



腾讯社交广告
Tencent Social Ads



高校算法大赛

我很难受

演讲人 李 强



腾讯社交广告
Tencent Social Ads



团队
介绍

框架
设计

特征
工程

模型
融合

2017-07-06

2

3

1

4

var

我很难受

time

if

*

01 团队介绍

2017-07-06

Your

var

time

if

<=Date()

freeze()

for

8

//

9

+

~

~

≡

李 强
吉林大学



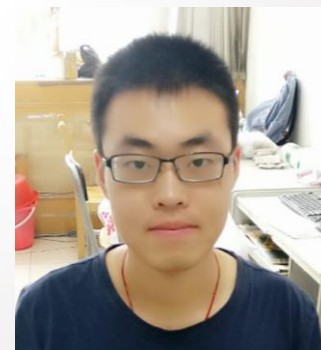
研一在读
计算机专业

李 智
北京航空航天大学
研二在读
电子科学与技术



我很难受

李博
北京邮电大学
研一在读
自动化专业



02 框架设计

2017-07-06



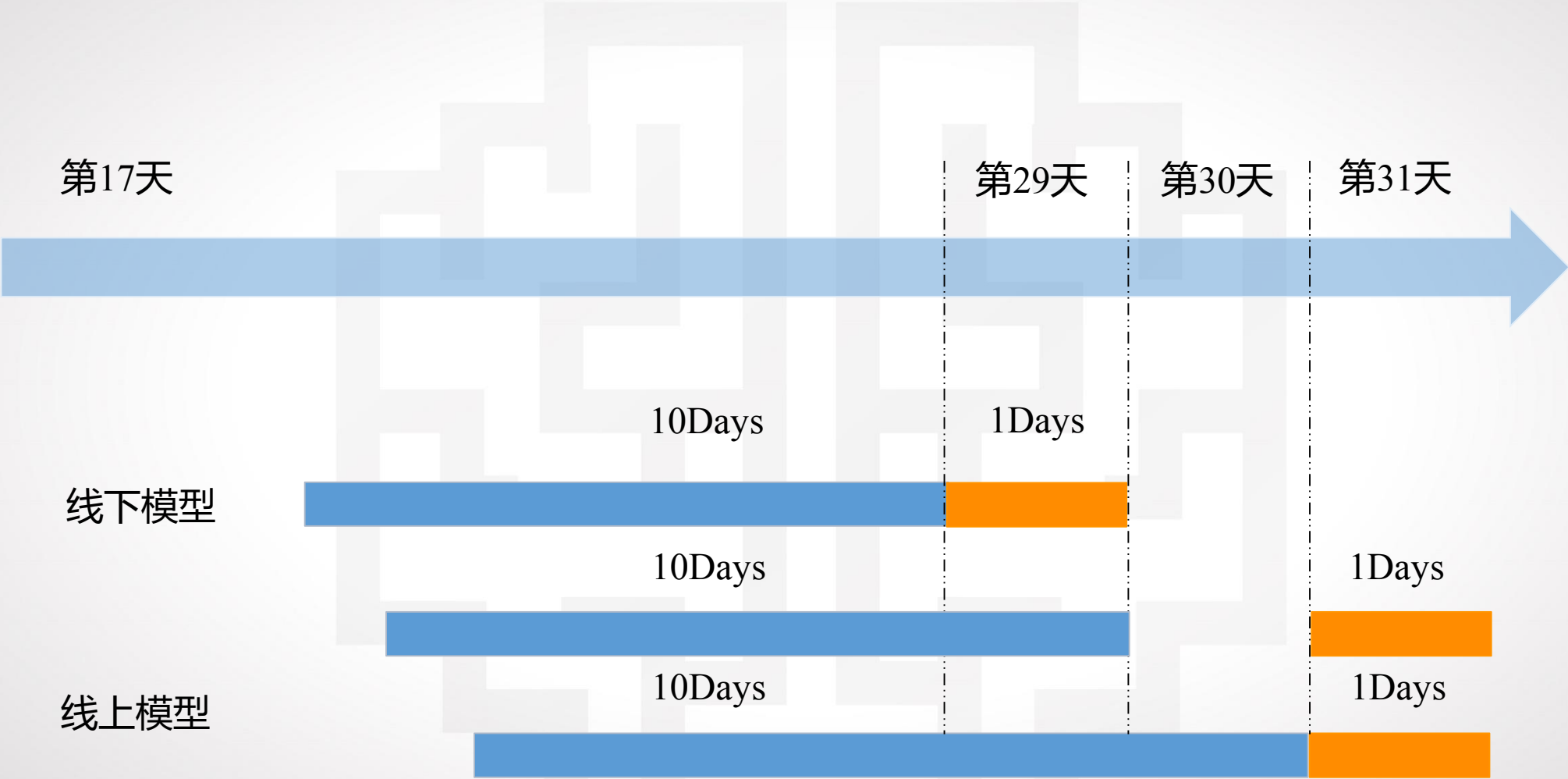
给定：用户信息表，app种类表，广告信息表，广告位信息表，第一天之前用户安装的app列表，第1天到第30天的app安装列表，第17天到30天的训练表。

