

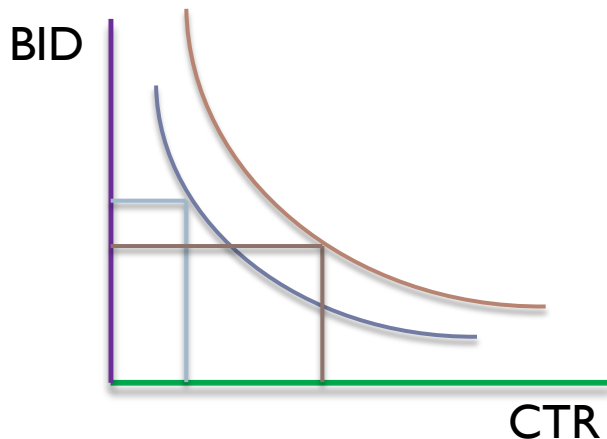
扶翼ECPM提升

回顾与计划

提升eCPM：问题描述

$$ad = \text{Max}_i \{ \text{ctr}_i \text{bid}_i \}$$

- ▶ eCPM：收益的数学期望
- ▶ “提升CTR” => 预估更准的CTR
 - ▶ 排序
 - ▶ 绝对值 vs 相对值
- ▶ “提升bid” => 充分竞价



提升eCPM：技术工作

- ▶ 预估更准的ctr
 - ▶ 引入更多有效特征
 - ▶ 采用更全面的数据
 - ▶ 应用更复杂的模型
- ▶ 转化效果为目标
 - ▶ 引入转化数据
 - ▶ Transfer learning
- ▶ 充分竞价
 - ▶ 预算配速
 - ▶ pacing策略



预估更准的ctr

$$p(\text{click} | \text{user}, \text{ad}, \text{page})$$

▶ 模型误差的来源

- ▶ Bias：模型建模能力不足，underfitting
- ▶ Variance：数据集不足以描述实际情况，overfitting

▶ 针对误差来源解决

- ▶ 降低bias：提升模型复杂度，提高特征描述能力
- ▶ 降低variance：应用更多的训练样本，增加数据的全面性



2015.12-2016.2 有效提升ctr/ecpm的工作

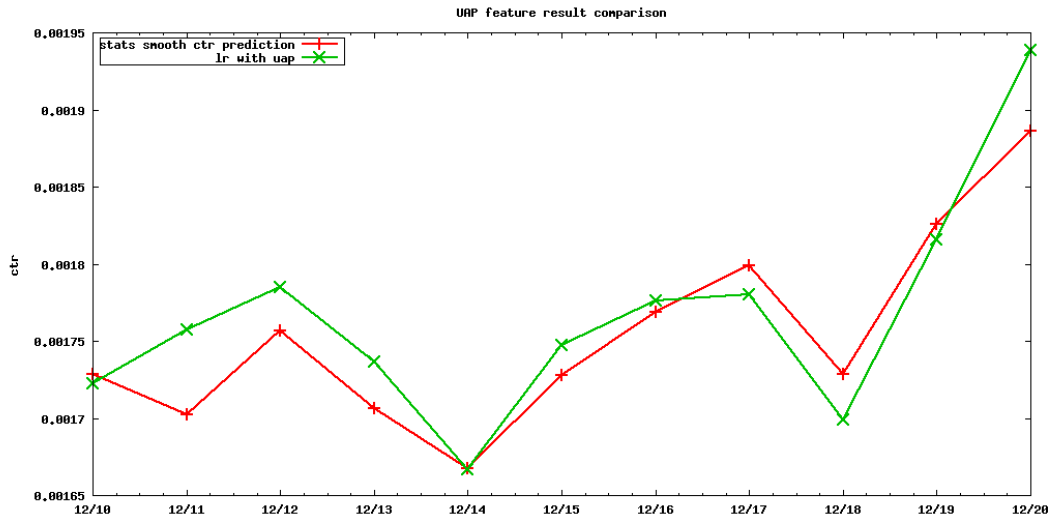
- ▶ 引入有效特征
 - ▶ User agent , pagechannel
- ▶ 增加数据的全面性
 - ▶ 单广告位 => 多广告位数据融合
- ▶ 增强模型建模能力
 - ▶ LR线性 => FM非线性
 - ▶ 引入特征层次关系
 - ▶ 行业-广告主-广告组-投放单-创意，解决新广告问题



引入有效特征

▶ 通过重构ETL

- ▶ User agent, user profile, pagechannel, ...



