

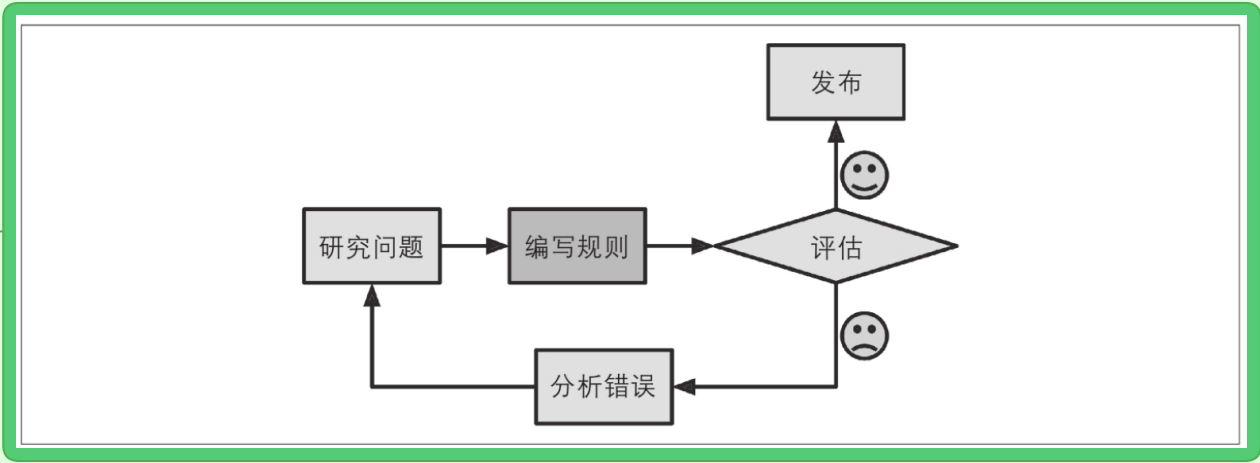
机器学习是什么？

机器学习是一门通过编程让计算机从数据中进行学习的科学(和艺术)

