

天猫推荐算法大挑战

第二赛季 总决赛



阿里巴巴大数据竞赛分享

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

目录

1 解题思路

2 特征提取

3 算法实现

4 总结回忆



目录

1 解题思路

2 特征提取

3 算法实现

4 总结回忆

1 解题思路

■ 赛题的任务

根据用户4个月在天猫的行为日志，建立用户的品牌偏好，并预测他们在将来一个月内会购买哪些品牌。

$$\text{Precision} = \frac{\sum_i^N \text{hitBrands}_i}{\sum_i^N \text{pBrands}_i}$$

$$\text{Recall} = \frac{\sum_i^M \text{hitBrands}_i}{\sum_i^M \text{pBrands}_i}$$

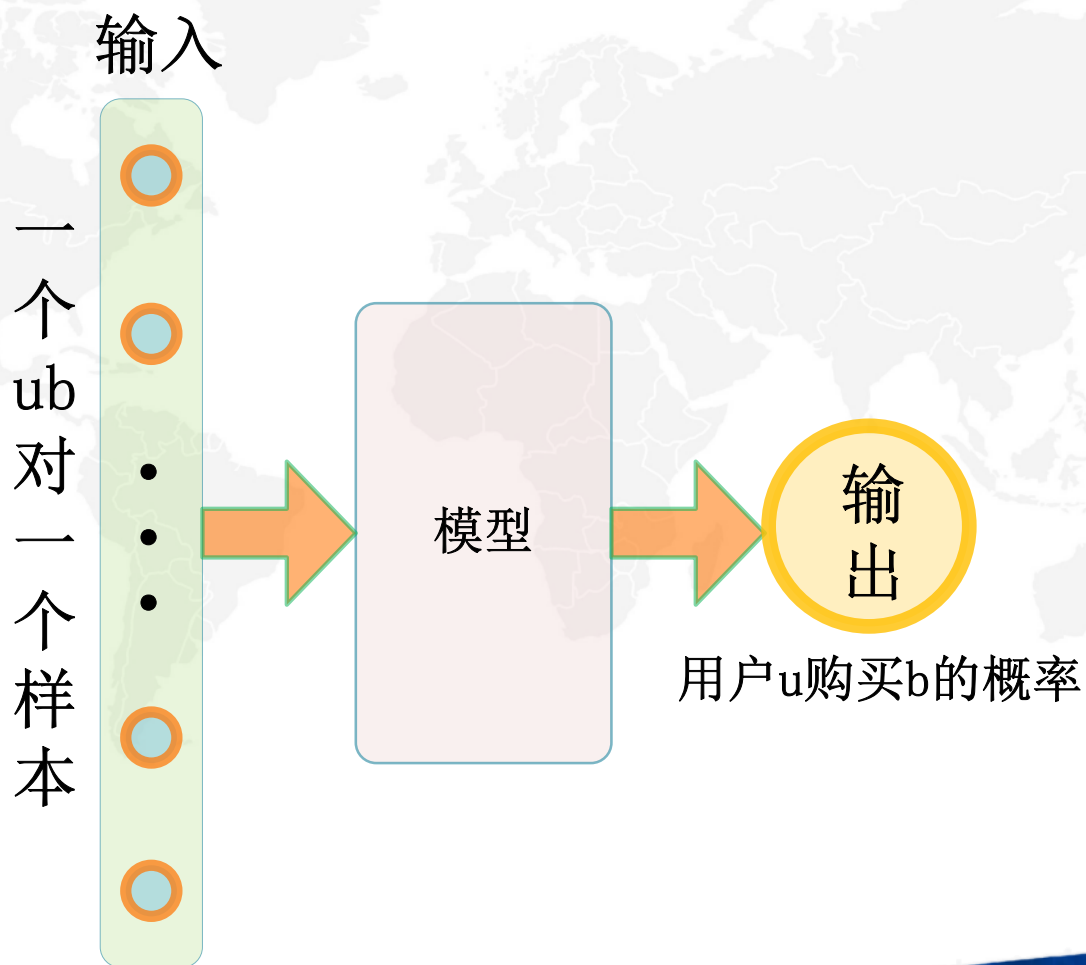
$$F_1 = \frac{2 * P * R}{P + R}$$



根据用户4个月在天猫的行为日志，预测用户u在将来一个月**是否**会购买这个品牌b——**分类问题**

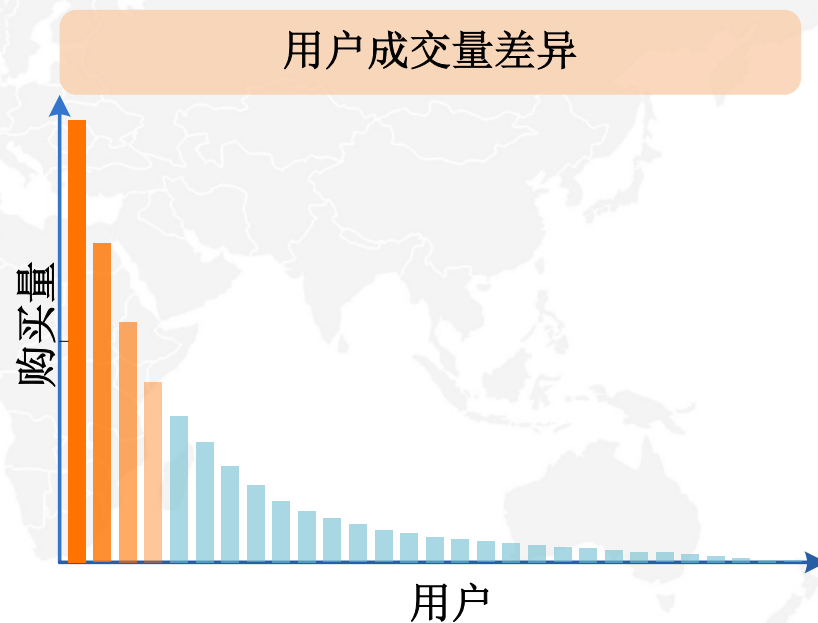
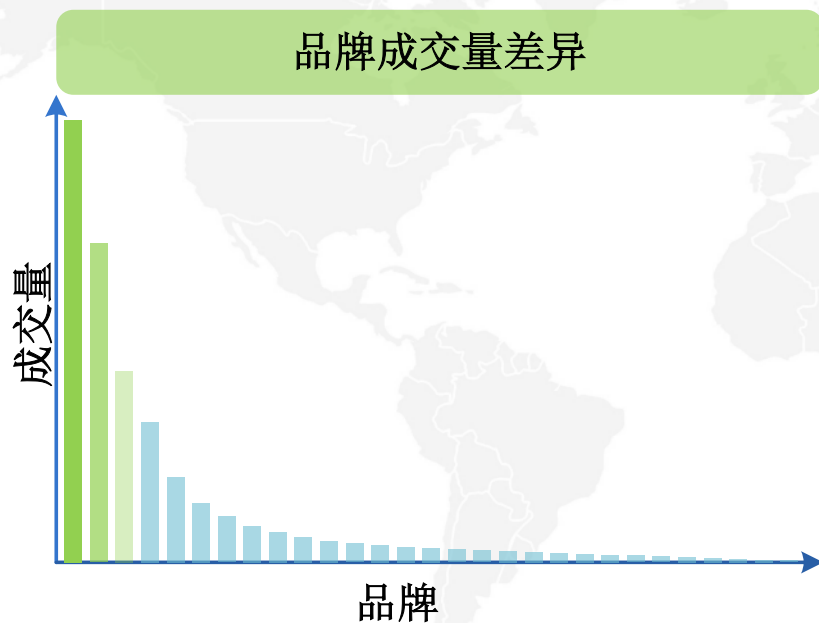
1 解题思路

