True Positive (TP)} + \text{False Positive (FP)}}$$

F1

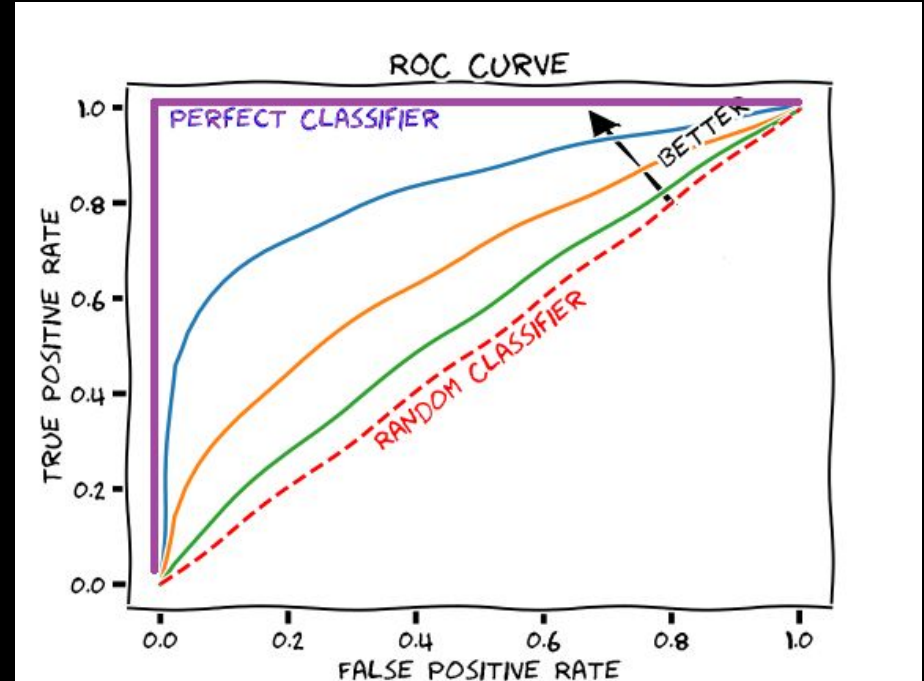
The harmonic mean encourages similar values for precision and recall.

That is, the more the precision and recall scores deviate from each other, the worse the harmonic mean.

$$F_1 = 2 * \frac{\textit{precision} * \textit{recall}}{\textit{precision} + \textit{recall}}$$

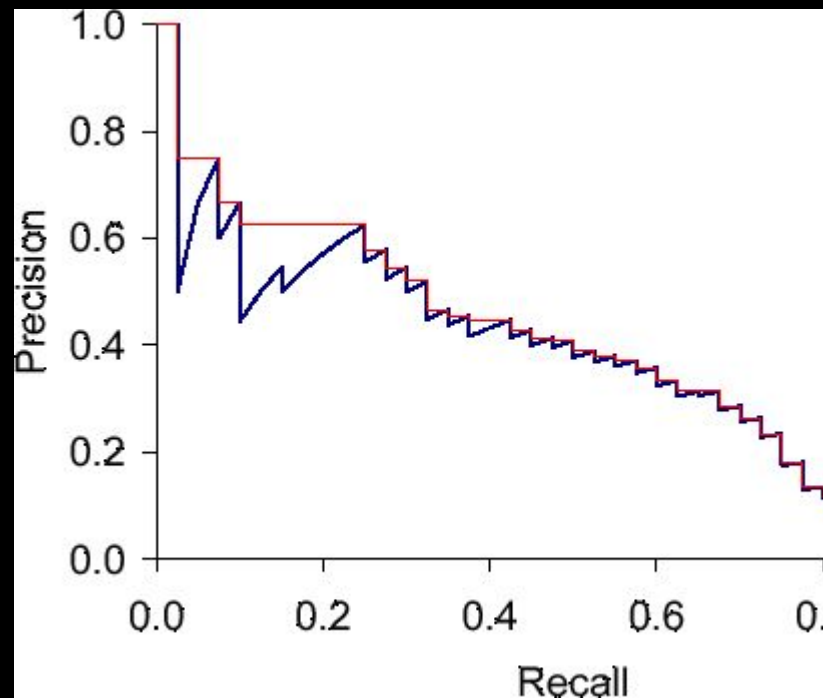
AUCROC

- Area under receiver operating characteristic curve or c-statistic or “concordance statistic.”
- ROC curves should be used when there are roughly equal numbers of observations for each class.

