

# 阿里巴巴大数据竞赛

天猫推荐算法大挑战

第二赛季 总决赛

## 大数据，小细节



Fly402 团队

# 参赛历程

**阿里巴巴大数据竞赛**  
天猫推荐算法 大挑战

第二赛季 总决赛



# 思路

预测购买

未交互



协同过滤、关联规则、隐语义，矩阵分解.....

交互



逻辑回归、贝叶斯、SVM、随机森林，GBRT.....

# 交互



# 核心要点

阿里巴巴大数据竞赛  
天猫推荐算法大挑战

第二赛季 总决赛

- 特征选择

- 模型选择

- 训练

- 评价

- 融合



# 特征选择

样本形式：

User<sub>i</sub>

Brand<sub>x</sub>

总分类 \ 时间维度	7天	15天	30天	60天	90天
User <sub>i</sub>					
Brand <sub>x</sub>					
User <sub>i</sub> -Brand <sub>x</sub>					

类型维度：

点击

购买

收藏

购物车

# 特征选择

业务指标

数据分析

特征提取

购买该品牌的用户数

b\_bought\_u\_count

品牌被购买次数

b\_bought\_times

回头客的数量

b\_rebought\_u\_count

平均每个回头客的回头次数

b\_avg\_rebought\_days

回头客数量/购买的用户数

b\_rebought\_ratio

品牌购买价值



# 特征选择

- 比值类特征
  - 部分基本特征相除
- 分类训练预测的启发
  - u买过b跟u没买过b的特征有区别
- 均值、方差等
  - 体现特征值的平均水平与波动情况
- 策略
  - Season1的策略转换为特征

- 特征选择

- 模型选择

- 训练

- 评价

- 融合



# 模型选择

逻辑回归

线性模型

正则化项

逐步逻辑回归

训练预测效率高

大数据量下精度下降

# 模型选择

随机森林

averaging方法

降低方差

基于分类决策树

抗噪能力强

训练快、预测慢



# 模型选择

GBRT

boosting方法

降低偏差

回归树

训练慢、预测快

- 特征选择

- 模型选择

- 训练

- 评价

- 融合



# 训练集构造

04-15

08-15

122

30

0

标记

预测集

92

0

时间尺度统一；树模型特征无需处理；

?

