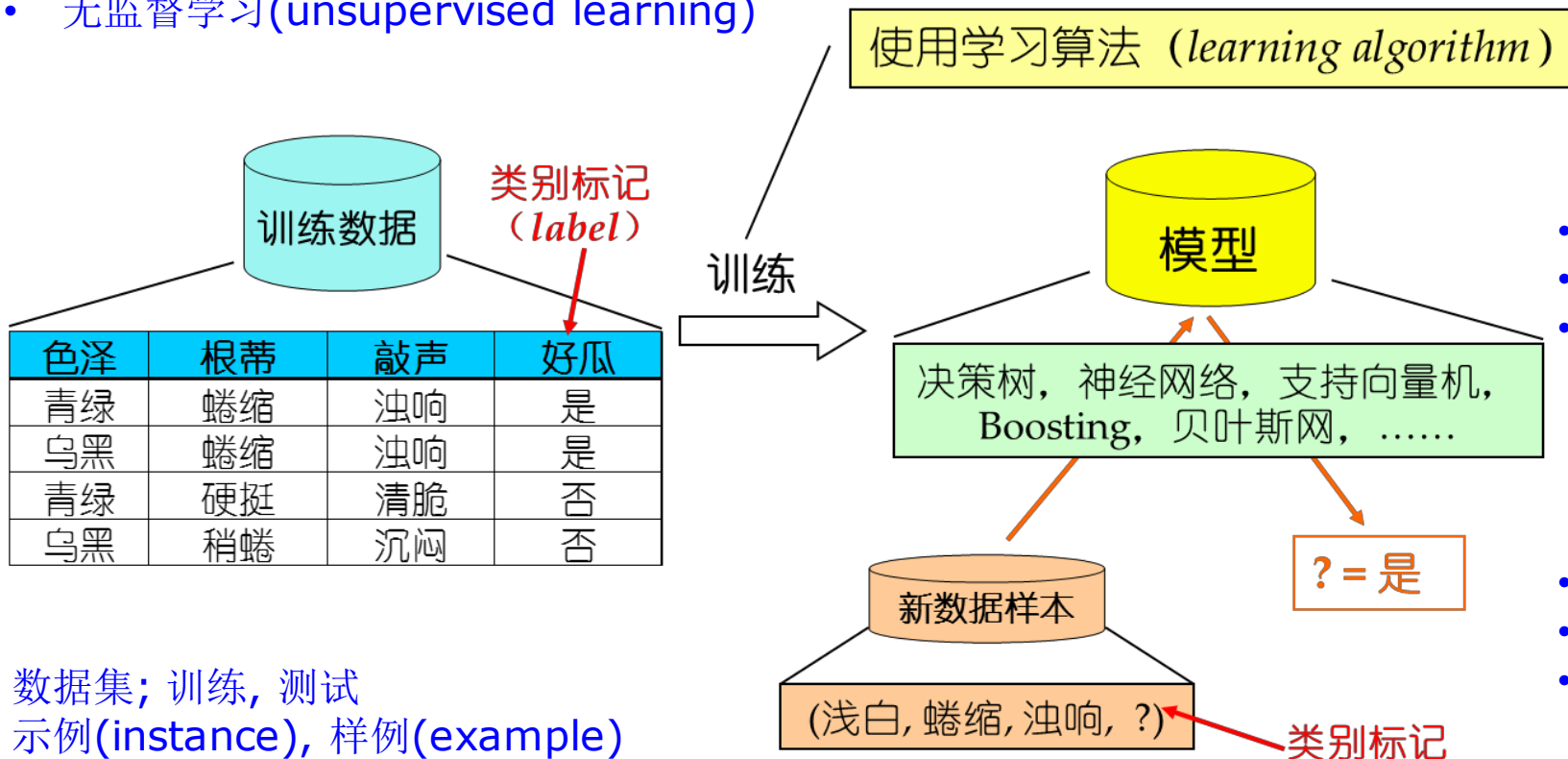


Spark机器学习

寒小阳
2018-01-13

基本术语与概念

- 监督学习(supervised learning)
- 无监督学习(unsupervised learning)

