



中国联通
应用商店运营中心



基于用户画像的大数据挖掘实践

杨步涛

2014年11月



纲要

- 1 沃商店定位
- 2 沃商店大数据体系架构
- 3 用户画像建设
- 4 个性化推荐
- 5 广告
- 6 用户画像的其他应用实例

渠道聚合

- 2013年中国手机应用分发总量快速上升，其中应用商店的分发量占比超过**80%**；
- TOP10渠道占总分发量的**90%**。

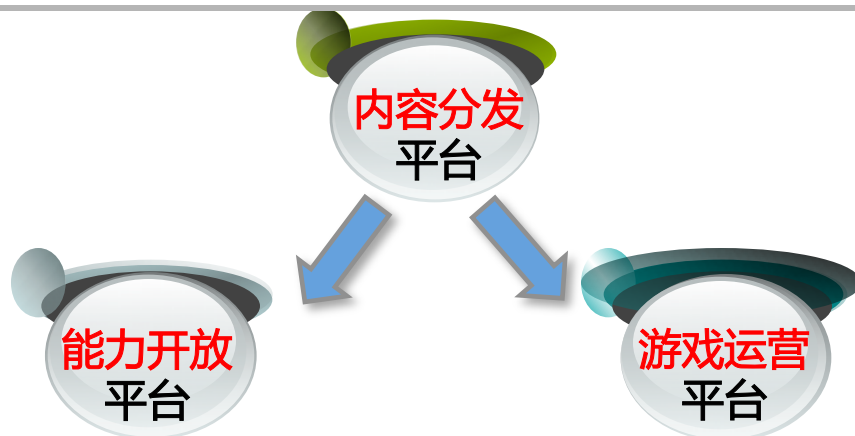
手游爆发

- 2013年中国手机游戏市场近**100亿元**，**2014年预计将达到180亿元**。
- 多款优质手游月流水超过**5000万元**。

支付变革

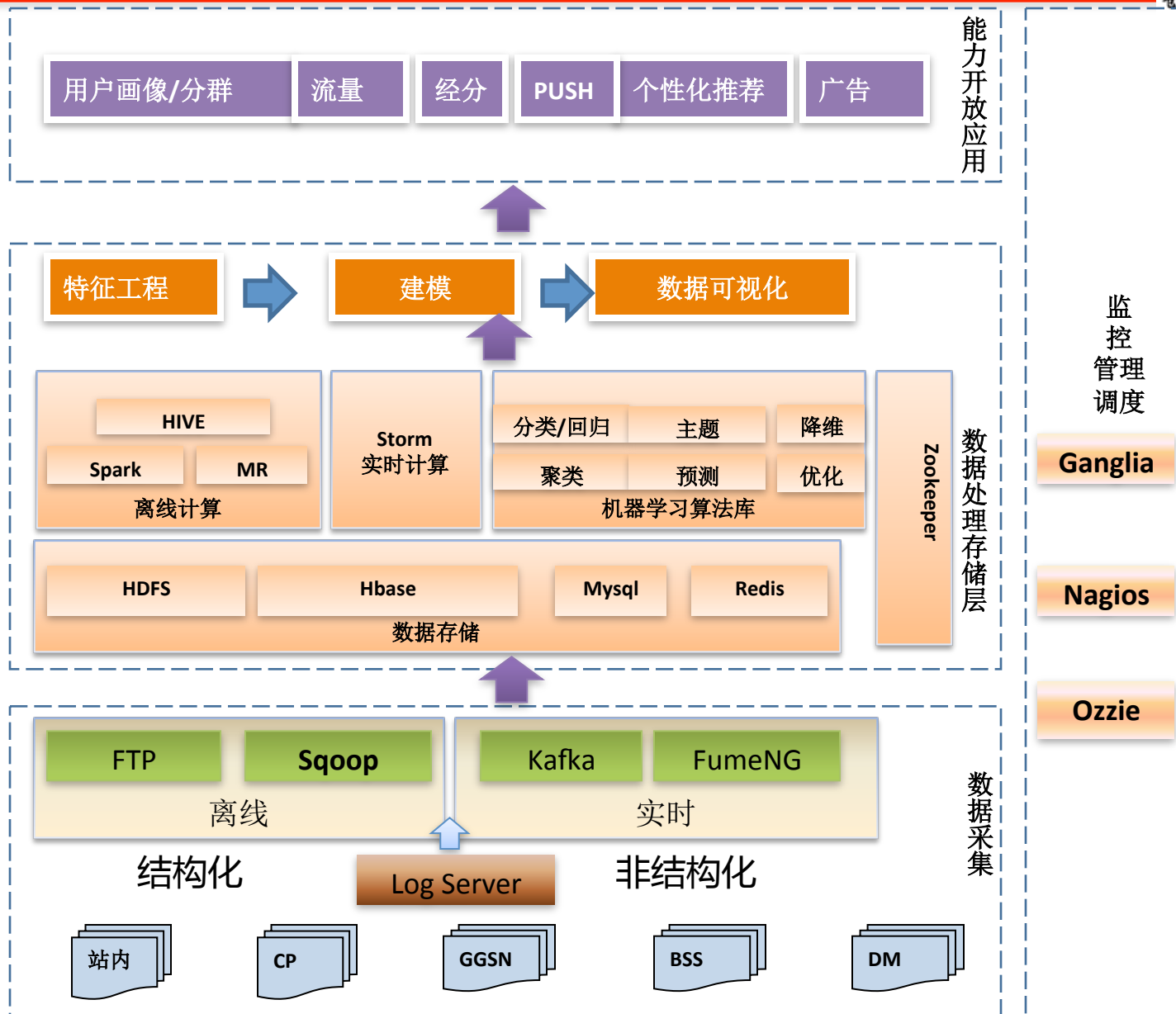
- 运营商通信账户SDK支持APP应用内付费，提升付费转换率，从不足5%提升至**20%**以上。
- 话费支付的便捷优势机遇期短暂仅**1-2年**，移动互联网支付的替代转瞬即至。

