

2018年 中国高校计算机大赛 ——大数据挑战赛

团队:今我来思



1 团队介绍

2 问题描述

3 框架设计

4 数据分析

5 特征工程

6 算法及模型融合

7 总结

01

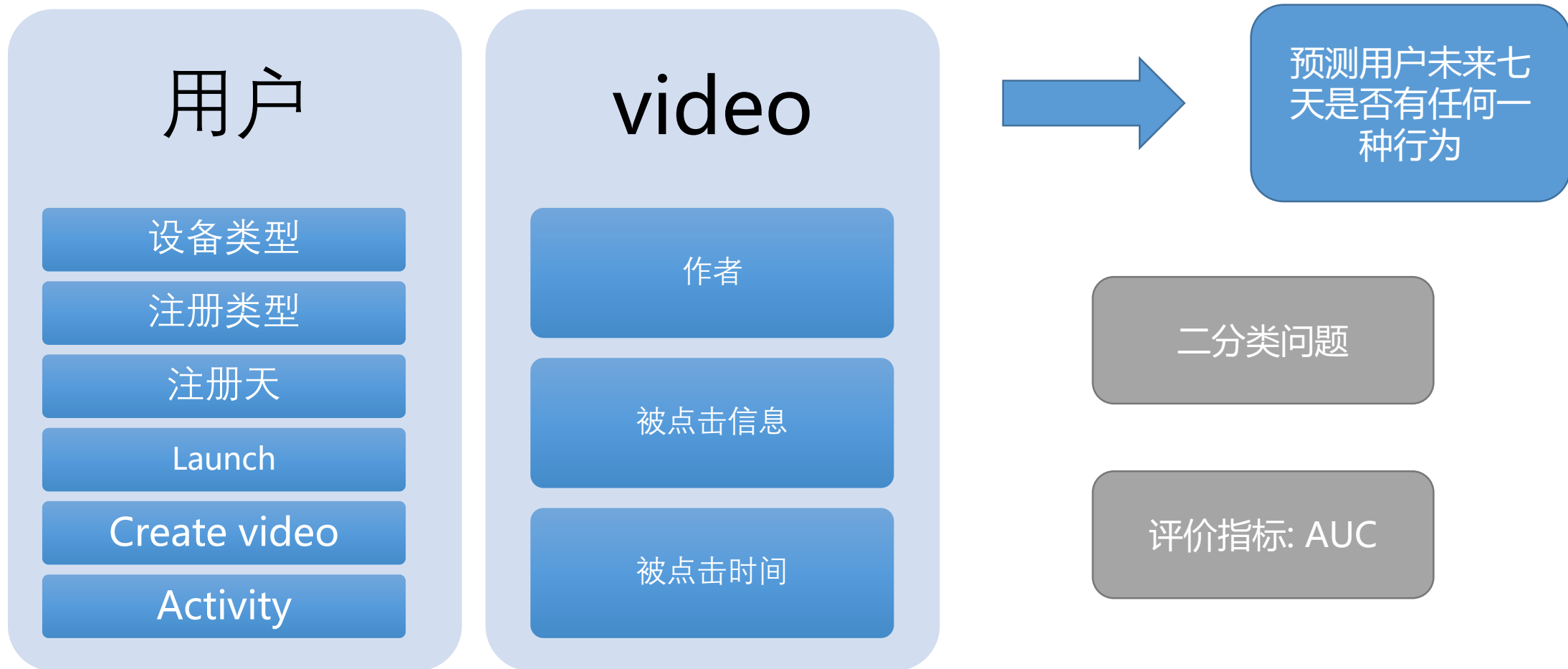
团队介绍

上分历程-破釜沉舟



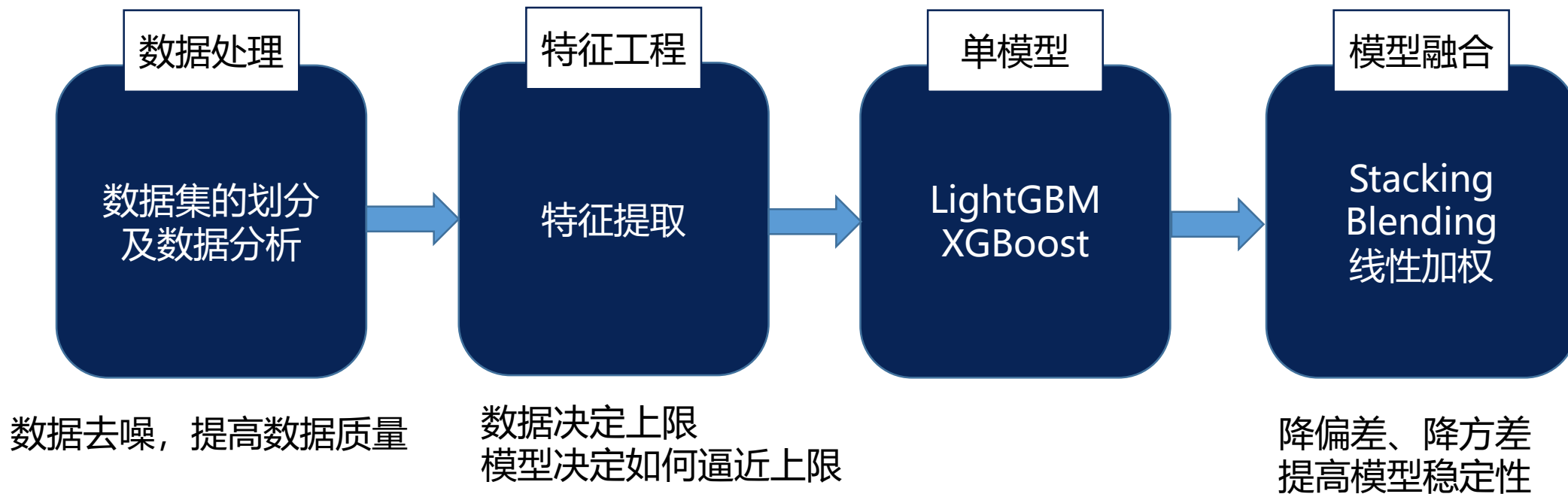
02

问题描述



03

解题思路