■ 解题思路



1 解题思路

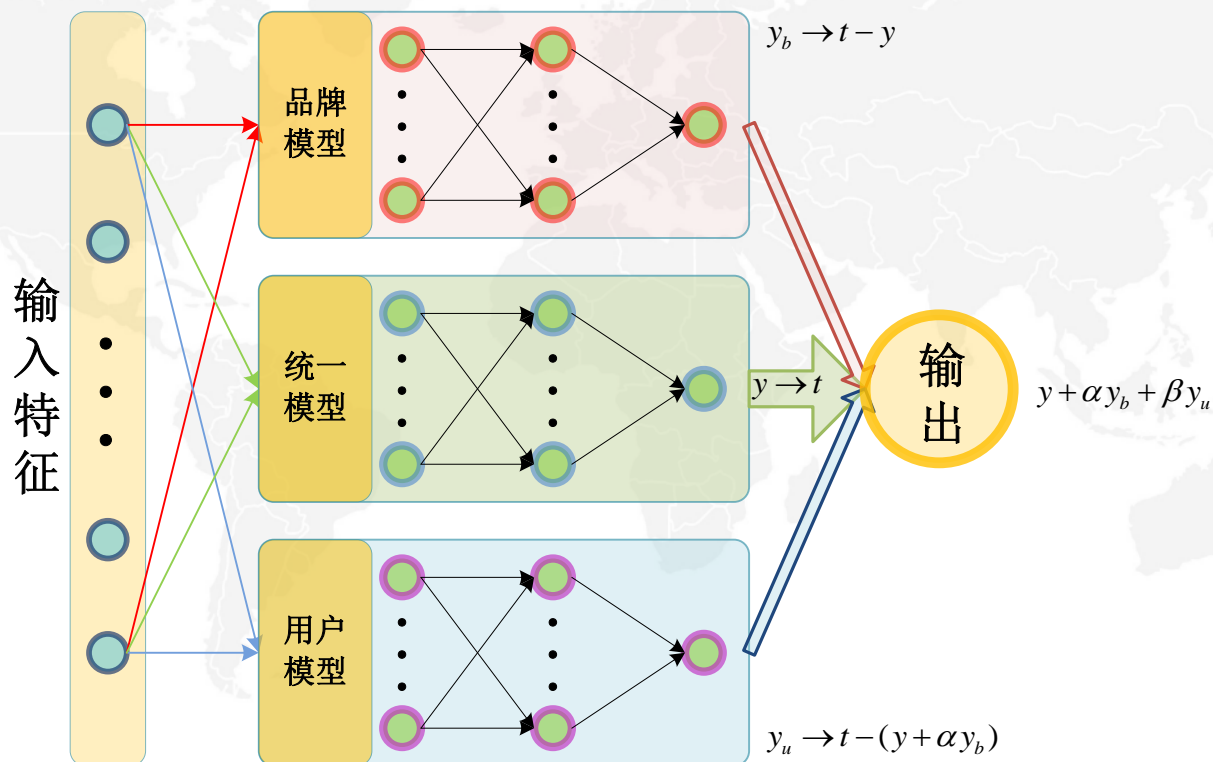
用户和品牌的类内差异



物以类聚，人以群分

1 解题思路

三位一体架构



目录

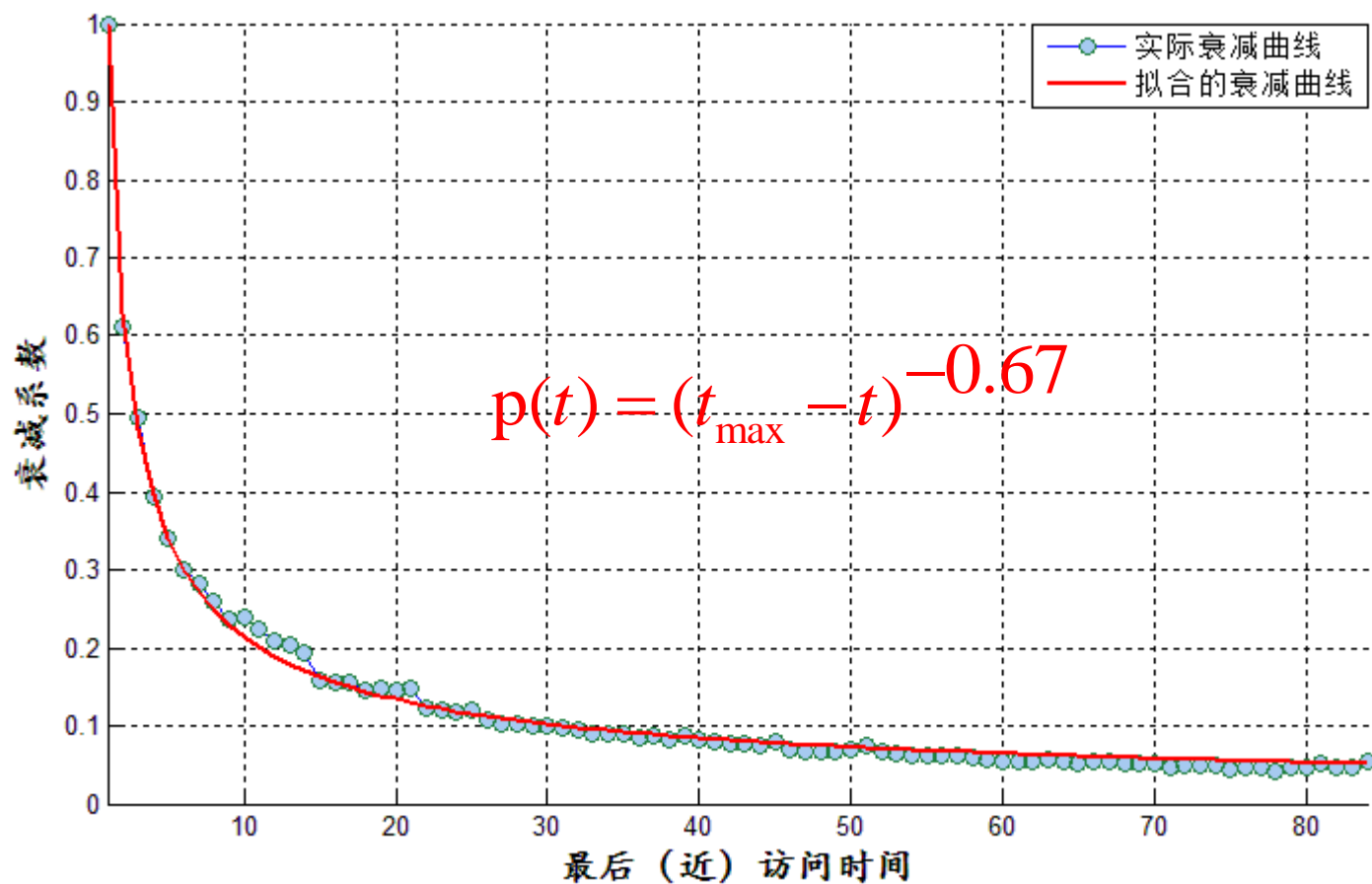
1 解题思路

2 特征提取

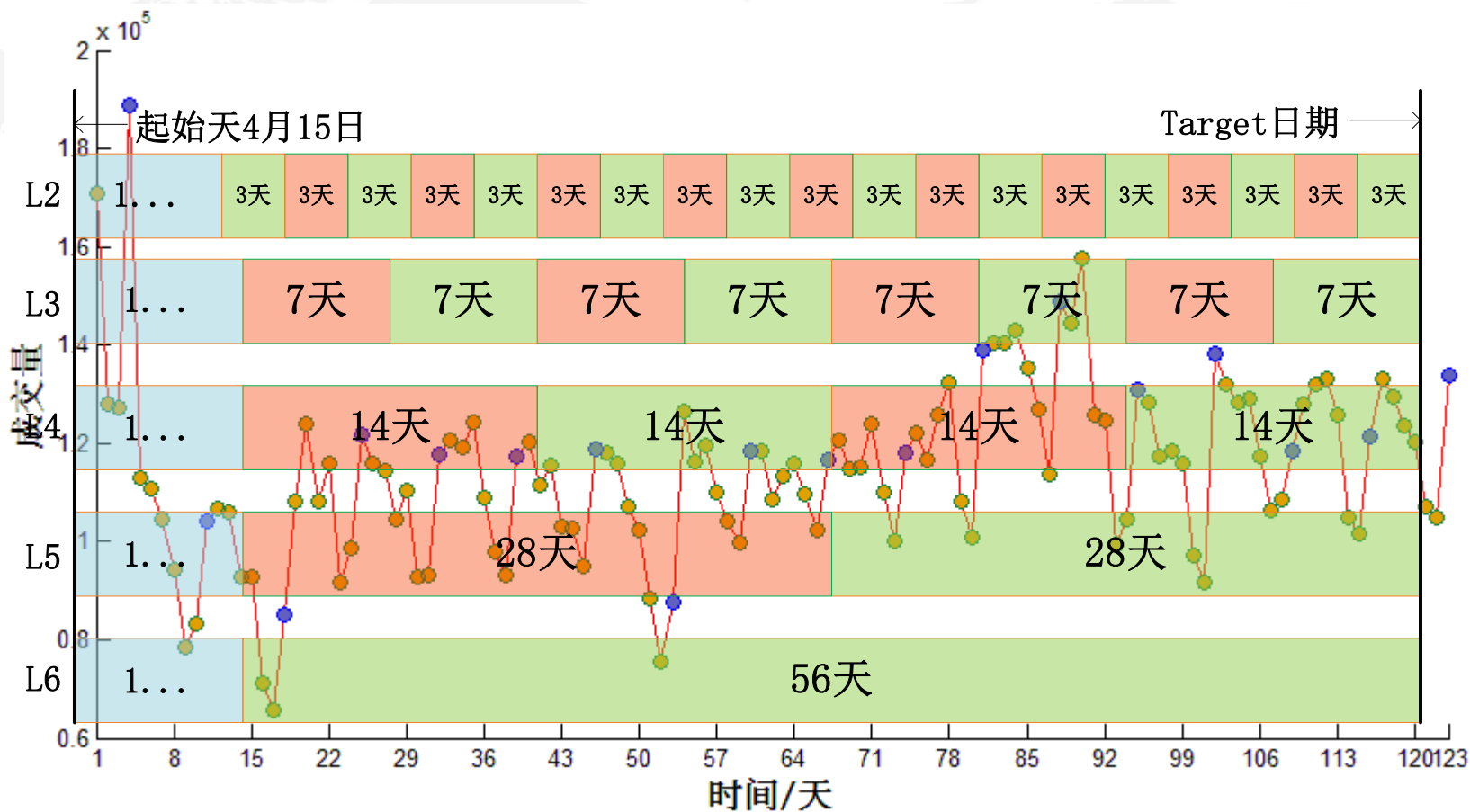
3 算法实现

4 总结回忆

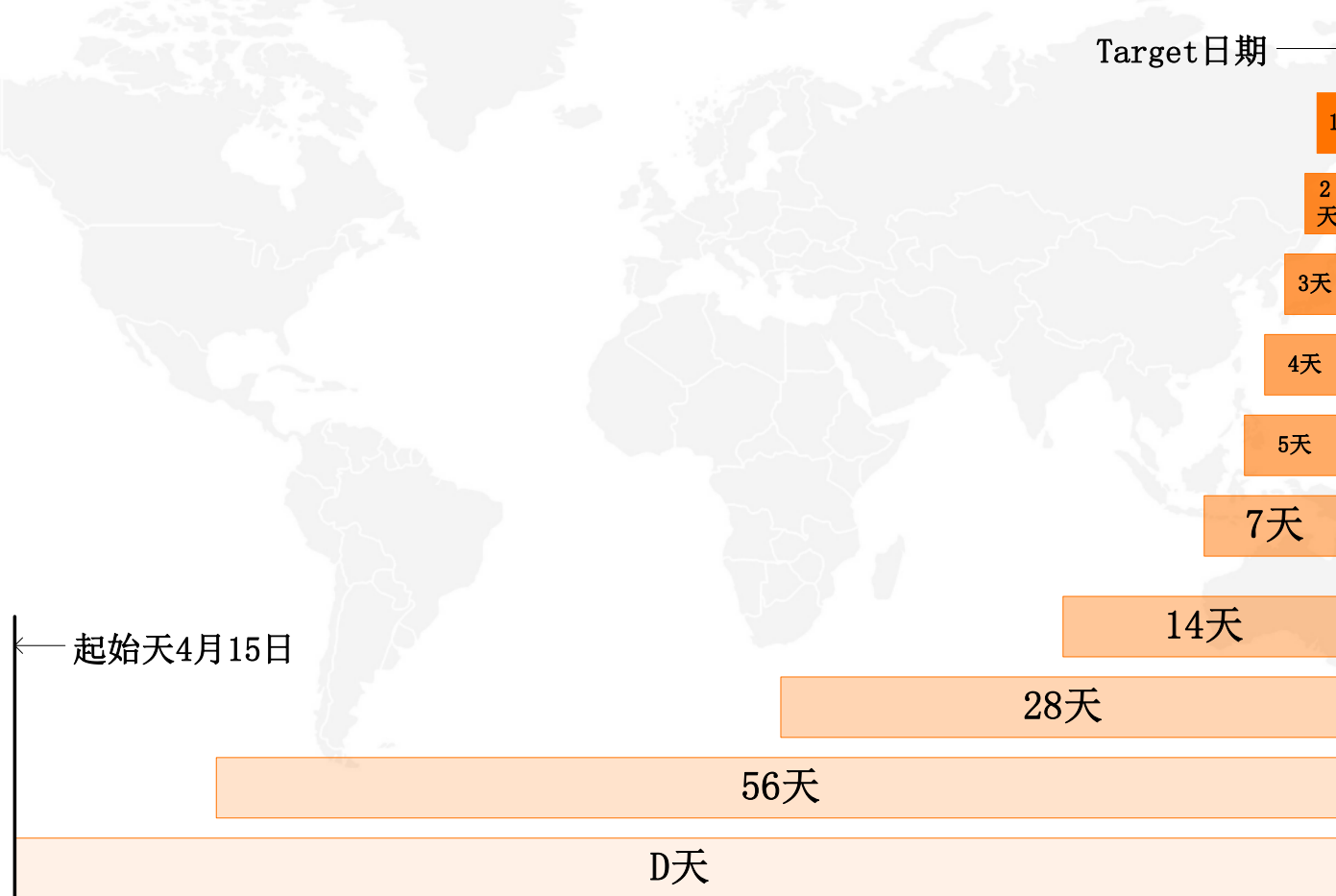