沃商店定位





沃商店大数据架构体系

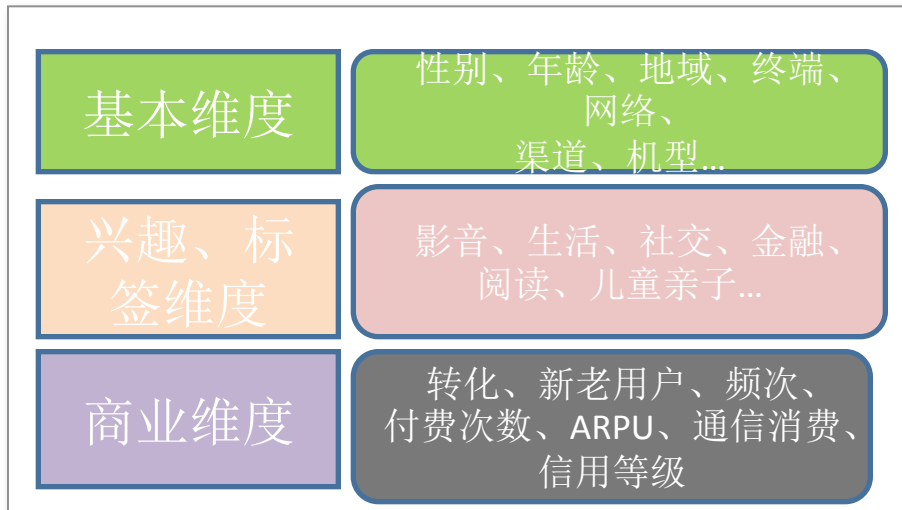
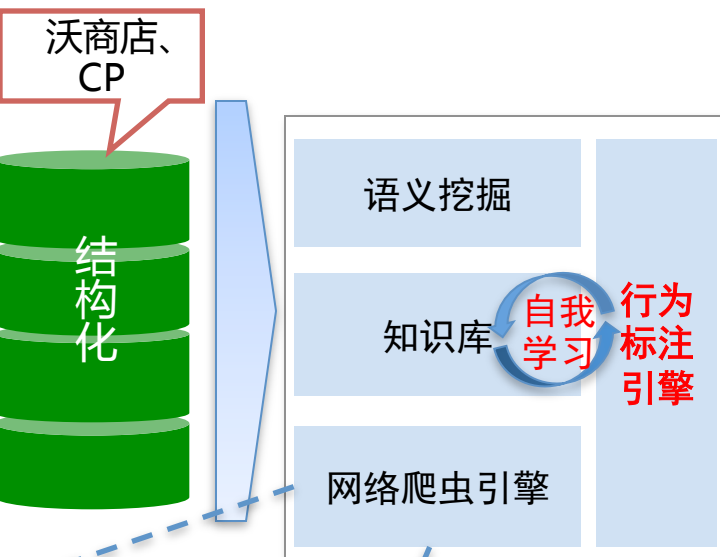