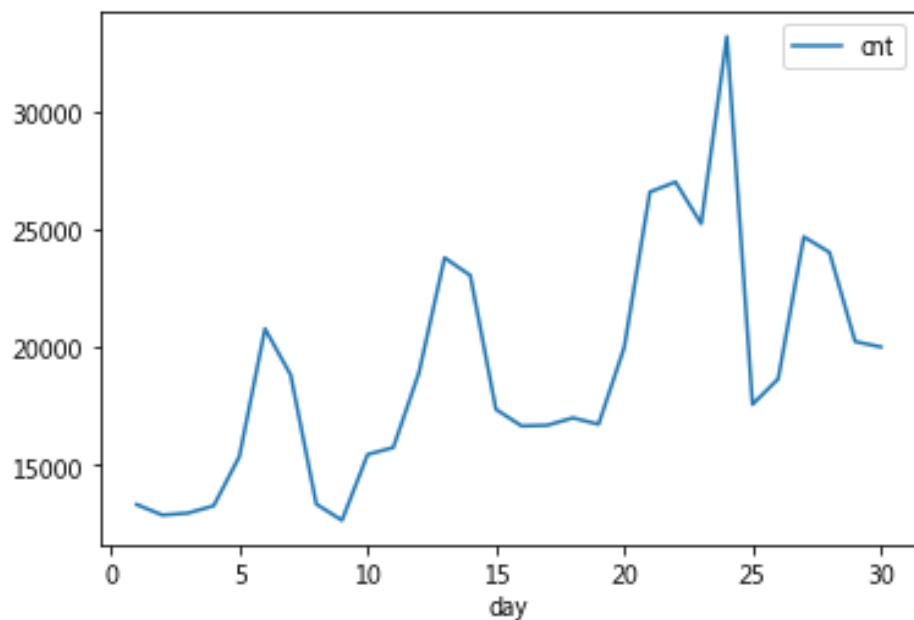


04

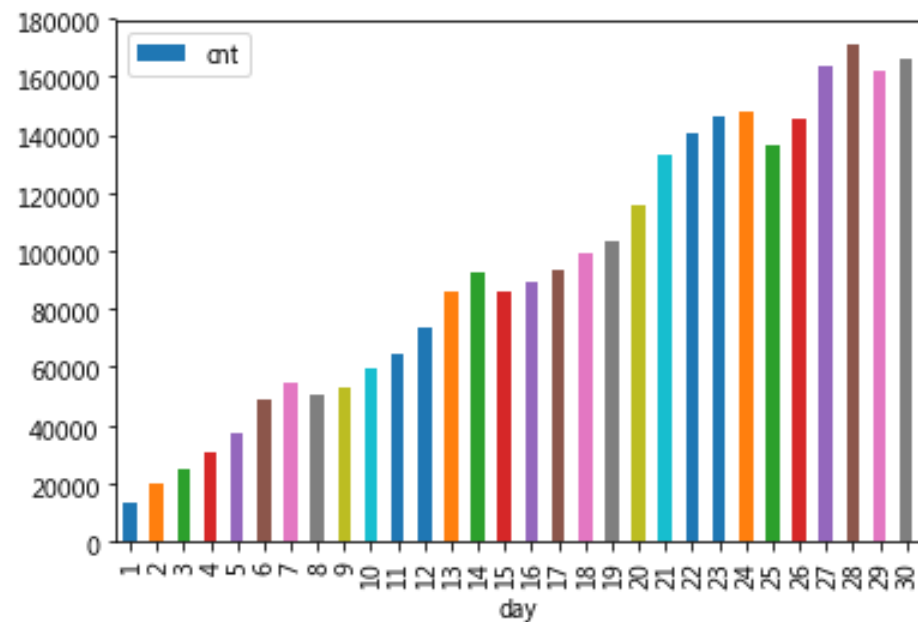
数据分析

4.1 数据清洗

- 有activity,必须有登录, 用activity表补全launch
修复了0.022%的数据
- Activity 数据去重
- 异常值处理

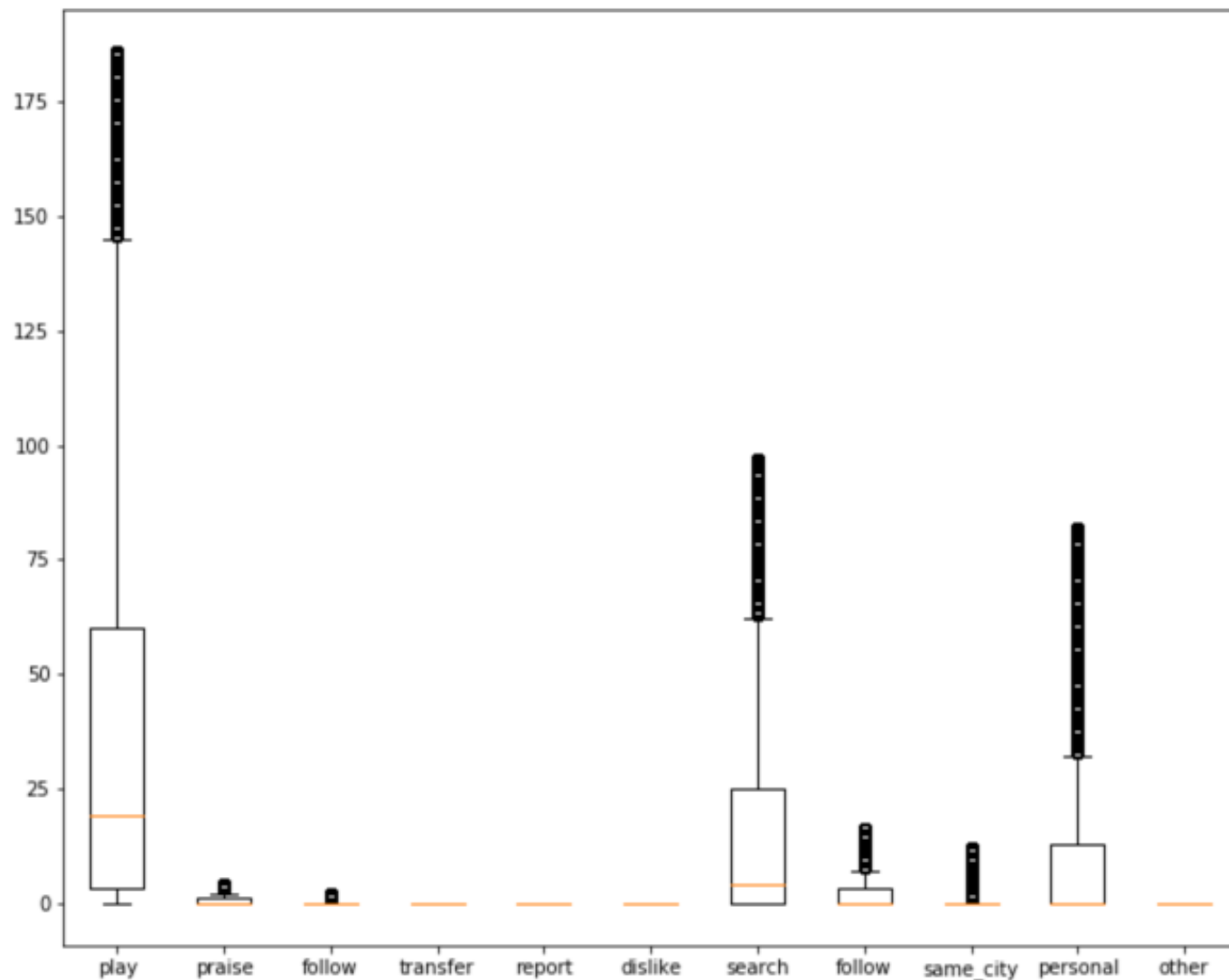


注册次数周期性



用户活动量统计图

activity	Cover_rate
play	0.8618
praise	0.3500
follow	0.2384
transfer	0.0668
report	0.0003
dislike	0.0019
search	0.6502
follow	0.4715
Same_city	0.3118
personal	0.5368
other	0.0838



Activity (include 0)

