

队名：levy

Contents

目录

1 团队介绍

2 赛题理解

2.1 评价指标

2.2 建模分析

3 特征工程

3.1 特征介绍

3.2 特征有效性

4 模型介绍

4.1 通用操作

4.2 旧广告模型

4.3 新广告模型

5 总结与思考



赛题理解

评测指标/赛题建模分析

大赛赛题：广告日曝光预估



1) 准确性指标

01

$$SMAPE = \frac{1}{n} \sum_{t=1}^n \frac{|F_t - A_t|}{(F_t + A_t)/2}$$

2) 出价单调相关性指标

02

$$\text{score} = \frac{1}{n} \sum_{k=1}^n \frac{(imp_0 - imp_k)(bid_0 - bid_k)}{|(imp_0 - imp_k)(bid_0 - bid_k)|}$$

3) 最终得分

03

$$\text{TotalScore} = w_1 * \left(1 - \frac{SMPAE}{2}\right) + w_2 * \frac{\text{MonoScore} + 1}{2}$$

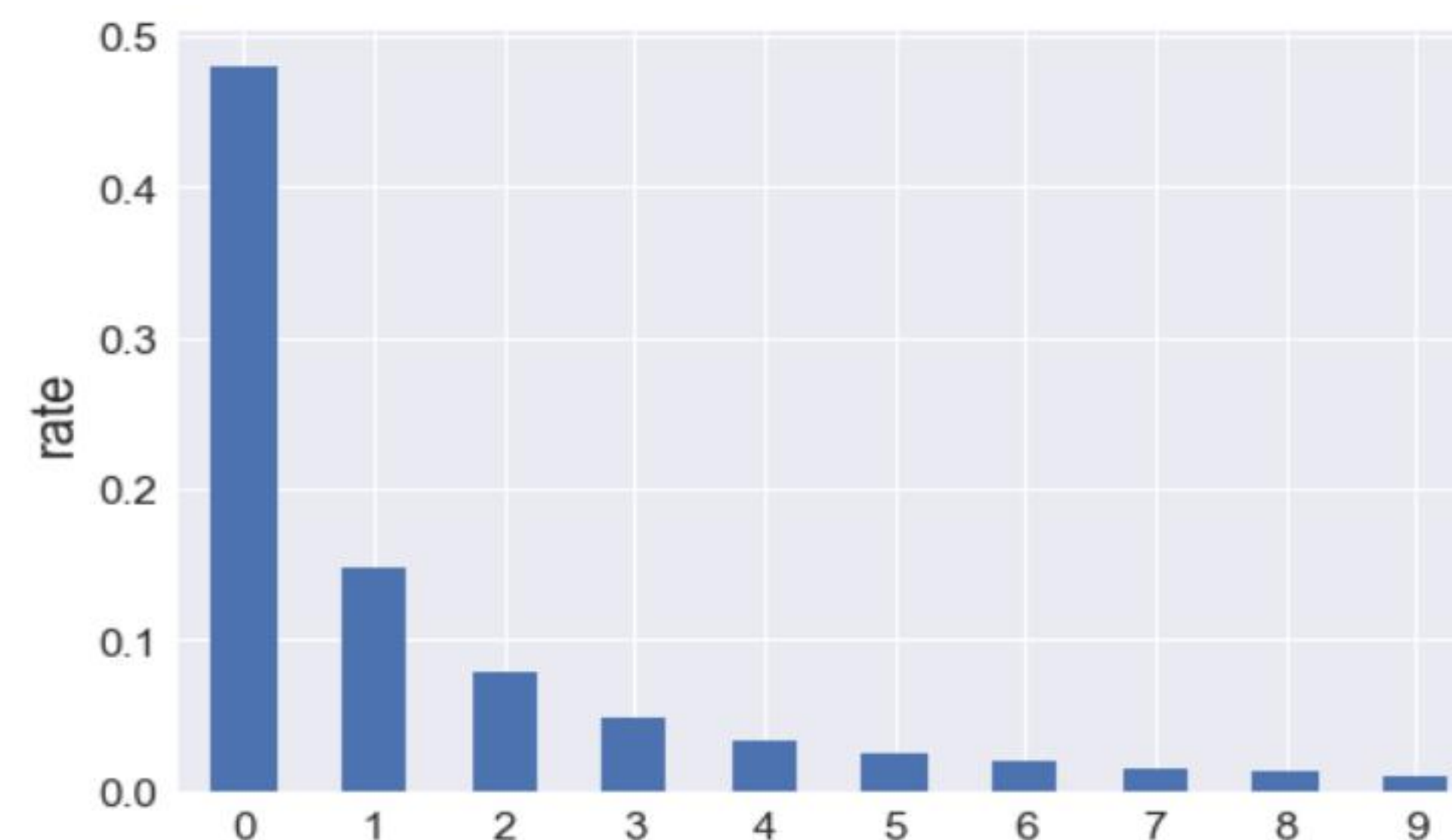
曝光量为非负整数

小曝光量样本占绝大多数

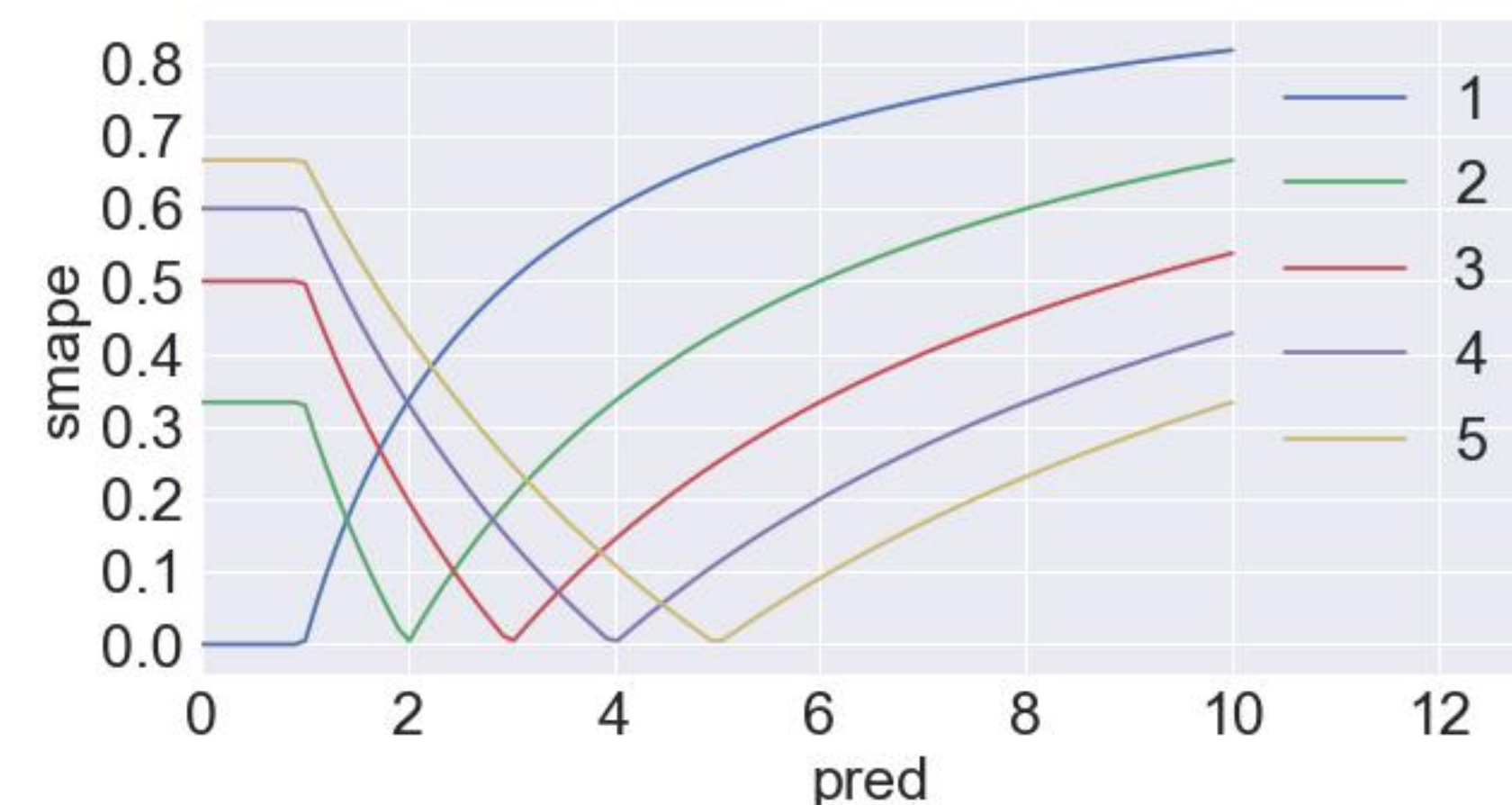
线上对0曝光样本的平滑处理

$$pred \leftarrow \max(pred, 1)$$