特征工程—用户画像



数据开放



用户画像

非结构化

手机搜狐

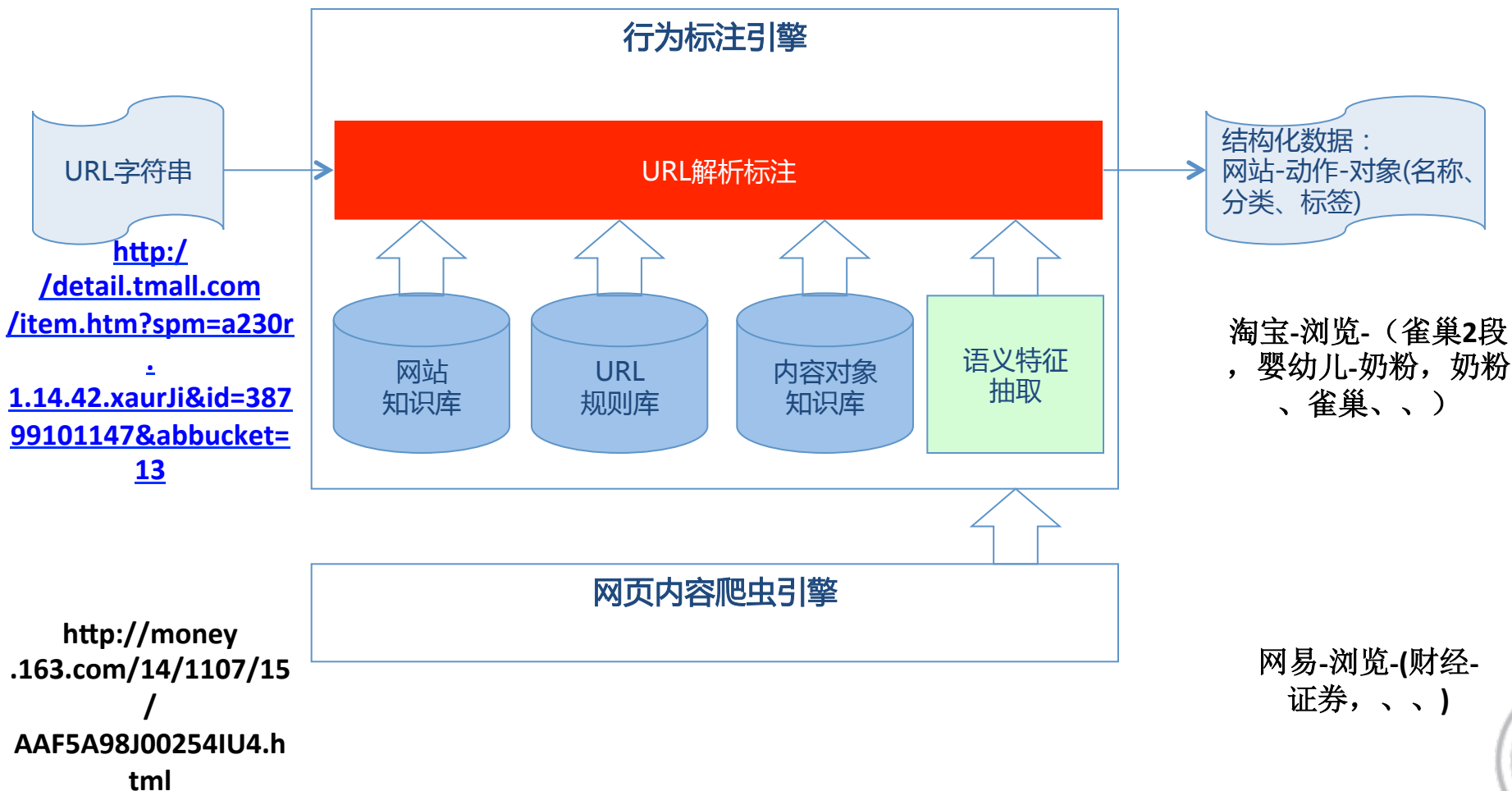


VANCL

淘宝网

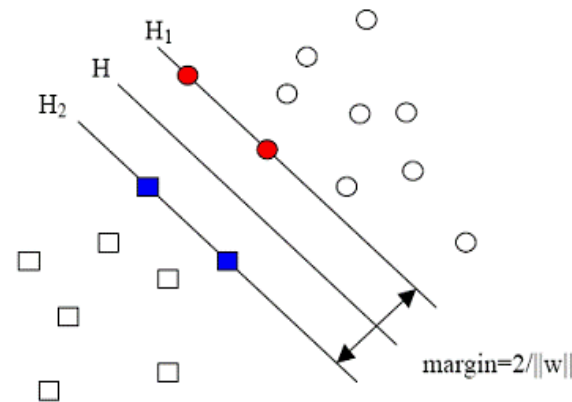


用户画像—行为标注





- 特征预处理、特征筛选(降维)
- 支持向量机SVM
 - 结构风险最优化
 - 非线性(核函数、松弛变量)
 - 1对1方式多分类支持
- 评估：准确率、召回率、F1

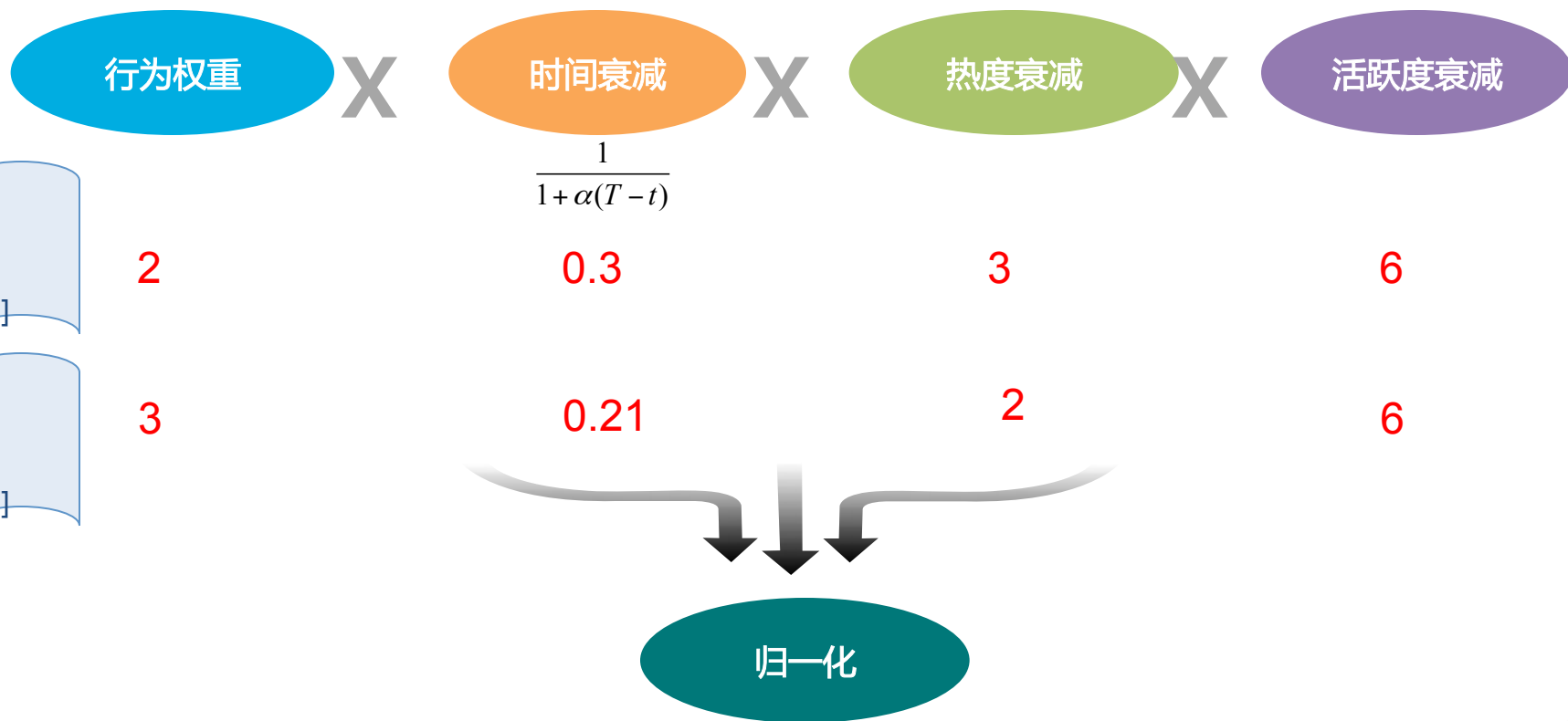


$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \zeta_i \\ \text{subject to} \quad & y_i[(w x_i) + b] \geq 1 - \zeta_i \quad (i=1, 2, \dots, l) \end{aligned}$$