增加数据的全面性

▶ 原有不足

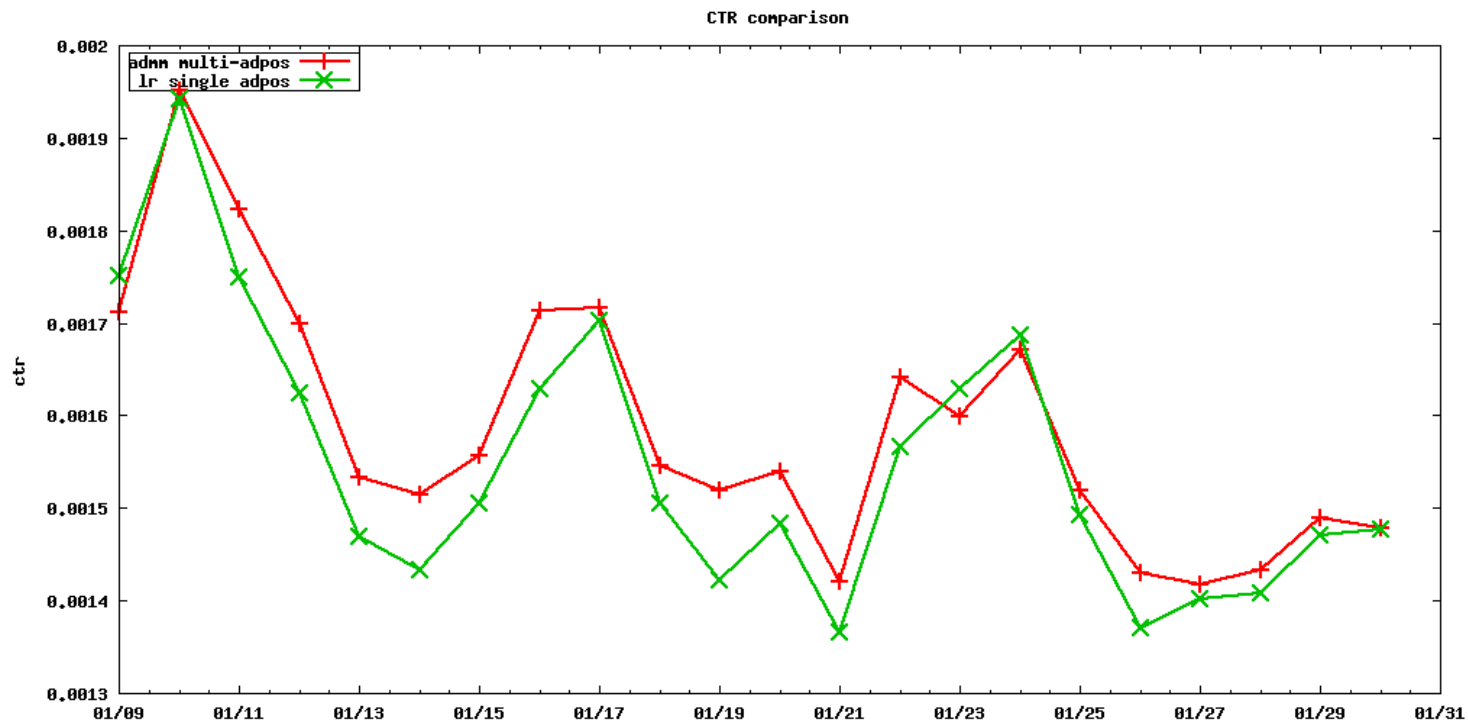
- ▶ 单广告位模型，数据不全面，稀疏性

▶ 解决

- ▶ 分解模型参数 $\vec{w}_t = \vec{w}_0 + \vec{v}_t$
 - ▶ \vec{w}_0 : 多广告位的整体因素, \vec{v}_t : 当前广告位的影响
- ▶ ADMM算法框架解决的多任务学习