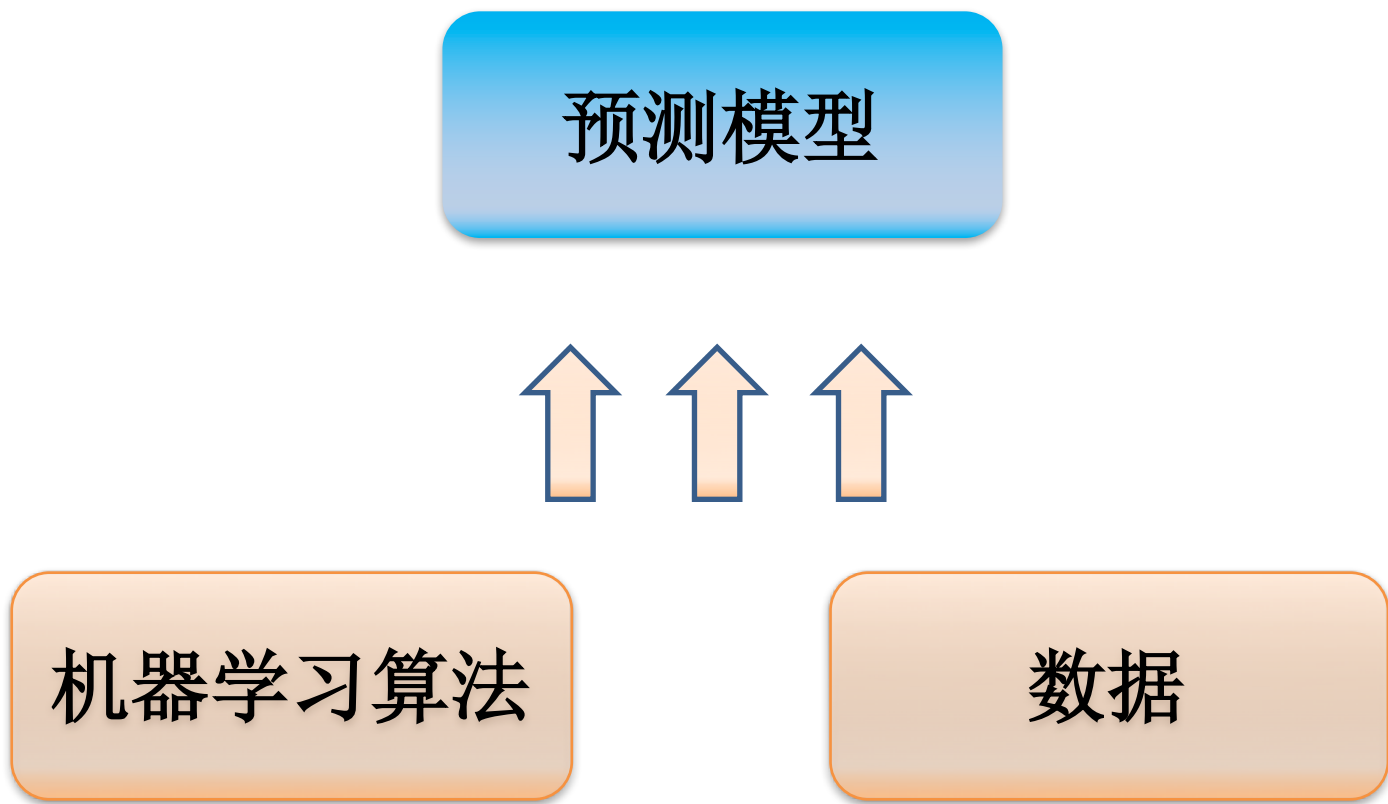


- 假设(hypothesis)
- 真相(ground-truth)
- 学习器(learner)