算法	准确率(P)	召回率(R)	F1
朴素Bayes	85%	86.2%	85.5
SVM	92%	93%	92.4



用户画像—兴趣建模





- 个性化推荐
- 广告
- 信用等级分群
- 用户流失预警
- 游戏潜在用户群体筛选
- 异常监控分析





推荐一应用场景

首页推荐



应用详情推荐



用户粘性

转化率

广告

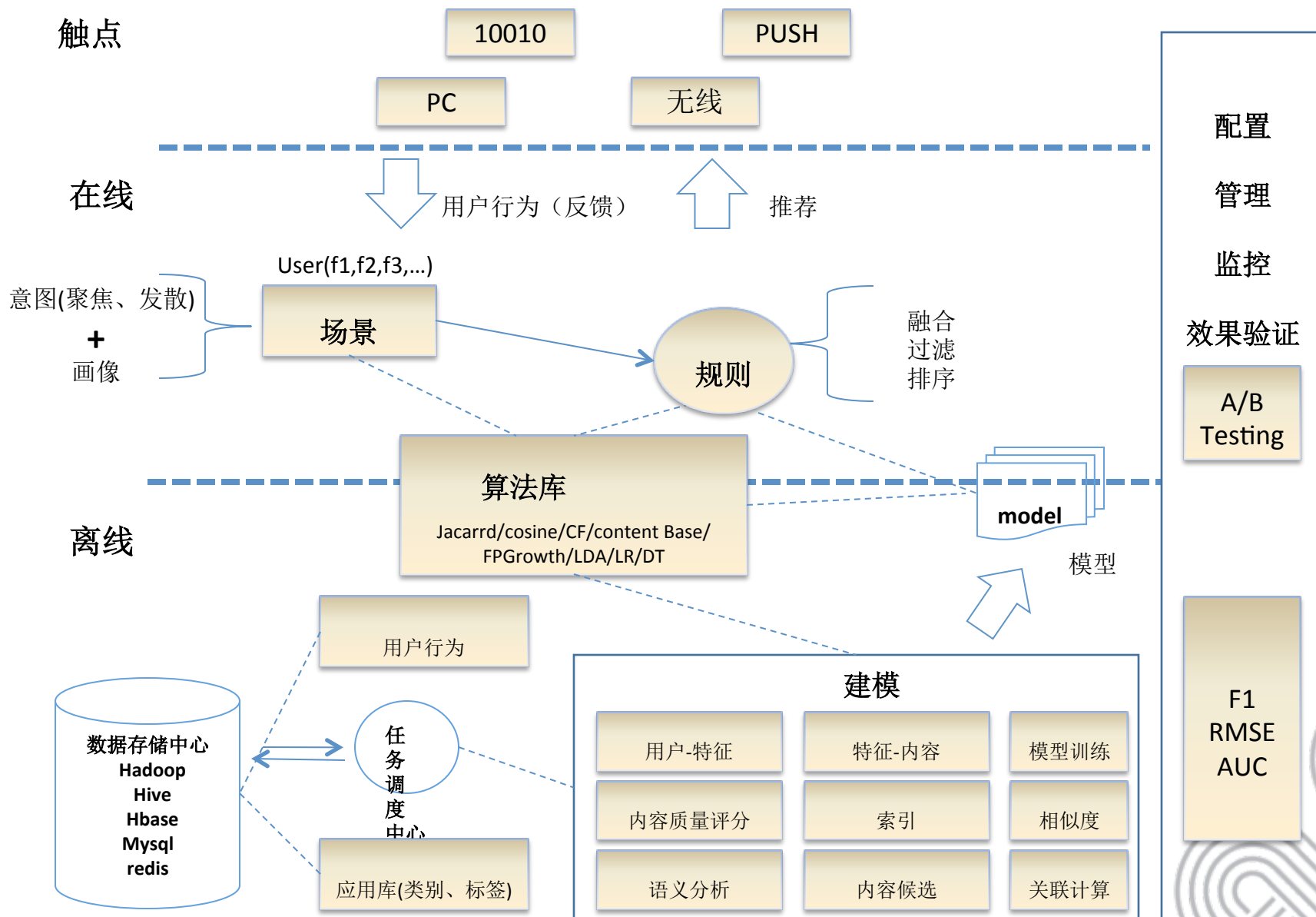


猜你喜欢





个性化推荐—平台架构





Content-Based

Category

冷启动

推荐精度

ItemBased-CF

$$p(u, j) = \sum_{i \in N(u) \cap S(j, K)} w_{ij} r_{ui}$$

新颖

$$w_{ij} = \frac{|N(i) \cap N(j)|}{\sqrt{|N(i)| |N(j)|}}$$

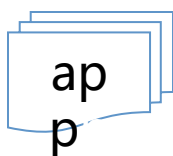
稀疏

$$w_{ij} = \frac{\sum_{u \in N(i) \cap N(j)} \frac{1}{\log 1 + |N(u)|}}{\sqrt{|N(i)| |N(j)|}}$$

活跃用户

聚类模型

Model-Based



语义分析

LDA

Topic分布

基于KL距离
推荐语义相关应用

来源融合

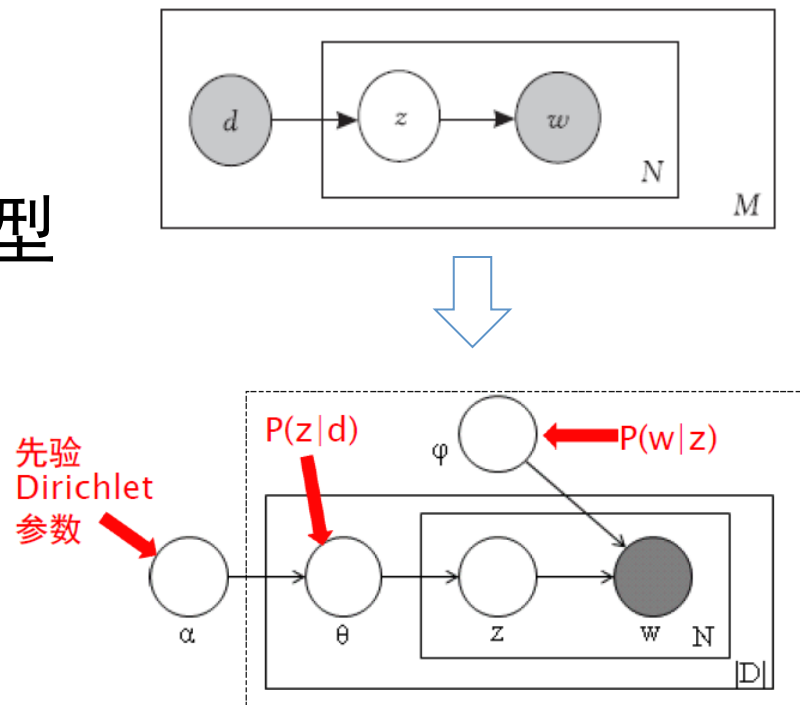
$$D_{KL}(P \| Q) = \sum_i \ln \left(\frac{P(i)}{Q(i)} \right) P(i).$$