增加数据的全面性



增强模型建模能力

▶ 原始特征数据

user

ad

page

▶ LR：线性模型

user

ad

page

user-ad

page-ad

User-page-ad

$$n = n_u + n_a + n_p + n_u n_a + n_p n_a + n_u n_p n_a$$

额外的特征工程与特征选择

▶ FM：非线性模型

user

latent

ad

latent

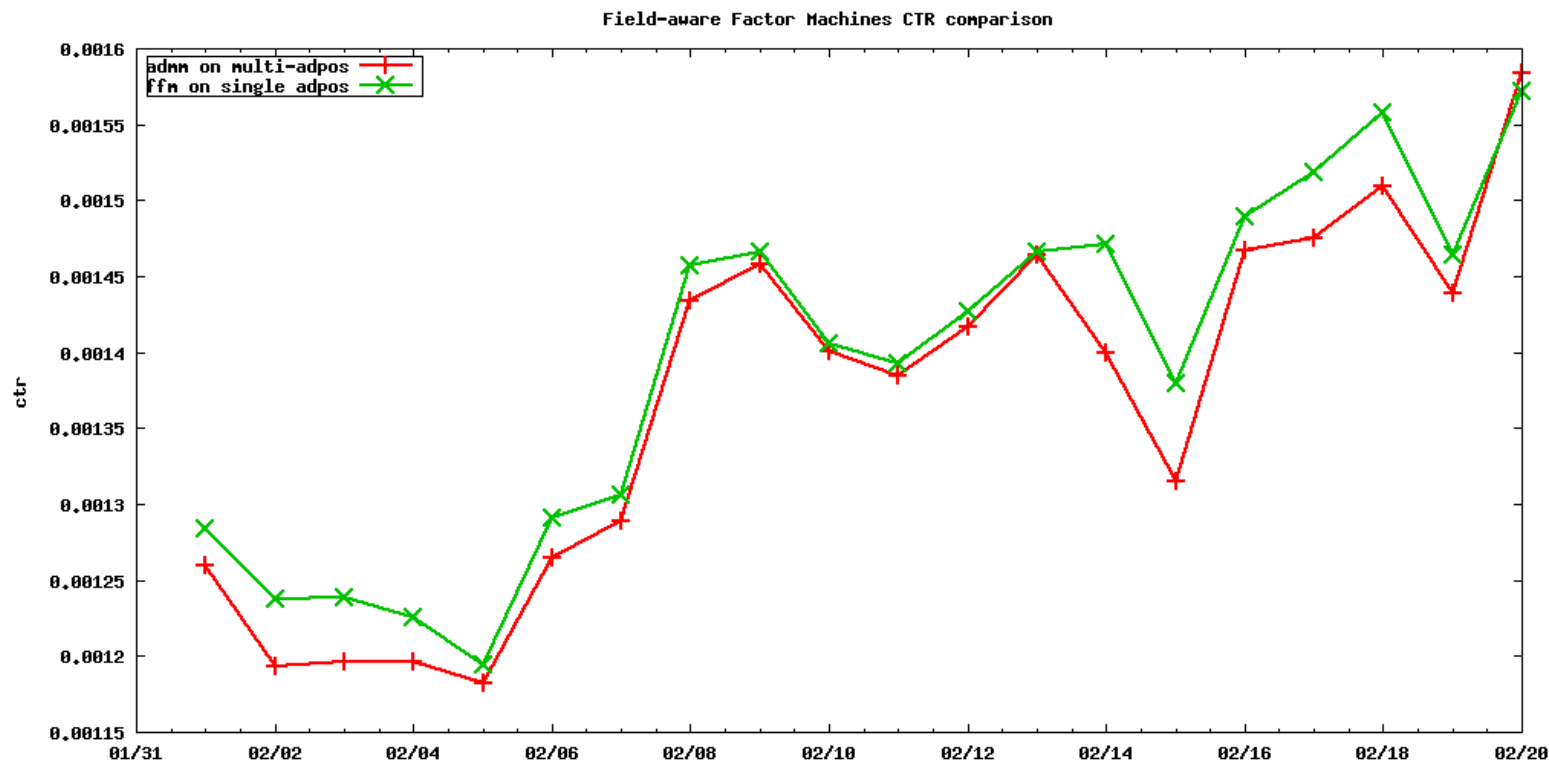
page

latent

$$m = (n_u + n_a + n_p)(1 + k)$$

模型自身的高次项完成特征组合

增强模型建模能力



预估更准的ctr：2016H1工作

▶ 特征引入

- ▶ 标签化的广告特征
- ▶ 创意特征挖掘

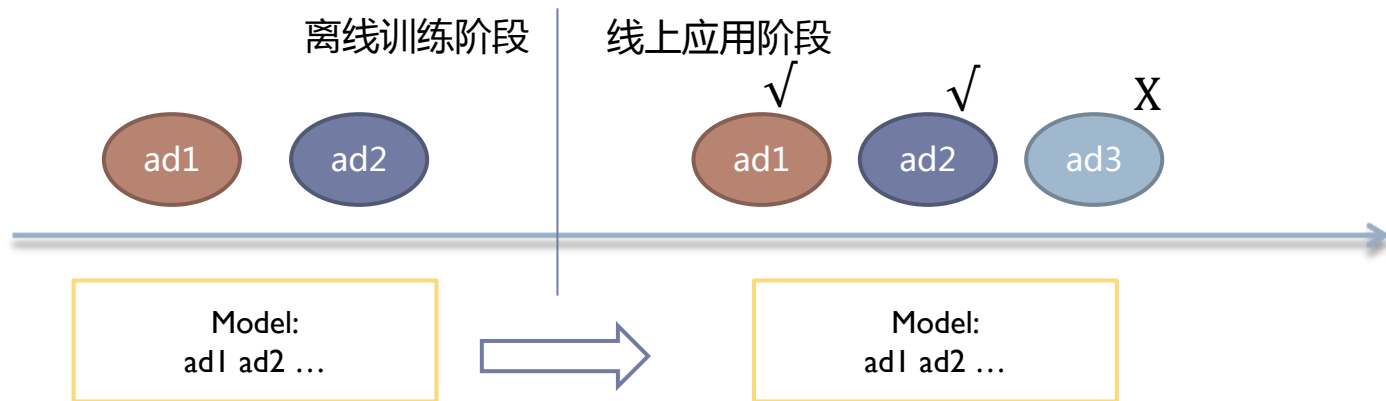
▶ 模型

- ▶ 增强 数据全面性+模型建模能力
- ▶ 建模目标：单纯点击率 => 点击+排序
- ▶ 支持转化数据的引入



2016H1工作与计划：引入有效特征

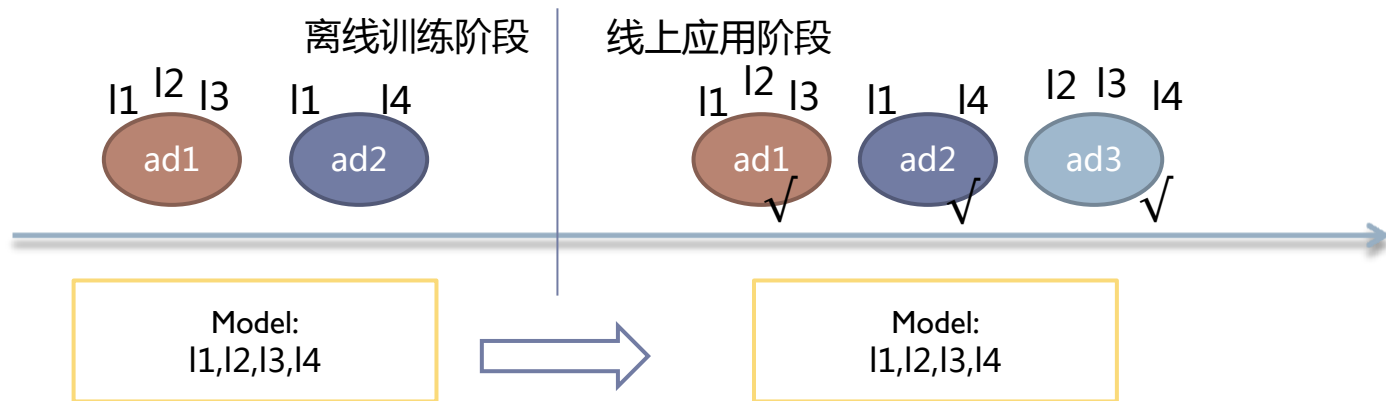
▶ 新广告冷启动问题



- ▶ 长期方案：实时在线模型训练
- ▶ 短期方案：标签化广告/广告主特征; 引入层次关系

2016H1工作与计划：引入有效特征

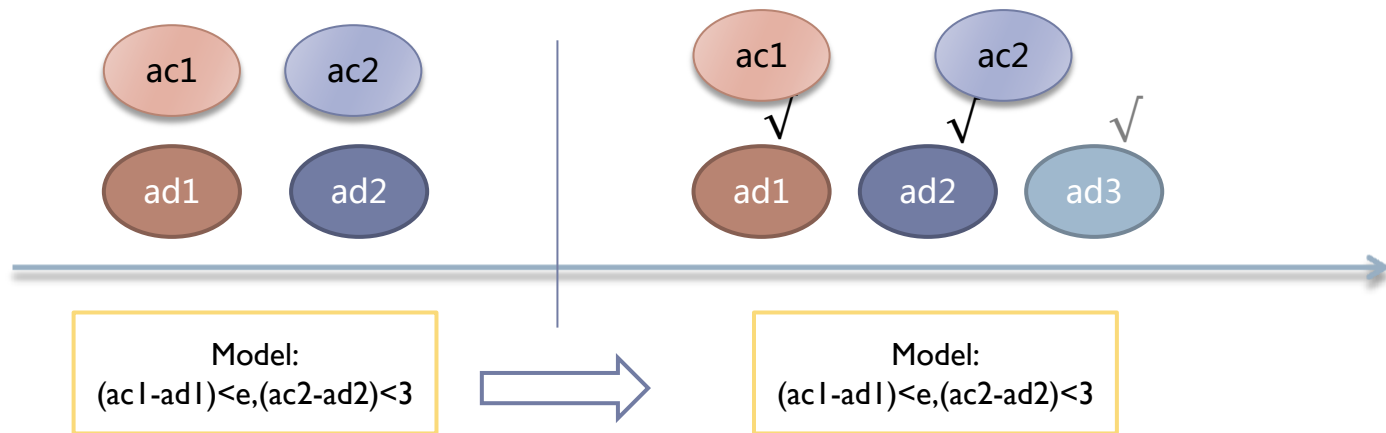
