

ML Algorithm	Algorithm Comparison
k-NN	<p>Type: Supervised Classification, Non-Parametric, lazy learner, discriminative.</p> <p>Feature Data: For majority of the scenario, Numeric/Quantitative sometimes categorical.</p> <p>Target Variable: Categorical</p> <p>Noisy Data: Is sensitive to noisy data and outliers are liable to get misclassified.</p> <p>Missing Data: k-NN cannot train on missing values therefore some reasonable imputation has to be performed or NA records must be deleted.</p> <p>Correlated Data: Does not affect the classification.</p> <p>Hyperparameters: K: number of neighbors, D: distance function like Euclidean, Cosine, Manhattan for numeric feature space and Hamming for categorical feature space.</p> <p>Evaluation Measures: Confusion Matrix, Specificity, Sensitivity etc.</p> <p>Package Implementation: class, caret.</p>
Naïve Bayes'	<p>Type: Supervised Classification, Parametric, eager learner, generative.</p> <p>Feature Data: For majority of the scenario Categorical</p> <p>Target Variable: Categorical</p> <p>Noisy Data: Averages out the noisy therefore not sensitive to noisy data. Its performance degrades more slowly as compared to other algorithms.</p> <p>Missing Data : The specific record is ignored for the frequency count. The R package implementation can also determine whether to ignore the example/record or the feature.</p> <p>Correlated Data: Does not work optimally due to class conditional independence assumption.</p> <p>Hyperparameters: Maximum Likelihood Estimation, Bayesian Estimation, optimization of loss criterion.</p> <p>Evaluation Measures: Confusion Matrix ,Sensitivity, Specificity.</p> <p>Package Implementation: e1071,naïve bayes, caret</p>
Classification Trees	<p>Type: Supervised Classification and Prediction(Dual), Non-Parametric, eager learner, discriminative, greedy learner, top down divide and conquer, bottom up pruning, recursive partitioning</p> <p>Feature Data: Numeric/Quantitative or Categorical</p> <p>Target Variable: Categorical or Numeric.</p> <p>Noisy Data: Decision trees are very sensitive to the noise in the data and even a small change in the data can throw the classification/prediction off therefore an ensemble of tree is better than just one tree.</p> <p>Missing Data : Is handled as per the package implementation for instance CART provides the functionality of assigning the missing data examples to the partition which has maximum examples already assigned to it.</p> <p>Correlated Data: Typically, it will not be affected by correlated variables and the split will occur as per the best feature in terms of info gain.</p>

	<p>Hyperparameters: Gini index, information gain, chi square ,pre pruning ,post pruning, assigning levels of trees</p> <p>Evaluation Measures: Confusion Matrix ,Sensitivity, Specificity, R squared, MSE.</p> <p>Package Implementation: CART,C5.0,ID 3</p>
Classification Rules	<p>Type: Supervised Classification , Non-Parametric, eager learner, discriminative, greedy learner, bottom up separate and conquer.</p> <p>Feature Data: In most scenarios Categorical.</p> <p>Target Variable: Mostly Categorical .</p> <p>Noisy Data: Decision rules are very sensitive to the noise in the data and even a small change in the data can throw the classification/prediction off .</p> <p>Missing Data : Is handled as per the package implementation for instance CART provides the functionality of assigning the missing data examples to the partition which has maximum examples already assigned to it.</p> <p>Correlated Data: Typically, it will not be affected by correlated variables and the split will occur as per the best feature in terms of info gain.</p> <p>Hyperparameters: Gini index, information gain, chi square ,pre pruning ,post pruning, assigning levels of trees</p> <p>Evaluation Measures: Confusion Matrix ,Sensitivity, Specificity</p> <p>Package Implementation: CART,C5.0,ID 3</p>