$$P(y=1 | x) = \frac{1}{1 + e^{-(\beta_0 + \sum \beta_i * x_i)}}$$

■ LDA

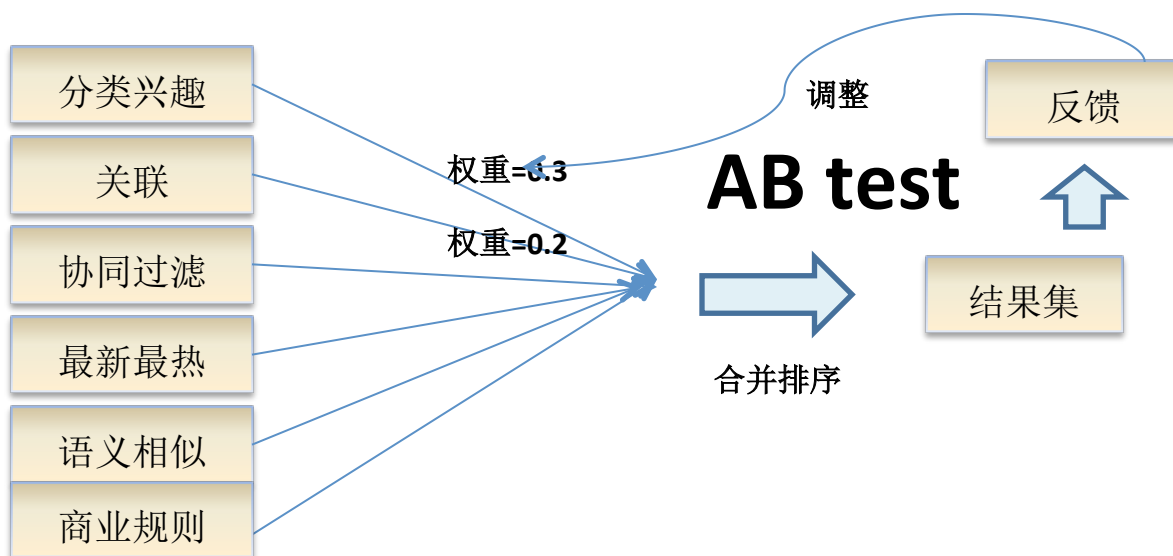
- EM
- 无 $p(z|d)$ 的生成概率模型
- 容易过拟合

- Gibbs
- 参数少，过拟合风险小
- 新文档处理能力强





如何确定各个模型、特征的权重？



人工对权重的调整，很难把控
新加入特征难以快速设置特征



Logic Regression

$$P(y=1|x) = \frac{1}{1+e^{-(\beta_0 + \sum \beta_i * x_i)}}$$

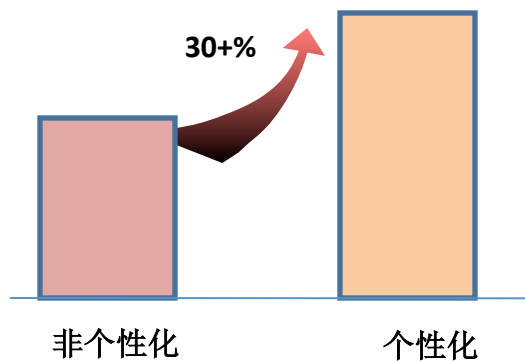
权重系数: $h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

候选集	model1	model2	model3			score
App1	0.2	0.54				0.7
App2		0.32	0.6			0.5

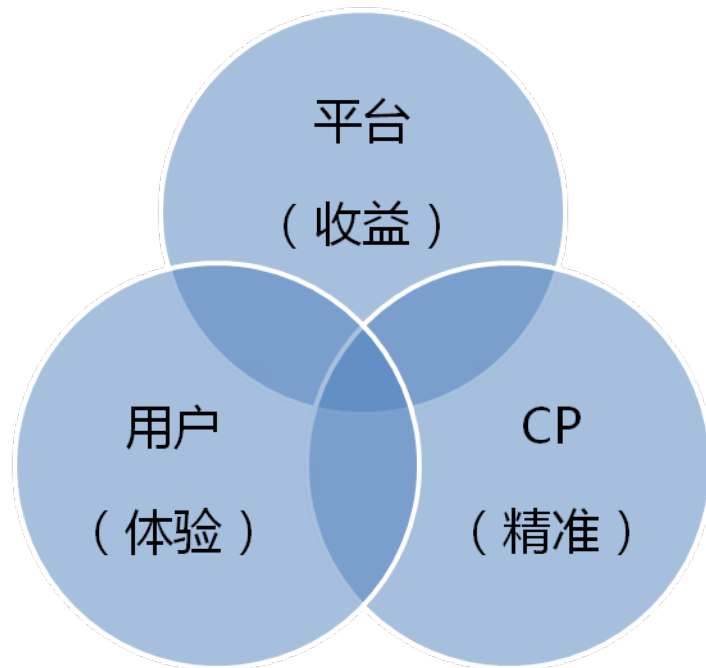
学习隐藏在用户群体行为背后的规律



PV转化率 (CTR*CVR)，效果提高30%



个性化推荐的下载量占比21%



■ 公式：ctr*Bid

$$\text{ctr} = \text{click} / \text{PV}$$

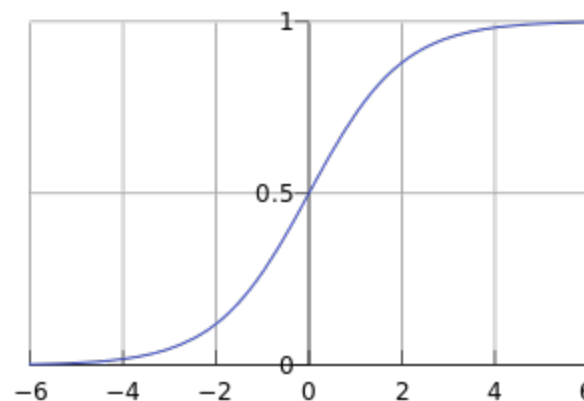
■ 缺点：冷启动、缺少个性化诉求



- 公式: $pCTR * Bid$
- $pCTR: p(\text{click} | \text{ad}, \text{user})$
- 基于LR的点击预估模型