▶ 新广告冷启动问题



- ▶ 长期方案：实时在线模型训练
- ▶ 短期方案：标签化广告/广告主特征; 引入层次关系

2016H1工作与计划：引入有效特征

▶ 新广告冷启动问题

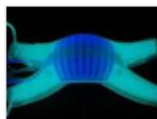


- ▶ 长期方案：实时在线模型训练
- ▶ 短期方案：标签化广告/广告主特征; 引入层次关系

2016H1工作与计划：引入有效特征

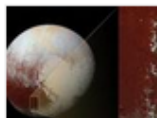
▶ 创意级别特征的挖掘

- ▶ 标题文字
- ▶ 图片
- ▶ 落地页内容
- ▶ 第三方数据
 - ▶ 百度搜索结果数等



科学家参照章鱼发明发光机器人
制造能变色和显示信息的软体机器人，将用于显示器。

3



NASA在冥王星山顶发现“积雪”
在冥王星一片巨大的深色区域内，奇异群山上覆盖着白雪。

8



原价1980智能手表 今天特价
已有3913人购买 先验货再付款

赞助 2



夏普或明日就归入鸿海旗下签约
郭台铭在日本逗留到4日，指挥对夏普财务状况的细查工作。

7



互联网金融迎来规范发展年
今年的政府工作报告提到互联网金融时用词是“规范发展”。

2

2016H1工作与计划：数据+模型提升（1）

▶ ADMM多广告位+FFM模型训练

- ▶ 已有的两项工作的合并
- ▶ 数据全面性+模型建模能力
- ▶ 目标eCPM提升1.5%



2016H1工作与计划：数据+模型提升（2）