为什么使用机器学习？

和传统学习方法的区别

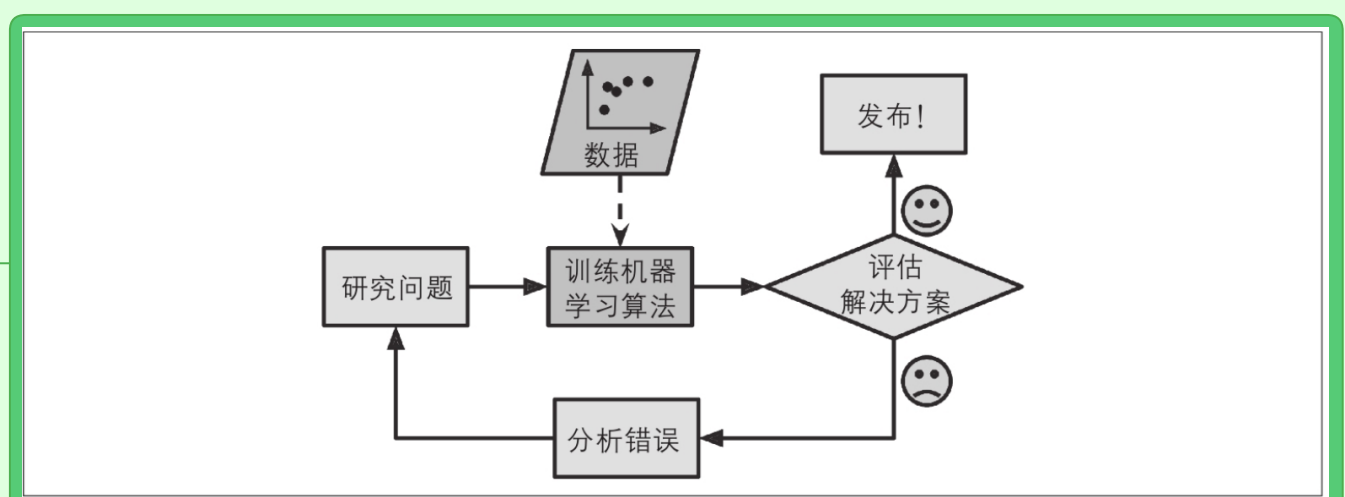
数据



ML适用于

- 有复杂的解决方案(需要大量的人工参与)和复杂的code
- 传统方法根本解决不了的问题
- 环境有波动, ML可以适应新数据
- 洞察复杂问题和大量数据

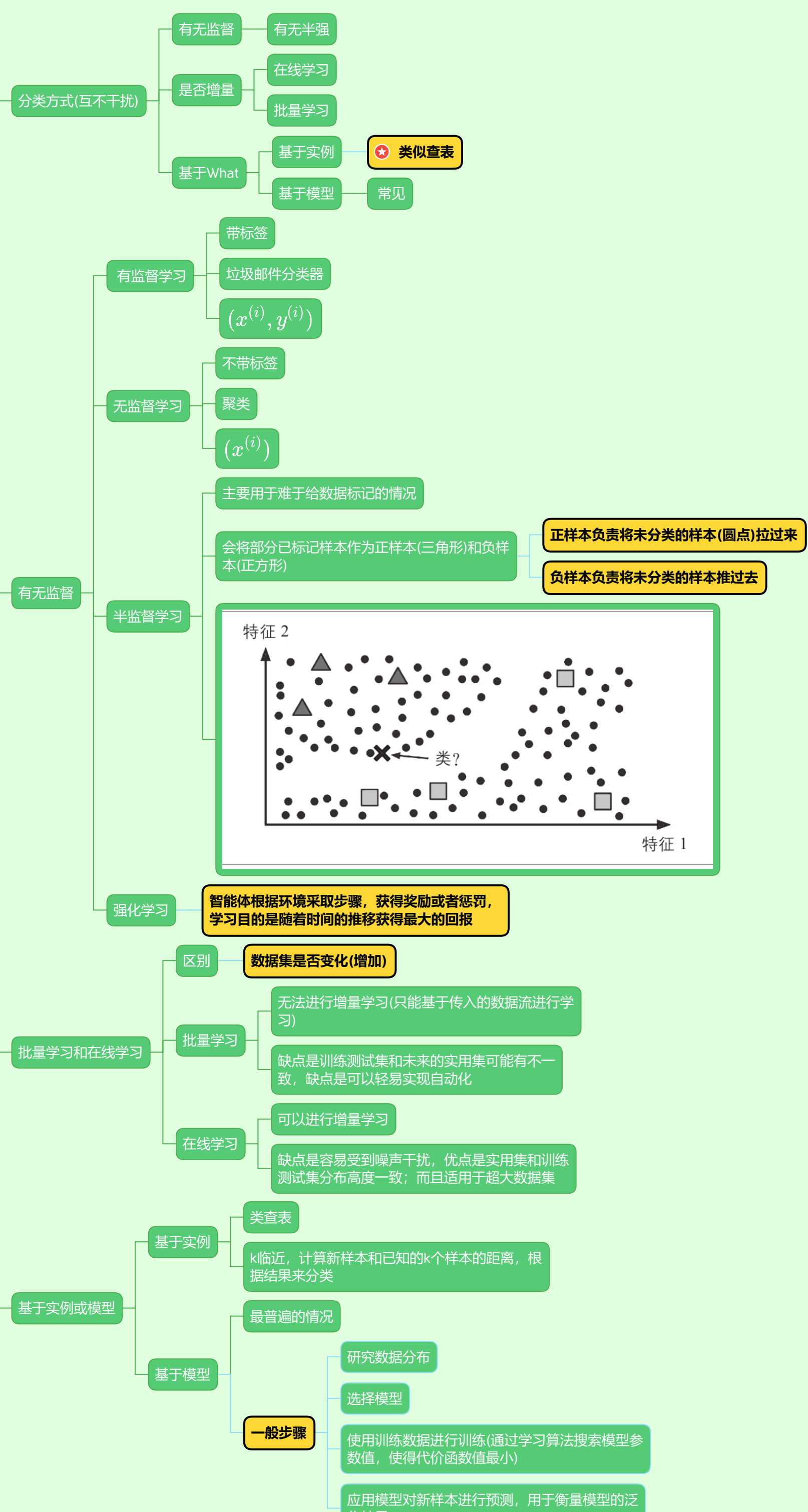
ML可以简化代码, 减少人力投入



机器学习的应用示例

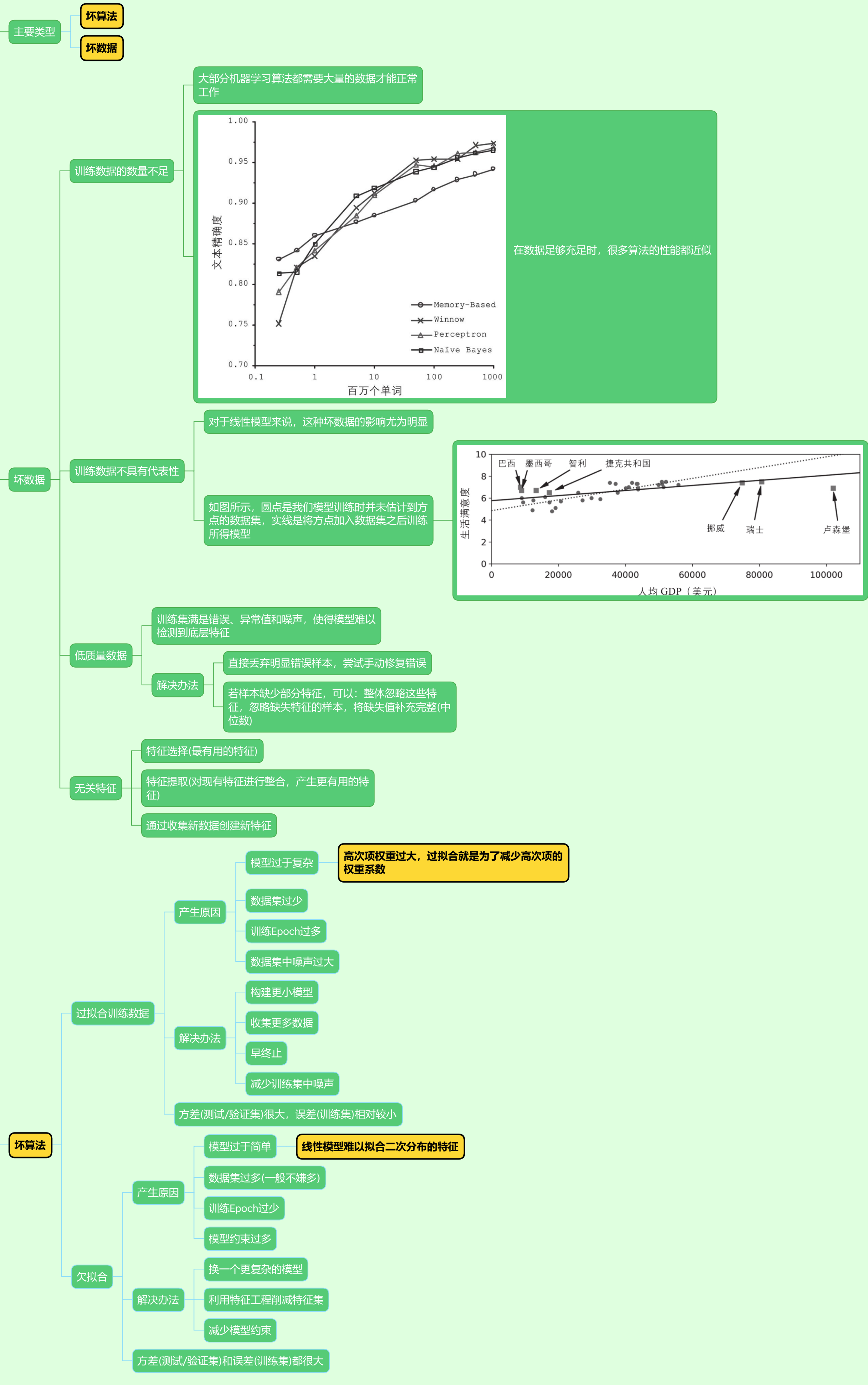
- CV
 - 图像分类 - 生产产品分类
 - 语义分割 - 医学影像处理
- NLP
 - 新闻分类
 - 标记不良发言
 - 对长文本进行总结
- NLU+问答模块
 - 创建小冰和客服
- 利用多个指标来预测下一年的收入
 - 回归问题
 - 线性、多项式、SVM、随机森林、人工神经网络
- 语音识别
- 异常检测 - 检测欺诈
- 用户分类 - 聚类
- 利用图表表示高维数据集 - 数据可视化
- 推荐产品 - 推荐系统
- 创造人机 - RL

机器学习系统的类型



第一章 机器学习概论

机器学习的主要挑战



测试和验证

超参数调整和模型选择

我们的模型最终会运用在新场景(新数据), 为了了解模型性能, 我们不能只通过基于训练集的表现来看(毕竟存在过拟合), 所以我们需要分割出部分数据来模拟“新场景”, 这便是训练集。

单个模型可以通过基于测试集上面的表现来衡量, 多个模型如何衡量呢? 我们可以将这多个模型都训练来看基于测试集的泛化效果。

但是在我们多次基于测试集挑选出效果最好的模型之后, 该模型可能对测试集过拟合而导致在实用集上面表现不好, 所以我们仍需要另外一个集合用于拟合在实用集上面的性能, 这个集合被称为开发集。

- 三个集合总结
 - 训练集 - 用于训练模型
 - 开发集 - 用于调整超参数和选择模型
 - 测试集 - 用于测试模型的最终性能

数据不匹配

实用集和训练测试集分布不同的情况

我们可以将训练集分出一部分用于开发(train-dev), 表现不好则来自数据不匹配