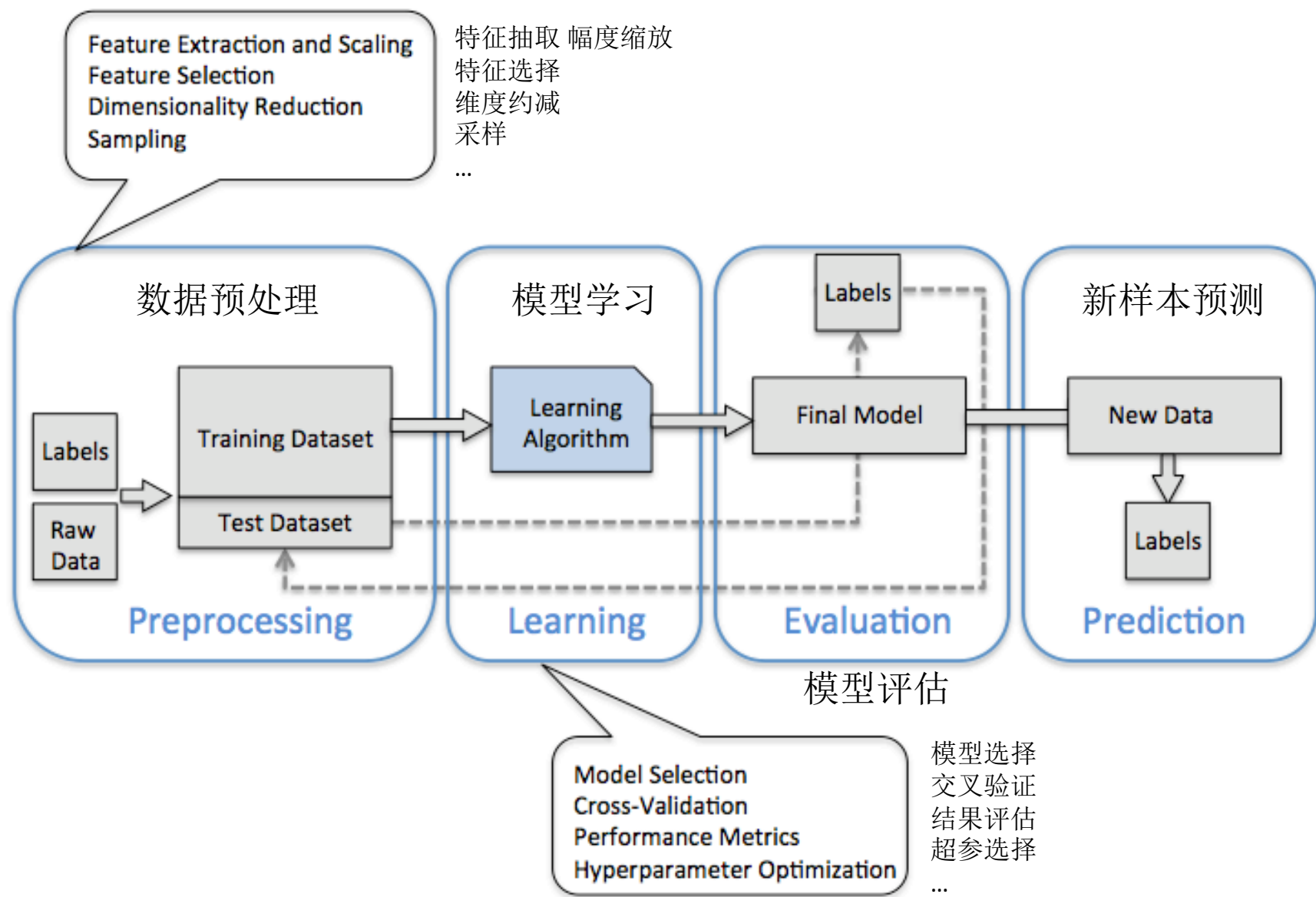
- 数据集; 训练, 测试
- 示例(instance), 样例(example)
- 样本(sample)
- 属性(attribute), 特征(feature); 属性值
- 属性空间, 样本空间, 输入空间
- 特征向量(feature vector)
- 标记空间, 输出空间