预测集

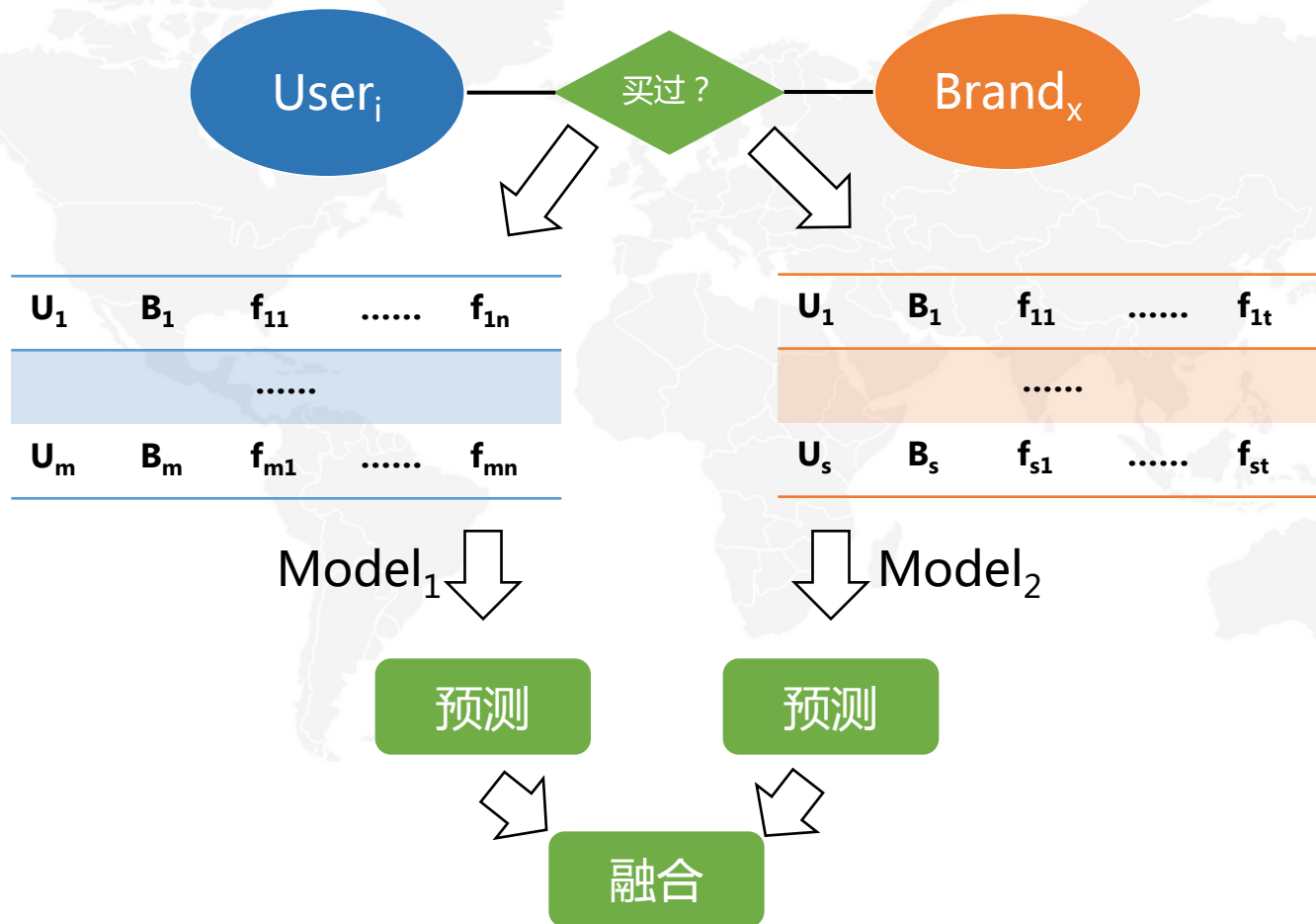
122

0

信息量更足；更多的用户品牌对；时间尺度不统一，需要特殊处理

?

