



机器学习在票房预估中的实战

李明辉

猫眼电影高级技术专家

WOTD

World Of Tech
2017年12月1-2日

全球软件开发技术峰会

[深圳站]

报名咨询：010-68478816

议题提交：wot@51cto.com

市场合作：yangxh@51cto.com

商务合作：songjc@51cto.com

媒体合作：yankk@51cto.com

在线咨询（微信）：18401576051

团·购·享·受·更·多·优·惠

5折 优惠（截止8月31日）
现在报名，立省1400元/张



内容提纲

- 行业背景
- 技术体系
- 预测算法
- 工作展望



内容提纲

- 行业背景
- 技术体系
- 预测算法
- 工作展望



行业背景

机遇与风险


- 中国电影市场高速发展
- 高风险、投资回报率难以预测
- 大投入未必有大产出
- 预测工具的缺失

票房预测的意义

- 投资：预估投资回报率，控制风险
- 宣发：营销成本控制
- 上映：智能排片，利益最大化





 行业背景





行业背景





行业背景

票房预测的难点

- 中国电影处于野蛮生长期
- 信息量很大，但是垃圾更多
- 可供参考或学习的样本量少
- 有些感性特征难以量化

影响票房的因素

- 题材（受众群范围）
- 卡司阵容（粉丝群范围）
- 影片质量（口碑效应）
- 档期（同行竞争）
- 宣传力度（主动传播）
- 非市场因素



行业背景

票房预测的难点

- 中国电影处于野蛮生长期
- 信息量很大，但是垃圾更多
- 可供参考或学习的样本量少
- 有些感性特征难以量化

影响票房的因素

- 题材（受众群范围）
- 卡司阵容（粉丝群范围）
- 影片质量（口碑效应）
- 档期（同行竞争）
- 宣传力度（主动传播）
- 非市场因素





行业背景

• 票房预测的发展阶段

萌芽阶段(1915~1960)

对变量（因素）的摸索

- 高昂的拷贝价格
- 制片需要了解影片的质量
- 盖洛普研究影响票房的因素

初级阶段(1980~2006)

复杂因素分析模型的建立

- 电视的发展、刺激电影行业控制风险
- 巴里·利特曼票房回归模型
- 斯格特·苏凯的竞争市场预测模型
- 机器学习模型BPNN

发展阶段(2006~2013)

采用单一数据源为核心的分析

- 基于博客的票房预测模型
- 基于新闻报道的票房预测模型
- 基于Twitter进行票房预测的模型
- 基于google搜索引擎的预测模型
- 基于维基百科的预测模型



行业背景

国内的票房预测服务

猫眼

淘宝

腾讯

百度

艺恩

时光



一起拍电影



票房透视镜



行业背景

- 票房预测服务的分布

| | 猫眼 | 淘宝 | 腾讯 | 百度 | 艺恩 | 时光 | 一起拍 | 犀牛娱乐 |
|--------|----|----|----|----|----|----|-----|------|
| 全国实时票房 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | |
| 影院实时票房 | ✓ | | | | | | | |
| 次日票房预测 | ✓ | | ✓ | ✓ | | | ✓ | ✓ |
| 总票房预测 | ✓ | | | ✓ | | | ✓ | ✓ |



内容提纲

- 行业背景
- 技术体系
- 预测算法
- 工作展望



技术体系

实时票房

- 每部影片在全国、各影院、城市的票房、排片、人次、平均票价等

天级票房

- 每部影片在全国次日票房预测
- 前一日17点、21点

总票房

- 每部影片在整个上映周期的全国票房预测
- 首映日、首周末时间节点

预售票房

天级票房预测结果



技术体系

合作影院售票数据

猫眼交易数据

基础数据

全国实时票房

影院实时票房

城市地区实时票房

天级票房预测

总票房预测

预测服务



猫眼专业版



技术体系

猫眼专业版-票房

排片 影院院线 **票房预测** 上座率 北美票房

日票房 2017年5月15日 含服务费

前 一天 今日实时 **6417.4万** 后 一天
北京时间 22:59:08

票房排名 更多指标 >

| 影片 | 综合票房 (万元) | 排片占比 | 场均人次 | 人次 (万人) |
|---------------------------|----------------|-------|------|---------|
| 摔跤吧! 爸爸 上映11天 4.49亿 | 3542.59 | 29.6% | 19 | 117.2 |
| 银河护卫队2 上映11天 5.66亿 | 1239.83 | 17.8% | 9 | 32.7 |
| 亚瑟王: 斗兽争霸 上映4天 4006.8万 | 432.50 | 12.3% | 5 | 12.6 |

票房 网播量 资料库 找合作 我的

实时票房

总票房

票房预测 票房指数

周二大盘报收 **5774.1万** 周三大盘预测 **5307.6万**
5月16日 周二 5月17日 周三

注: 每晚22点更新今日快报及明日预测

日票房排名 周二快报 (万元) 周三预测 (万元)

| 影片 | 预测总票房 | 周二快报 (万元) | 周三预测 (万元) |
|--|---------|---------------|---------------|
| 摔跤吧! 爸爸 预测总票房 9.49亿 上映12天 4.80亿 | 9.49亿 | 3190.6 | 2934.2 |
| 银河护卫队2 预测总票房 6.95亿 上映12天 5.78亿 | 6.95亿 | 1164.1 | 1032.2 |
| 亚瑟王: 斗兽争霸 预测总票房 6388.2万 上映5天 4407.1万 | 6388.2万 | 391.4 | 330.2 |
| 毒。诚 预测总票房 5286.8万 上映5天 3345.1万 | 5286.8万 | 264.7 | 211.1 |
| F8 速度与激情8 预测总票房 27.11亿 | 27.11亿 | 164.3 | 154.3 |