指导思想——源于对数据分析和业务理解



提取方式——层级提取



■ 提取方式——最近K天0/1提取



■ 特征类型



行为特征

1 解题思路 2 特征提取算法实现 3 总结回忆

行为特征——描述用户行为

最后访问时间

$$g(t, \sigma) = \exp\left(-\frac{(t_{\max} - t)^2}{2\sigma^2}\right)$$

$$\sum_{a=0}^3 w_a \sum_{t=1}^D p(t) \log(x_a(t) + 1)$$

平均访问间隔

层级行为次数

层级行为天数

对数特征

$$\sum_{t=1}^D p(t) x_a(t)$$

品牌特征

1 解题思路 2 特征提取算法实现 3 特征工程 4 总结回忆

品牌特征——描述品牌属性

层级四种行为次数

重复购买率

最近K天频率特征

转化率

老顾客比例

购买前访问次数

成交量变化趋势

用户特征

1 解题思路 2 特征提取算法实现 3 模型训练 4 总结回忆

用户特征——描述用户属性

四种行为次数

购买频率

用户行为天数

用户转化率

用户等级

总购买品牌数

平均购买前访问天数

交叉特征

1 解题思路 2 特征提取算法实现 3 特征工程 4 总结回忆

交叉特征

行为品牌交叉

访问次数 \times 用户转化率

行为用户交叉

访问次数 \times 品牌转化率

品牌圈特征

U-B对点击次数/用户总点击次数

品牌用户交叉

目录

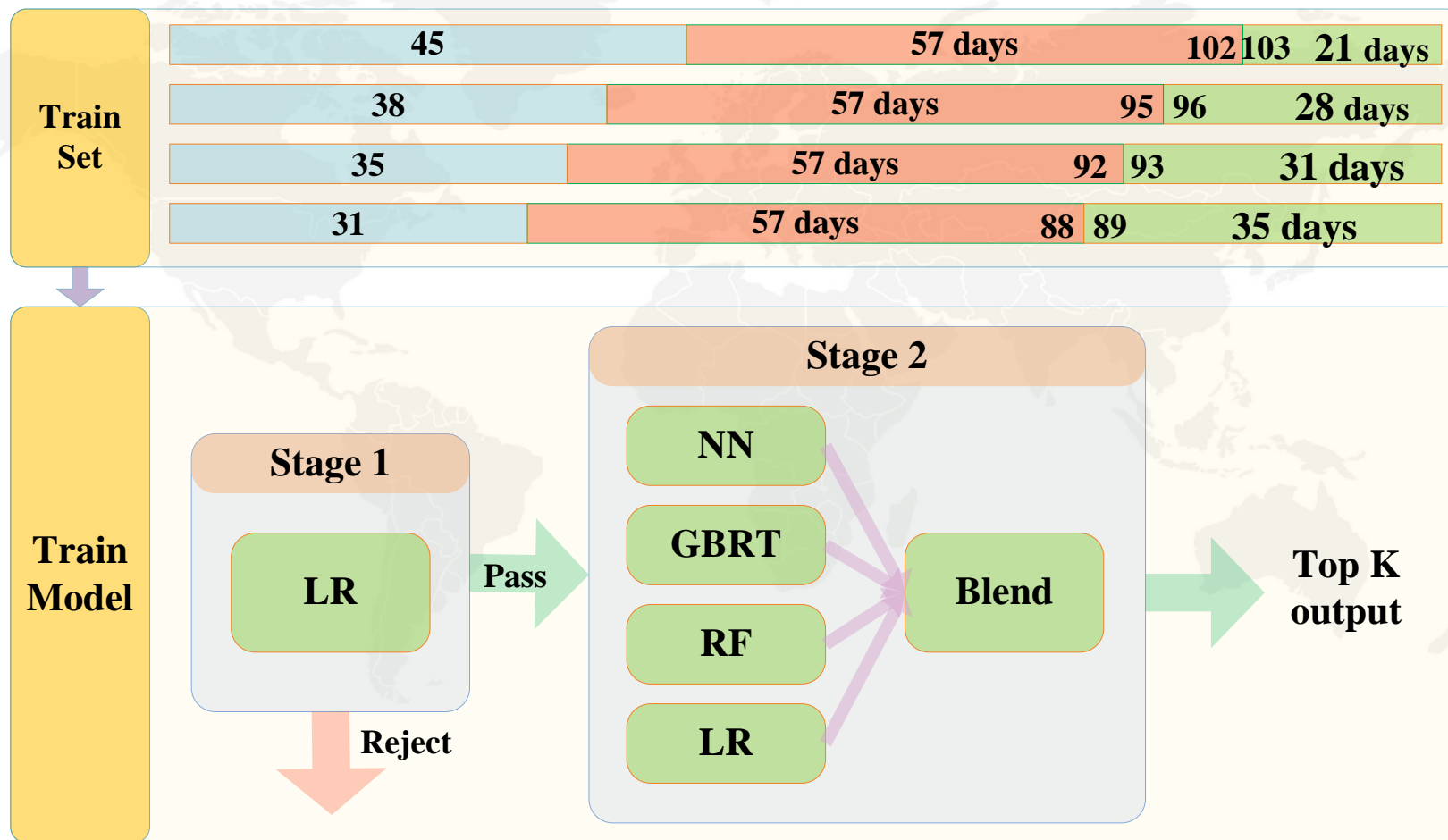
1 解题思路

2 特征提取

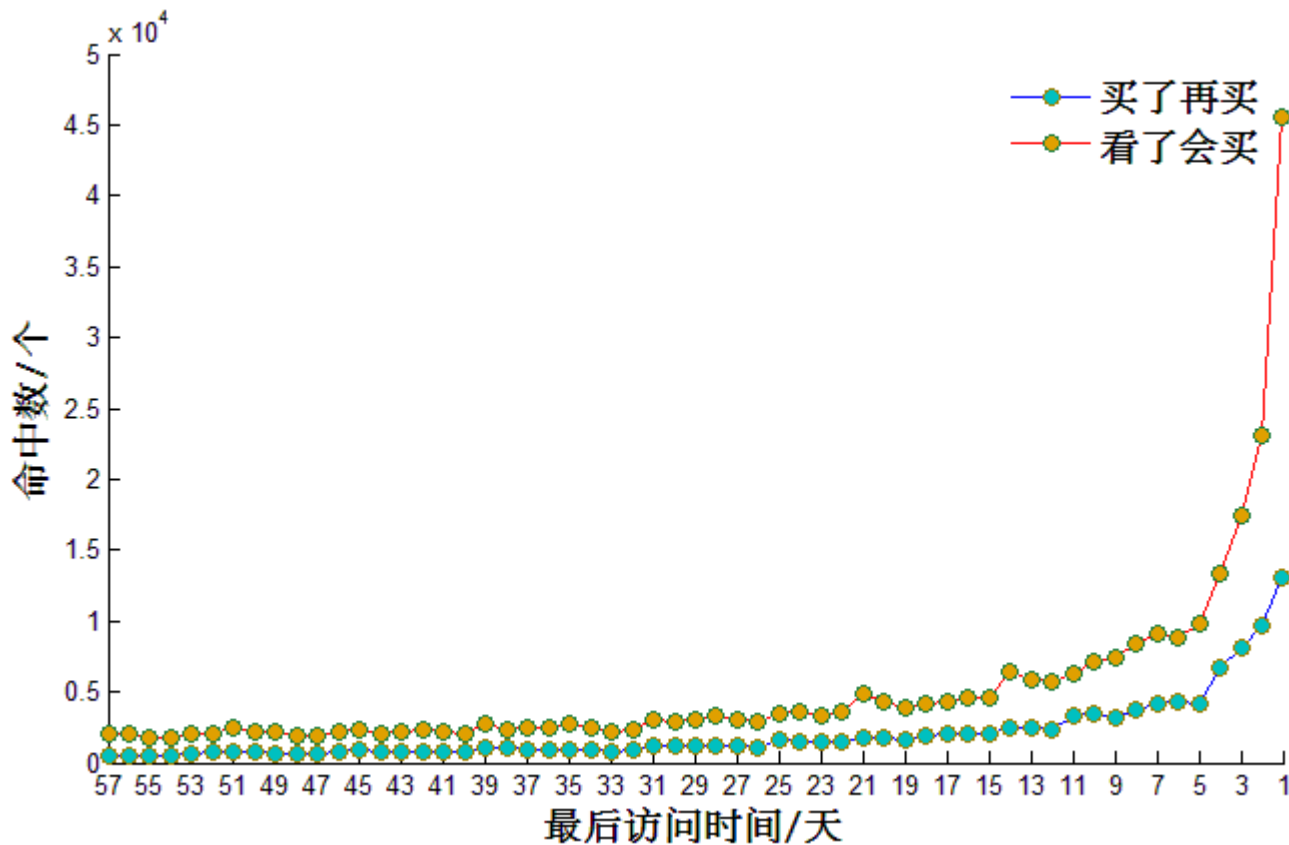