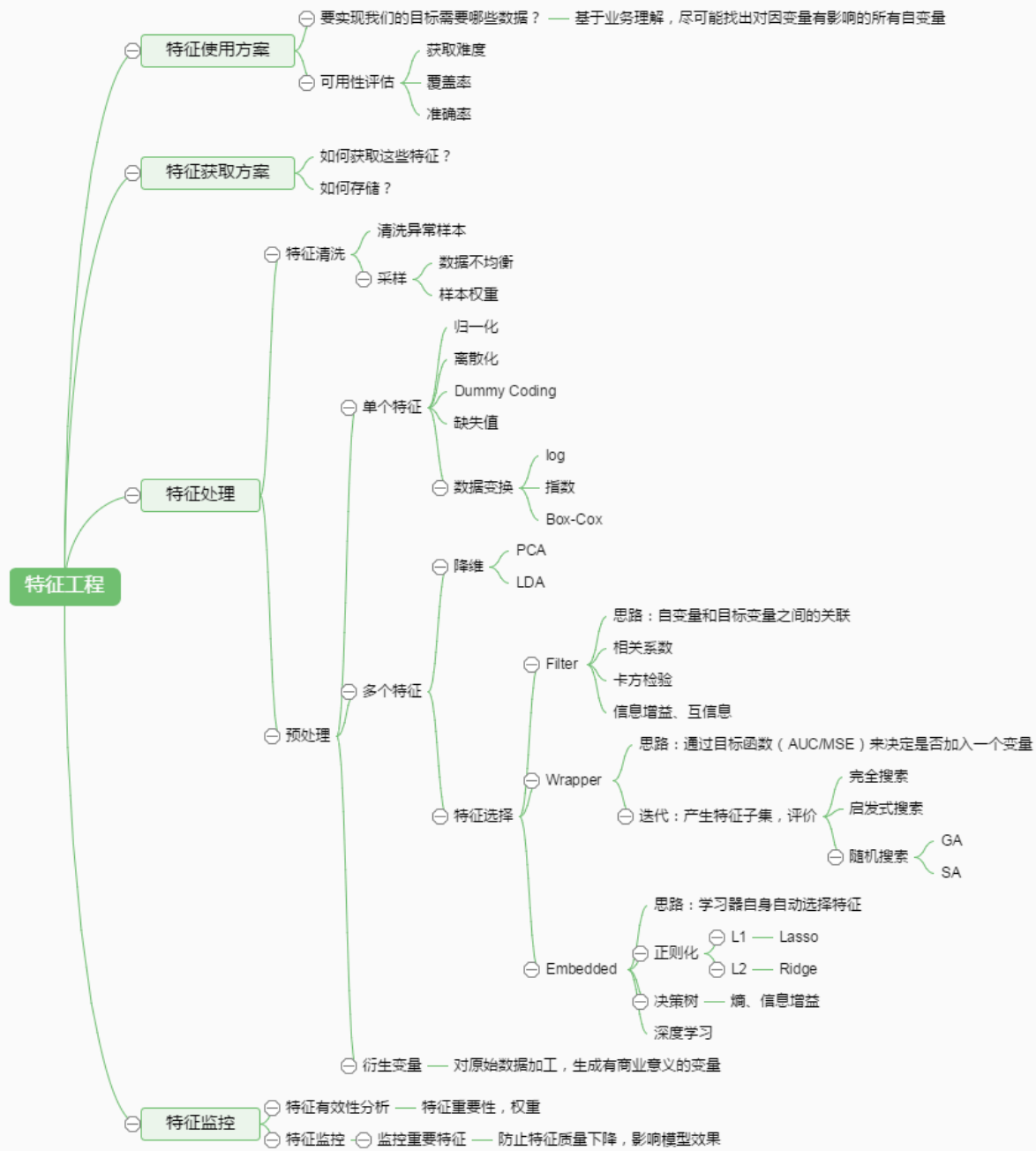
- 分类, 回归
- 二分类, 多分类
- 正类, 反类

- 未见样本(unseen instance)
- 未知“分布”
- 独立同分布(i.i.d.)
- 泛化(**generalization**)