预测：用户在第31天点击某个广告后会发生转化的概率。

评价：

$$\text{logloss} = -\frac{1}{N} \sum_{i=1}^N (y_i \log(p_i) + (1 - y_i) \log(1 - p_i))$$





03 特征工程

2017-07-06

Your

var

time

if

<=Date()

freeze()

9

+

//

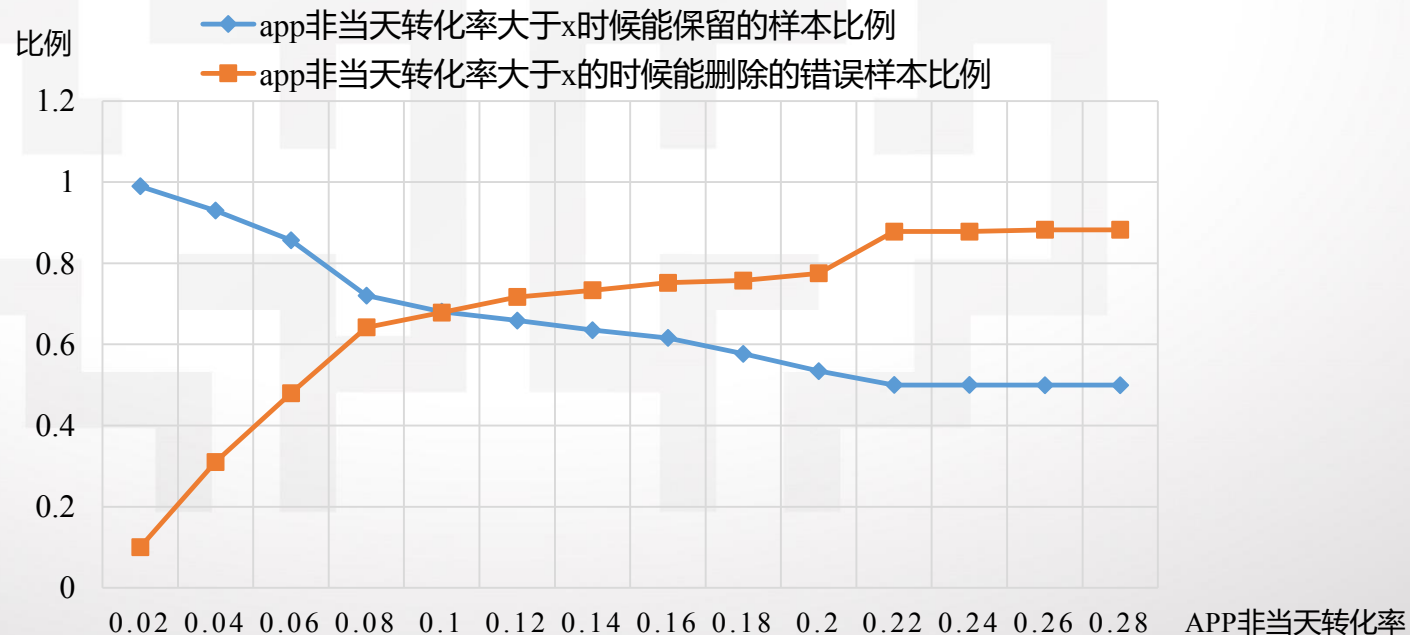
~

~

≡

训练数据的末尾几天由于转化时间的滞后性，存在错误标签的样本，越靠近第30天，错误样本的比例越大，线上成绩反馈表明，直接将第30天放入训练样本中训练，成绩将会大幅下降。为此我们定义了app非当天转化率来对数据做清洗：