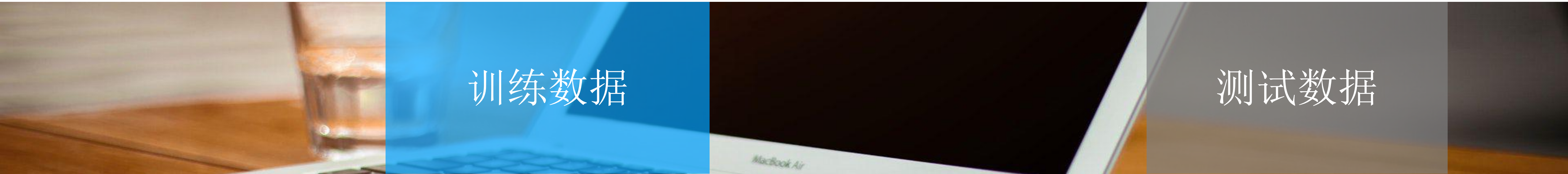
$$label \leftarrow \max(label, 1)$$



不同曝光量样本的训练集占比（长尾）



不同曝光量下smape与预测日曝光量的关系



线下验证集构造

4/10 4/11 4/12 4/15 4/16 4/17

4/10 - 4/15的竞价日志中竞价数>35的广告

4/16 - 4/17的竞价日志(按竞价数分布采样挖去5000+广告
的竞价日志做广告采样)

线上训练数据集

4/10 4/11 4/12 4/22 4/23 4/24

4/10 - 4/22的竞价日志中竞价数>35的广告

4/24参与的竞价队列的广告
当日竞价数>35

建模设想1: Point-wise rank

Basic_ecpm = cpm_bid
= 1000 * cpc_bid * pctr
= 1000 * cpa_bid * pctr * pcvr

Quality_ecpm ≈ 20*pctr

Total_ecpm = Basic_ecpm + Quality_ecpm

曝光



广告id	Rank of Total_ecpm	是否过滤
广告a	1	True
广告b	2	True
广告c	3	True
广告d	4	False
广告e	5	False
广告f	6	False

建模设想1: Point-wise rank

存在的问题:

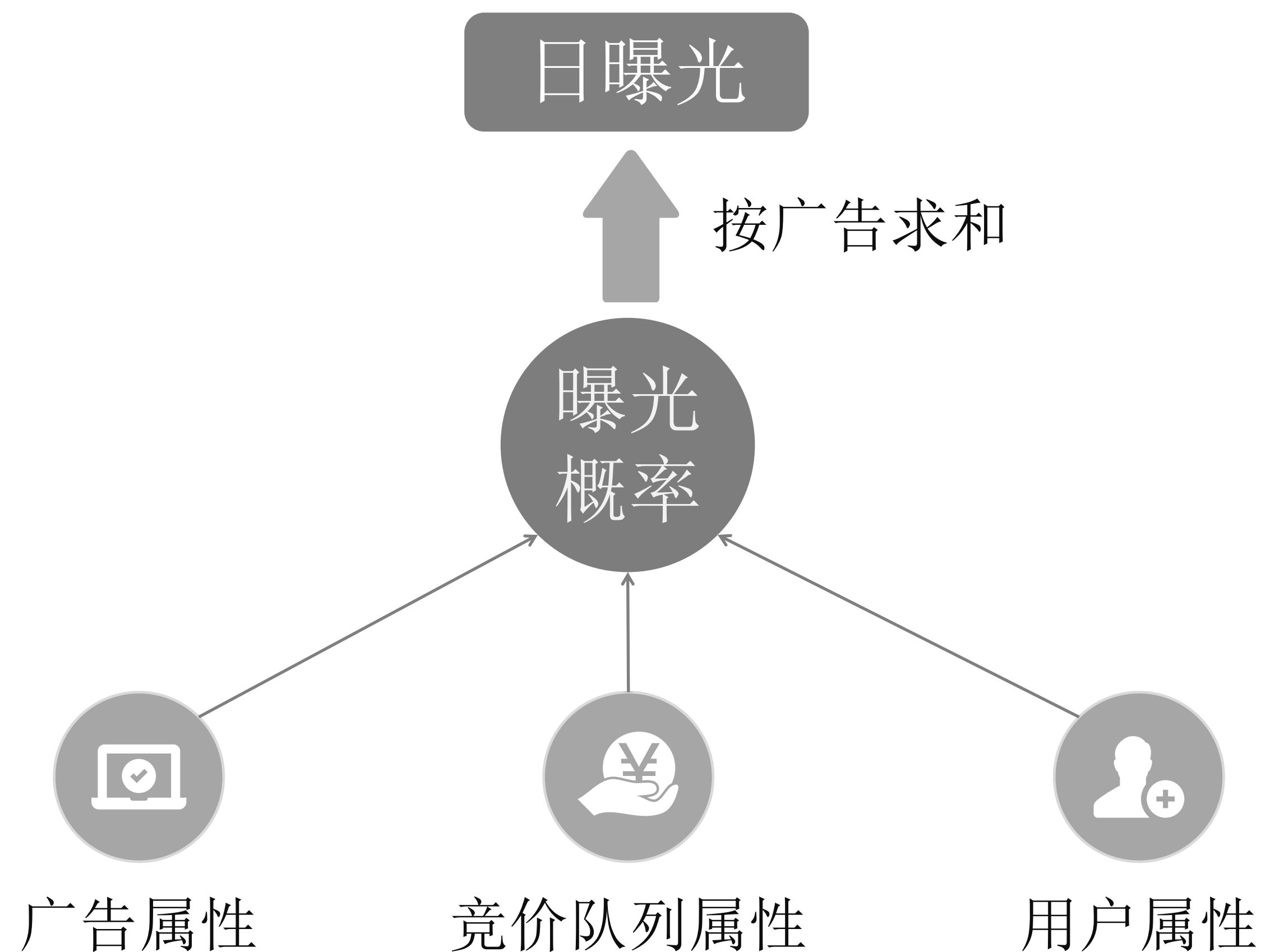
- 1. pctr, pcvr的预测偏差
- 2. 每条测试广告的真实出价隐藏在多条不同的虚拟出价中

广告id	Rank of Total_ecpm	是否过滤
广告a	1	True
广告b	2	True
广告c	3	True
测试广告1	4	True
测试广告2	4	False
测试广告3	4	False
广告d	5	False
广告e	6	False
广告f	7	False

建模设想2：预测某广告的某次 竞价是否曝光（2分类）

存在的问题：

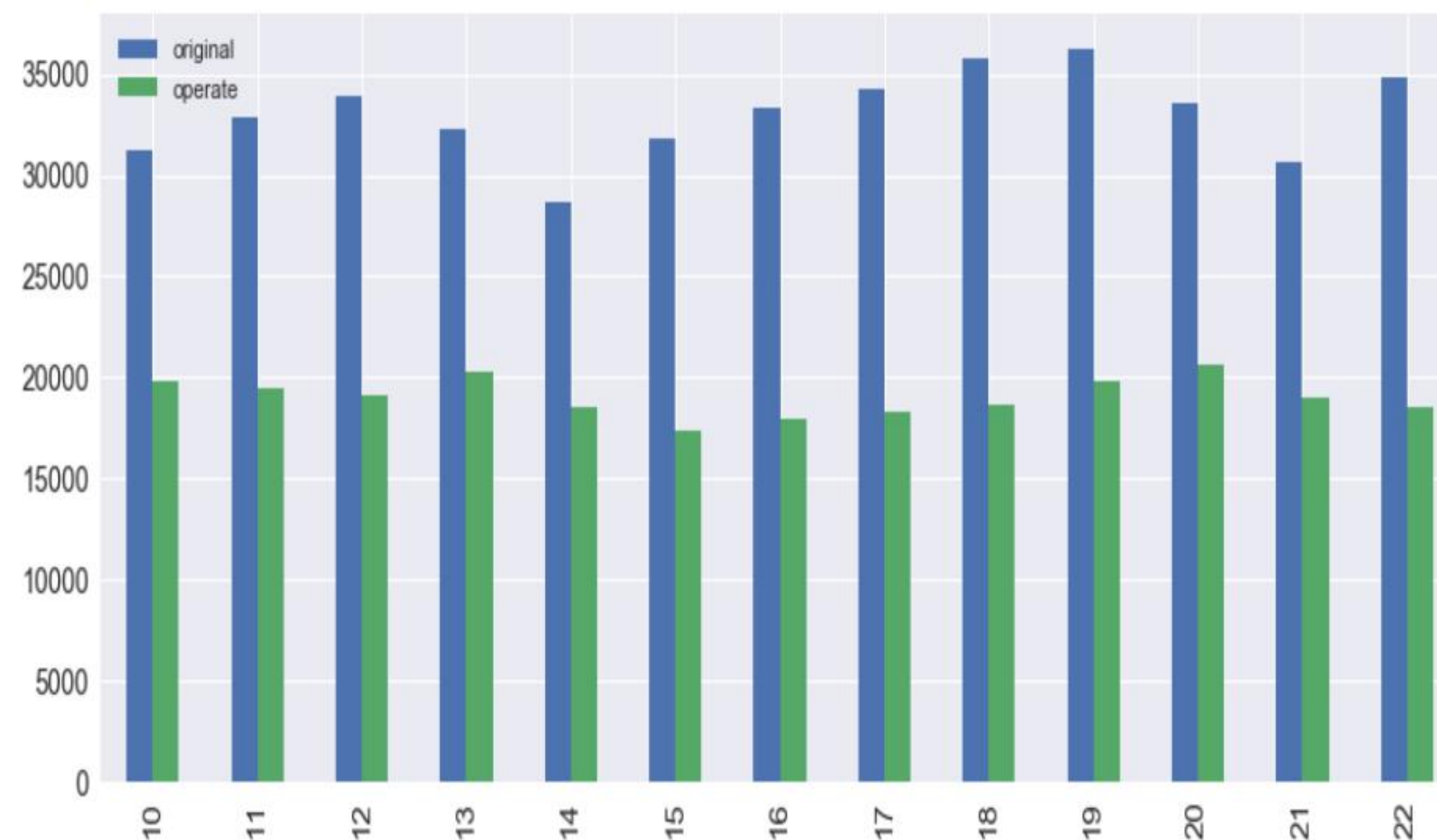
1. 正负样本极度不均衡
2. 预测的概率和真实的曝光期望之间存在偏差
3. 未充分考虑广告过滤的影响
4. 数据量很大



