



5

# Modèles Boosting Comparaison

# 5. Modèles – Boosting - Comparaison



Function	XGBoost	CatBoost	Light GBM
Important parameters which control overfitting	<ul style="list-style-type: none"> <li>1. <b>learning_rate or eta</b> – optimal values lie between 0.01-0.2</li> <li>2. <b>max_depth</b></li> <li>3. <b>min_child_weight</b>: similar to <b>min_child_leaf</b>; default is 1</li> </ul>	<ul style="list-style-type: none"> <li>1. <b>Learning_rate</b></li> <li>2. <b>Depth</b> - value can be any integer up to 16. Recommended - [1 to 10]</li> <li>3. No such feature like <b>min_child_weight</b></li> <li>4. <b>L2-leaf-reg</b>: L2 regularization coefficient. Used for leaf value calculation (any positive integer allowed)</li> </ul>	<ul style="list-style-type: none"> <li>1. <b>learning_rate</b></li> <li>2. <b>max_depth</b>: default is 20. Important to note that tree still grows leaf-wise. Hence it is important to tune <b>num_leaves</b> (number of leaves in a tree) which should be smaller than <math>2^{\text{max\_depth}}</math>. It is a very important parameter for LGBM</li> <li>3. <b>min_data_in_leaf</b>: default=20, alias= <b>min_data</b>, <b>min_child_samples</b></li> </ul>
Parameters for categorical values	Not Available	<ul style="list-style-type: none"> <li>1. <b>cat_features</b>: It denotes the index of categorical features</li> <li>2. <b>one_hot_max_size</b>: Use one-hot encoding for all features with number of different values less than or equal to the given parameter value (max – 255)</li> </ul>	<ul style="list-style-type: none"> <li>1. <b>categorical_feature</b>: specify the categorical features we want to use for training our model</li> </ul>
Parameters for controlling speed	<ul style="list-style-type: none"> <li>1. <b>colsample_bytree</b>: subsample ratio of columns</li> <li>2. <b>subsample</b>: subsample ratio of the training instance</li> <li>3. <b>n_estimators</b>: maximum number of decision trees; high value can lead to overfitting</li> </ul>	<ul style="list-style-type: none"> <li>1. <b>rsm</b>: Random subspace method. The percentage of features to use at each split selection</li> <li>2. No such parameter to subset data</li> <li>3. <b>iterations</b>: maximum number of trees that can be built; high value can lead to overfitting</li> </ul>	<ul style="list-style-type: none"> <li>1. <b>feature_fraction</b>: fraction of features to be taken for each iteration</li> <li>2. <b>bagging_fraction</b>: data to be used for each iteration and is generally used to speed up the training and avoid overfitting</li> <li>3. <b>num_iterations</b>: number of boosting iterations to be performed; default=100</li> </ul>