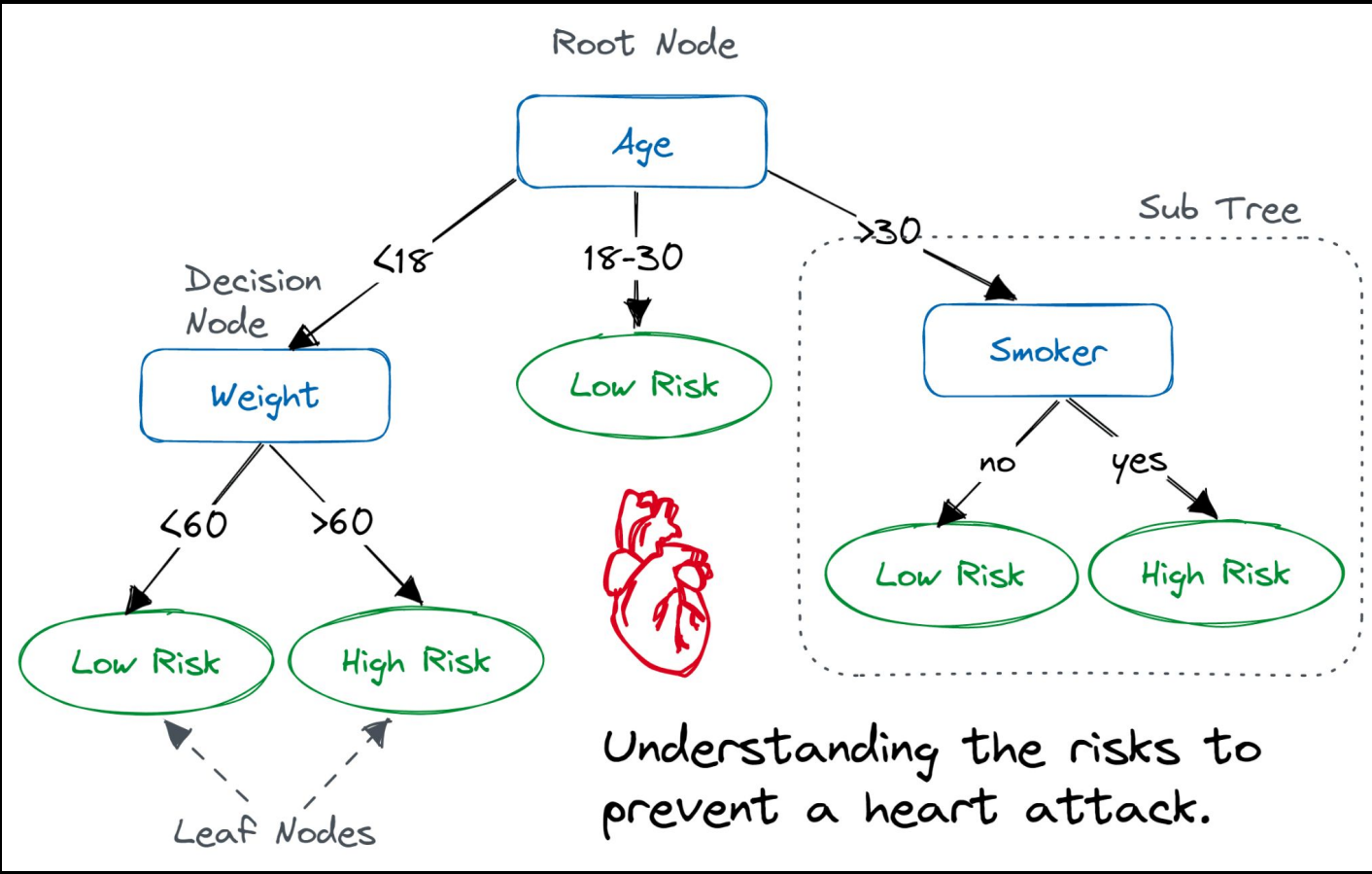


Precision-Recall Curve

Precision-Recall curves should be used when there is a moderate to large class imbalance.