建模设想3：回归日曝光/出价

优势：
利用了出价信息，模型输出即已满足
出价单调性

缺点：
可用样本数减少近一半



建模设想4：直接回归日曝光（最终方案）



缺点

- 上限较低
- 难以充分利用用户信息和竞争队列信息
- 输入label的处理和loss function的选择对结果影响较大

优势

- 模型简单
- 易于快速迭代
- 鲁棒性好
- 能充分利用训练集样本





特征工程

特征介绍/特征有效性

广告属性特征



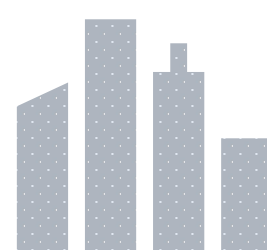
- 广告静态属性（广告主、素材尺寸、广告领域、.....）
- 广告动态属性（计费类型、转化类型）

竞价日志特征



- 竞价数及其时空信息（不同广告位和时间段的竞价数分布）
- 用户重复曝光情况
- 针对旧广告，统计对应的历史 n 天的曝光量以及竞价次数

统计特征



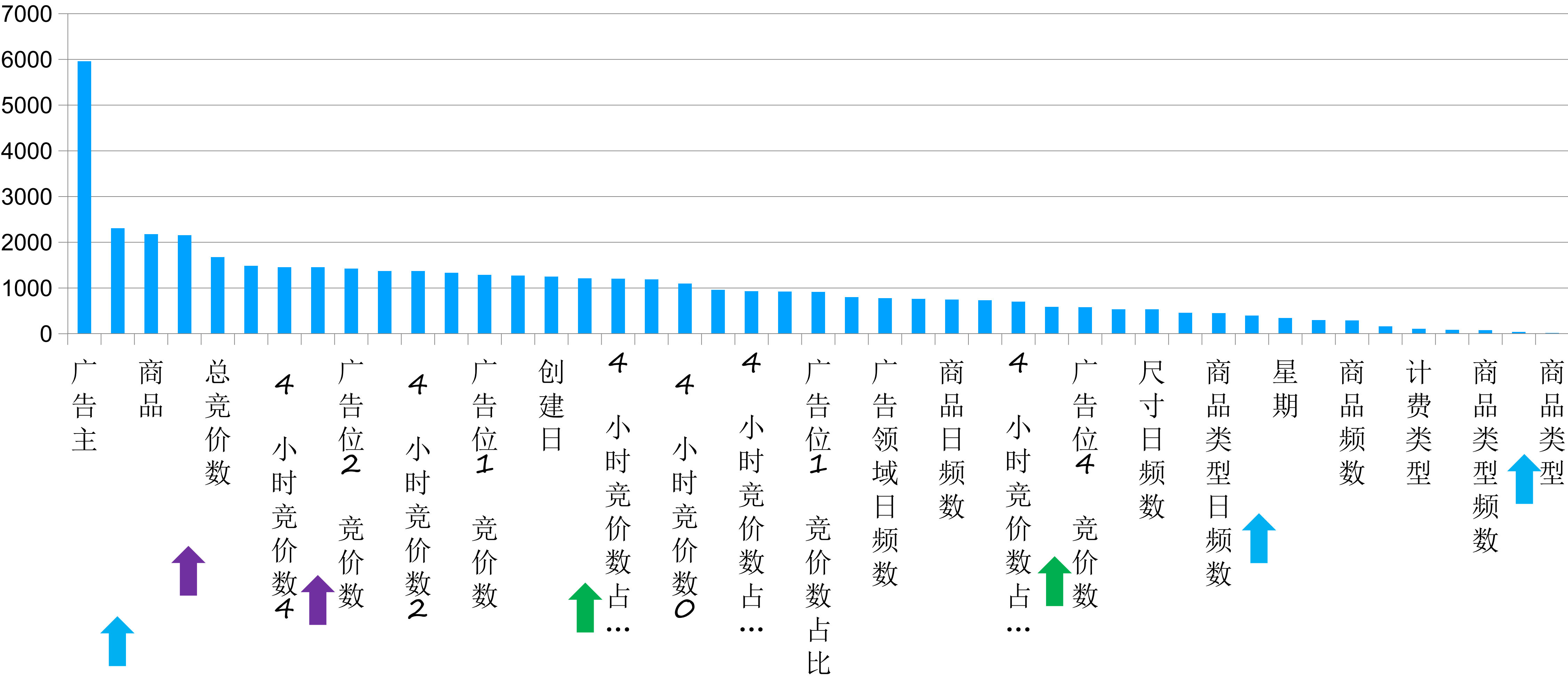
- 广告属性频数（总体）
- 广告属性频数（当天）

时间特征



- 预测日
- 星期
- 广告生存周期（预测日-创建日）

Lgb中特征重要性





模型介绍

通用处理/旧广告/新广告

STEP1



Label处理

- 用0-1之间的给定值REPLACE_ZERO替换0
曝光样本的标签并取log

Why not log1p?

$$\left\| \log \frac{e_1}{e_0} \right\|^2 \quad \left\| \log \left(\frac{1+e_1}{1+e_0} \right) \right\|^2$$