天级票房



内容提纲

- 行业背景
- 技术体系
- 预测算法
- 工作展望



实时票房



❖ 案例：《我不是潘金莲》



11.18首映：万达排片低于其他院线

冯小刚《潘金莲至王健林的一封信》

王思聪反击

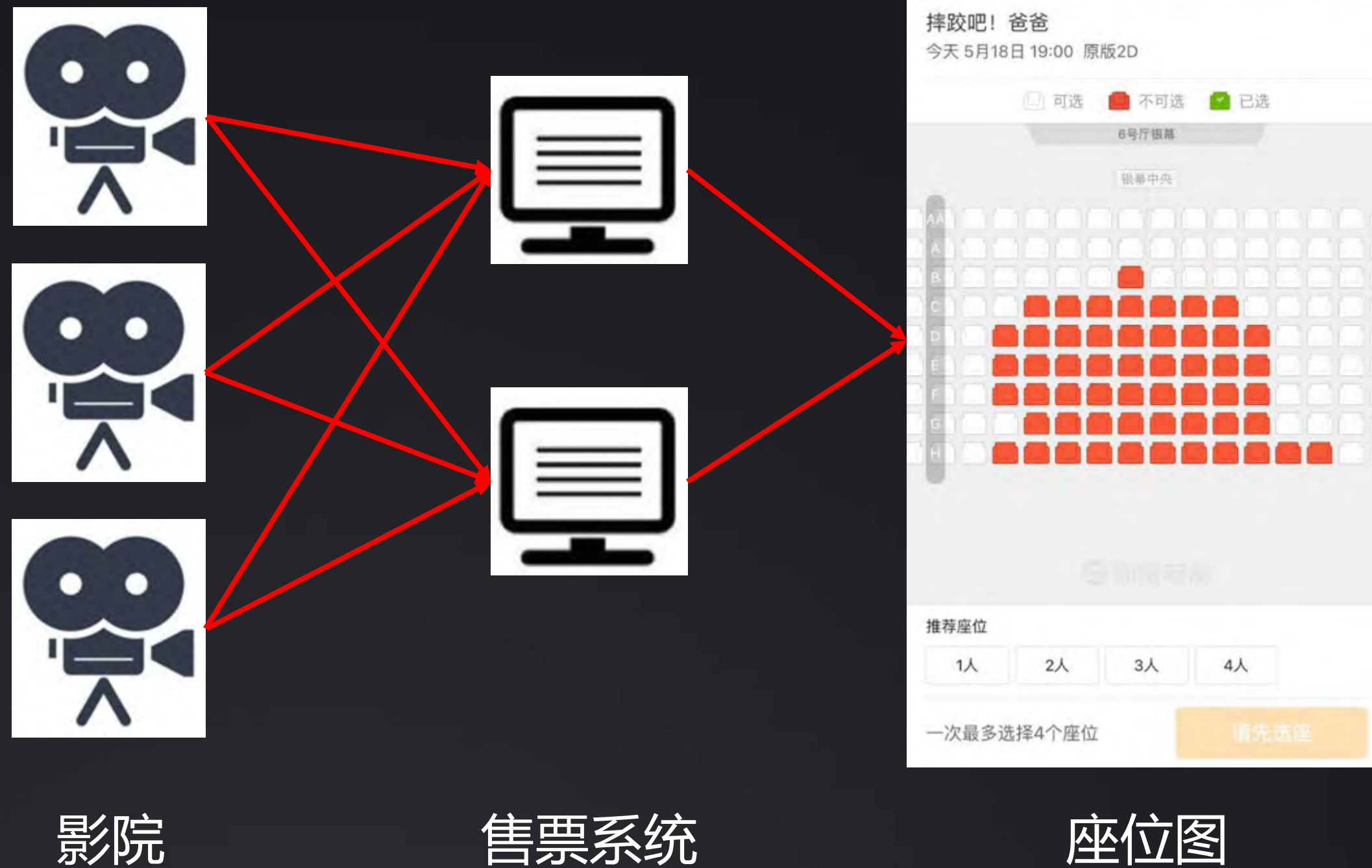
曝光度激增

票房



实时票房

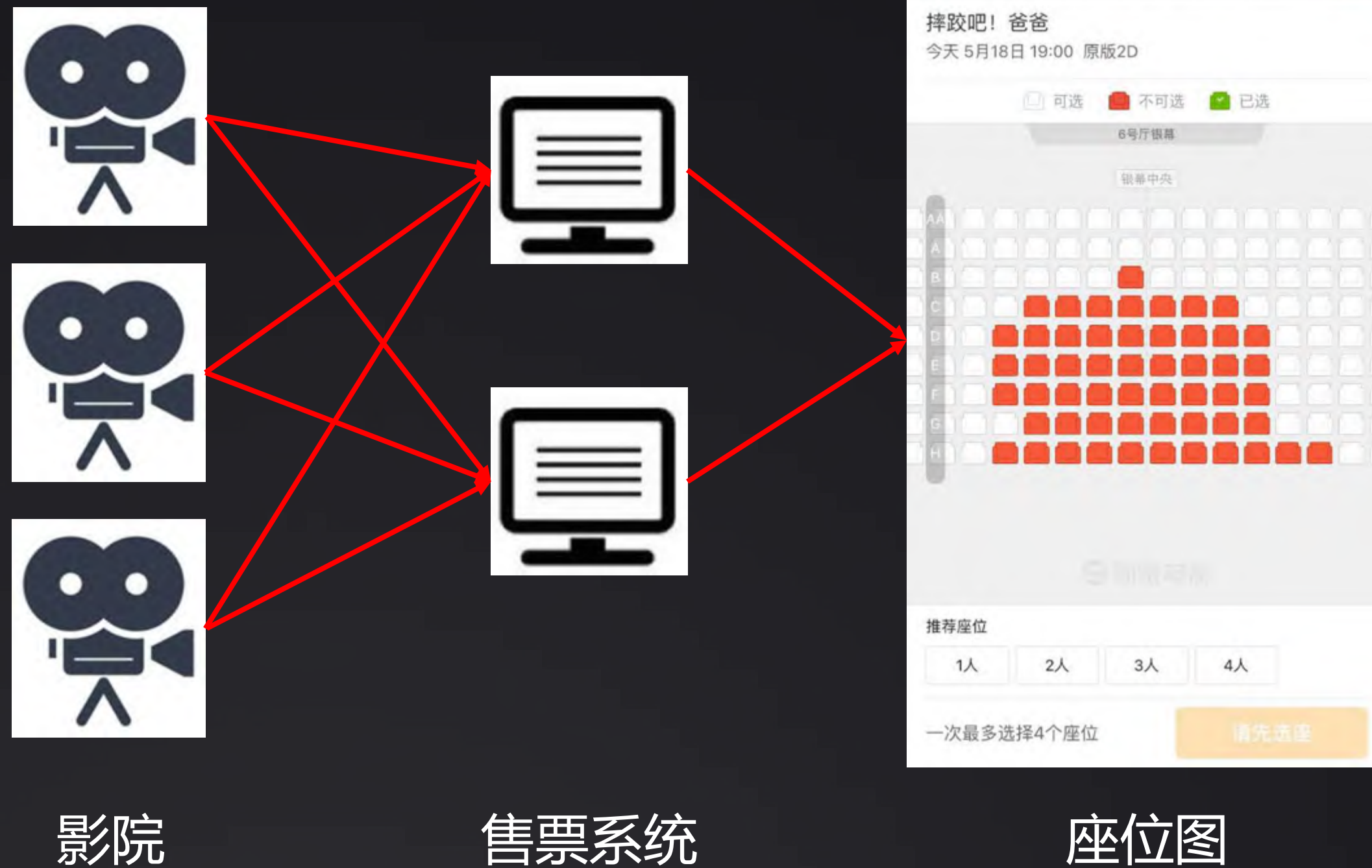
- 猫眼的数据优势
 - 与90%以上的影院合作，对接18种售票系统
 - 拥有详细的影院场次及座位状态





实时票房

- 猫眼的数据优势
 - 与90%以上的影院合作，对接18种售票系统
 - 拥有详细的影院场次及座位状态



问题：不可选 ≠ 已售



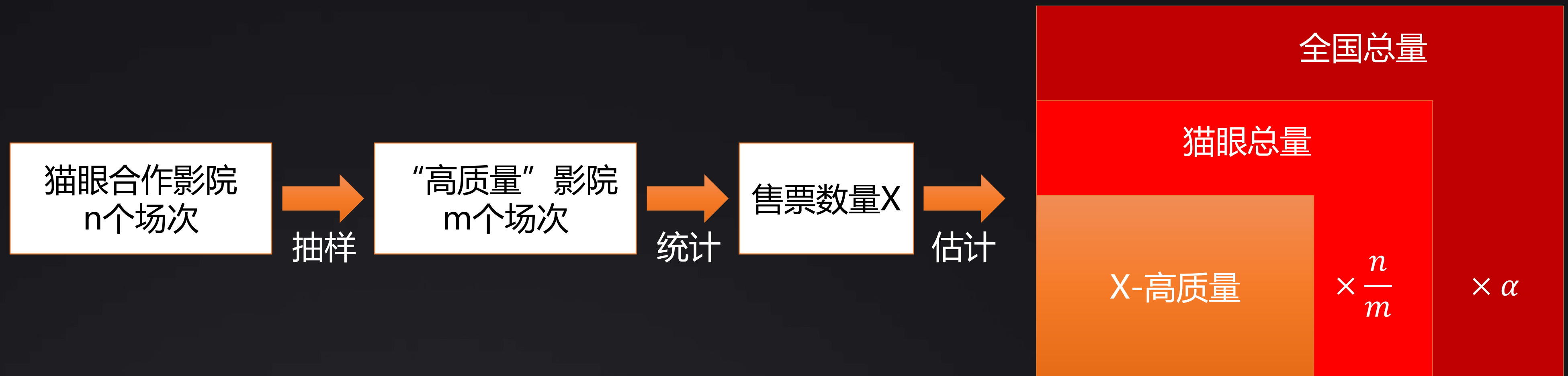
实时票房

- 问题定义
 - 已知猫眼合作影院（占全国约90%）每部电影、每个场次的实时座位图数据
 - 预测每部电影在全国全部影院的实时票房总和
- 难点
 - 数据噪音：座位图中状态不明确
 - 数据不完备：非全部影院
- 解决方案
 - 数据抽样：以部分样本为基础，估计全量