建模目标的调整 目标ecpm提升 2.5%

▶ 点击概率的预估

$$p(\text{click} \mid \text{user}, \text{ad}, \text{page})$$

▶ 行业惯例，标准方案

▶ 要求概率模型足够精准



▶ 广告间的pk

$$p(\text{ad}_1 > \text{ad}_2 \mid \text{user}, \text{page}, \text{ctx})$$

▶ 与线上竞价一致

▶ 增加对排序准确性的建模



$$w = \operatorname{argmin}_w \eta \text{Loss}(w, D) + (1 - \eta) \text{Loss}(w, P)$$

2016H1工作与计划:支持转化数据的引入



▶ 模型上支持引入

- ▶ 作为特征：过于稀疏
- ▶ 作为样本：分布与点击不同
- ▶ Transfer learning方案



提升eCPM：技术工作

- ▶ 预估更准的ctr
 - ▶ 引入更多有效特征
 - ▶ 采用更全面的数据
 - ▶ 应用更复杂的模型
- ▶ 转化效果为目标
 - ▶ 引入转化数据
 - ▶ Transfer learning
- ▶ 充分竞价
 - ▶ 预算配速
 - ▶ pacing策略



保证充分竞价：问题与目标

▶ 问题

- ▶ 优质的广告会快速用光预算，并提前退出
 - ▶ 广告的流程倾斜，客户体验不佳
 - ▶ 竞价程度降低，GSP竞价下扣费会降低
 - ▶ 预算余额到投放控制之间有延迟，造成0扣费点击