Dislike、Report行为很少，如何处理？

我们认为，除了播放行为，其他行为都是带有感情色彩的，因此可以直接相加。

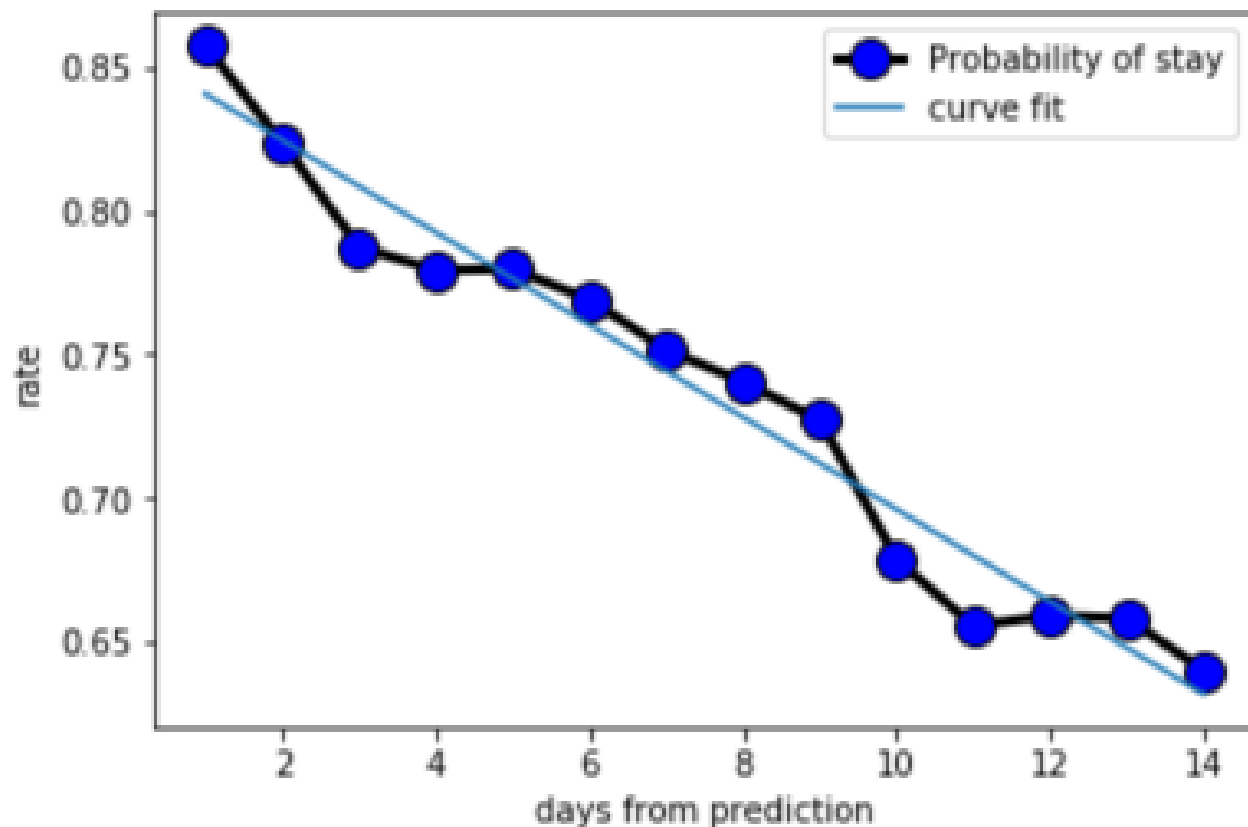
从感情色彩的强度分析：

举报 > 转发 > 喜欢 == 不喜欢 > 播放

$\text{Emotion_activity} = \text{点赞} + \text{不喜欢} + 