实时票房

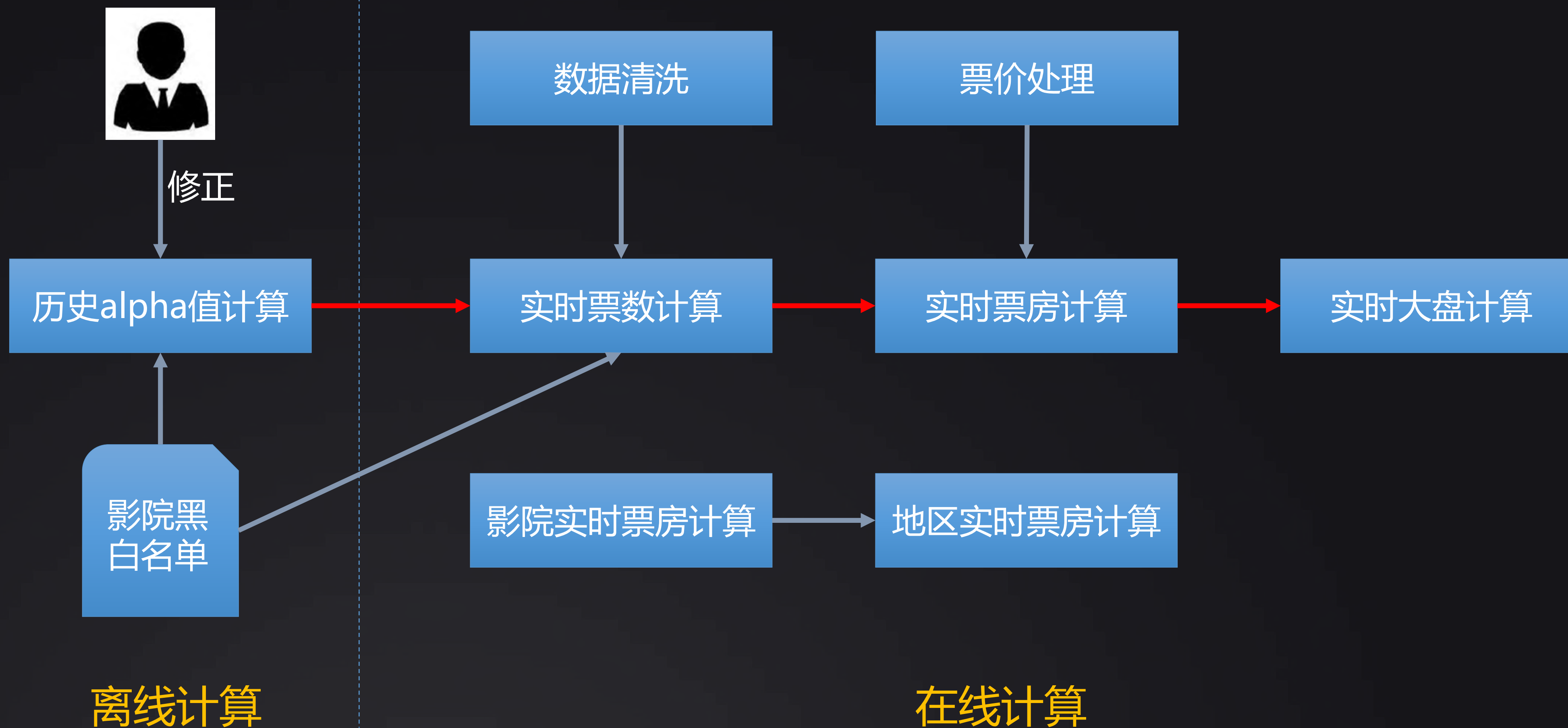
- 统计模型



$$g_t = X_t * \frac{n_t}{m_t} * \alpha_{t-1}$$



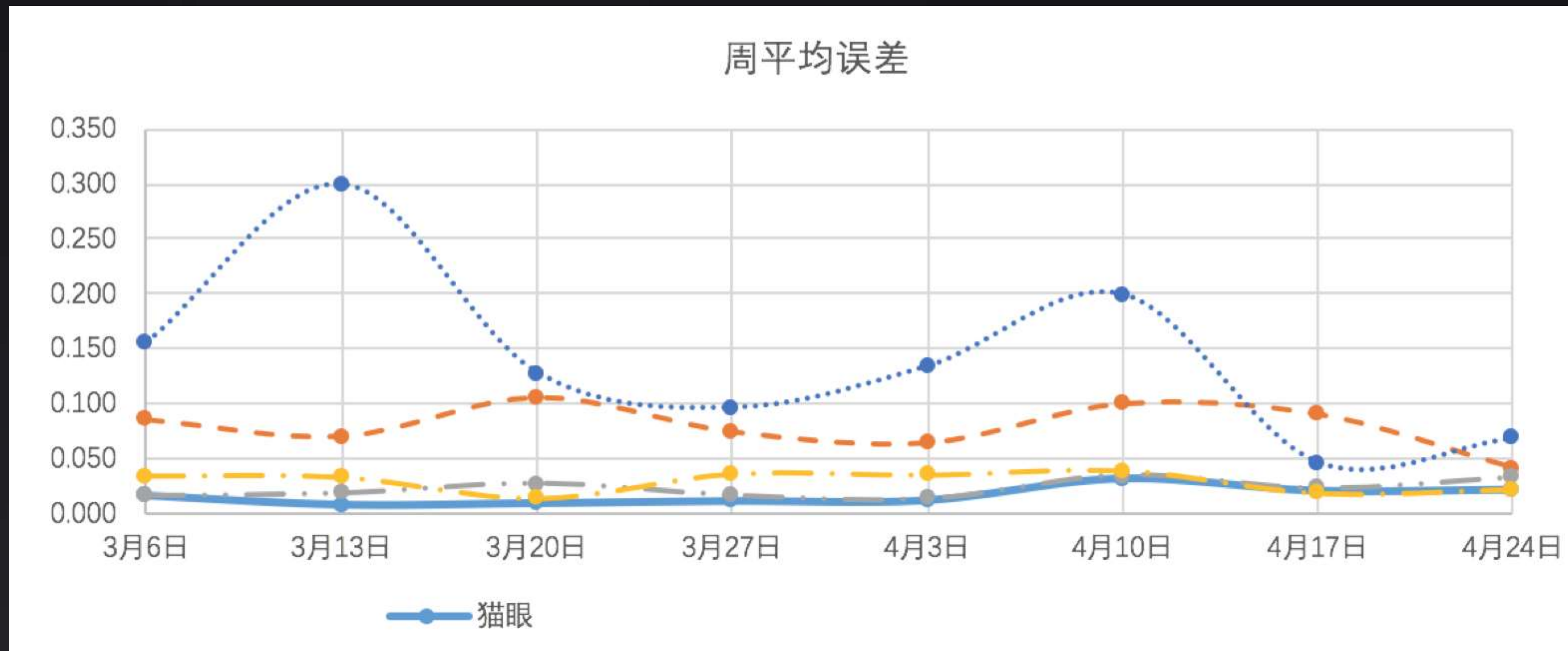
实时票房





实时票房

- 效果对比



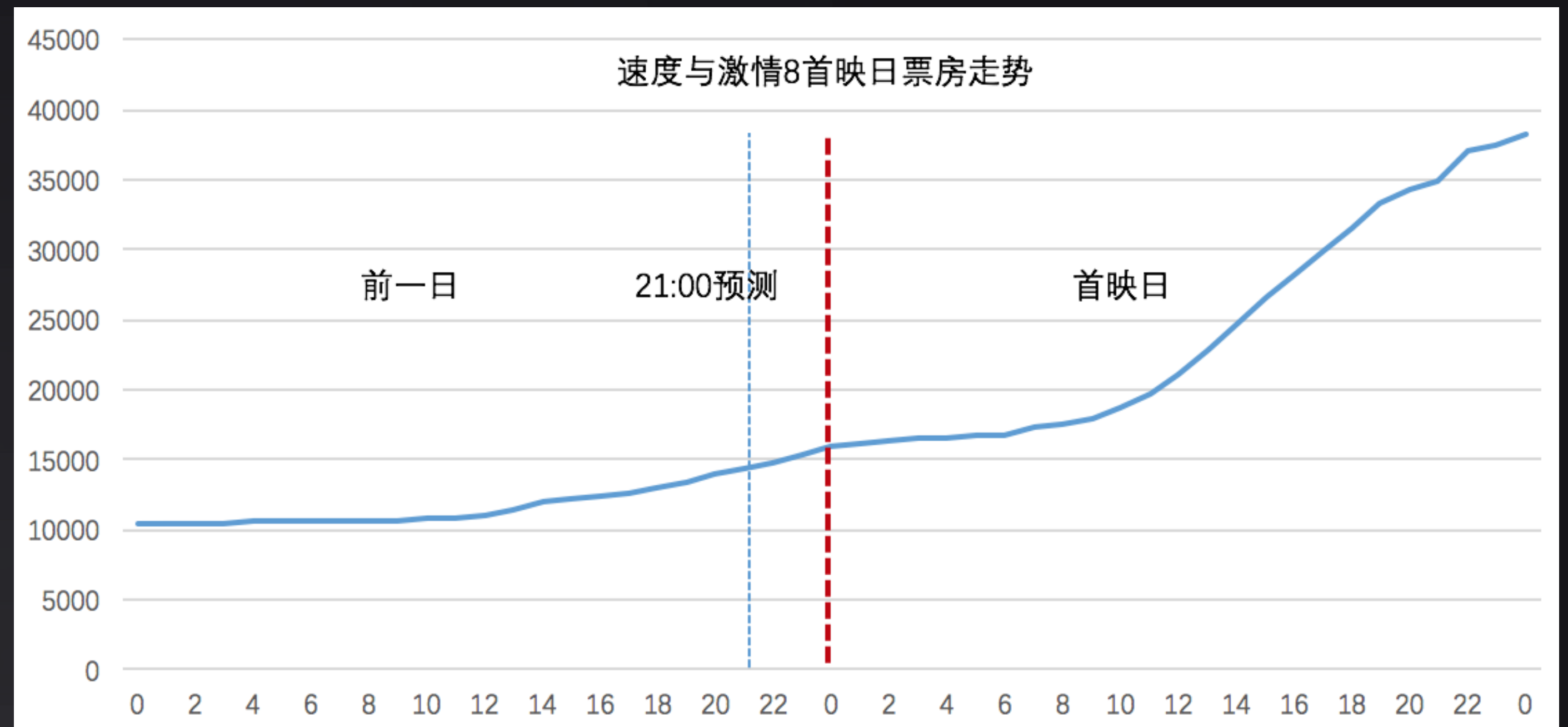


天级票房预测



天级票房预测

- 猫眼的优势
 - 在线售票系统：直接的交易数据，预售票房、排片场次等
 - 精准的实时票房为天级票房提供可靠的特征
- 问题定义
 - 次日票房预测
 - 前一日21:00
 - 前一日17:00
- 难点
 - 影片数量少
 - 异常数据





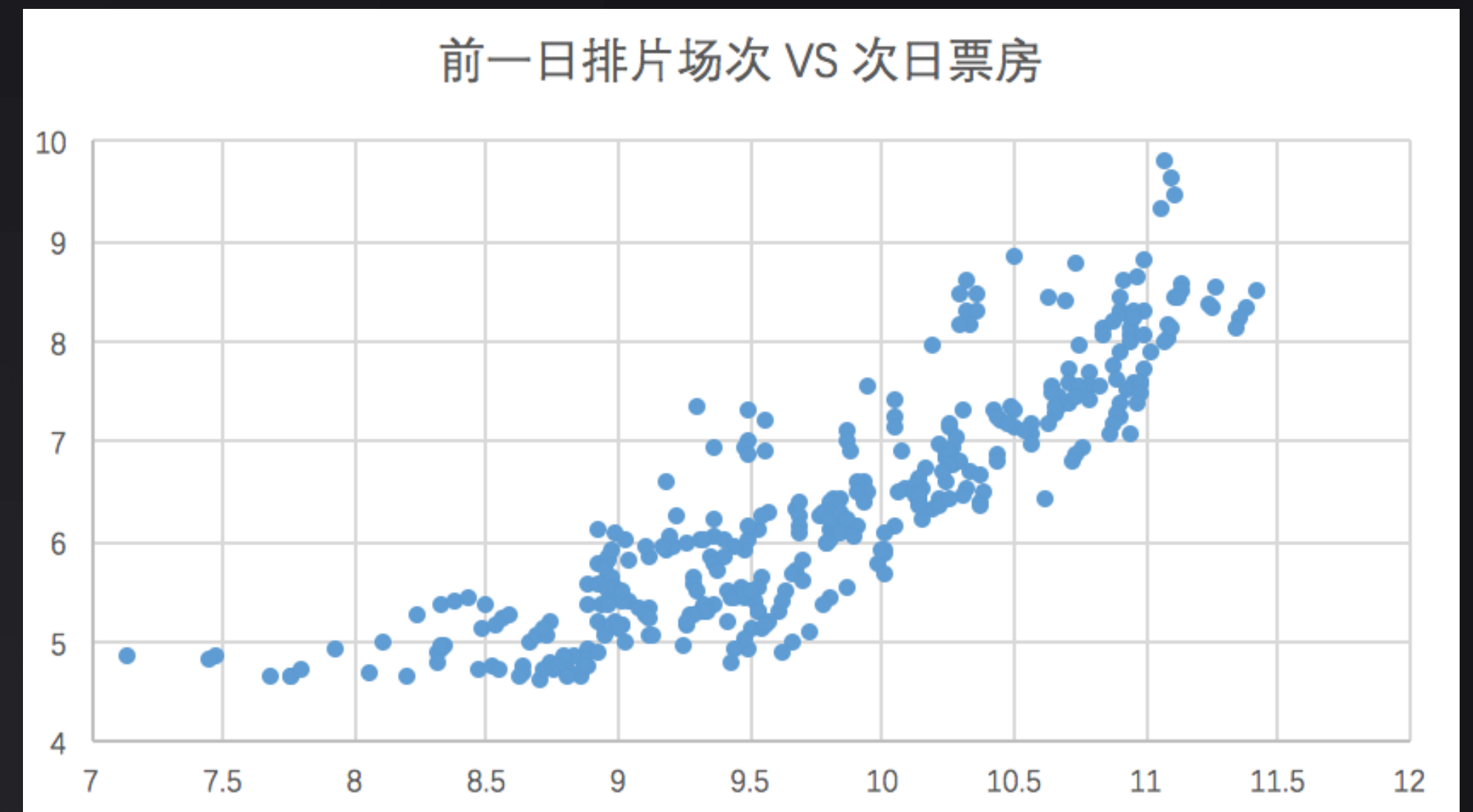
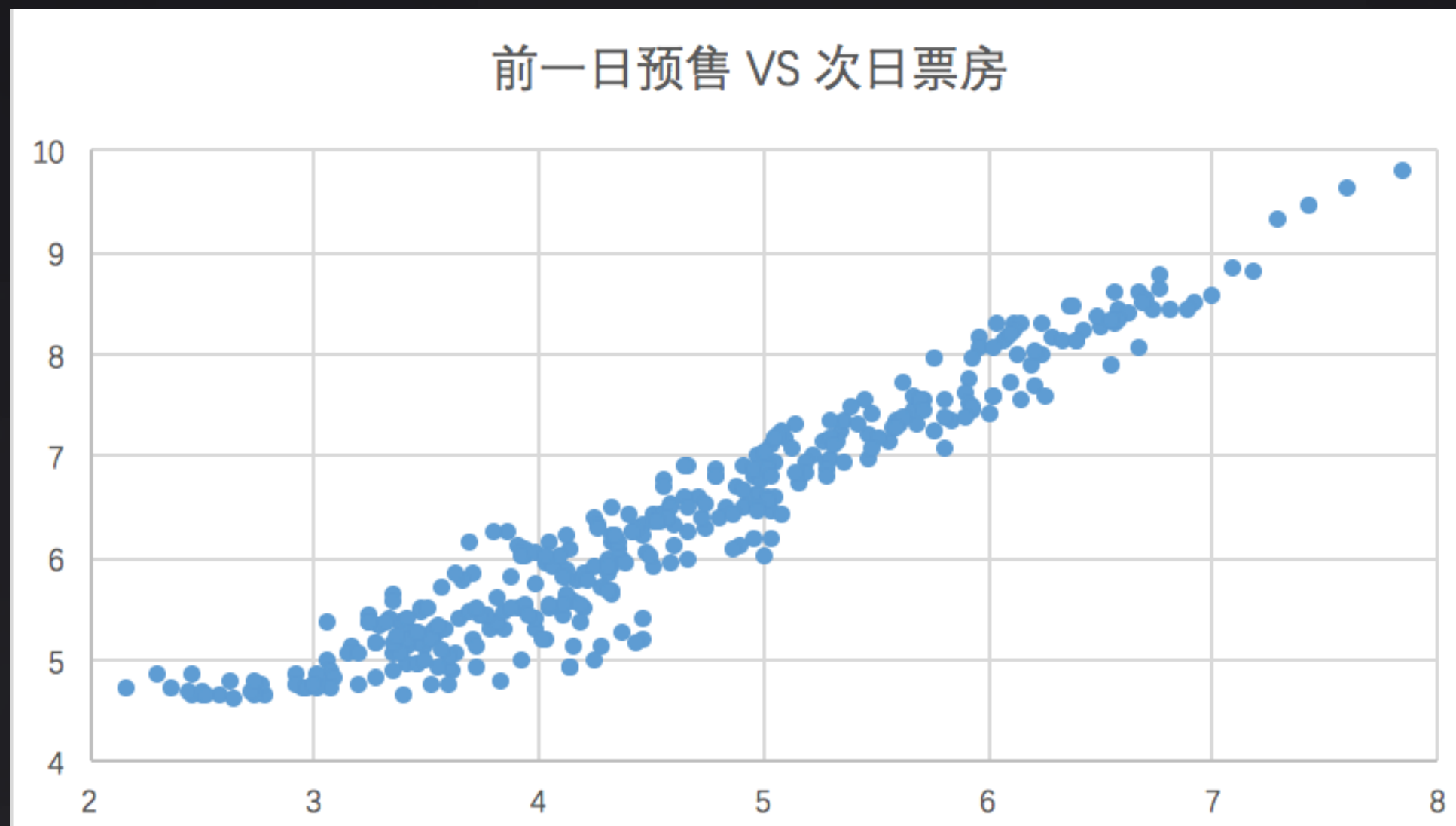
天级票房预测