$$\text{app非当天转化率} = \text{app非当天转化总数} / \text{app转化总数}$$





基础id特征

广告类id

广告位id

app类id

用户id

统计特征

点击次数

转化次数

转化率

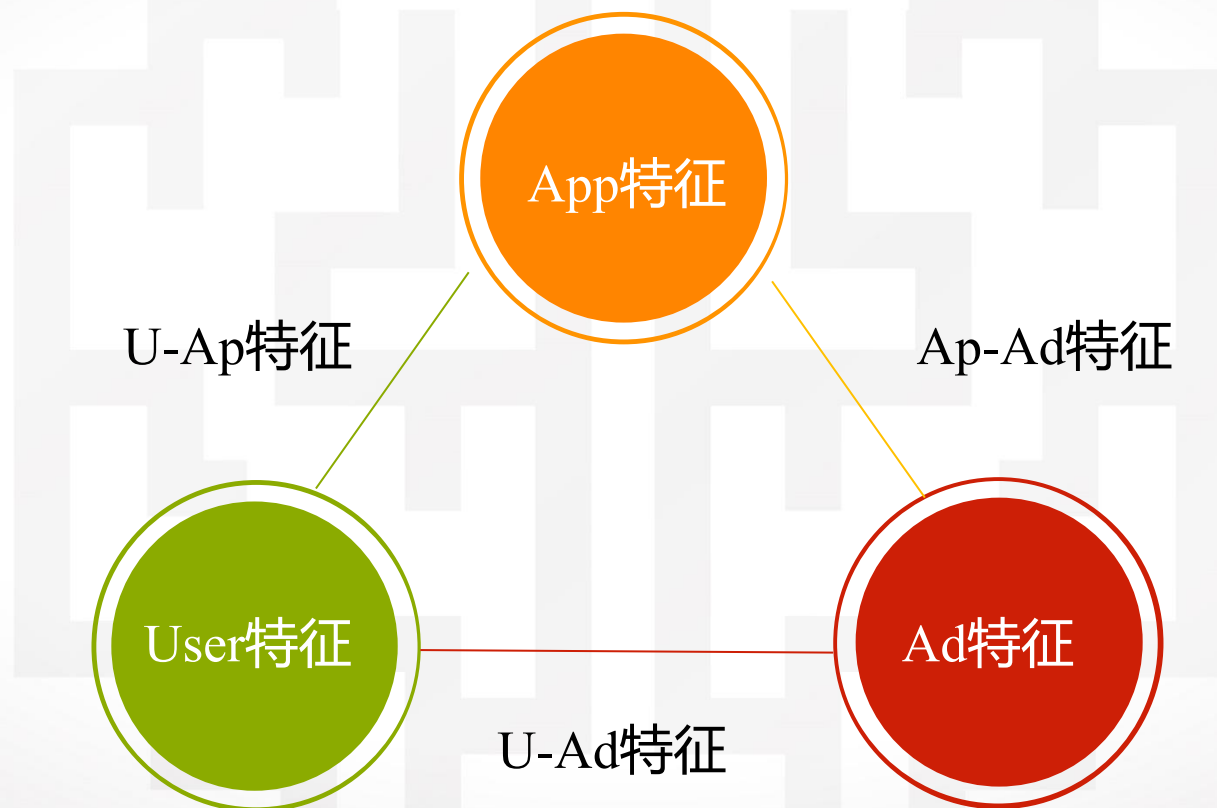
时间特征

时间间隔

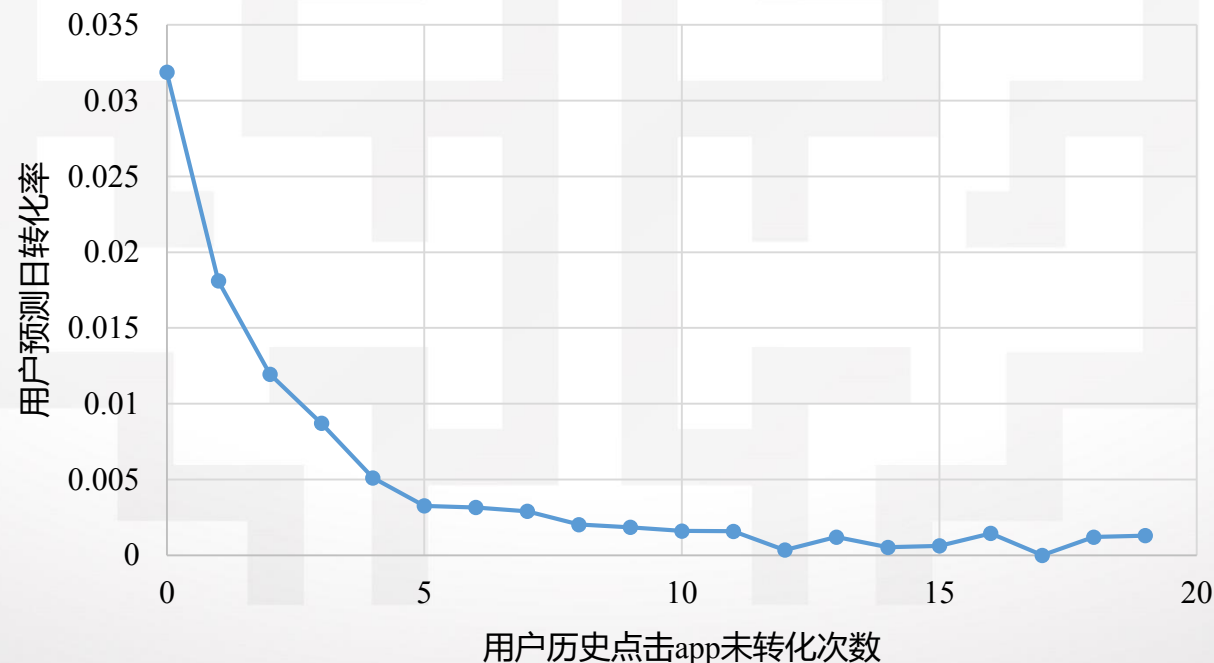
用户活跃度

最早最晚时间





因为数据量庞大，盲目的通过模型测试特征变得十分浪费时间，因此我们在每次测试特征前都会对当前特征做分析，当特征在不同取值下转化率有较大差别的情况下才会用模型测试。



- 广告，广告平台，app的历史转化率作为特征能极大提升预估的准确性，但是当数据量比较少的时候，历史转化率与实际转化率相差很大，极端情况下，广告就被点击1次并被转化了，直接计算转化率为1与实际转化率相差很大。

引入贝叶斯平滑对历史转化率做出修正：

$$r = (C + \alpha) / (I + \alpha + \beta)$$

- 通过划窗来增加样本势必会带来时间跨度不统一而带来的特征量级不同的问题，我们用日均点击量代替历史点击量，日均转化量代替历史转化量，使得特征量级归一化。



04 模型融合



2017-07-06



freeze()



<=Date()