2 * \text{转发} + 3 * \text{举报}$

activity	Cover_rate
Emotion_activity	0.4436





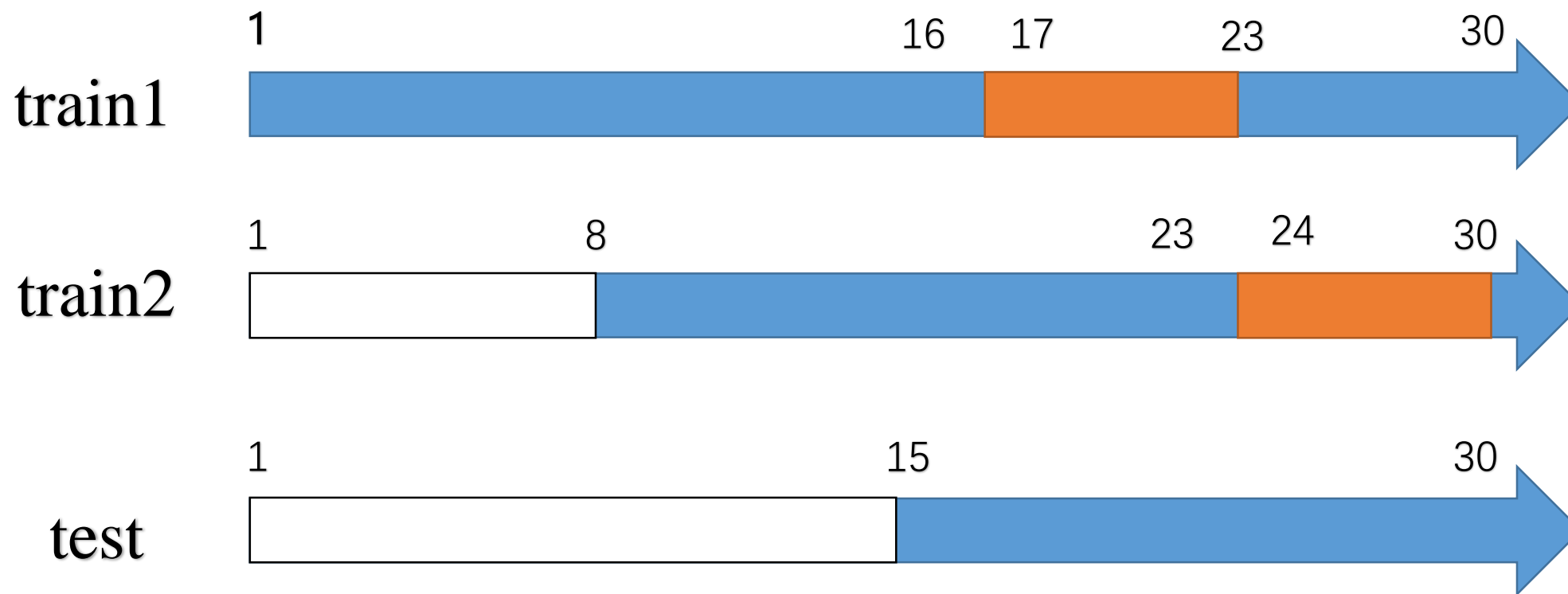
$$W(day) = -0.01606279 * DFP + 0.85668832$$