点击=1, 不点击=0
点击的概率

$$P(y=1|x) = \frac{1}{1 + e^{-(\beta_0 + \sum \beta_i * x_i)}}$$





■ 特征

用户profile(活跃度、性别、年龄、标签)、
广告(广告质量、历史点击率、新颖性)、
CP、
用户和广告交叉主题特征

■ 样本选择

- 去噪、样本抽样

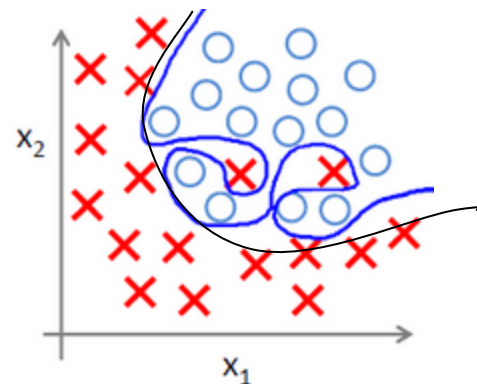
■ 特征处理

- 归一化
- 离散化、交叉
- 泛化能力

正则化(惩罚)

L1, 使得大量无效特征权重为0

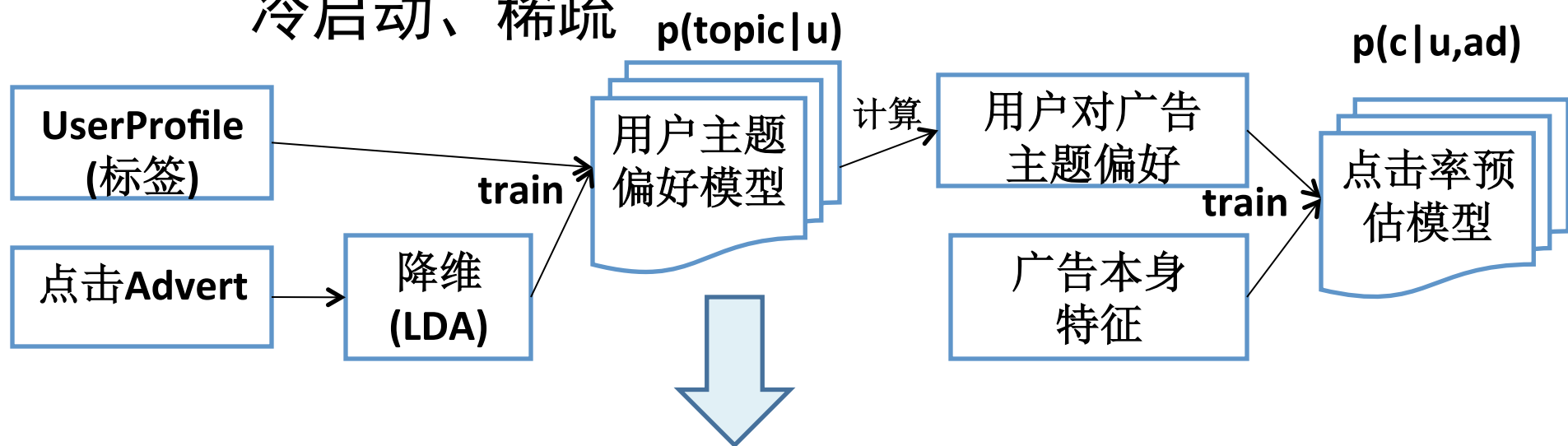
L2



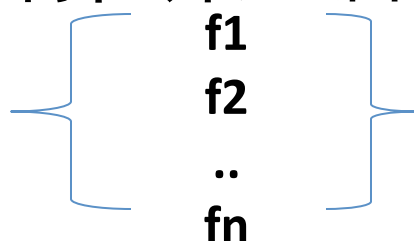


■ 特征处理

冷启动、稀疏



用户标签特征筛选(降维)



Click, PV
选取对点击PV贡献最大的特征TopN



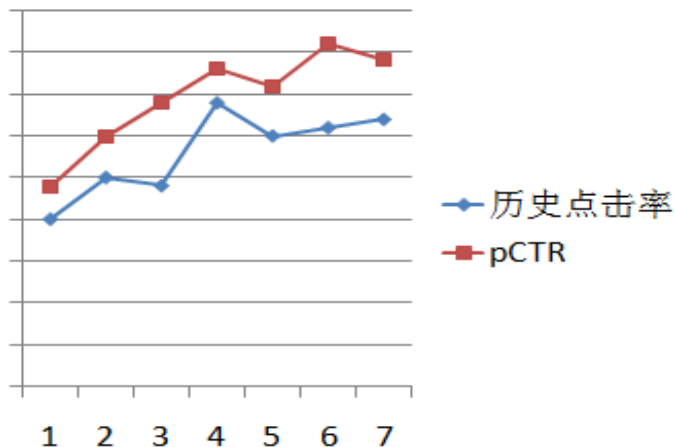
■ 参数估计优化：

L-BFGS

line search确定步长，无须手动选择，
利用有限内存近似BFGS，
利用历史值和梯度寻找当前方向(Two loop)，
实现快速迭代



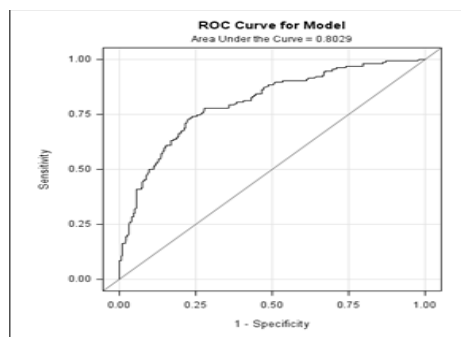
■ 产品层面



■ 算法层面

AUC

对于CTR高的广告，
预测的是否也高？



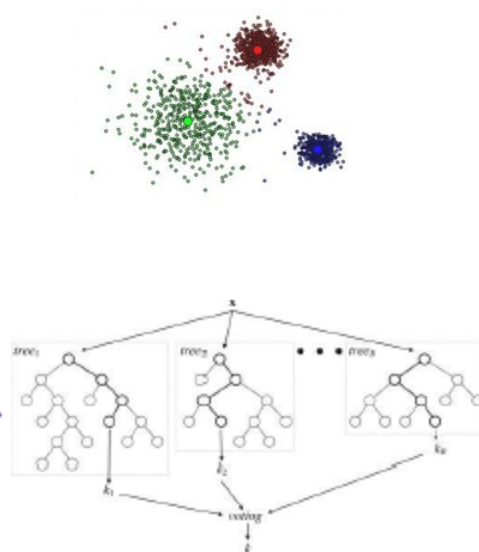


算法、数据、人机交互





用户信用等级分群



L1

L2

L3

...