- 基本假设
 - 票房 = 票价*人次 = 票价*场次*场均人次
 - $\log(\text{票房}) = \log(\text{票价}) + \log(\text{场次}) + \log(\text{场均人次})$
 - $y = w_1 * x_1 + w_2 * x_2 + \dots$



天级票房预测

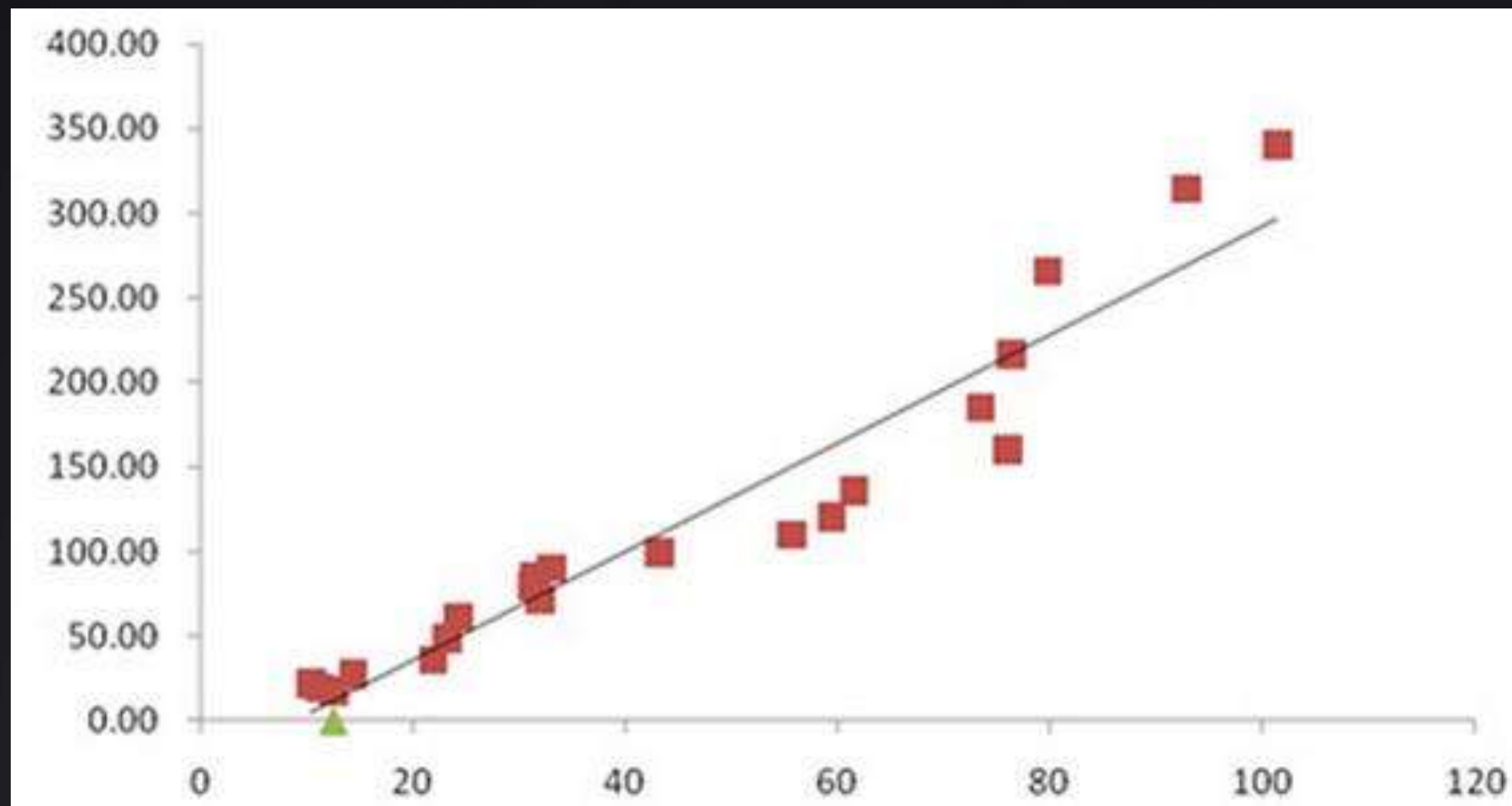
- 基本假设
 - 票房 = 票价*人次 = 票价*场次*场均人次
 - $\log(\text{票房}) = \log(\text{票价}) + \log(\text{场次}) + \log(\text{场均人次})$
 - $y = w_1 * x_1 + w_2 * x_2 + \dots$



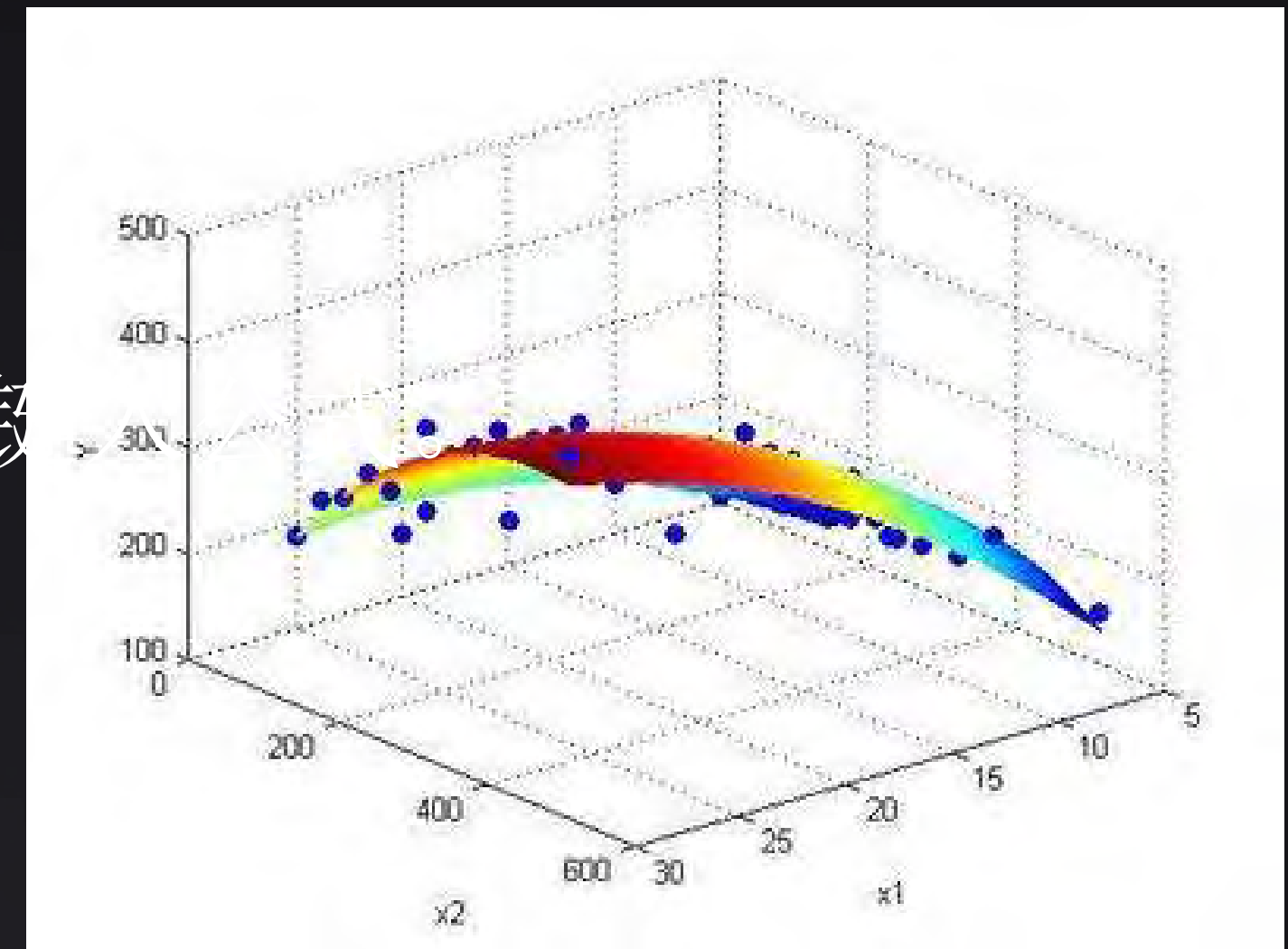


天级票房预测

- 线性回归模型



此处



$$f(X) = W^T X + b = w_1 * x_1 + w_2 * x_2 + \dots + w_m * x_m + b$$

$$\text{loss} = \frac{1}{n} \|Y - W^T X\|^2$$



天级票房预测

- 模型算法
 - 多元线性回归 (LR)
 - GA特征选择算法

| 特征大类 | 特征小类 | 特征 | | | |
|-------|------|----------------------|----------------------|------------------------|---------------------|
| 实时票房 | 日票房 | dailyBox(t-1, t) | dailybox(t-2, t) | dailybox(t-2, t-1) | dailybox(t-1, t-1) |
| | 平均票价 | avgPrice(t-1, t) | avgPrice(t-2, t) | avgPrice(t-2, t-1) | |
| | 上映场次 | totalShow(t-1, t) | totalShow(t-2, t) | totalShow(t-2, t-1) | totalShow(t-1, t-1) |
| | 人次 | totalView(t-1, t) | totalView(t-2, t) | totalView(t-2, t-1) | totalView(t-1, t-1) |
| | 大盘票房 | sumDailyBox(t-1, t) | sumDailyBox(t-2, t) | sumDailyBox(t-2, t-1) | |
| | 大盘人次 | sumTotalView(t-1, t) | sumTotalView(t-2, t) | sumTotalView(t-2, t-1) | |
| 猫眼订单 | 日票房 | maoyanOrder(t-1,t) | maoyanOrder(t-1,t-1) | | |
| 黄金时间 | 上映上次 | hotShow(t-1,t) | | | |
| | 座位数 | hotSeat(t-1,t) | | | |
| 节假日特征 | 节假日 | holiday(t-1) | holiday(t) | holiday(t+1) | |
| 组合特征 | ... | | | | |