9



本次比赛，我们分别训练了三种模型，分别是xgboost，lightGBM以及FFM。

xgb

lgb

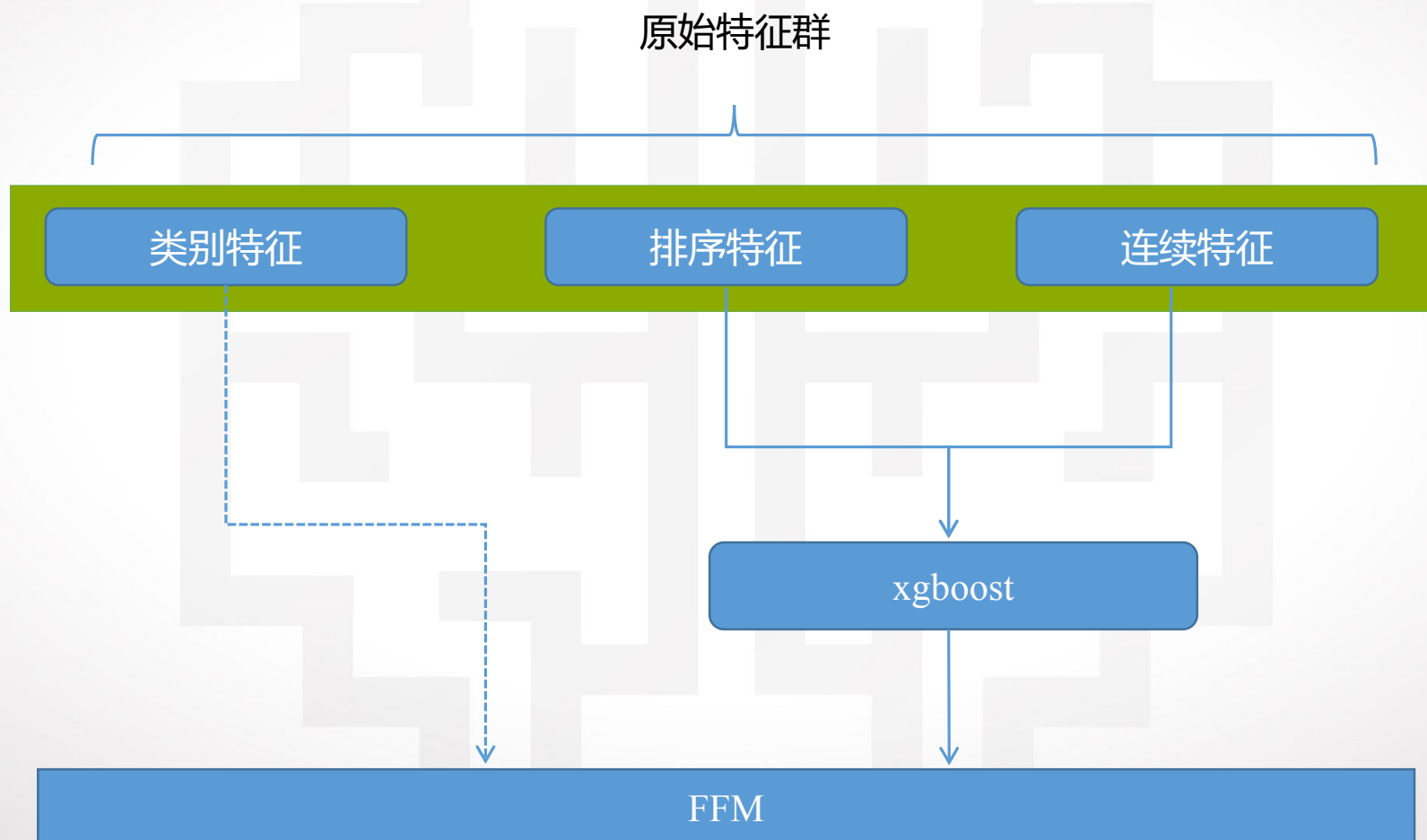