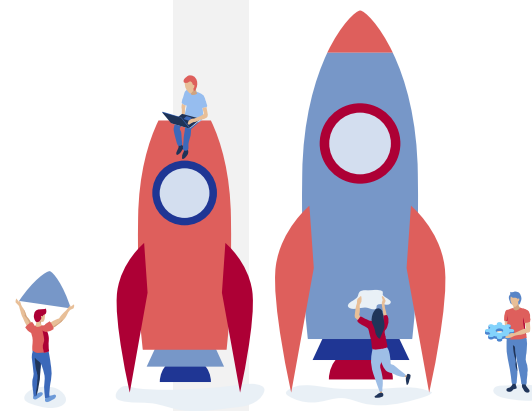
STEP2



梯度修正

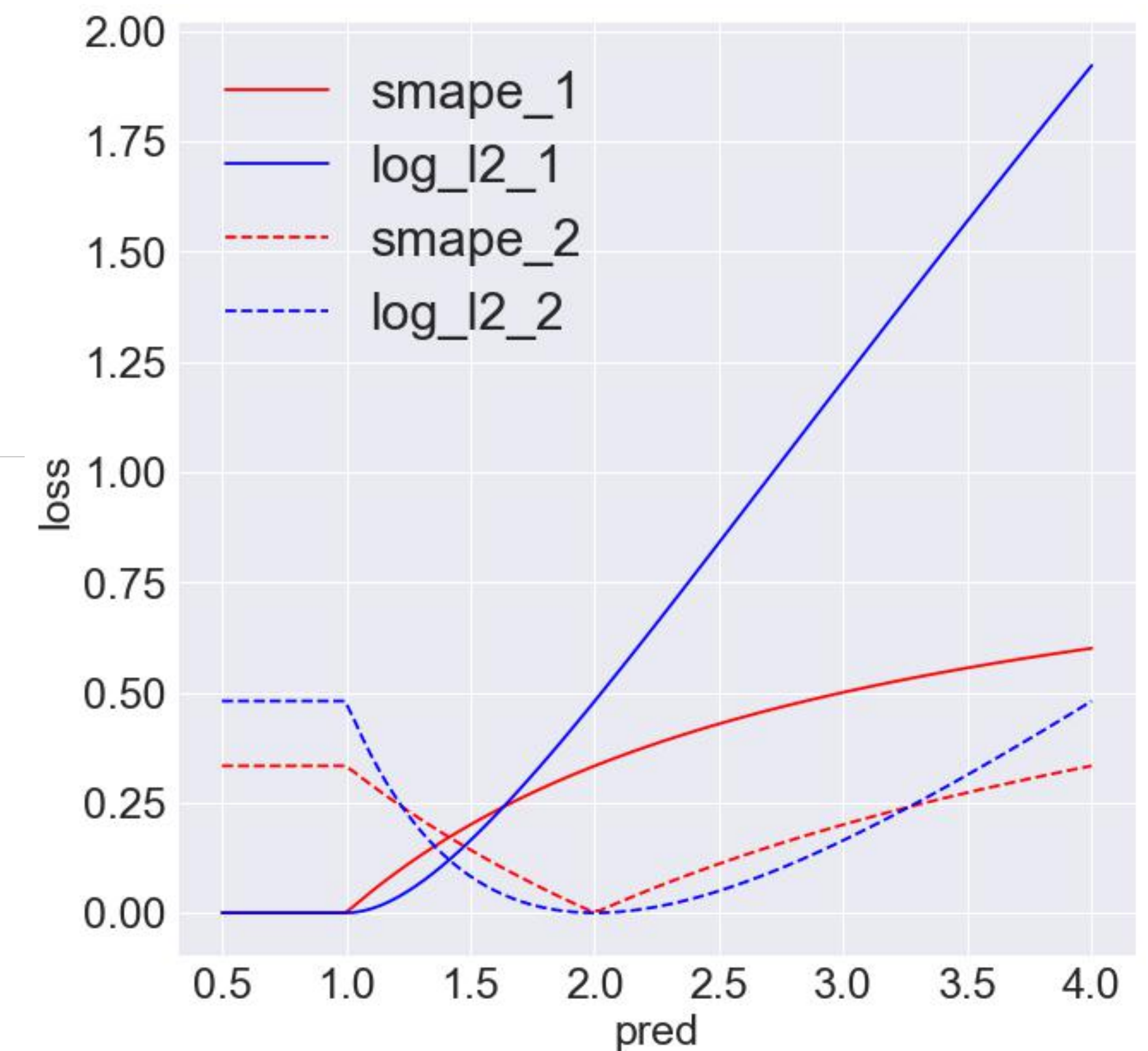
- 对于预测值偏小的0&1曝光样本不反传梯度

STEP3



组合loss

- L2-loss与smape的加权组合



旧广告预测

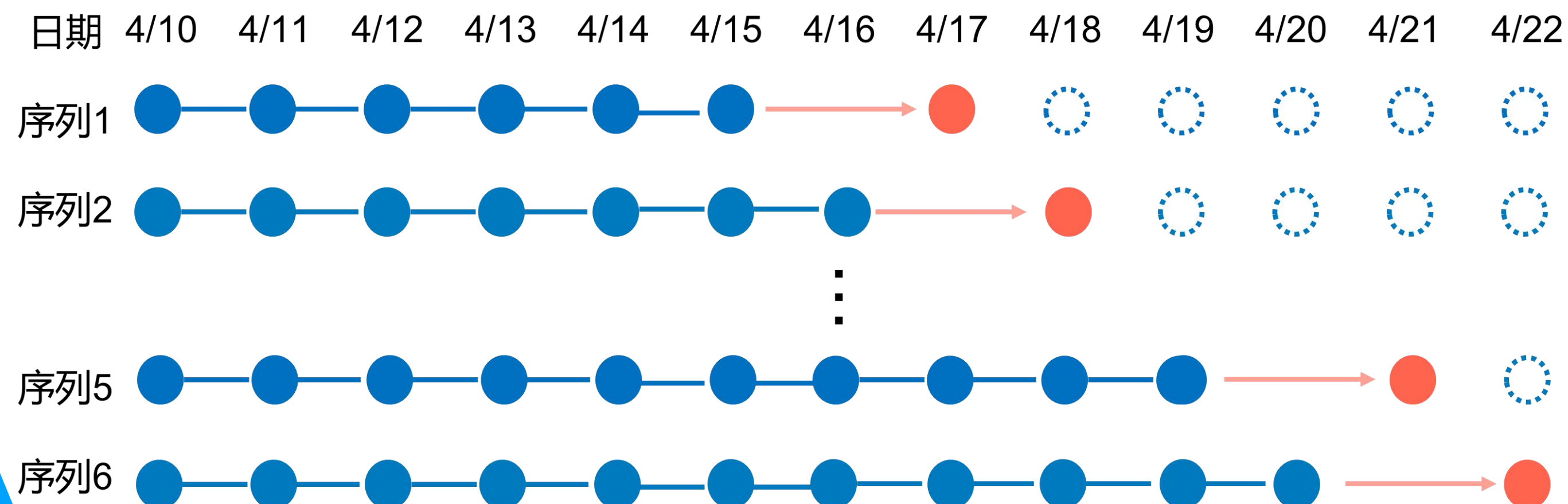