Decision Tree Classifier

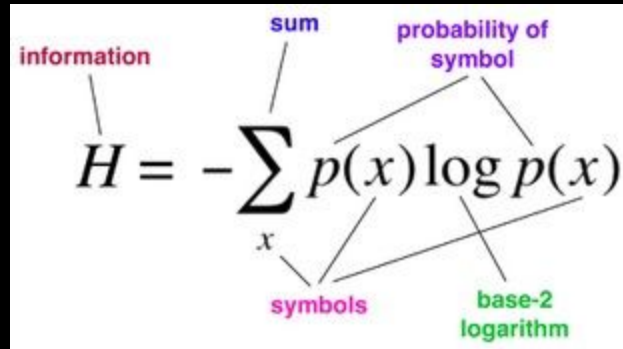


Understanding the risks to prevent a heart attack.

How do we pick which feature becomes the Node?

How do we pick which feature becomes the Node?

Entropy : Measure of information contained in a random variable



The diagram shows the entropy formula $H = -\sum_x p(x) \log p(x)$ with several annotations. A red line connects the word 'information' to the variable H . A purple line connects the word 'sum' to the summation symbol \sum . A purple line connects the words 'probability of symbol' to the variable $p(x)$. A purple line connects the word 'symbols' to the variable x . A green line connects the words 'base-2 logarithm' to the \log function. A purple line connects the word 'symbols' to the variable $p(x)$ in the logarithm's argument.

$$H = -\sum_x p(x) \log p(x)$$

Lower Probability , Higher Information

Higher Probability, Lower Information

How do we pick which feature becomes the Node?

Information gain calculates the reduction in entropy parent node to child node.

$$\text{Information Gain} = E_{\text{parent}} - \text{Avg}E_{\text{child}}$$

	Entropy Node	Average Entropy	Information Gain
Parent	0.9968		
working	0.9183	0.8110	0.1858
Not_work	0.6500		
Bkgd_Ma	0.9852	0.4598	0.5370
Bkgd_CS	0.0000		
Bkgd_oth	0.0000		
online_co	0.9544	0.9688	0.0280
online_no	0.9852		

How do we pick which feature becomes the Node?

Gain Index :

The probability for a random instance being misclassified when chosen randomly. The lower the Gini Index, the better the lower the likelihood of misclassification.

CART (Classification and Regression Tree) uses the Gini method to create split points.

$$Gini = 1 - \sum_{i=1}^j P(i)^2$$

Please [cite us](#) if you use the software.

1.10. Decision Trees

1.10.1. Classification

1.10.2. Regression

1.10.3. Multi-output problems

1.10.4. Complexity

1.10.5. Tips on practical use

1.10.6. Tree algorithms: ID3, C4.5, C5.0 and CART

1.10.7. Mathematical formulation

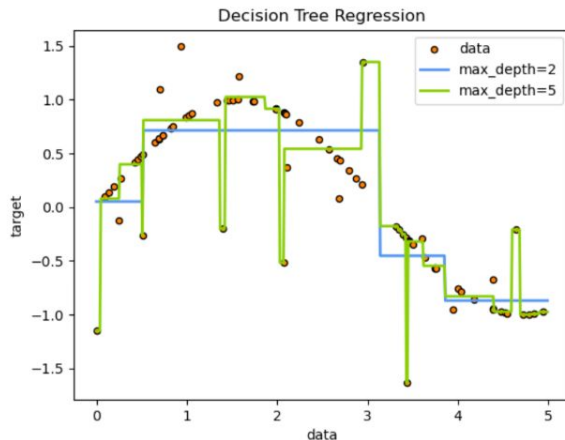
1.10.8. Missing Values Support

1.10.9. Minimal Cost-Complexity Pruning

1.10. Decision Trees

Decision Trees (DTs) are a non-parametric supervised learning method used for [classification](#) and [regression](#). The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features. A tree can be seen as a piecewise constant approximation.

For instance, in the example below, decision trees learn from data to approximate a sine curve with a set of if-then-else decision rules. The deeper the tree, the more complex the decision rules and the fitter the model.