3 算法实现

4 总结回忆

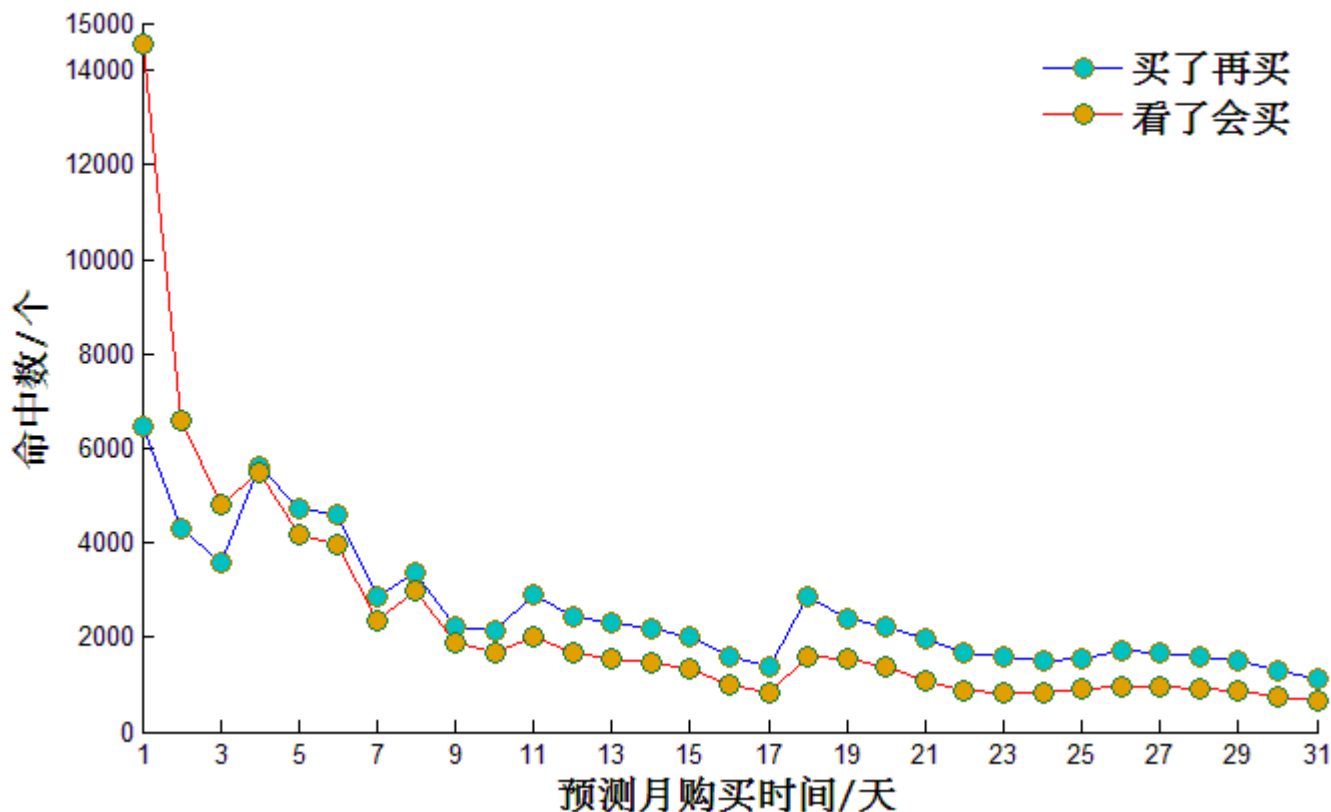
■ 算法整体框架



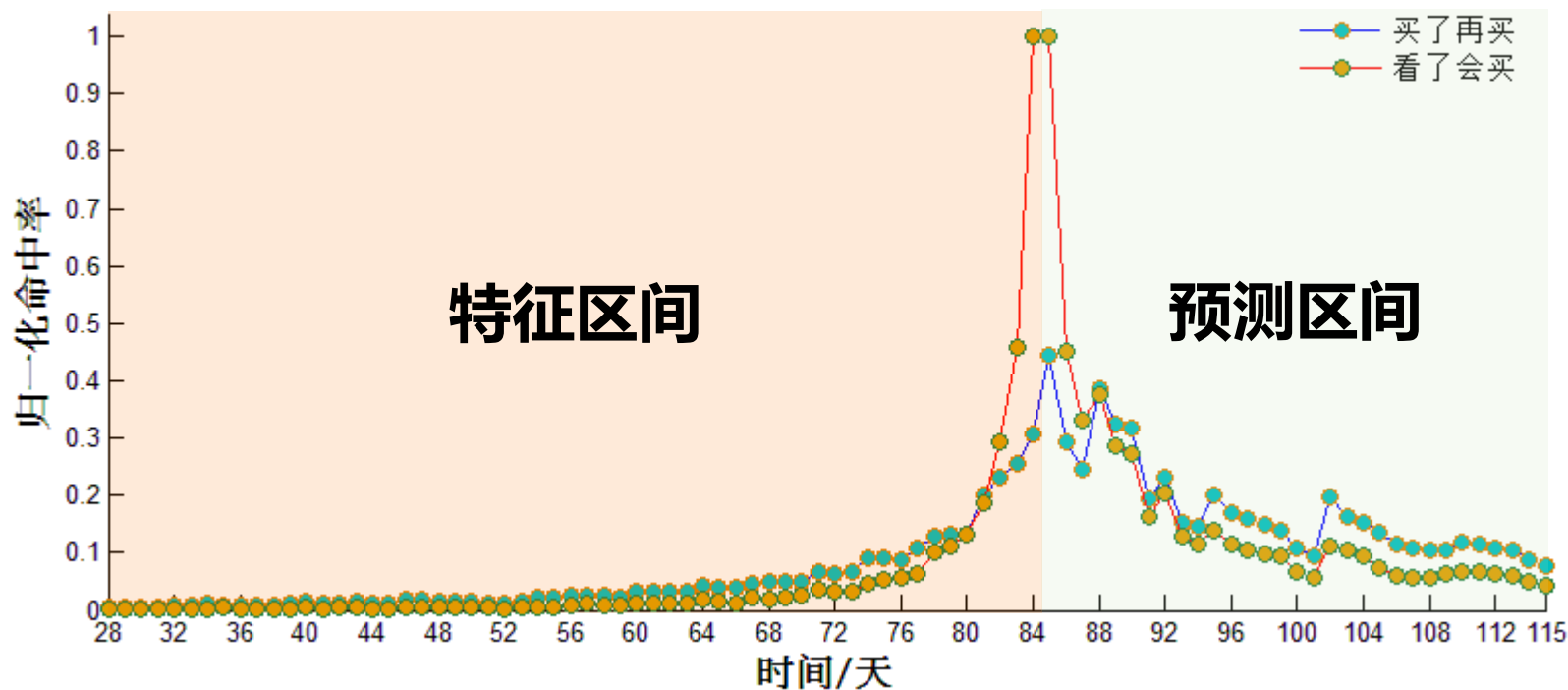
不同行为命中的时间分布



不同行为 模型命中时间分布



模型的预测命中分布



不同行为的差异，算法怎么解决？

模型命中量与时间节点相关，如何选择最优？

■ 不同行为单独训练

不同行为的差异，算法怎么解决？

买了再买

看了会买