未利用预测日竞价
数信息

历史日曝光？

加权历史日曝光率

历史日曝光率？

剥离地利用不同天
的曝光率信息



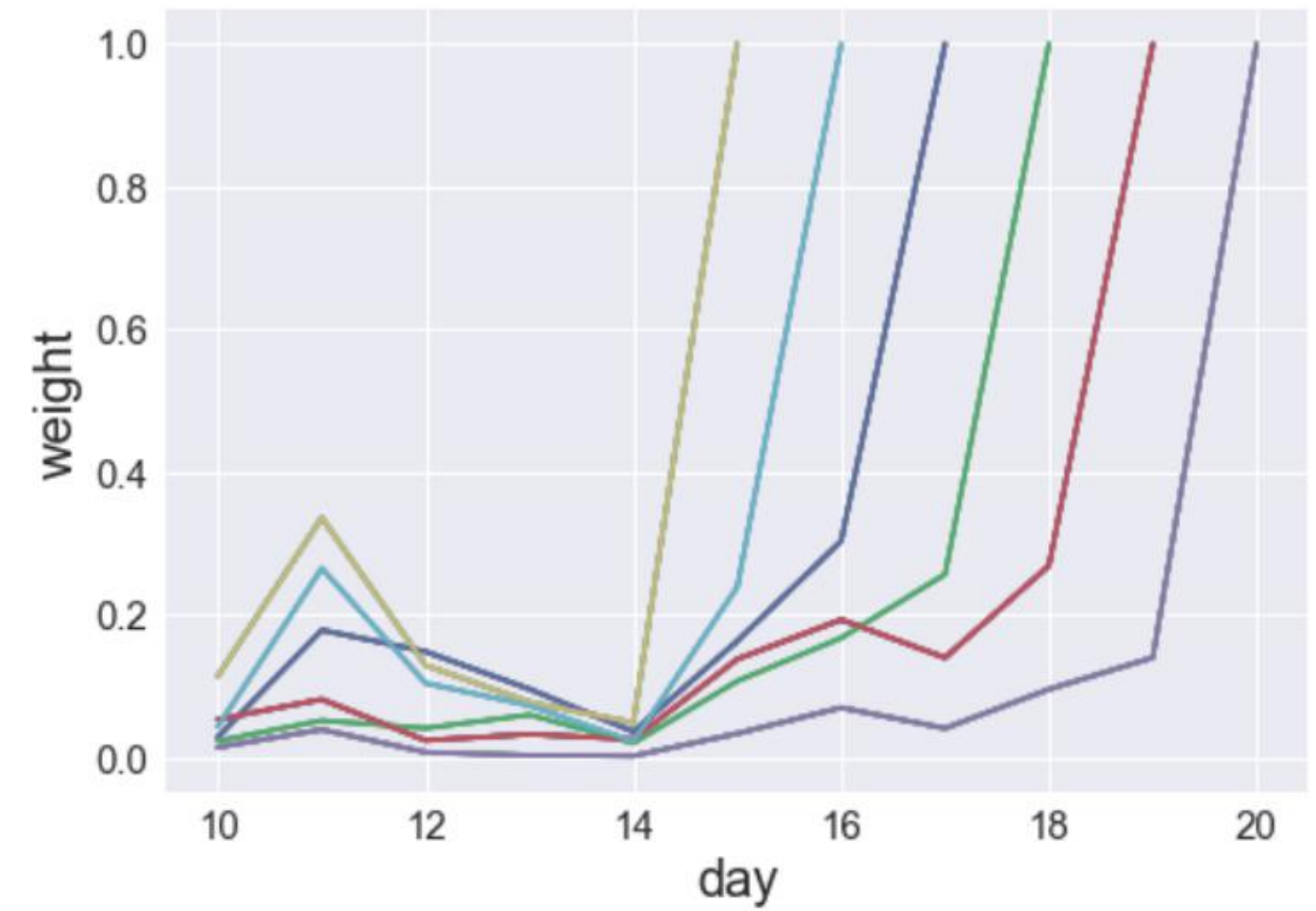
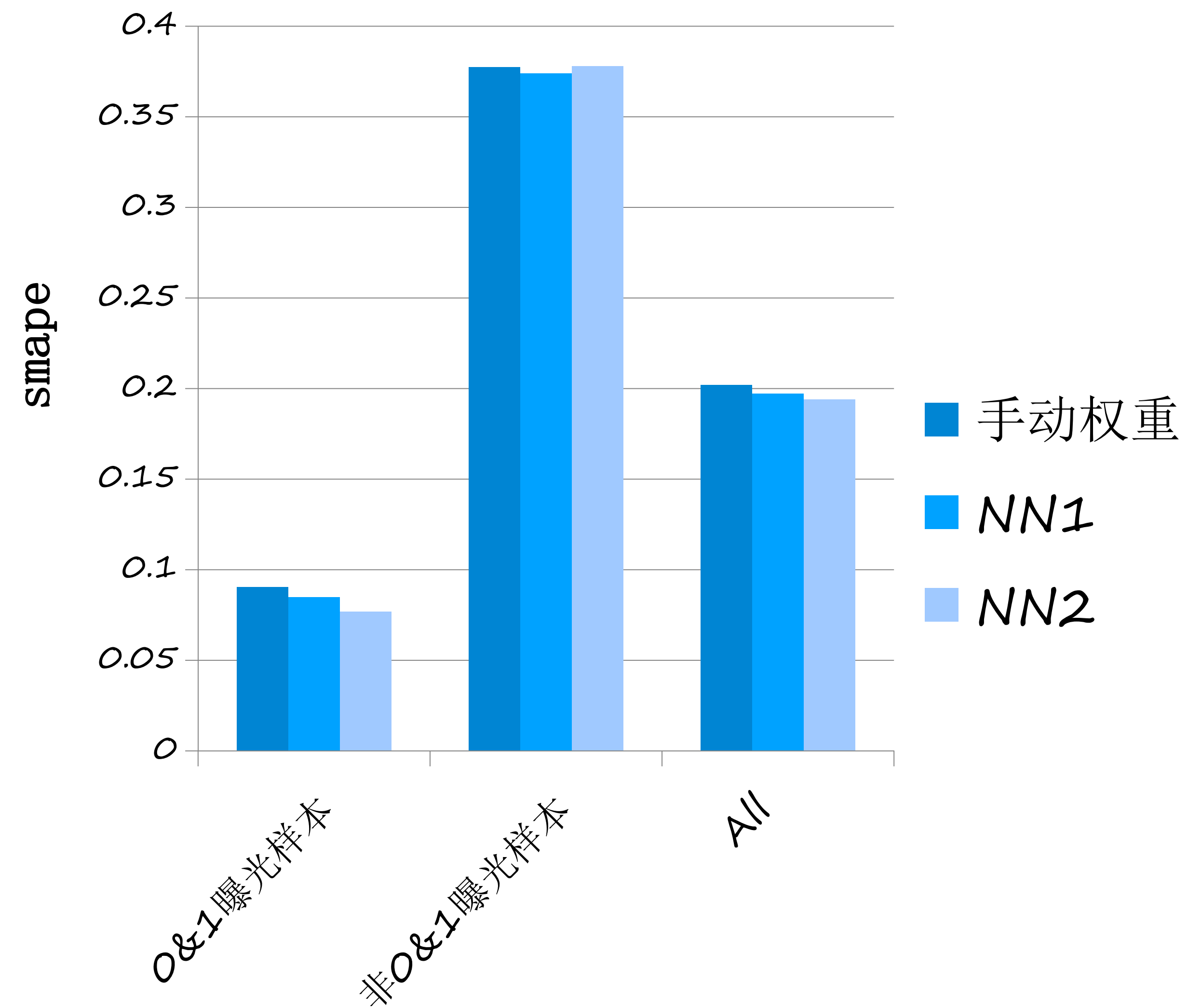
$$e_{pred} = c_{test} \times \frac{\sum w_i e_i}{\sum w_i c_i}$$