▶ 优化目标

- ▶ 平滑预算+提升广告效果，“预算配速”
- ▶ 避免调整bid值



预算配速

$\mathbf{B} = \{B^{(1)}, \dots B^{(T)}\}$ T个时段的预算分配, $\mathbf{C} = \{C^{(1)}, \dots C^{(T)}\}$ T个时段的实际消耗;
给定广告投放单, 对于一次广告展示机会*i*:

设定的参与竞价的概率 r_i

$s_i \sim \text{Bern}(r_i)$ 参与竞价; w_i 竞价胜出; c_i cost; p_i 点击率; $q_i \sim \text{Bern}(p_i)$;

$$C = \sum_i s_i \cdot w_i \cdot c_i ; P = C / \sum_i s_i \cdot w_i \cdot q_i$$

预算控制+收益最大化:

$$\begin{aligned} & \text{minimize} \quad P \\ & \text{subject to} \quad C = B, \sqrt{\frac{1}{T} \sum_{t=1}^T (C^{(t)} - B^{(t)})^2} \leq \varepsilon \end{aligned}$$

线上应用: 实时根据 $\{B^{(1)}, \dots B^{(t)}\}$ 和 $\mathbf{C} = \{C^{(1)}, \dots C^{(t)}\}$ 来**调整** r_i ,
相当于online的求解这一最优化问题