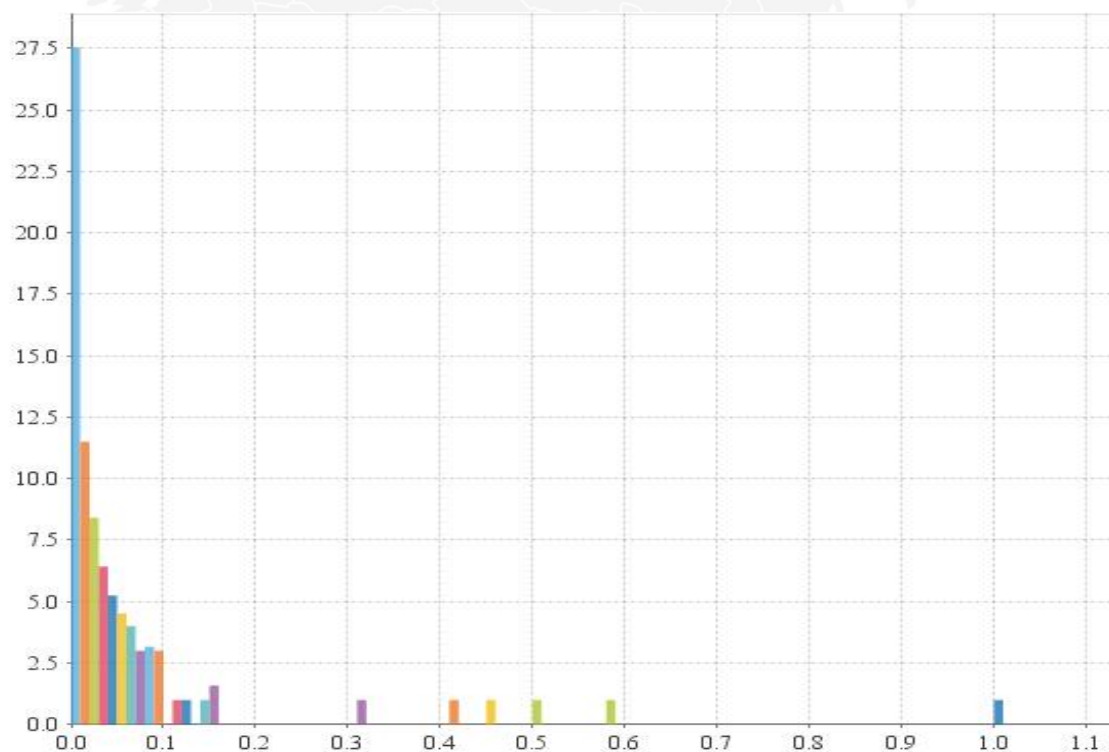
# 分类训练预测



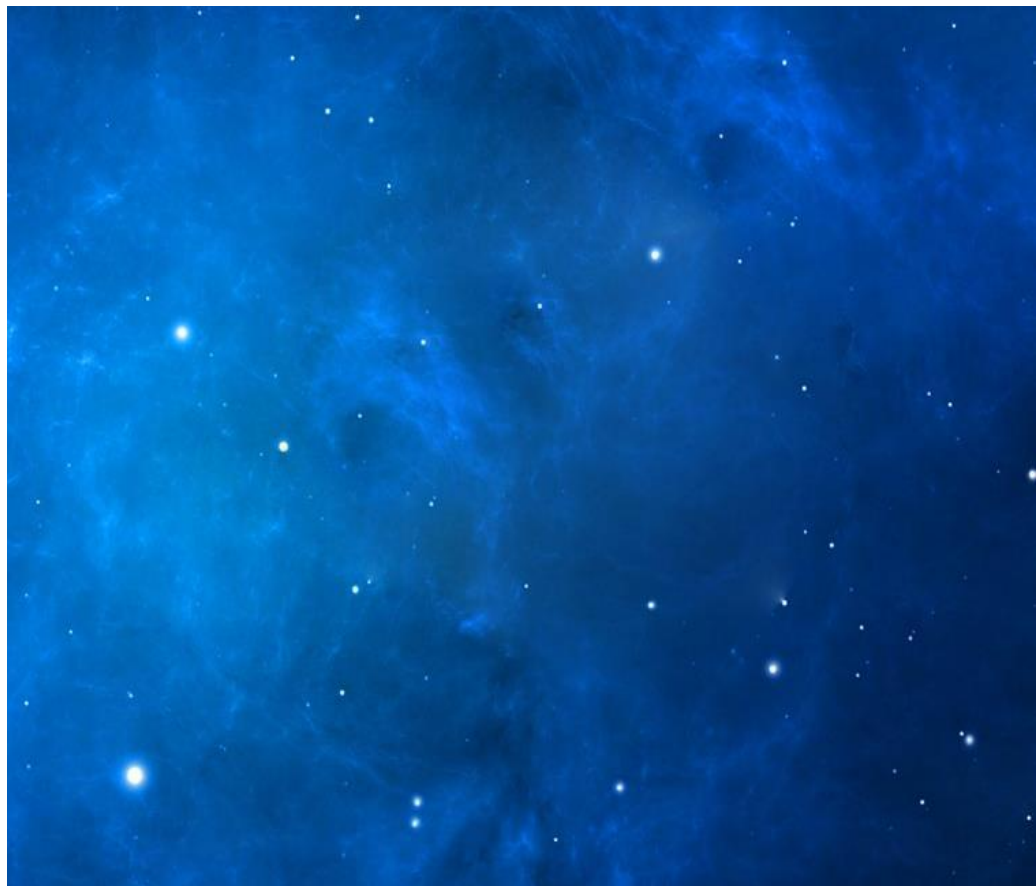