$$Score = \sum_{i=1}^t W(i) * info(i)$$

05

特征工程

5.1 数据集划分



5.2 特征挖掘

全局数据特征群

全局角度

从全局不同角度描述用户习惯。
包含用户信息特征、相关统计特征等。

行为序列特征群

用户行为角度

由用户出发，用细粒度的时间窗口描述行为的变化情况

Feat1

Feat3

Feat2

Feat4

时间衰减特征群

业务角度

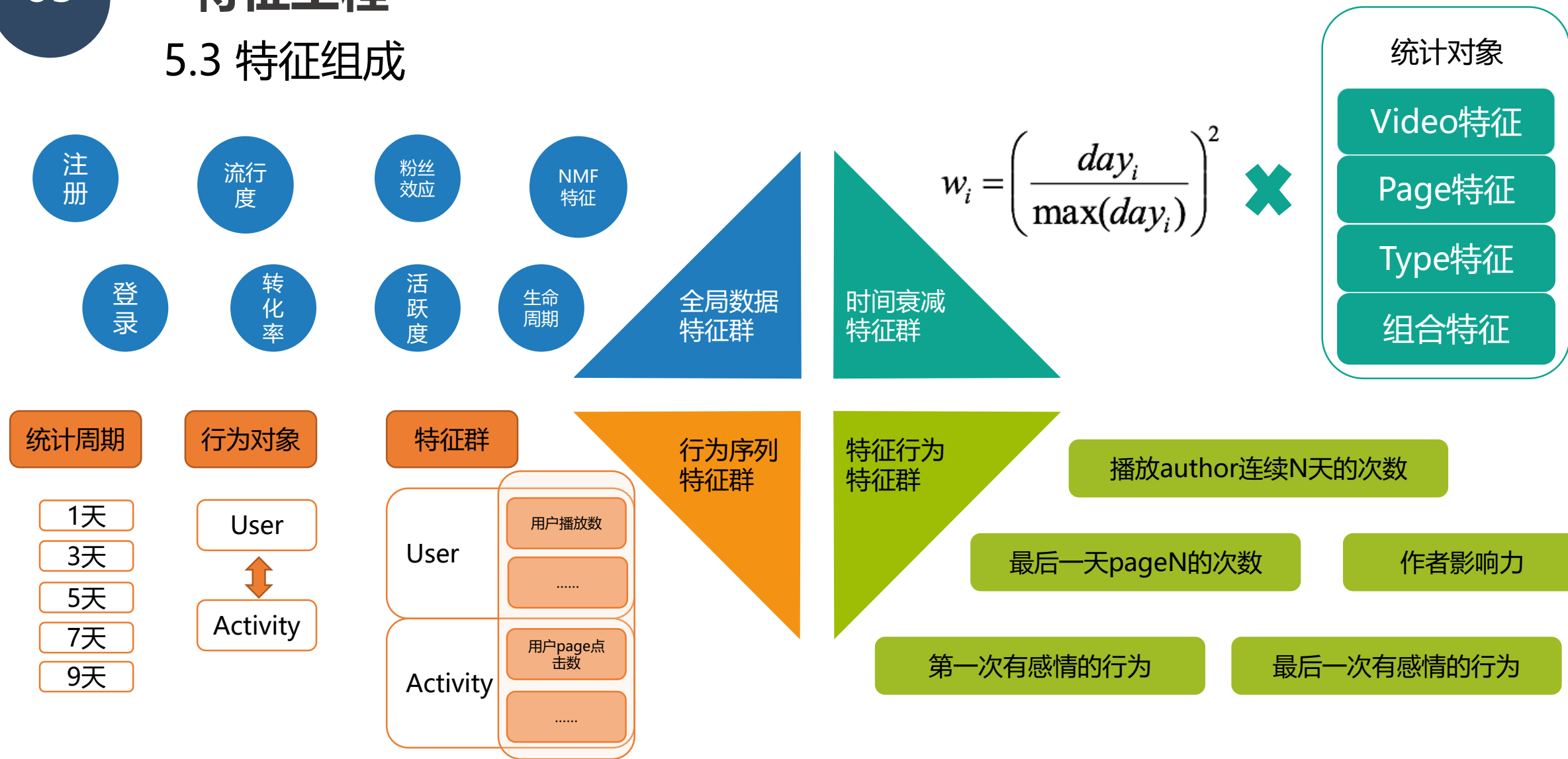
考虑用户在不同时间上相同行为的不同含义，给予不同的权重表示。

特殊行为特征群

用户业务角度

考虑用户在特征行为，如粉丝属性，看某个author的连续天数、次数等。

5.3 特征组成



标签最后七天的权重真的一样吗？

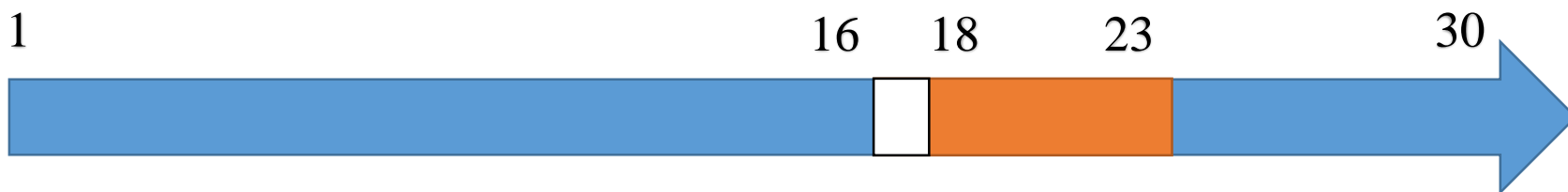