模型命中量与时间节点相关，如何选择最优？

线上
预测集

66

57天

训练集

45

57天

102

103

21天

38

57天

95

96

28天

35

57天

92

93

31天

31

57天

88

89

35天

线下
测试集

Test
set

27

57天

84

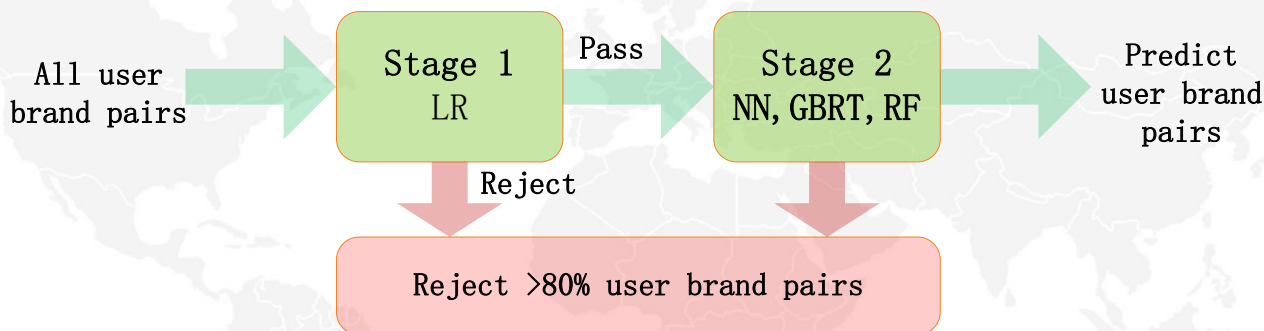
85

39天

target提取区间

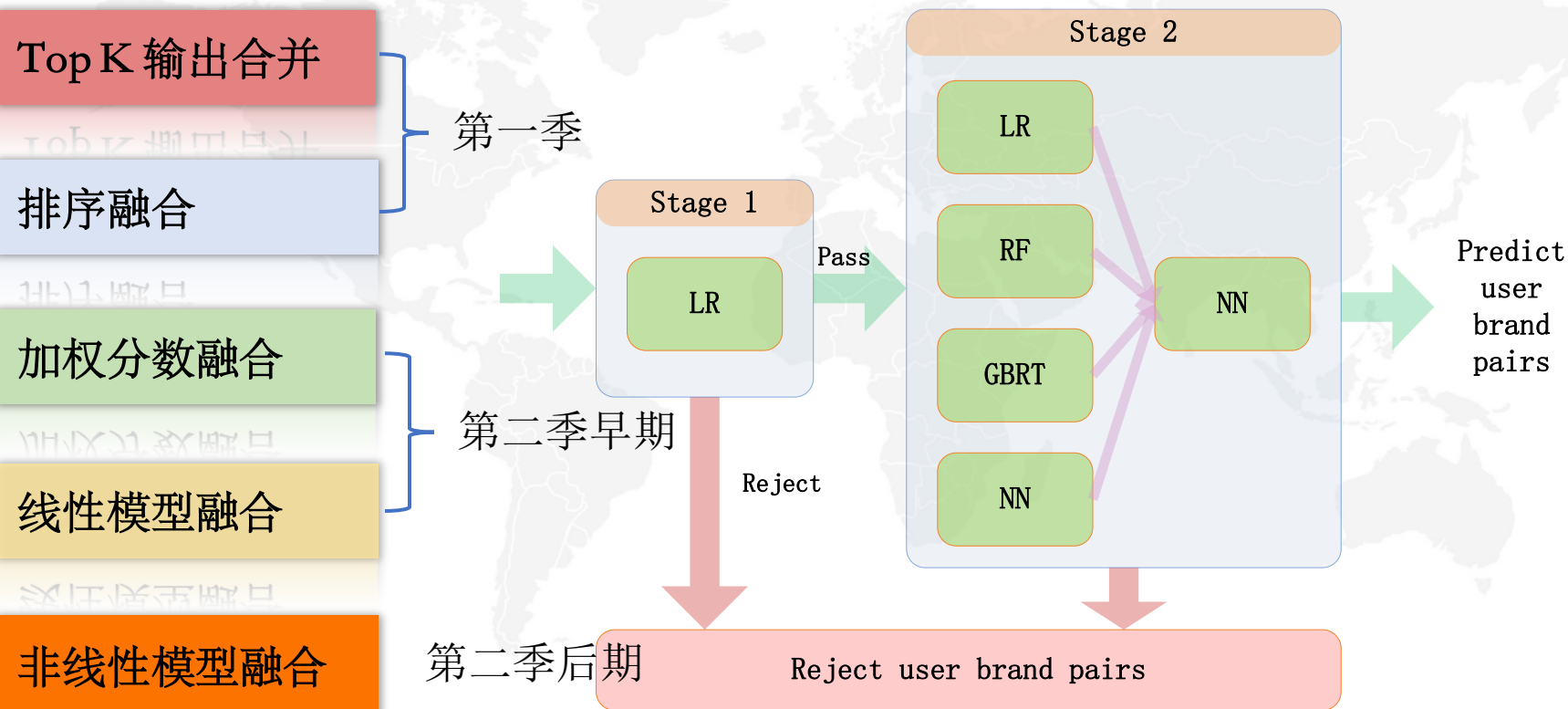