ffm

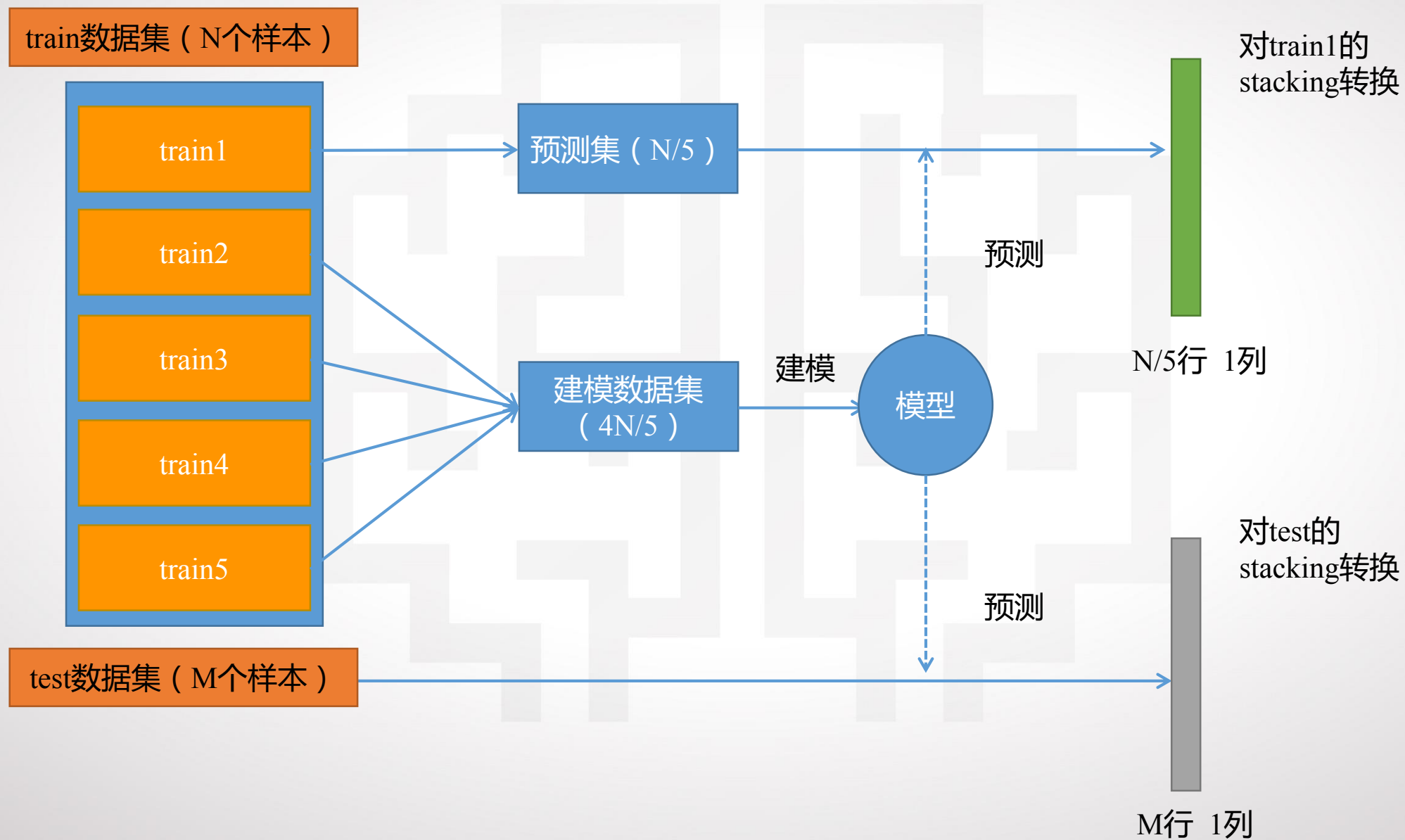
线上B榜0.101435

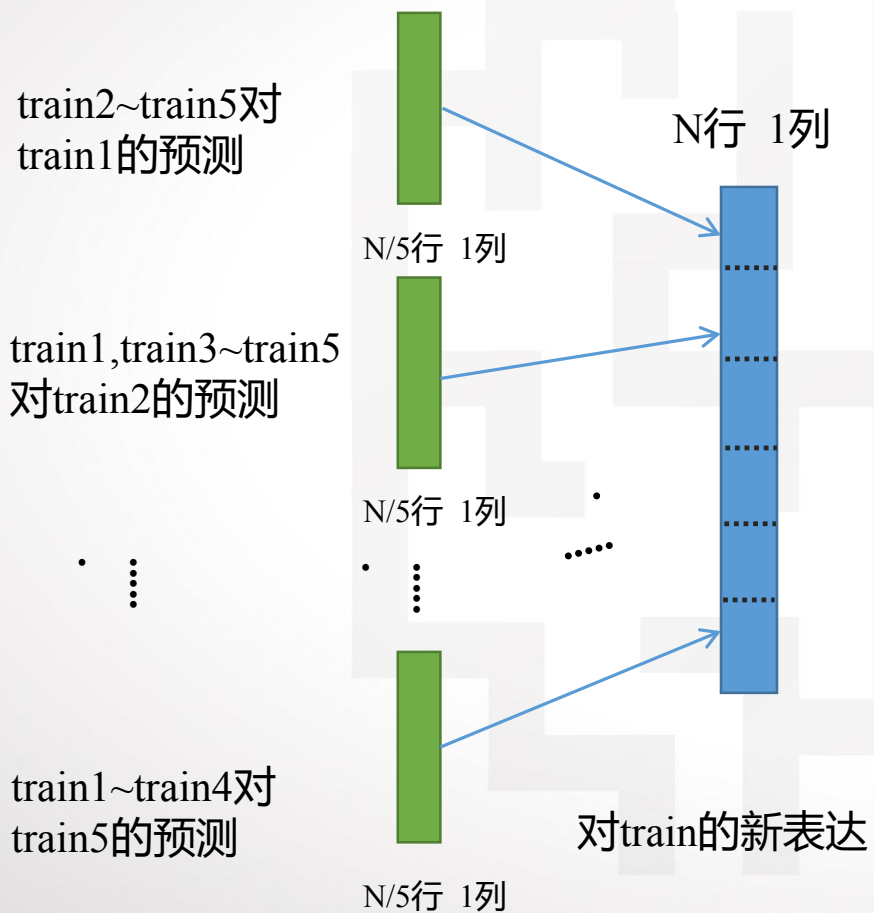