LN

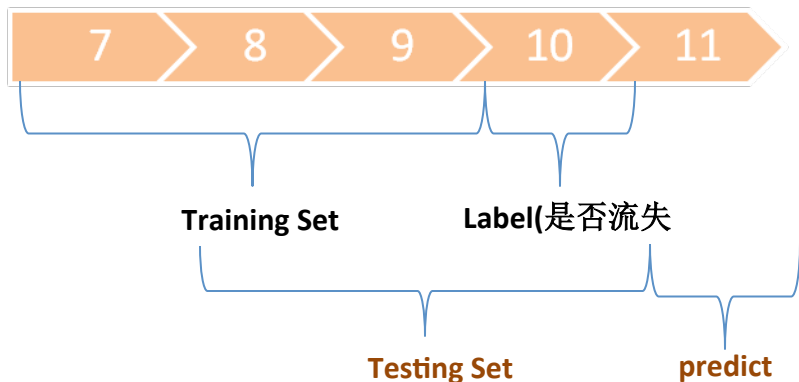
基本算法

● K-means

● RF



■ 模型



10月份流失的用户，
分析其前3个月的行为数据

■ 特征

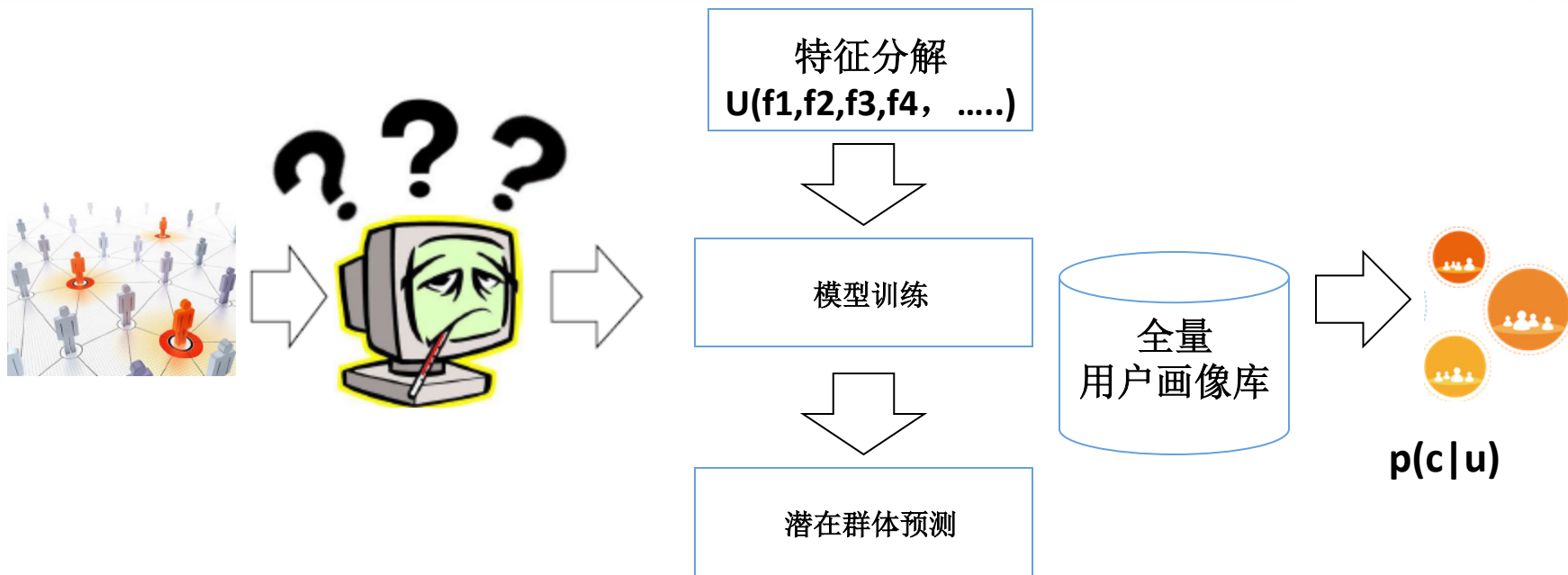
活跃度、登陆情况、下载情况、预装机情况、机型、....

