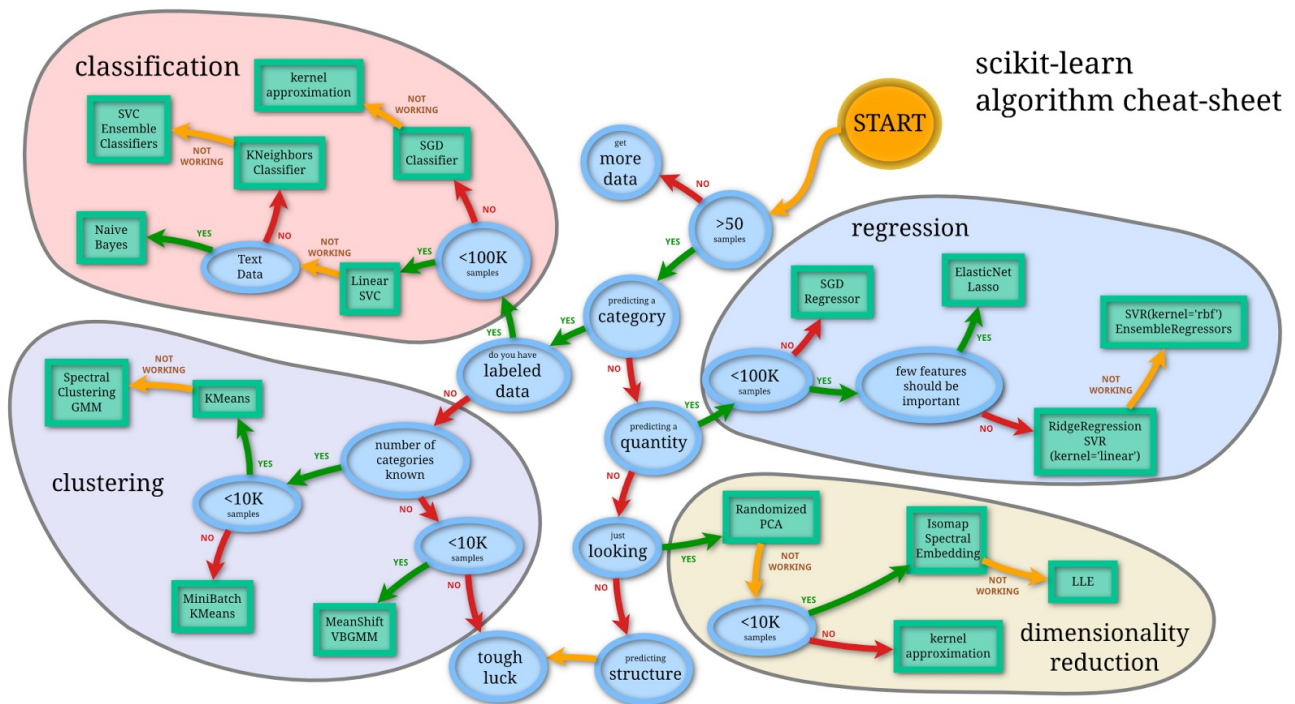


17.08.01(机器学习算法使用)



179种分类模型在UCI所有的121个数据上的性能，发现Random Forests 和 SVM 性能最好。

	Decision Trees	Neural Networks	Naïve Bayes	kNN	SVM	Rule-learners
Accuracy in general	**	***	*	**	****	**
Speed of learning with respect to number of attributes and the number of instances	***	*	****	****	*	**
Speed of classification	****	****	****	*	****	****
Tolerance to missing values	***	*	****	*	**	**
Tolerance to irrelevant attributes	***	*	**	**	****	**
Tolerance to redundant attributes	**	**	*	**	***	**
Tolerance to highly interdependent attributes (e.g. parity problems)	**	***	*	*	***	**
Dealing with discrete/binary/continuous attributes	****	***(not discrete)	***(not continuous)	***(not directly discrete)	** (not discrete)	*** (not directly continuous)
Tolerance to noise	**	**	***	*	**	*
Dealing with danger of overfitting	**	*	***	***	**	**
Attempts for incremental learning	**	***	****	****	**	*
Explanation ability/transparency of knowledge/classifications	****	*	****	**	*	****
Model parameter handling	***	*	****	***	*	***

Table 4. Comparing learning algorithms (**** stars represent the best and * star the worst performance)

各算法比较

但是直接应用算法后，一般精度都不是很理想，这个时候需要调节参数，最干货的问题来了，什么模型需要调节什么参数呢？

Model	Parameters to optimize	Good range of values
Linear Regression	<ul style="list-style-type: none"> fit_intercept normalize 	<ul style="list-style-type: none"> True / False True / False
Ridge	<ul style="list-style-type: none"> alpha Fit_intercept Normalize 	<ul style="list-style-type: none"> 0.01, 0.1, 1.0, 10, 100 True/False True/False
k-neighbors	<ul style="list-style-type: none"> N_neighbors p 	<ul style="list-style-type: none"> 2, 4, 8, 16 2, 3
SVM	<ul style="list-style-type: none"> C Gamma class_weight 	<ul style="list-style-type: none"> 0.001, 0.01.....10...100...1000 'Auto', RS* 'Balanced' , None
Logistic Regression	<ul style="list-style-type: none"> Penalty C 	<ul style="list-style-type: none"> L1 or L2 0.001, 0.01.....10...100
Naive Bayes (all variations)	NONE	NONE
Lasso	<ul style="list-style-type: none"> Alpha Normalize 	<ul style="list-style-type: none"> 0.1, 1.0, 10 True/False
Random Forest	<ul style="list-style-type: none"> N_estimators Max_depth Min_samples_split Min_samples_leaf Max features 	<ul style="list-style-type: none"> 120, 300, 500, 800, 1200 5, 8, 15, 25, 30, None 1, 2, 5, 10, 15, 100 1, 2, 5, 10 Log2, sqrt, None
Xgboost	<ul style="list-style-type: none"> Eta Gamma Max_depth Min_child_weight Subsample Colsample_bytree Lambda alpha 	<ul style="list-style-type: none"> 0.01,0.015, 0.025, 0.05, 0.1 0.05-0.1,0.3,0.5,0.7,0.9,1.0 3, 5, 7, 9, 12, 15, 17, 25 1, 3, 5, 7 0.6, 0.7, 0.8, 0.9, 1.0 0.6, 0.7, 0.8, 0.9, 1.0 0.01-0.1, 1.0 , RS* 0, 0.1, 0.5, 1.0 RS*

<http://blog.csdn.net/aliceyangxi1987/article/details/71079448>