

ML model	Assumptions	Advantages	Disadvantages	Feature Scaling	Missing Data	Outliers	Suitable for	Learning	Example Use
Naïve Bayes Classifier	Features are independent	<ul style="list-style-type: none"> <li>Performs well with categorical variables</li> <li>Converges faster: less training time</li> <li>Good with moderate to large training data sets</li> <li>Good when dataset contains several features</li> </ul>	<ul style="list-style-type: none"> <li>Correlated features affect performance</li> </ul>	No	Can handle missing data (it ignores missing data)	Robust to outliers	<ul style="list-style-type: none"> <li>Classification</li> <li>Multiclass classification</li> </ul>	Supervised	<ul style="list-style-type: none"> <li>Sentiment Analysis</li> <li>Document categorisation</li> <li>Email Spam Filtering</li> </ul>
Support Vector Machine (SVM)	None	<ul style="list-style-type: none"> <li>Good for datasets with more variables than observations</li> <li>Good performance</li> <li>Good of-the-shelf model in general for several scenarios</li> <li>Can approximate complex non-linear functions</li> </ul>	<ul style="list-style-type: none"> <li>Long training time required</li> <li>Tuning is required to determine which kernel is optimal for non-linear SVMs</li> </ul>	Yes	Sensitive	Robust to outliers	<ul style="list-style-type: none"> <li>Classification</li> <li>Regression</li> </ul>	Supervised	<ul style="list-style-type: none"> <li>Stock market forecasting</li> <li>Value at risk determination</li> </ul>
Linear Regression	<ul style="list-style-type: none"> <li>Linear relation between features and target</li> <li>Residuals are normally distributed</li> <li>Homoscedasticity</li> </ul>	<ul style="list-style-type: none"> <li>Interpretability</li> <li>Little tuning</li> </ul>	<ul style="list-style-type: none"> <li>Correlated features may affect performance</li> <li>Extensive feature engineering required</li> </ul>	Yes	Sensitive	Sensitive	Regression	Supervised	<ul style="list-style-type: none"> <li>Sales forecasting</li> <li>House pricing</li> </ul>
Logistic Regression	<ul style="list-style-type: none"> <li>Linear relation between features and the log odds</li> <li>Residuals are normally distributed</li> <li>Homoscedasticity</li> </ul>	<ul style="list-style-type: none"> <li>Interpretability</li> <li>Little tuning</li> </ul>	<ul style="list-style-type: none"> <li>Correlated features may affect performance</li> <li>Extensive feature engineering required</li> </ul>	Yes	Sensitive	Potentially sensitive	Classification	Supervised	<ul style="list-style-type: none"> <li>Risk Assessment</li> <li>Fraud Prevention</li> </ul>
Classification and Regression Trees	None	<ul style="list-style-type: none"> <li>Interpretability</li> <li>Render feature importance</li> <li>Less data pre-processing required</li> </ul>	<ul style="list-style-type: none"> <li>Do not predict a continuous output (for regression)</li> <li>It does not predict beyond the range of the response values in the training data.</li> <li>Overfits</li> </ul>	No	Some implementations do not need missing data imputation. The one in Scikit-learn does	Robust to outliers	<ul style="list-style-type: none"> <li>Classification</li> <li>Regression</li> </ul>	Supervised	<ul style="list-style-type: none"> <li>Risk Assessment</li> <li>Fraud Prevention</li> </ul>
Random Forests	None	<ul style="list-style-type: none"> <li>Interpretability</li> <li>Render feature importance</li> <li>Less data pre-processing required</li> <li>Do not overfit (in theory)</li> <li>Good performance /accuracy</li> <li>Robust to noise</li> <li>Little if any parameter tuning required</li> <li>Apt for almost any machine learning problem</li> </ul>	<ul style="list-style-type: none"> <li>Do not predict a continuous output (for regression)</li> <li>It does not predict beyond the range of the response values in the training data</li> <li>Biased towards categorical variables with several categories</li> <li>Biased in multiclass problems toward more frequent classes</li> </ul>	No	Some implementations do not need missing data imputation. The one in Scikit-learn does.	Robust to outliers	<ul style="list-style-type: none"> <li>Classification</li> <li>Regression</li> </ul>	Supervised	<ul style="list-style-type: none"> <li>Credit Risk Assessment</li> <li>Predict breakdown of mechanical parts (automobile industry).</li> <li>Assess probability of developing a chronic disease (healthcare)</li> <li>Predicting the average number of social media shares</li> </ul>