预算配速

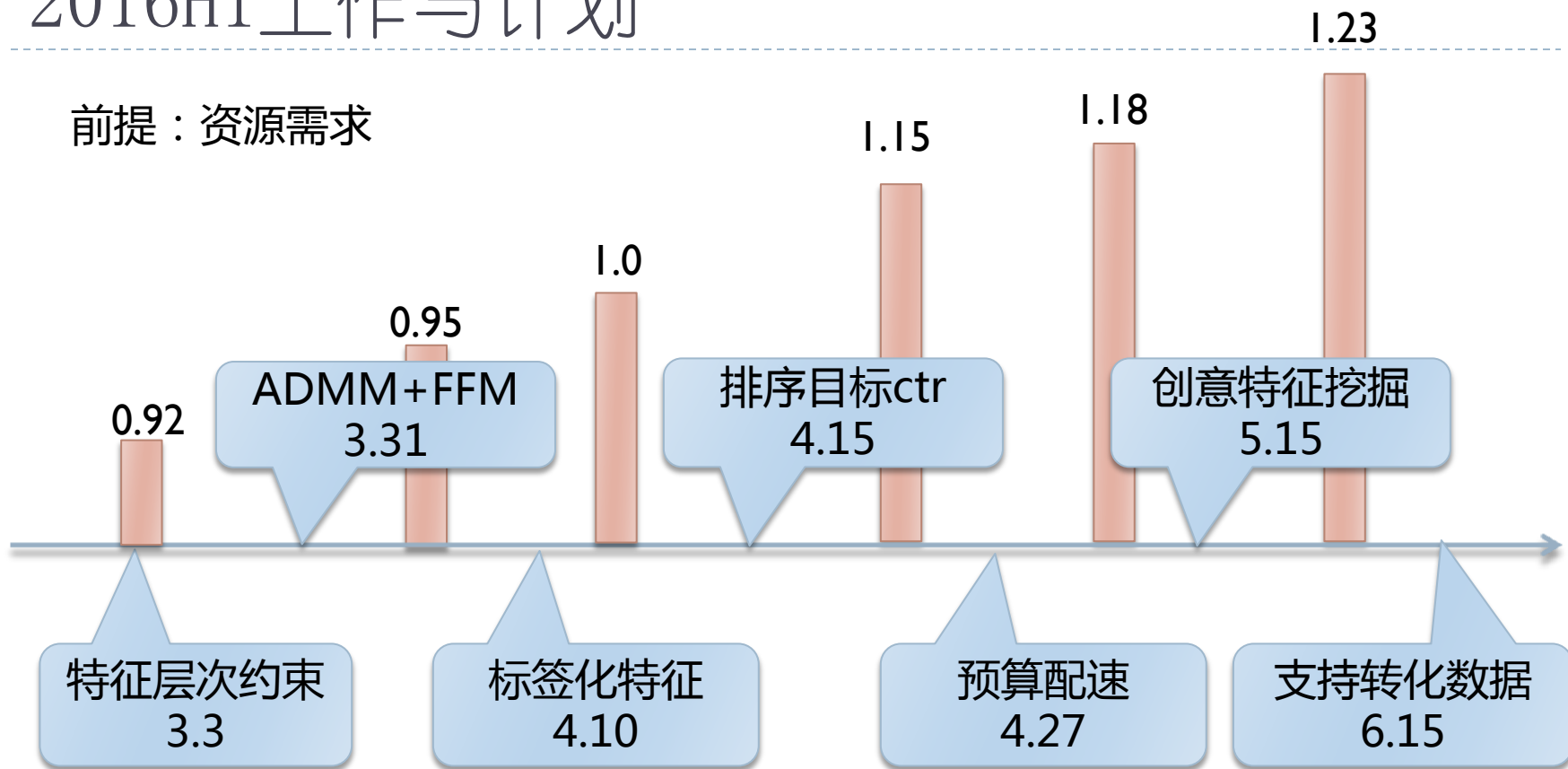
▶ 依赖

- ▶ **数据实时流**统计扣费信息，更新pacing factor
- ▶ **引擎**读取更新的pacing factor r 向量
- ▶ 算法逻辑中应用pacing factor



2016H1工作与计划

前提：资源需求



2016H1工作与资源需求

▶ 人

- ▶ 2~3人，能力≈现有人员，4月底前

▶ 机器

- ▶ 数据中心计算节点扩容
- ▶ 自有parameter server：现有4台，需要~10台

▶ 配合

- ▶ 引擎，数据实时流，日志，产品推动
 - ▶ 前端曝光日志数据
- ▶ 算法服务与算法插件