■ 如何解决大海捞针问题？

1亿多 vs 40多万

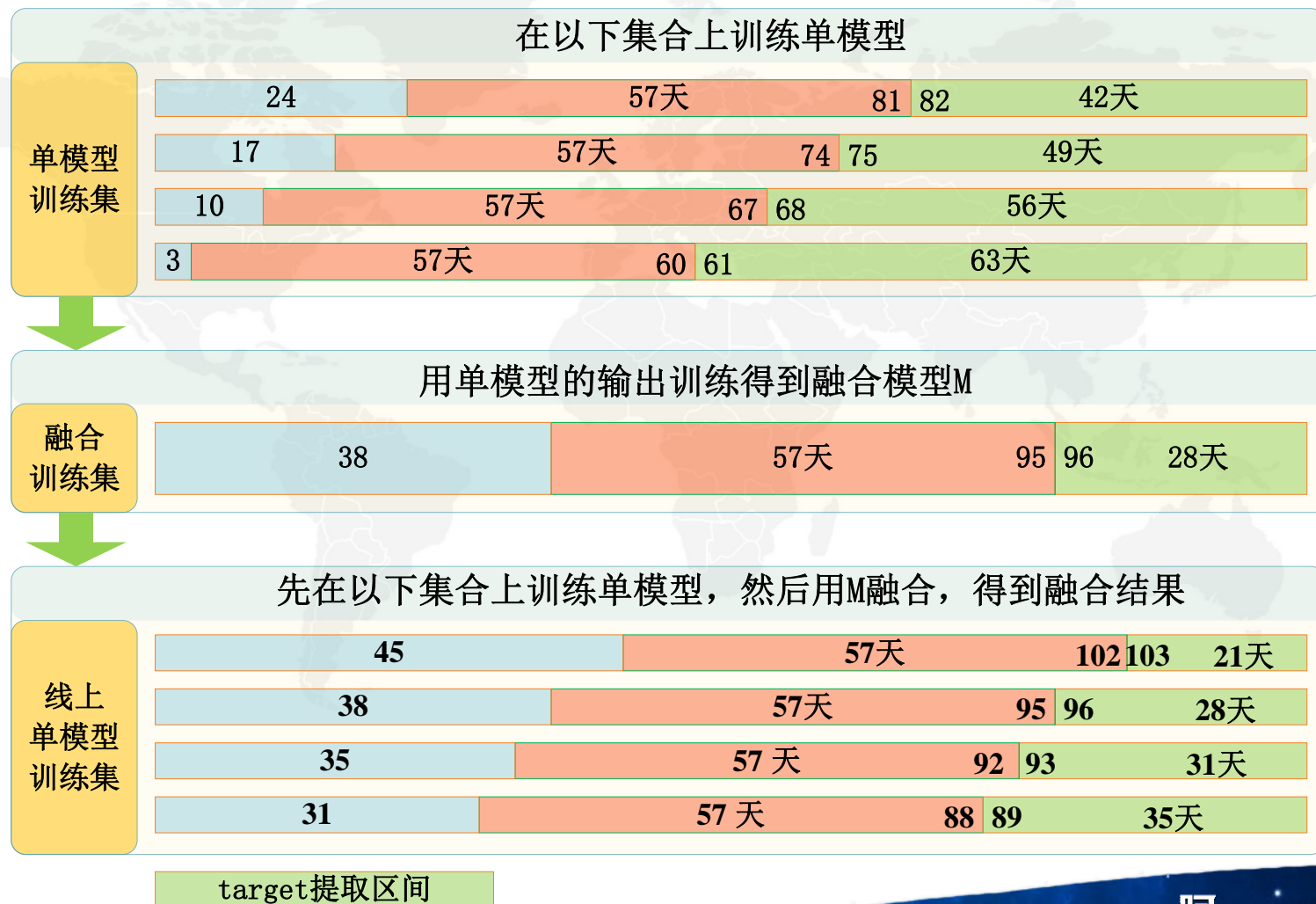


- 通过快速的快速的LR可以滤除>80% 的样本，节约训练预测时间
- 还有助于提高F1分数

模型融合方法



多集合非线性融合



Thank you!

please give some questions

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