天级票房预测

- 模型拆分

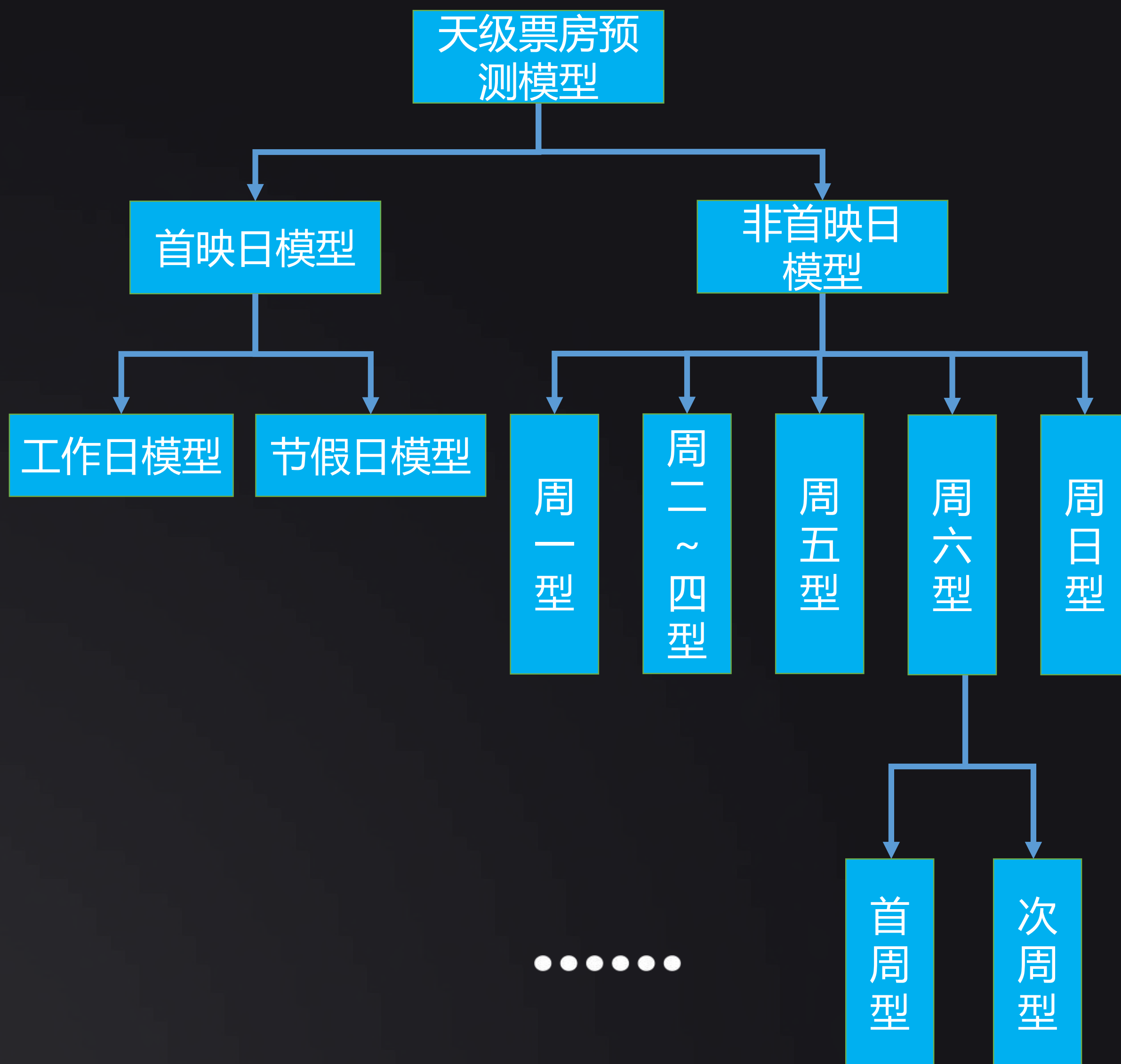




天级票房预测

模型拆分

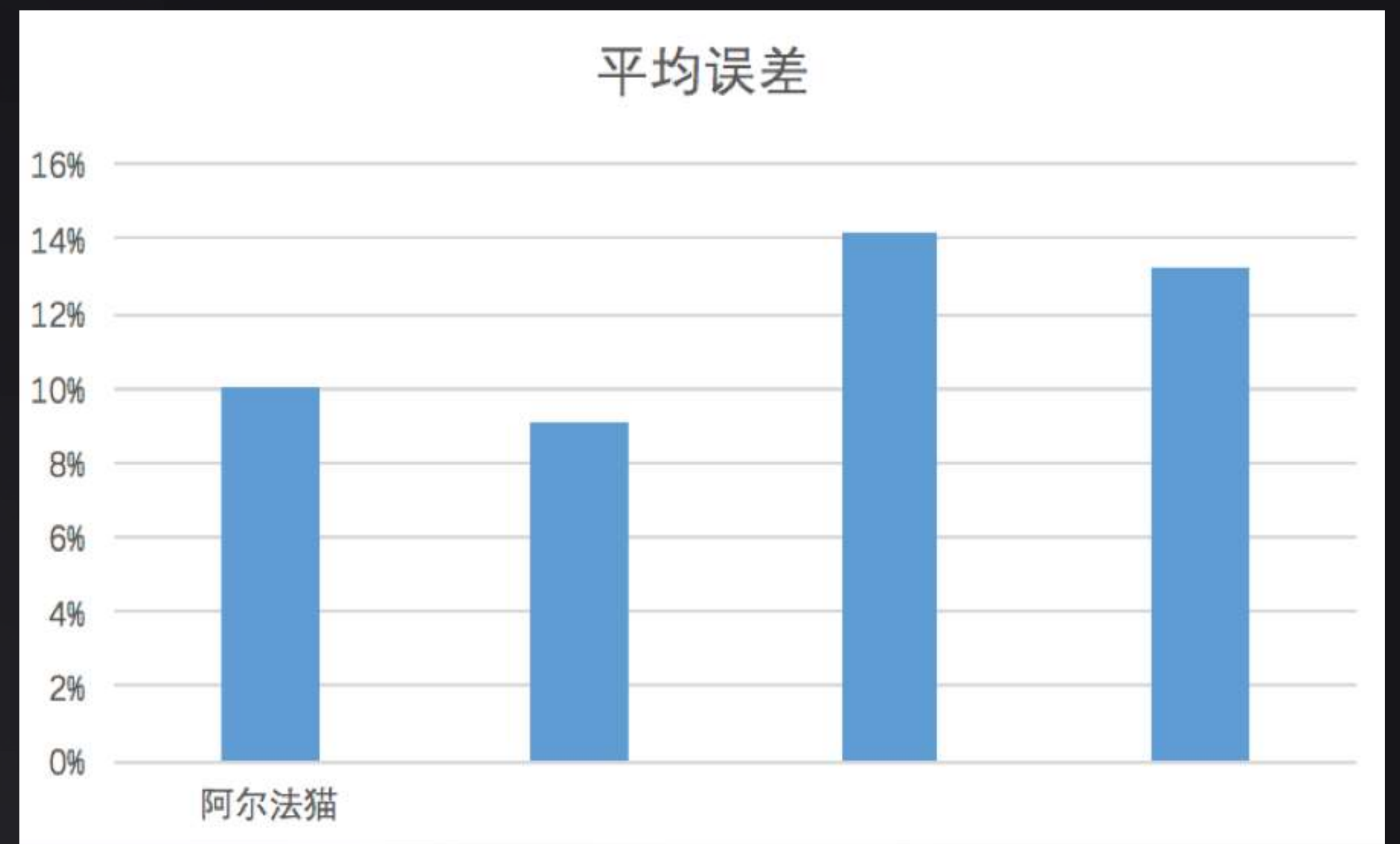
票房衰减规律





天级票房预测

- 效果对比



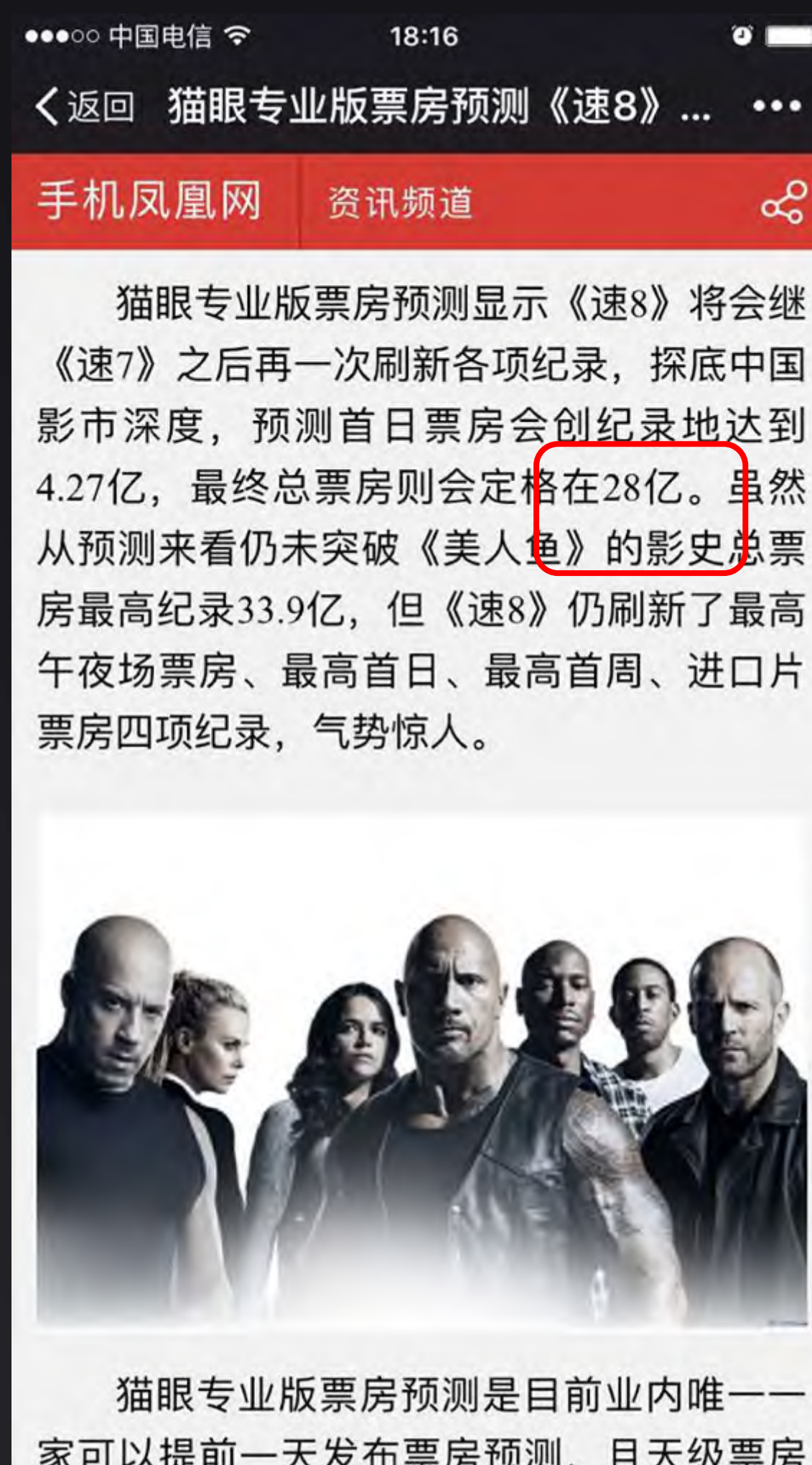
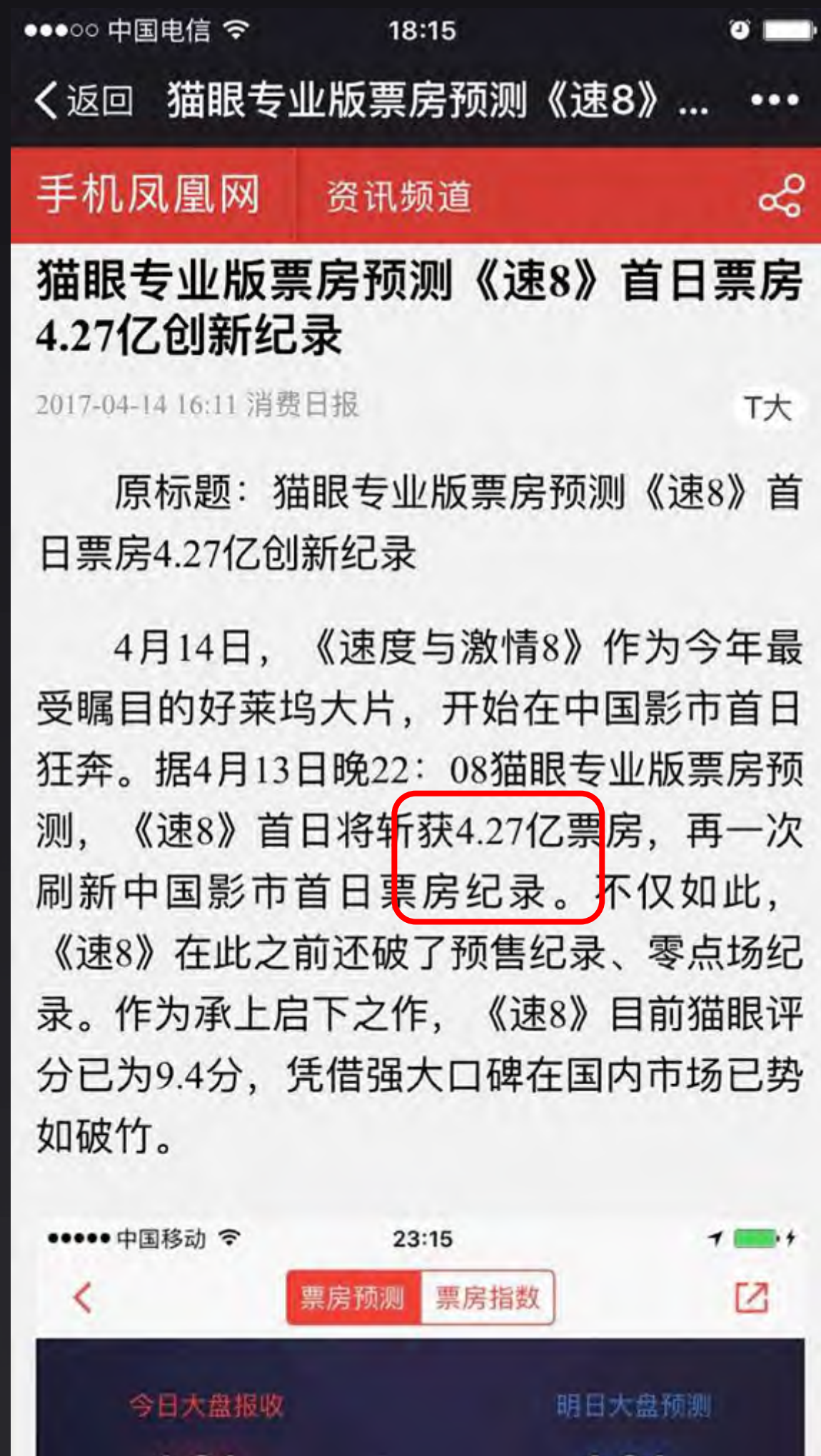
猫眼比竞对提前发布12个小时

R10=预测误差在10%以内的样本占比



天级票房预测

首映日准
确度98%



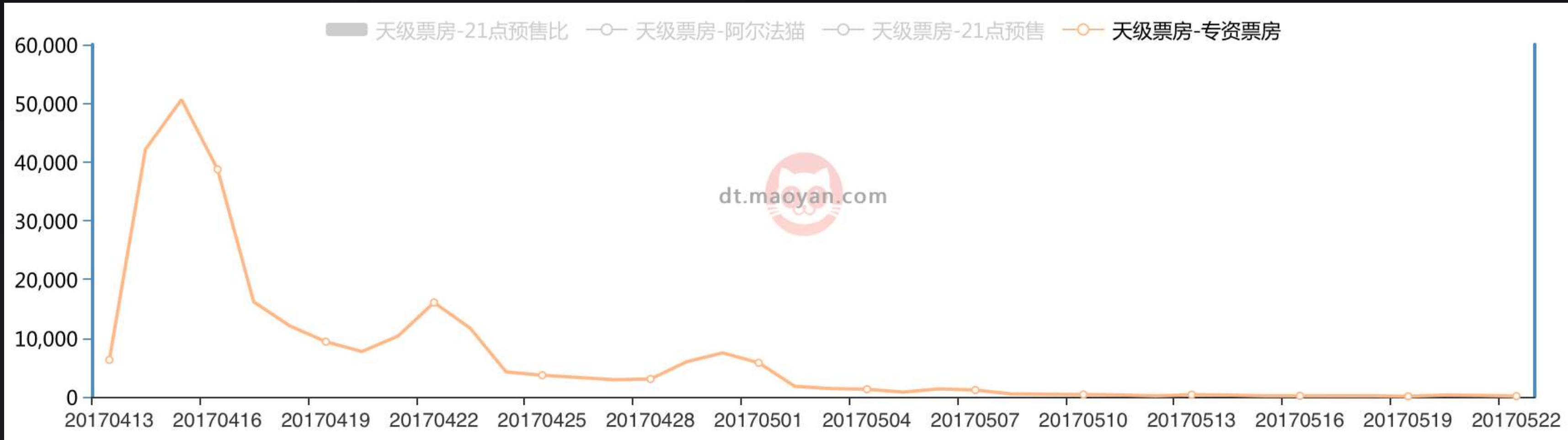
总票房准
确度95%



总票房预测



总票房预测





总票房预测

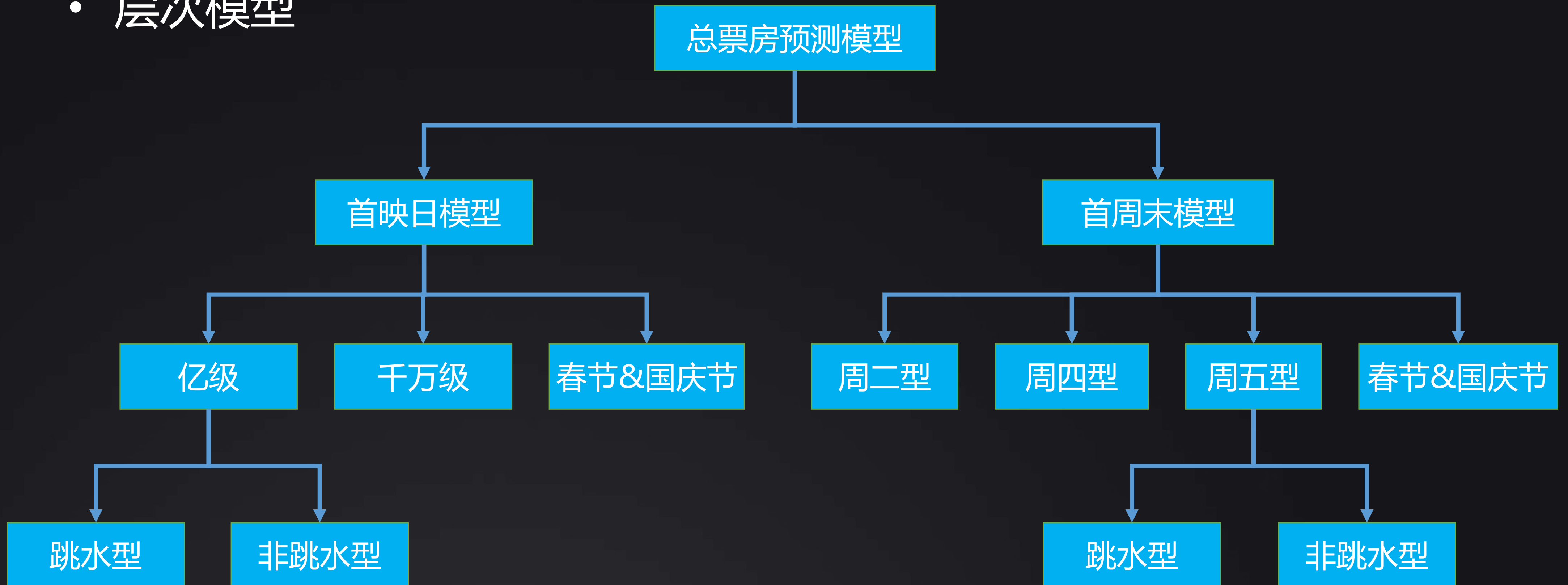
- 问题定义





总票房预测

- 层次模型





总票房预测

• 模型算法

| 特征大类 | 特征小类 | 特征 | | | | 影响因素 |
|------|---------|---|---|---|----------------|---------|
| 天级票房 | 票房 | box(t=1) | box(t=2) | box(t=3) | box(t=4) | 票房的体量 |
| | 票房走势 | $(\text{box}(t=3) - \text{box}(t=4)) / \text{box}(t=3)$ | $(\text{box}(t=1) - \text{box}(t=4)) / \text{box}(t=1)$ | $(\text{box}(t=3) - \text{box}(t=1)) / \text{box}(t=1)$ | | 口碑 |
| | 排片比 | showRate(t=1) | showRate(t=2) | showRate(t=3) | showRate(t=4) | 票房的体量 |
| | 票房比 | boxRate(t=1) | boxRate(t=2) | boxRate(t=3) | boxRate(t=4) | 票房的体量 |
| | 场均人次 | avgViewer(t=1) | | | avgViewer(t=4) | 票房的体量 |
| 影片属性 | 猫眼想看数 | wishNum | | | | 口碑 |
| | 猫眼评分 | score | | | | 口碑 |
| | 影片类型 | type | | | | |
| 档期 | 上映年月 | month | year | | | 档期 |
| 竞争 | 竞争影片想看数 | rivalWishNum | | | | 其他影片的影响 |
| | 竞争影片数量 | rivalNum | | | | 其他影片的影响 |