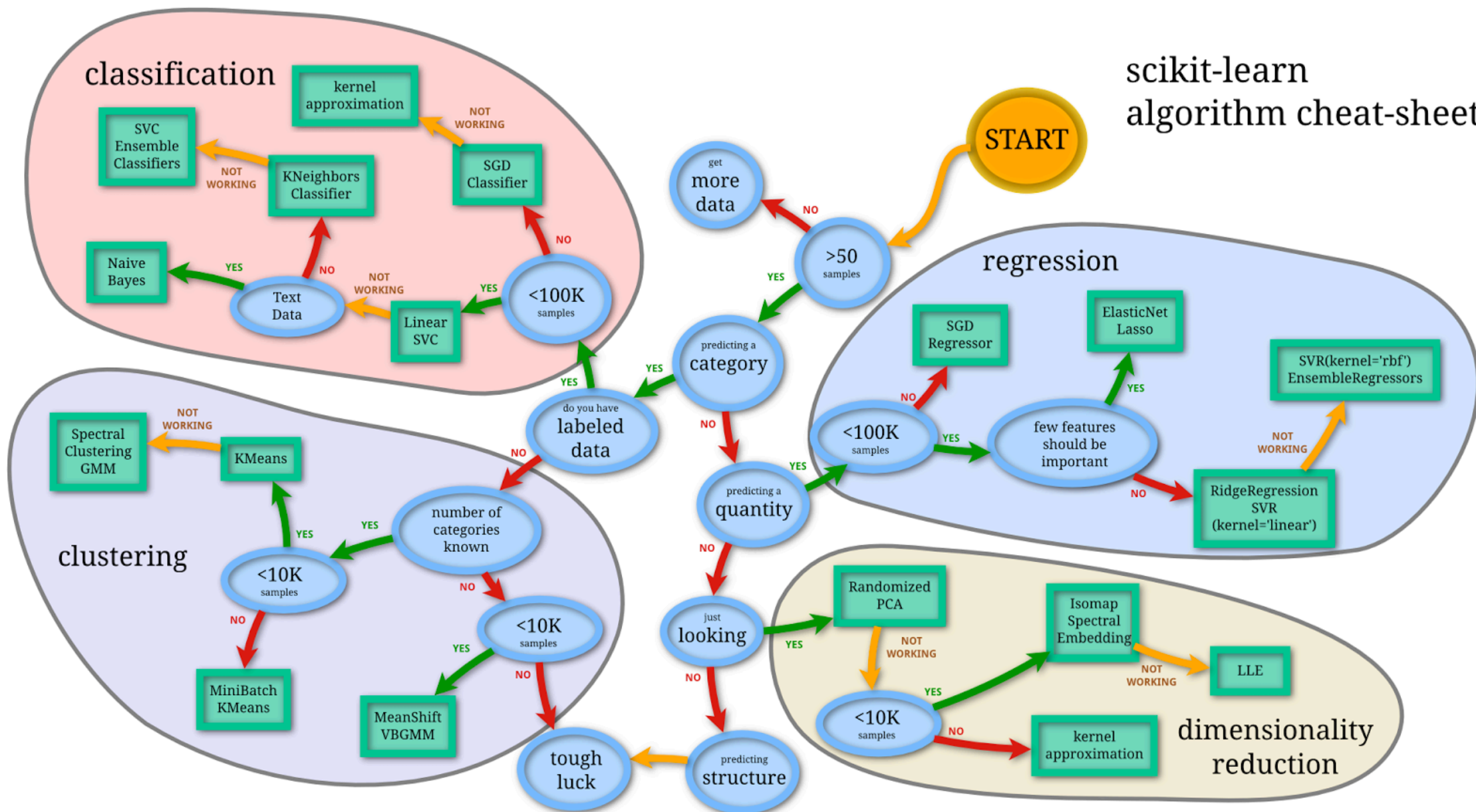


机器学习的应用工作是围绕着数据与算法展开的





样本调权
归一化
离散化
独热向量编码
Log/exp变换
PCA



感谢大家！

恳请大家批评指正！