# 数据处理

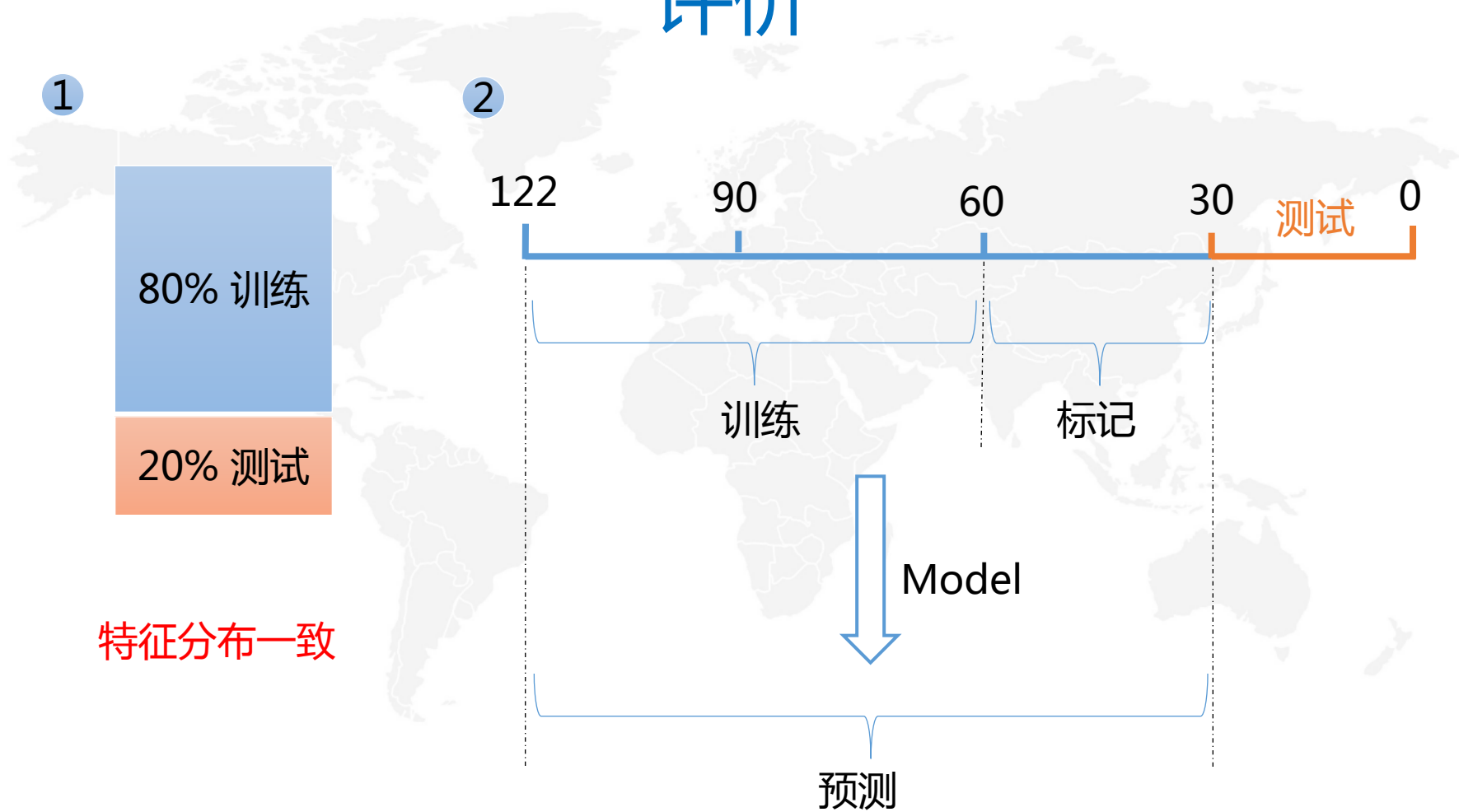
- 归一化
  - $(X - \min) / (\max - \min)$
- 分箱
  - 等量、等距
- 去噪