1 0 0 0 0 0 0



0 0 0 0 0 0 1

train1



train2



06

算法及融合模型

stacking

Xgb/Lgb/GBDT/RF/LR

Train data
5 Fold

Build New Model

New train
data

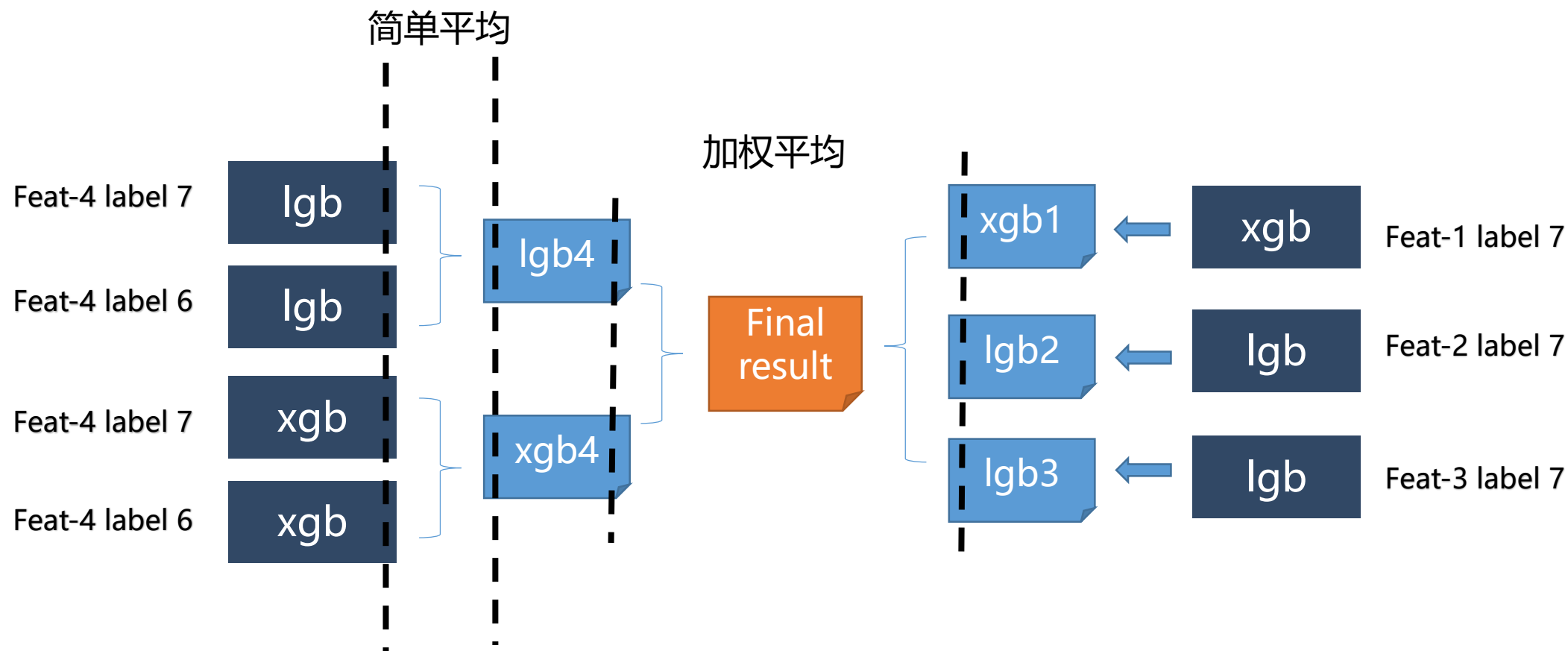
Label



X 5



算法及融合模型



为什么线性融合效果好？

	xgb1	lgb2	xgb4		xgb1	lgb1
xgb1	1	0.9905	0.9933	<	0.9982	
lgb2	0.9905	1	0.9922			
xgb4	0.9933	0.9922	1			

$$\rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y} = \frac{E((X-\mu_X)(Y-\mu_Y))}{\sigma_X \sigma_Y} = \frac{E(XY) - E(X)E(Y)}{\sqrt{E(X^2) - E^2(X)} \sqrt{E(Y^2) - E^2(Y)}}$$

07

总结

总结

- ✓ 对数据进行了充分的分析
- ✓ 模型多样性、特征多样性
- ✓ 融合方法适当

遗憾：

- I. 融合方式还有很多没尝试。如rank融合，Blending等。
- II. 没有充分挖掘用户行为 可以把user_id与authorid/videoid矩阵进行矩阵分解，获得用户兴趣爱好向量
- III. 神经网络没有尝试 可以把用户每天都行为embedding，进而用LSTM/CNN进行训练。

- ◆ 感谢所有参赛队伍；
- ◆ 感谢各周的周星星的分享；
- ◆ 感谢快手与清华一起举办这样成功的比赛，给了我们学习、锻炼和展示的机会。

谢谢聆听
THANKS!!