Gradient Boosted Trees	None	<ul style="list-style-type: none"> <li>Great performance</li> <li>Apt for almost any machine learning problem</li> <li>It can approximate most non-linear functions</li> </ul>	<ul style="list-style-type: none"> <li>Prone to overfit</li> <li>Needs some parameter tuning</li> </ul>	No	Some implementations do not need missing data imputation (e.g. xgboost). The one in Scikit-learn does.	Robust to outliers	<ul style="list-style-type: none"> <li>Classification</li> <li>Regression</li> </ul>	Supervised	<ul style="list-style-type: none"> <li>Same as Random Forests</li> </ul>
K-nearest neighbours	None	<ul style="list-style-type: none"> <li>Good performance</li> </ul>	<ul style="list-style-type: none"> <li>Slow when predicting</li> <li>Susceptible to high dimension (lots of features)</li> </ul>	Yes	Sensitive	Robust to outliers	<ul style="list-style-type: none"> <li>Classification</li> <li>Regression</li> </ul>	Supervised	<ul style="list-style-type: none"> <li>Gene expression</li> <li>Protein-protein interaction</li> <li>Content retrieval (of webpages for example)</li> </ul>
AdaBoost	None	<ul style="list-style-type: none"> <li>It doesn't overfit easily</li> <li>Few parameters to tune</li> </ul>	<ul style="list-style-type: none"> <li>Can be sensitive to noise and outliers</li> </ul>	No	Can handle	Sensitive	<ul style="list-style-type: none"> <li>Classification</li> <li>Regression</li> </ul>	Supervised	<ul style="list-style-type: none"> <li>Same as Random Forests, less used however, as xgboost and lightGBMs are more popular implementations of gradient boosted machines</li> </ul>
Neural Networks	None	<ul style="list-style-type: none"> <li>Can approximate any function</li> <li>Great Performance</li> </ul>	<ul style="list-style-type: none"> <li>Long training time</li> <li>Several parameters to tune, including neuronal architecture</li> <li>Prone to overfit</li> <li>Little interpretability</li> </ul>	Yes	Sensitive	Can handle outliers, and it affects performance if they are too many	<ul style="list-style-type: none"> <li>Classification</li> <li>Regression</li> </ul>	Supervised	<ul style="list-style-type: none"> <li>Image analysis</li> <li>Forecasting</li> <li>Text analysis</li> </ul>
K-Means Clustering	<ul style="list-style-type: none"> <li>clusters are spherical</li> <li>clusters are of similar size</li> </ul>	<ul style="list-style-type: none"> <li>Fast training</li> </ul>	<ul style="list-style-type: none"> <li>Need to determine k, the number of clusters</li> <li>Sensitive to initial points and local optima</li> </ul>	Yes	In the Scikit-learn implementation, missing data needs to be imputed	Sensitive	<ul style="list-style-type: none"> <li>Segmentation</li> </ul>	Unsupervised	<ul style="list-style-type: none"> <li>Customer segmentation</li> <li>Outlier detection</li> </ul>
Hierarchical clustering		<ul style="list-style-type: none"> <li>No a priori information about the number of clusters required</li> </ul>	<ul style="list-style-type: none"> <li>Final number of clusters to be decided by the scientist</li> <li>Slow training</li> </ul>	Yes	Sensitive	Sensitive	<ul style="list-style-type: none"> <li>Segmentation</li> </ul>	Unsupervised	<ul style="list-style-type: none"> <li>Customer segmentation</li> <li>Gene analyses</li> </ul>
PCA	<ul style="list-style-type: none"> <li>Correlation among features</li> </ul>	<ul style="list-style-type: none"> <li>Captures most of the variance in a smaller number of features</li> </ul>	<ul style="list-style-type: none"> <li>Number of principal components that explain most of the variance to be determined by the user</li> </ul>	Yes	Sensitive	Sensitive	Reducing feature space to train machine learning models	Unsupervised	<ul style="list-style-type: none"> <li>Creating few, informative, variables from tons of data</li> </ul>