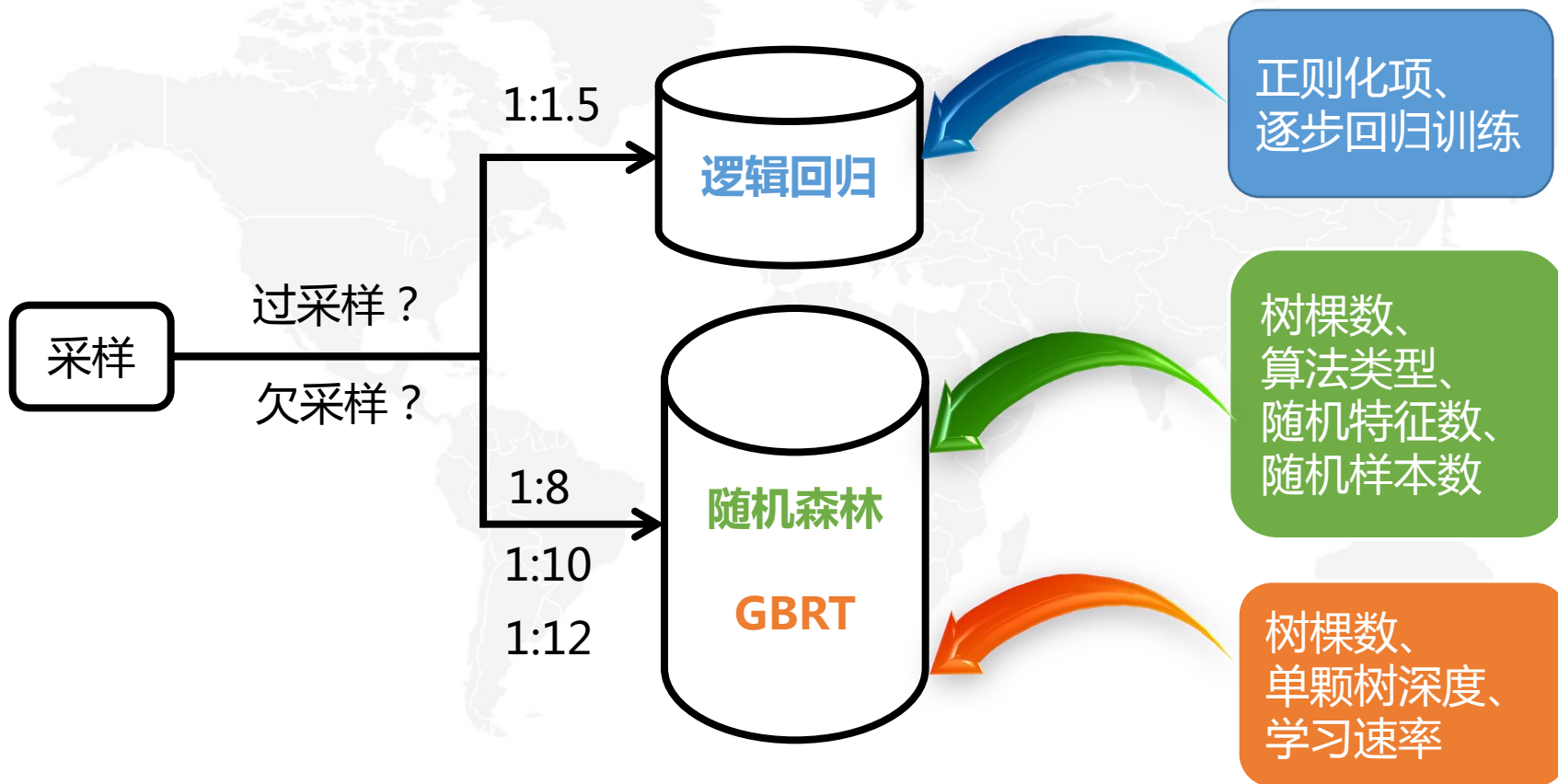
- 特征选择
- 模型选择
- 训练
- 评价
- 融合



# 评价



# 采样与调参





# 特征选择

- 逐步回归
- 相关矩阵
- 属性选择
  - 信息增益、基尼指数、信息增益率
- 编写脚本
  - 模型反馈
- 分批迭代测试

# 过拟合

- 健壮的本地图试
- 鲁棒性强的特征体系
- 数据分布
- 模型参数

- 特征选择
- 模型选择
- 训练
- 评价
- 融合