线上B榜0.101416

线上B榜0.1024左右

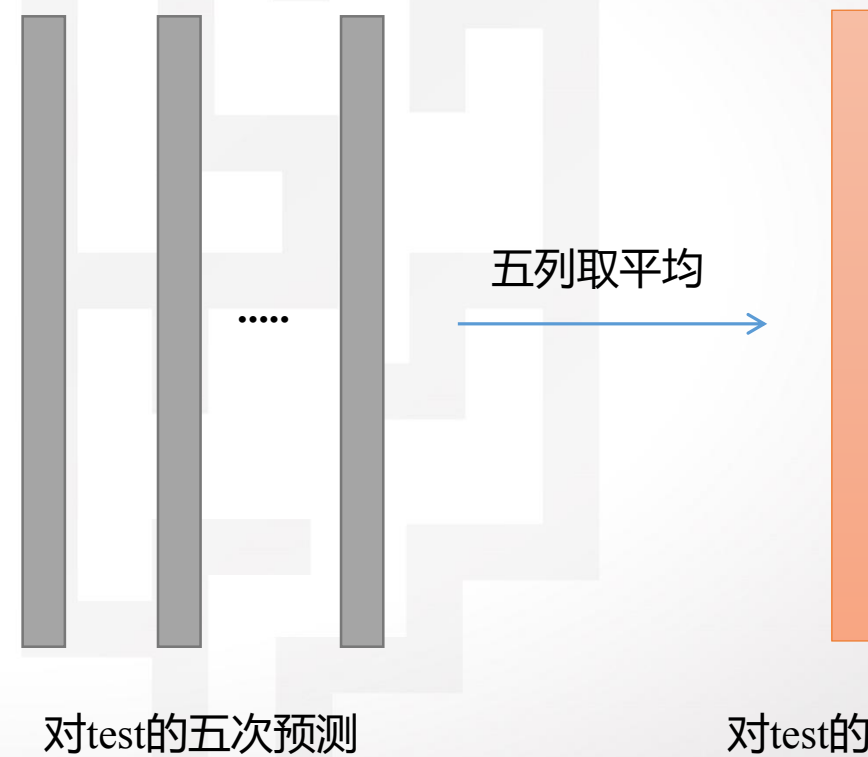


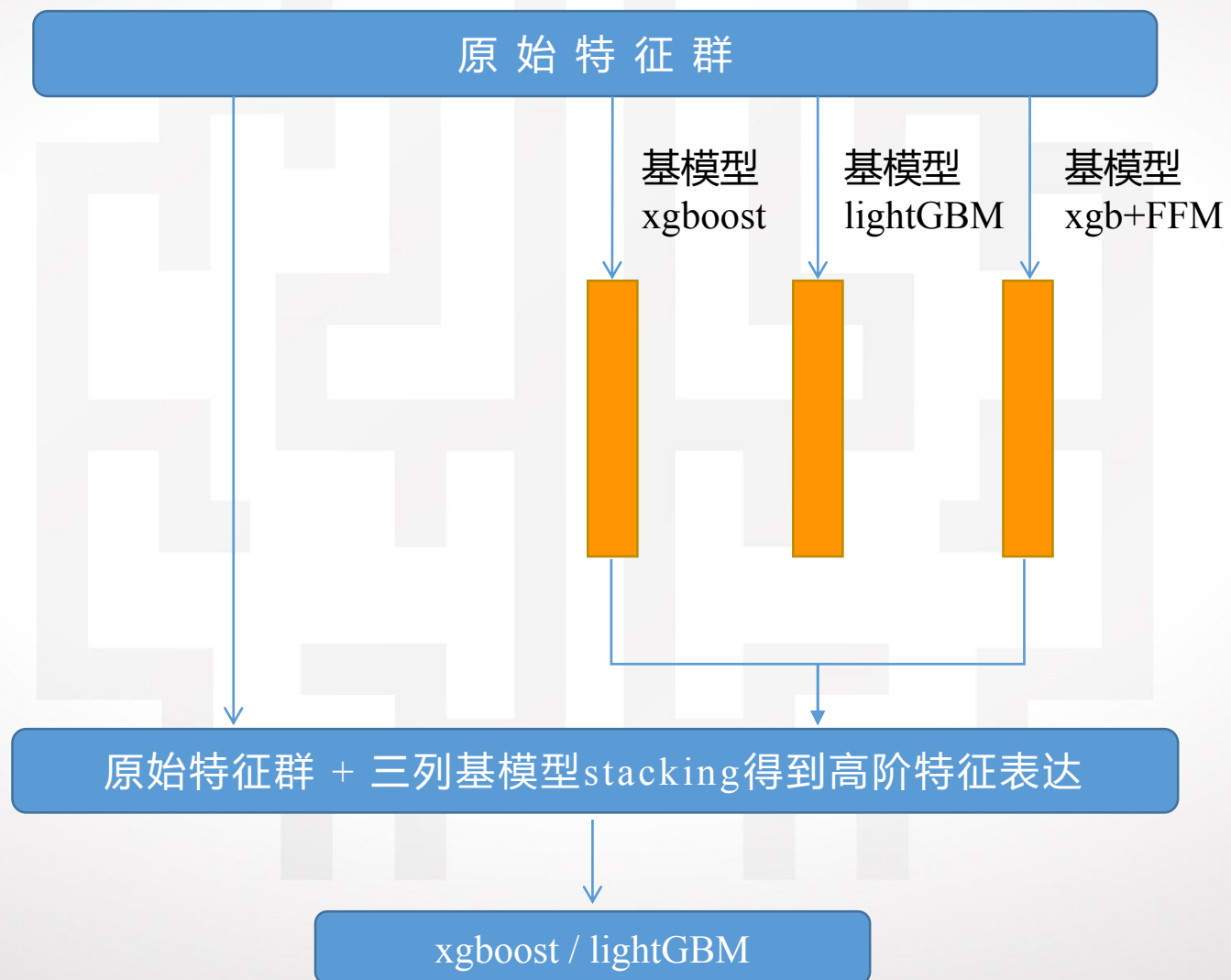






五部分均是 M行 1列





THANKS

freeze()

for

<=Date()

2017-07-06

Your

var

time

if