优势：

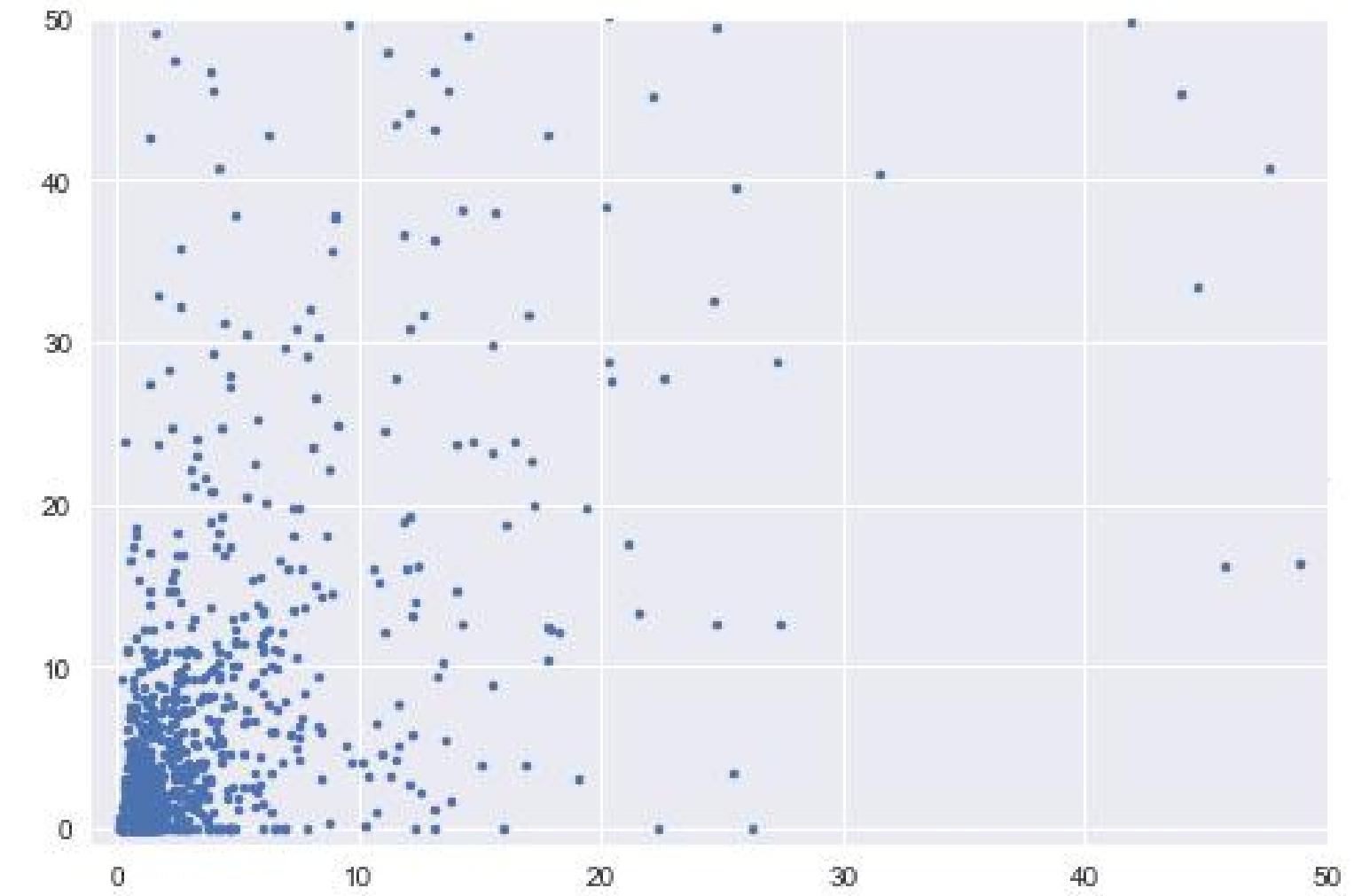
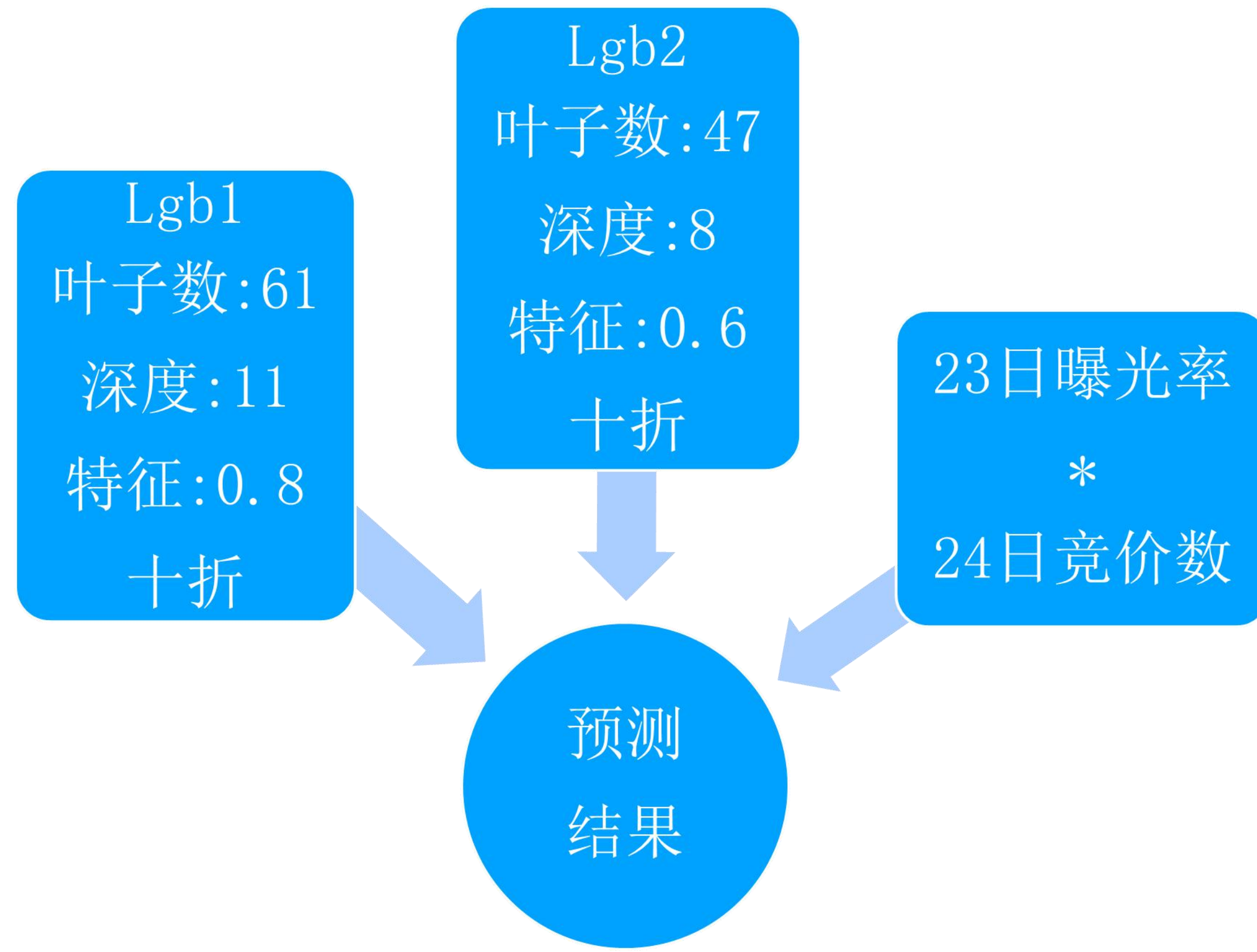
- 充分利用广告的所有历史竞价 e_i 和曝光 c_i
- 通过权重 w_i 综合考虑竞价数和与预测日的间隔天数