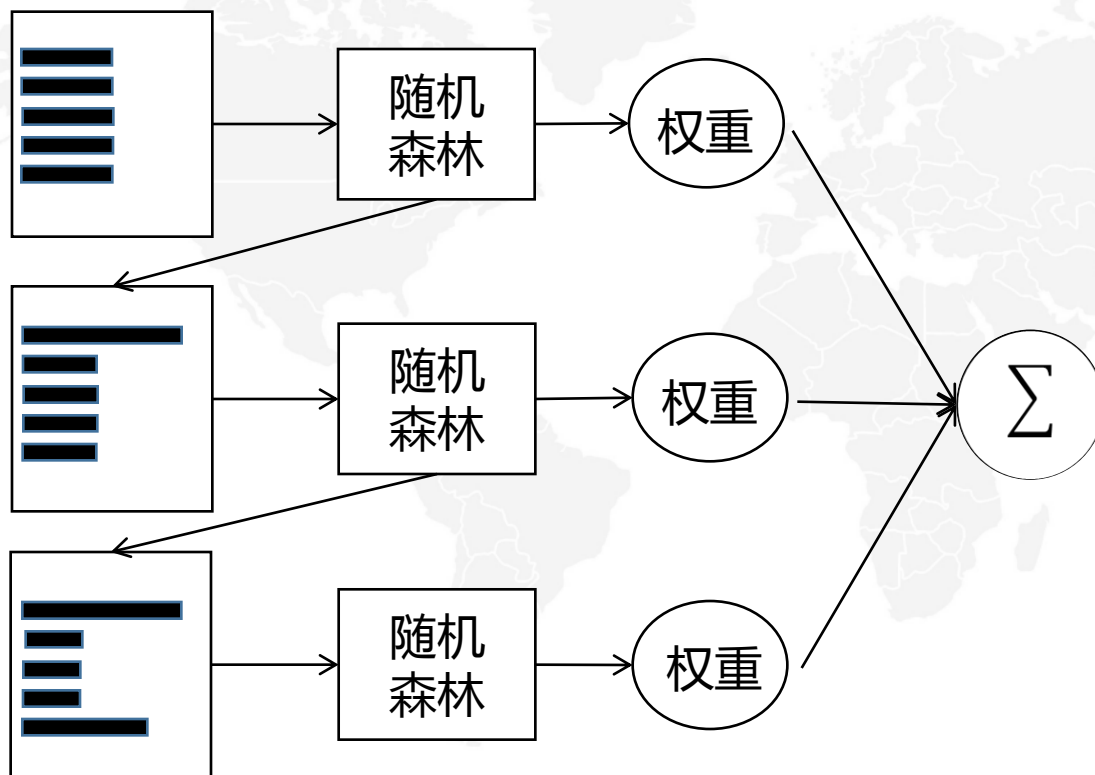


# 简单融合

- 不同时间尺度融合
  - 4个月+3个月+2个月
- 同模型融合
  - 不同采样比例、不同参数
- 不同模型融合
  - 逻辑回归、随机森林、GBRT
- 分开与全集



# 复杂融合



没效果原因：

- 随机森林是强分类器
- 迭代次数不够
- 错误率、权重调整不当

改进



- 使用弱分类器：LR或者参数简单的GBRT
- 不用权重列，使用有放回加权采样

谢谢！

阿里巴巴大数据竞赛  
天猫推荐算法大挑战

第二赛季 总决赛