总票房预测

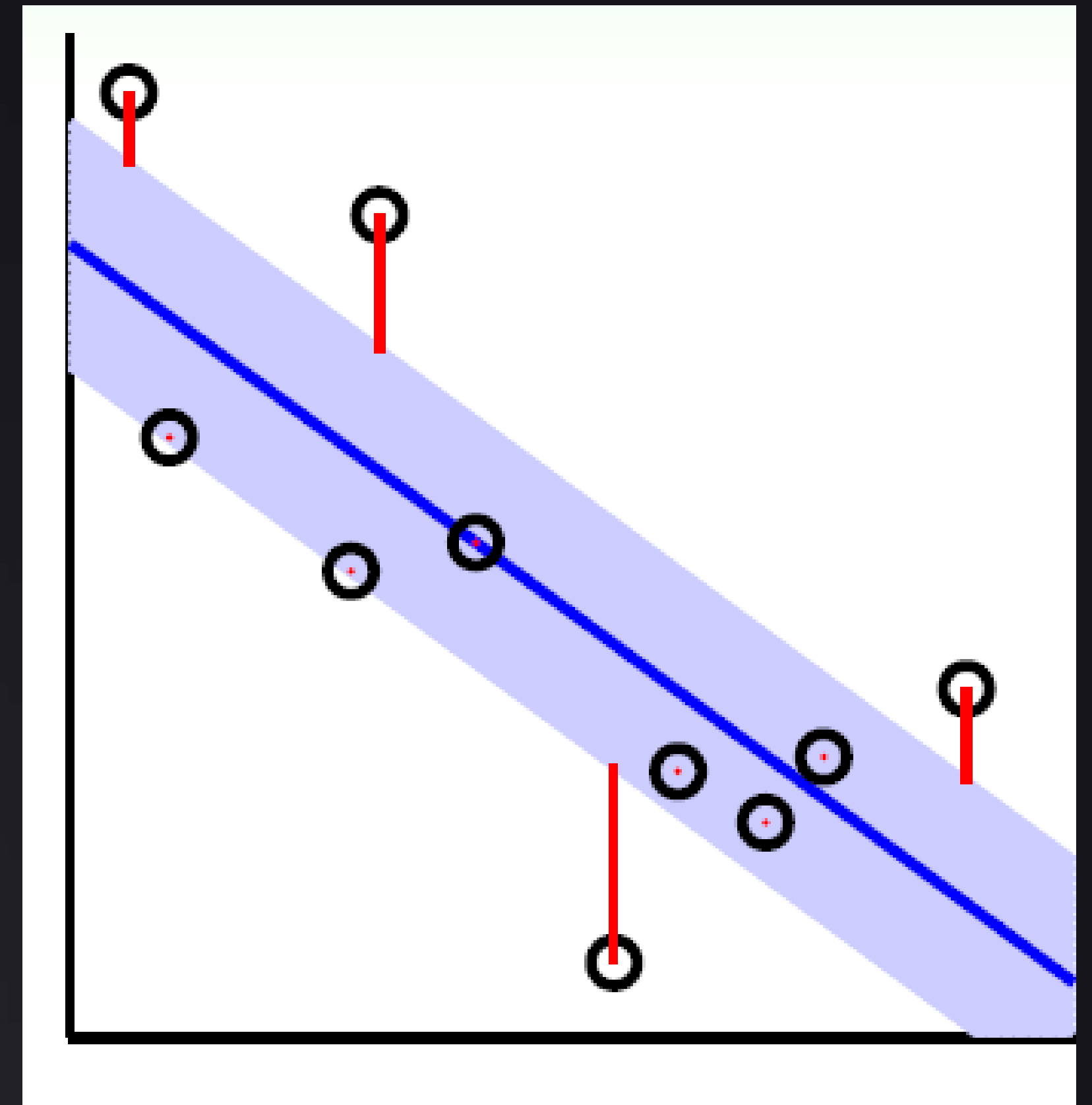
- 支持向量回归模型 (SVR)

- 损失函数：

$$err(x_i, y_i) = \begin{cases} 0 & |y_i - w \bullet \phi(x_i) - b| \leq \epsilon \\ |y_i - w \bullet \phi(x_i) - b| - \epsilon & |y_i - w \bullet \phi(x_i) - b| > \epsilon \end{cases}$$

- 目标函数：

$$\min \frac{1}{2} \|w\|_2^2 \quad s.t. \quad |y_i - w \bullet \phi(x_i) - b| \leq \epsilon (i = 1, 2, \dots, m)$$





总票房预测

• 支持向量回归模型 (SVR)

– 损失函数：

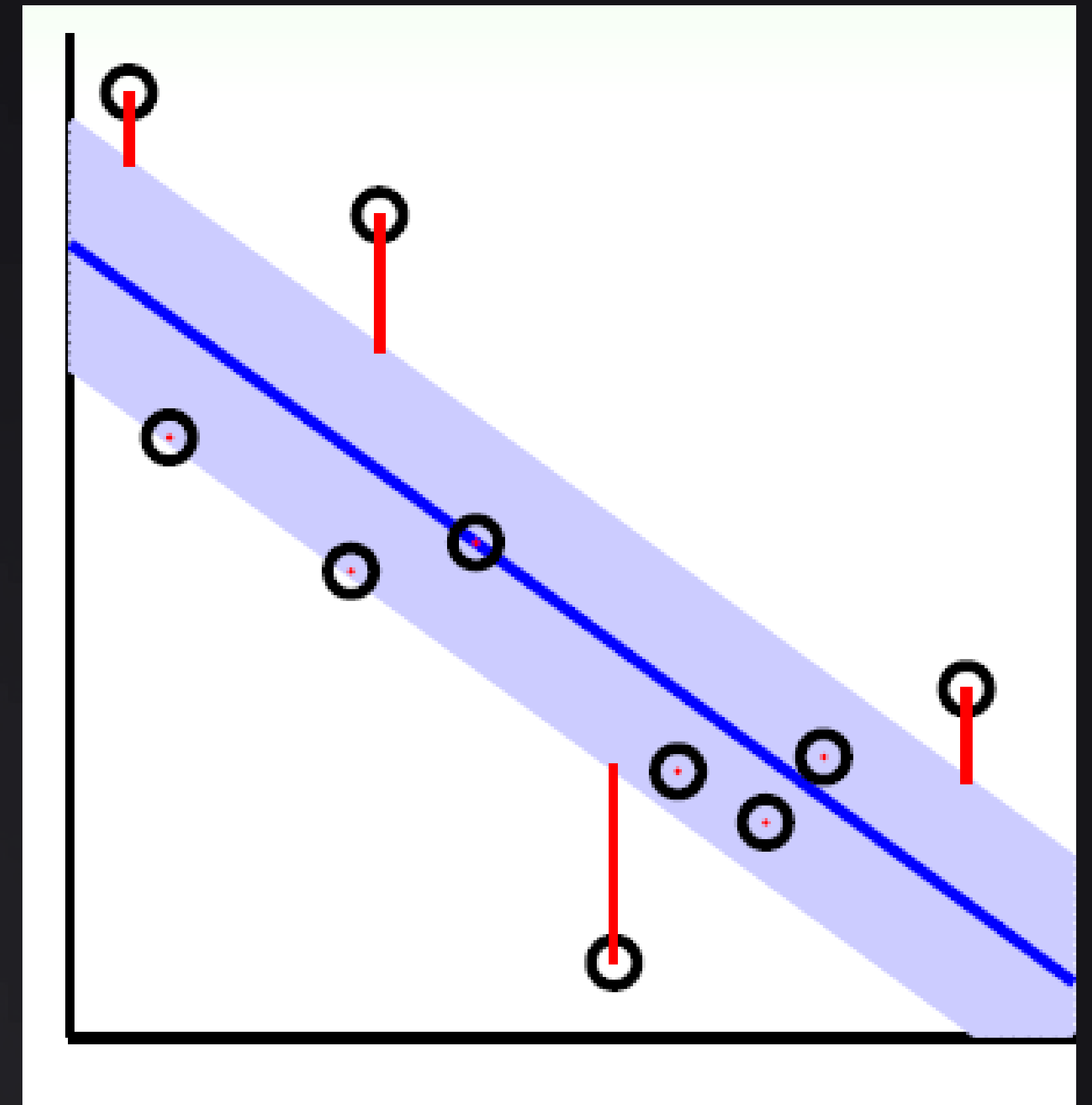
$$\text{err}(x_i, y_i) = \begin{cases} 0 & |y_i - w \cdot \phi(x_i) - b| \leq \epsilon \\ |y_i - w \cdot \phi(x_i) - b| - \epsilon & |y_i - w \cdot \phi(x_i) - b| > \epsilon \end{cases}$$

– 目标函数：

$$\min \frac{1}{2} \|w\|_2^2 \quad \text{s.t.} \quad |y_i - w \cdot \phi(x_i) - b| \leq \epsilon (i = 1, 2, \dots, m)$$