简单NN预测权重





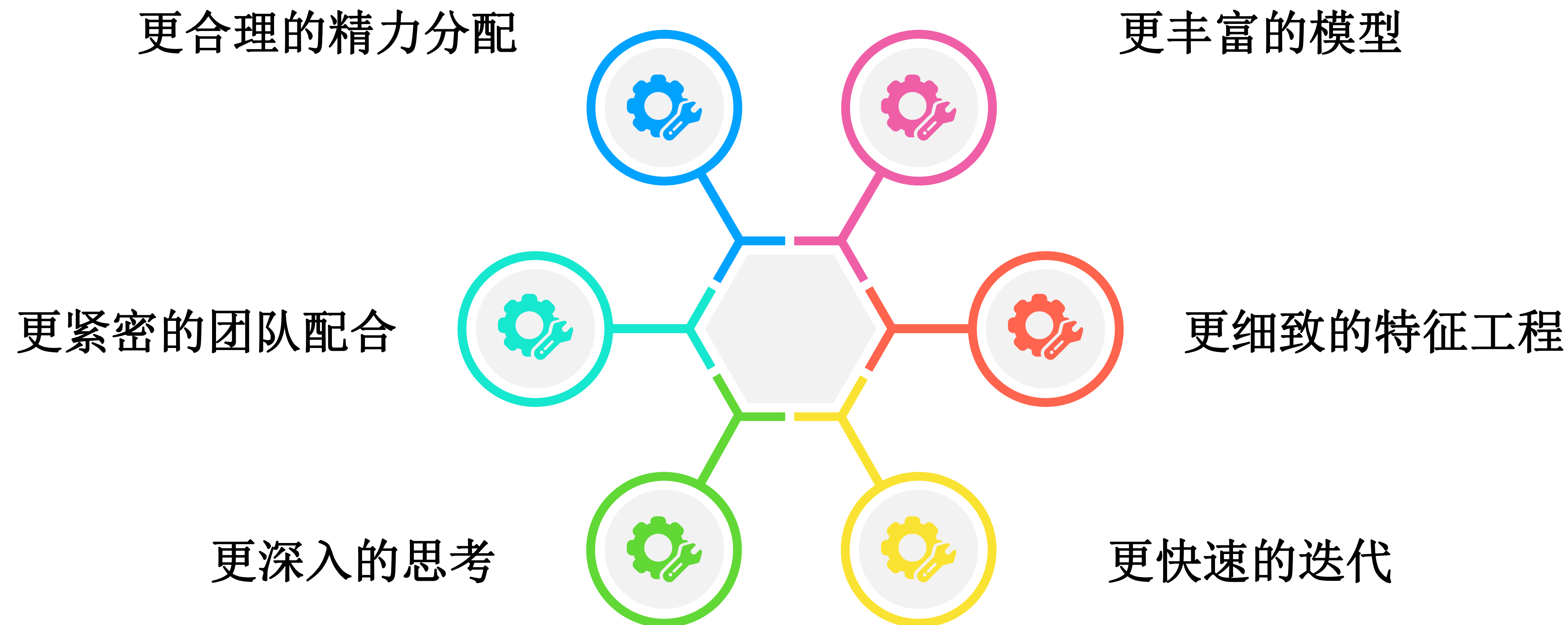
新广告预测





总结与思考

思考/感谢





感谢主办方的认真负责



感谢评委老师的倾听



感谢队友的一起拼搏

THANKS