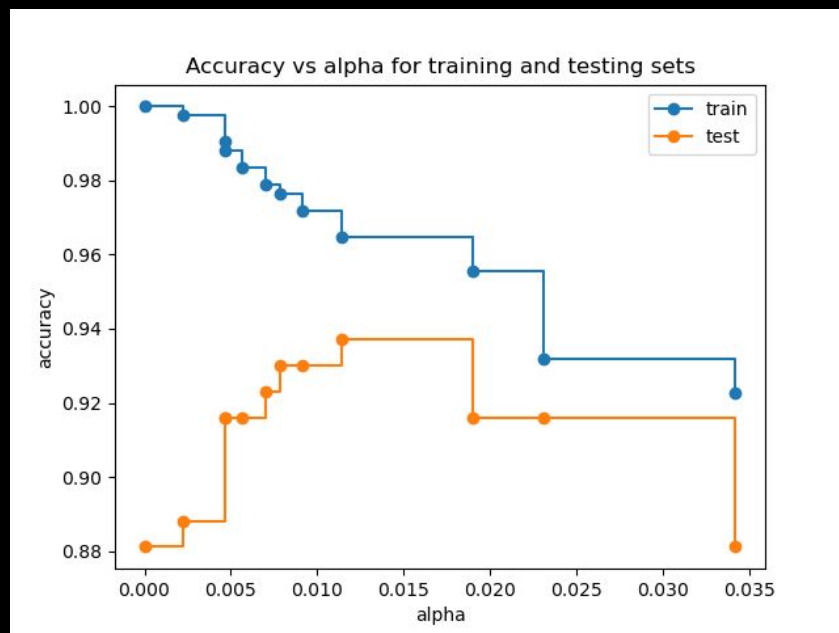
Preventing Overfitting

- Pre pruning or Early stopping: Preventing the tree from growing too big or deep
- Post Pruning: Allowing a tree to grow to its full depth and then getting rid of various branches based on various criteria
- Ensembling or using averages of multiple models such as Random Forest

Minimal Cost-Complexity Pruning

ccp_alpha

Complexity parameter used



K Nearest Neighbours

Birds of a feather flock together

Blog -

<https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761>

Link - <https://www.youtube.com/watch?v=UR2ag4lbBtc>

Time - 7m 17s

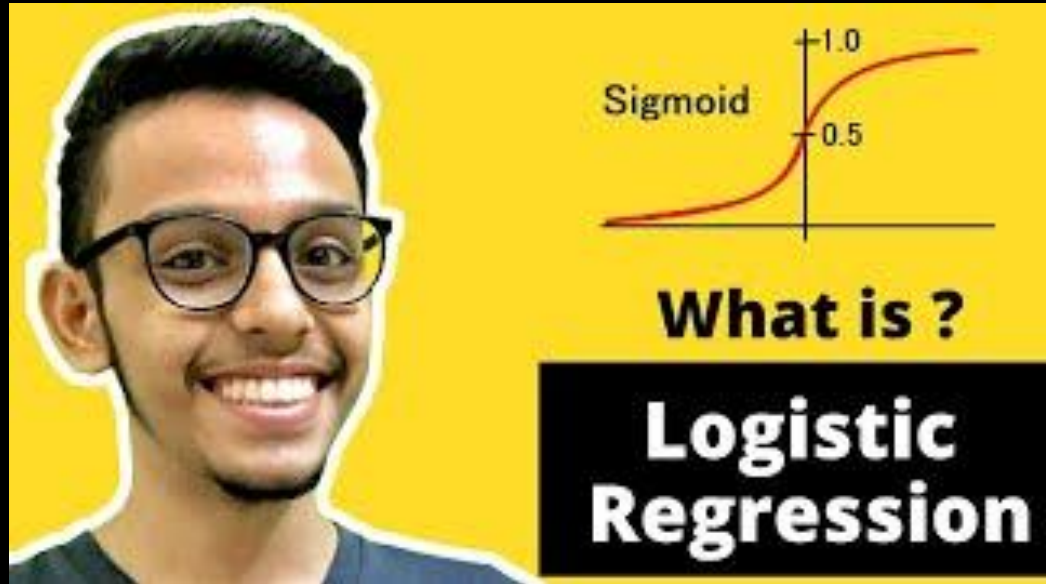


K
Nearest
Neighbor

Logistic Regression

Blog - <https://www.geeksforgeeks.org/understanding-logistic-regression/>

Video - <https://www.youtube.com/watch?v=U1omz0B9FTw>



Conclusion