➔潜在的流失用户

针对可能流失的用户做**PUSH**推广活动

基本算法

●GBDT



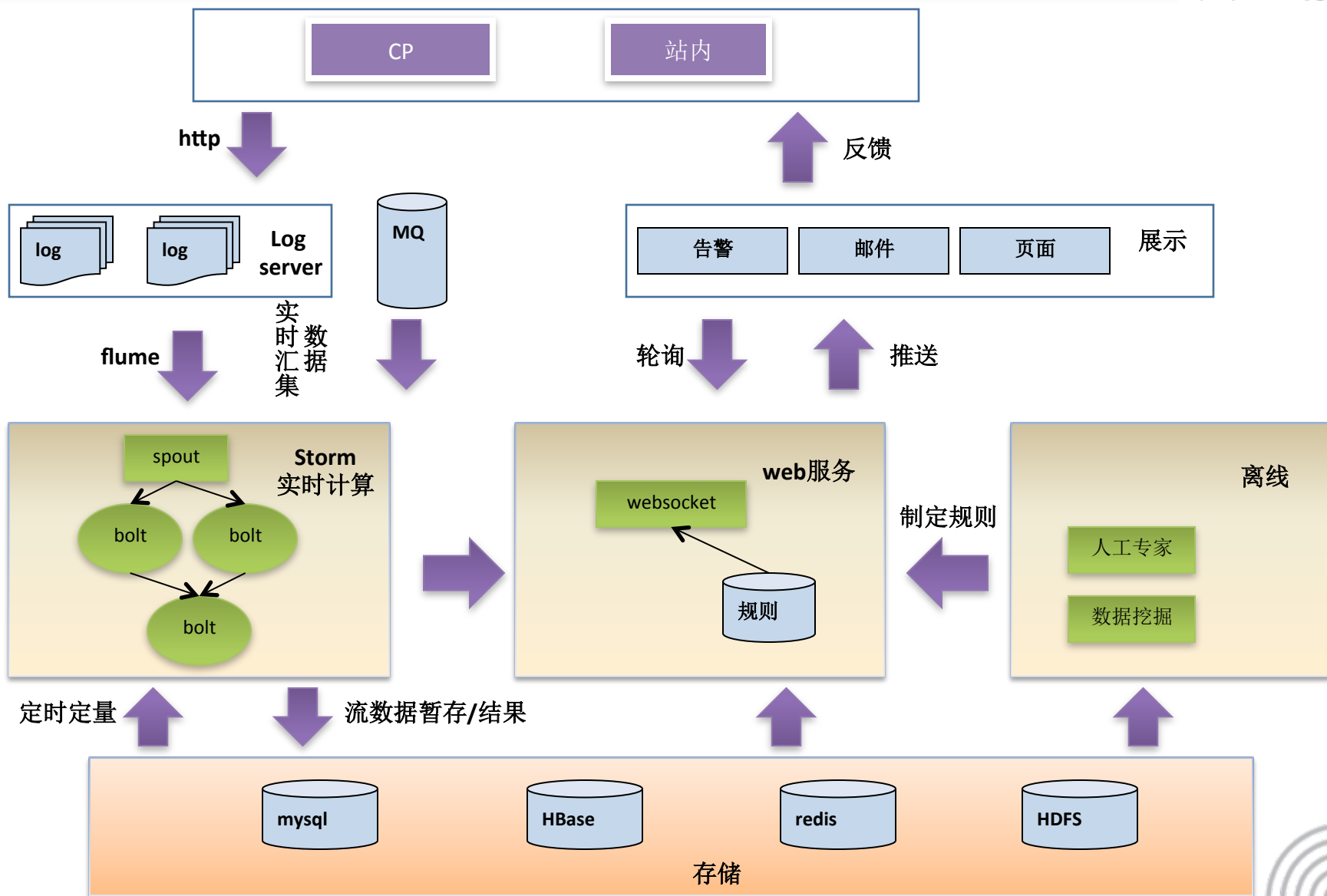
以在应用或者游戏中产生转化(注册、付费)的这些用户作为训练正样本，结合用户特征进行模型训练，从用户画像库中筛选出潜在的用户群体，推荐给CP，通过PUSH做相关的营销活动

基本算法

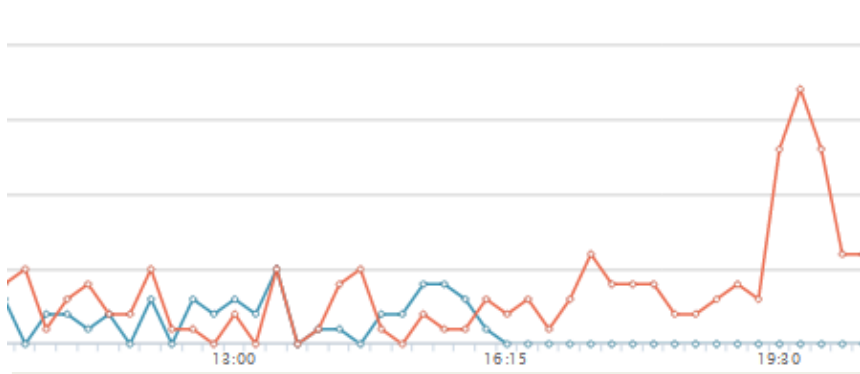
● Logic Regression



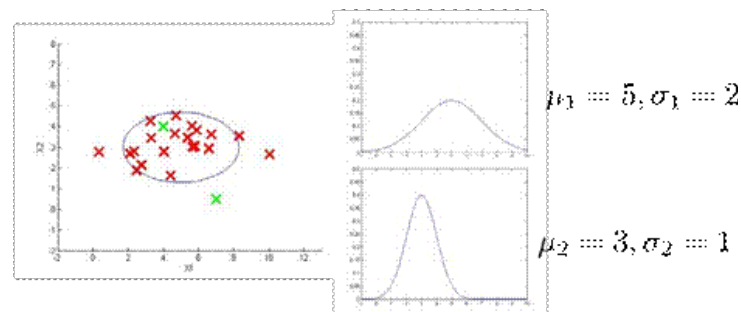
异常监控—Storm流计算



■ 异常检测



$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

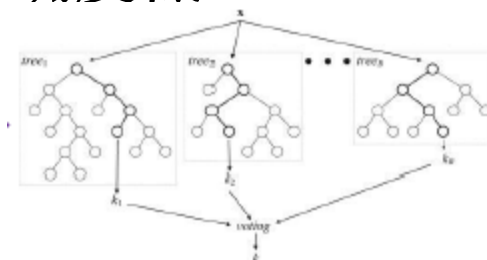


■ 刷机诊断

基于规则(IF ELSE): 依赖经验, 调整繁琐, 准确度低

基于模型:

利用用户刷机的一些行为特征和数据
进行模型训练, 结合模型来判断当前是否刷机



我的blog:

<http://blog.csdn.net/yangbutao>

我们招聘:

Hadoop/Hbase/Spark开发

算法工程师

数据挖掘工程师

...



敬请期待：
2015中华数据库大会
时间：2015.05.16
报名时间：2015.02.14
报名网址：meeting.zhdba.com

联系我们：
联系人：朱小姐
联系电话：136 5197 9898
联系QQ：378091820