– 优点：

- 特征维度大于样本数时，仍然适用
- 小样本情况下，模型泛化性强
- 非线性核函数，可解决非线性的回归问题

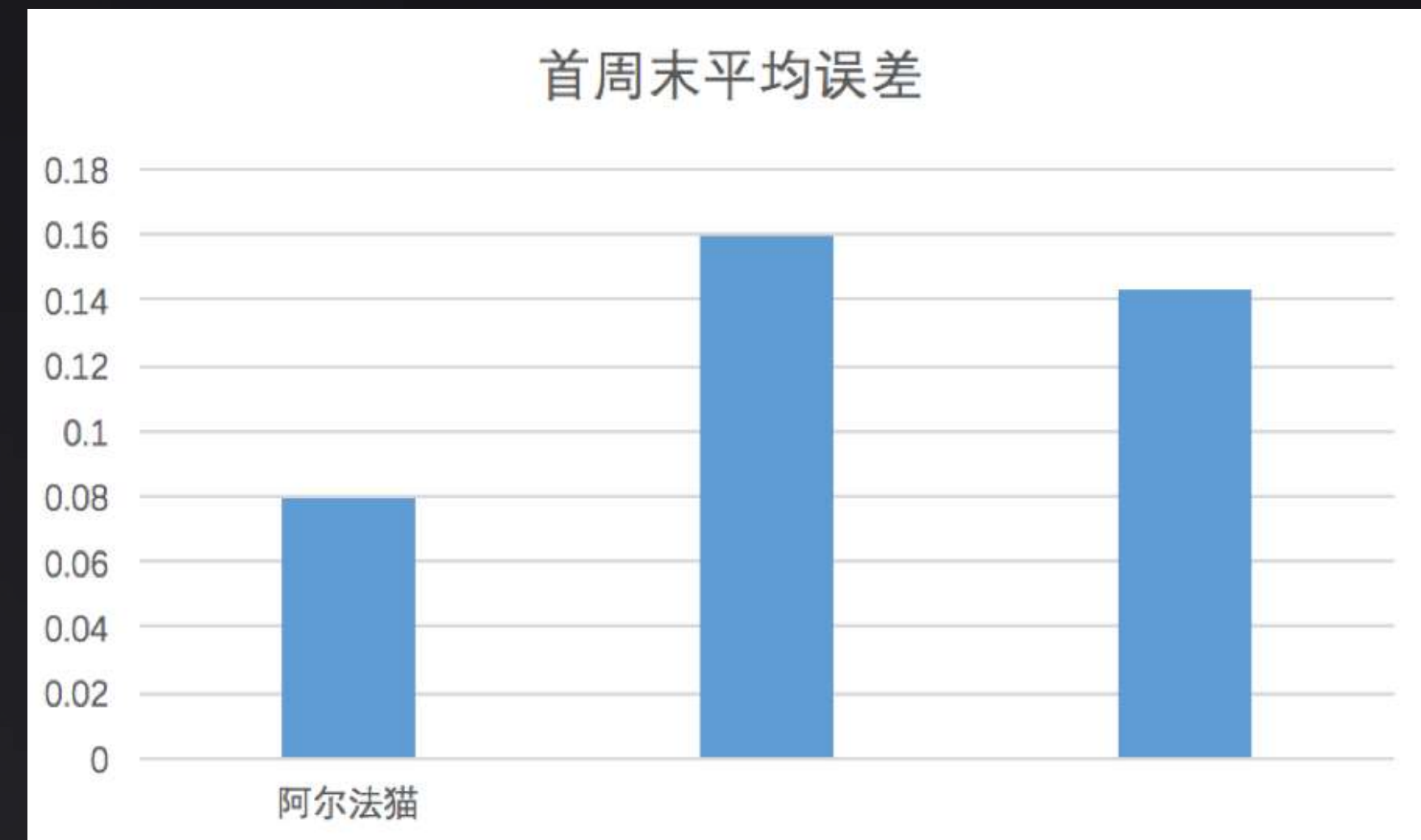
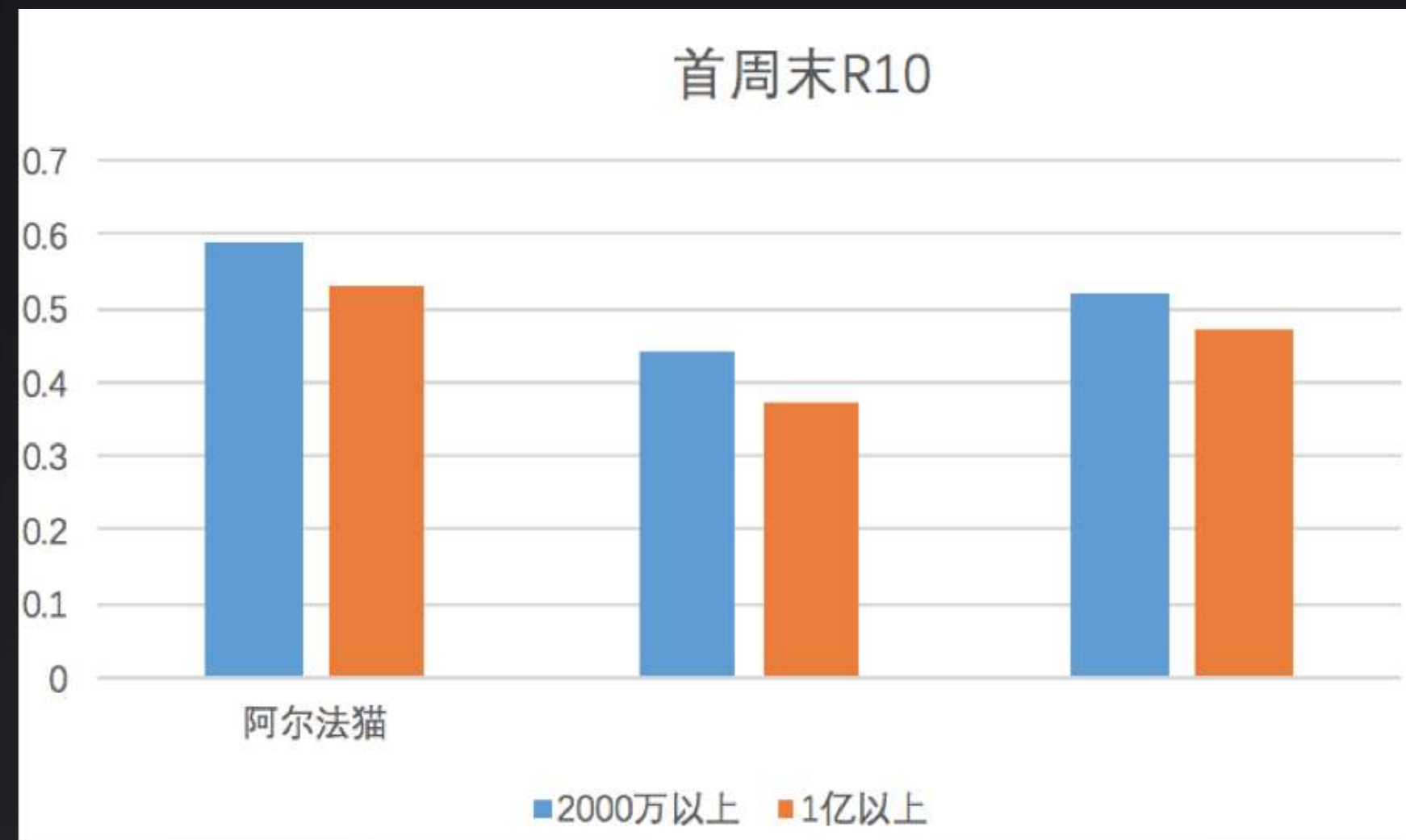




总票房预测

• 效果对比

– 数据统计时间：2017春节~2017.4月底下线，共计27部过2000万，19部过亿

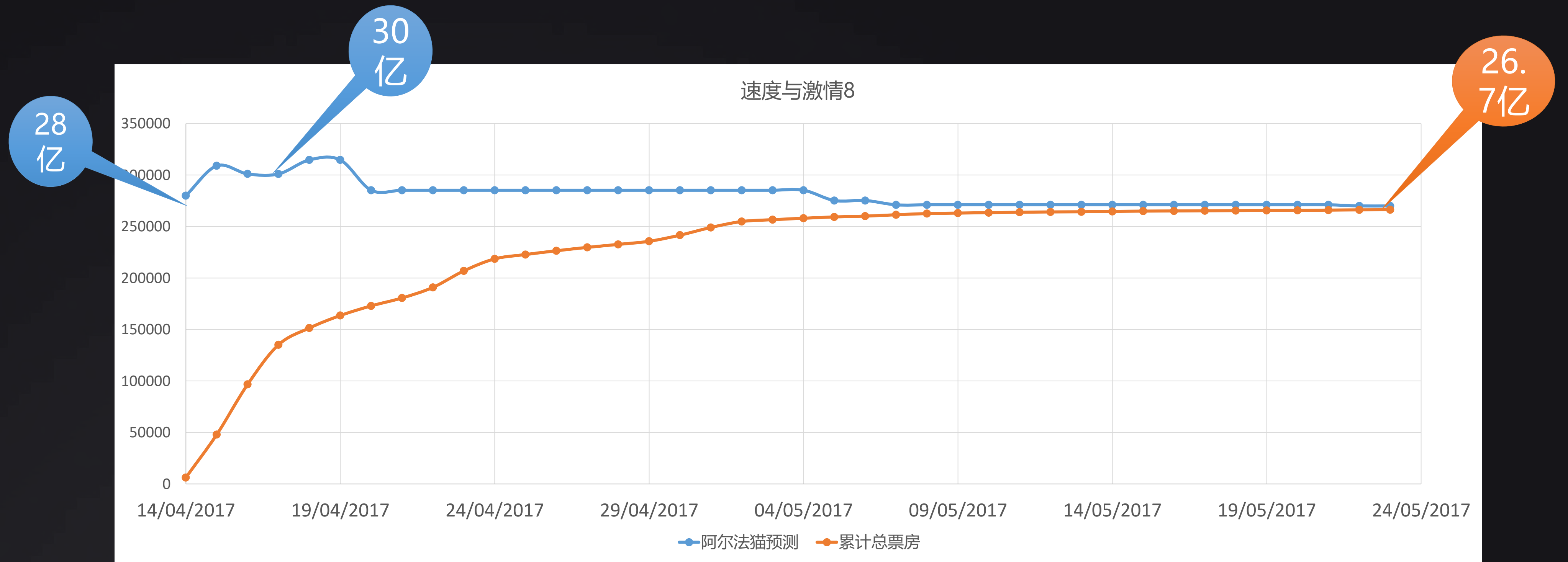


猫眼比竞对提前发布12个小时



总票房预测

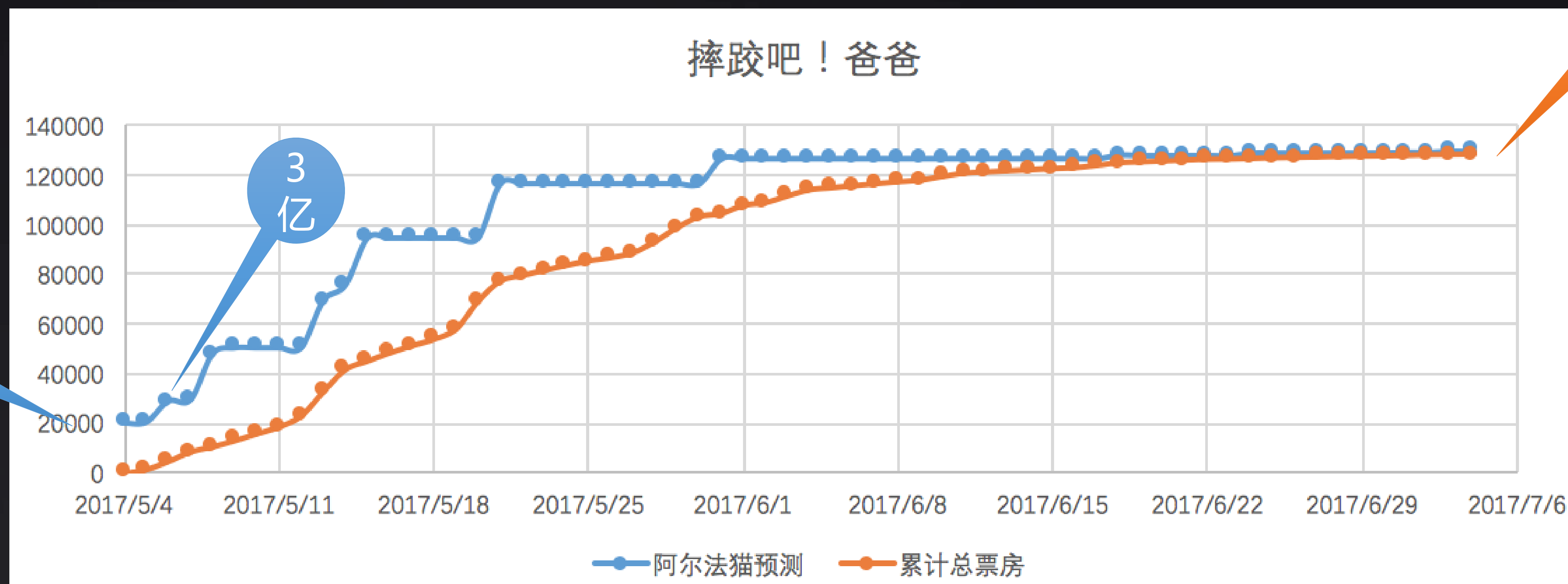
- 效果示例





总票房预测

- 效果示例





内容提纲

- 行业背景
- 技术体系
- 预测算法
- 工作展望



工作展望

票房预测的基础：提高效率、提前时间点、方法创新

- 洞察票房内在规律，不断探索新的方法
- 正在尝试预测票房走势、预售比
- 交互式预测系统：总票房、天级票房
- 天级票房再提前1~2天的小目标
- 总票房预测提前1个月的大目标

应用拓展

- 影片排片：排片助手=>智能排片
- 发行营销：参与营销计划，票补的优化
- 树立票房预测的行业标杆

