



# 阿里广告中的机器学习平台



晏宗

- 蒋龙
- 花名：昙宗
- 阿里巴巴高级技术专家
- 2006-2011，微软亚洲研究院自然语言计算组
  - 微软对联，机器翻译，实时搜索
- 2011-，阿里巴巴广告部（阿里妈妈）
  - 研发广告中的机器学习和NLP技术
- Email: [tanzong.jl@taobao.com](mailto:tanzong.jl@taobao.com)
- Weibo: 蒋龙\_淘昙宗



- 广告中的机器学习
- 阿里广告机器学习平台
  - 特征平台
  - 算法平台
    - MPI并行计算集群
    - 算法平台
    - 评测平台
  - 应用现状

s.taobao.com/search?q=%D1%A9%B7%C4+%C1%AC%D2%C2%C8%B9&commend=all&ssid=s5-e&search\_type=item&sourceId=tb.index&initiative

0-algo 5-生活服务 a-技术新闻 iphone 技术 淘宝网内 新浪微博 淘宝网 - 淘!我喜欢 TechWeb.com.cn ... W806手机内部购机... 资产管理系统

好评推荐: 黄钻买家最爱买 回头客最多 服务态度最好 发货速度最快 描述评分最高  
 袖长: 短袖 无袖/背心裙 长袖 吊带裙 五分袖/中袖 七分袖 +多选  
 图案: 纯色 花色 圆点 碎花 条纹 卡通 +多选  
 ▶女装: 连衣裙 (519.2万) 雪纺衫 (6.7万)  
 大码服装 (7.5万) 半身裙 (9052)  
 ▶童装/童鞋/亲子装: 儿童裙子 (7.0万)

你是不是想找: 大码雪纺连衣裙 碎花雪纺连衣裙 蕾丝雪纺连衣裙 雪纺连衣裙

所有宝贝 人气 天猫 二手 逛街  
 雪纺 连衣裙 确定 海外商品 货到付款 消费者保障 7天退换 正品保障 旺旺在线 亲, 合并同款换位置啦  
 默认 销量 信用 价格 总价 所在地 款式 店铺 列表 大图  
 在逛街中找到“雪纺 连衣裙”相关宝贝: 应季新品 特价商品  
 您好 rainbowgbh, 点击查看“雪纺 连衣裙”相关宝贝在收藏过店铺、购买过店铺、免邮费、同城内搜索的结果。 立即查看

排序分数:  
f(点击率, 竞价)  
=>点击率预估  
p(click | ad, query)

薰衣草 De 期許



今日8折特價  
RMB: 230

¥ 288.00

2012夏季新款碎花波西米亚长裙雪纺连衣裙子

最近成交1198笔

風潮來襲



仅限今日  
159 包邮

¥ 199.00

2012 女装 碎花 雪纺 连衣裙 夏季 新款

最近成交11685笔

凤色蝶服旗舰店



限时折扣  
7折包邮

红纱雀



淘宝第一

¥128 包邮 腰带

Unisonfm 果尚美





em.taobao.com/item.htm?id=17363372220



热卖单品 精品凉鞋 热卖女鞋 新款凉鞋 坡跟女凉鞋 休闲凉鞋 女鞋凉鞋 美鞋 拖鞋 性感鱼嘴 舒适单鞋 [更多热卖](#)



疯抢!漫步云端 女士包臀收腰蕾  
丝花边连衣裙  
¥598.00



经典精选  
8款口味  
79  
包邮  
Amovo纯可可脂巧克力8口味礼  
盒 零食 包邮  
¥109.00



2012夏季新款  
100%里外全皮  
168  
元  
包邮  
冰点促销  
思加图 专柜正品 2012夏季新款  
拖鞋凉拖鞋  
¥238.00



2件装  
2折  
0利润大牌高档退换包邮真丝连  
衣裙送彩票  
¥800.00

## 您可能对这些宝贝感兴趣

缀好



排序分数:  
f(点击率, 竞价)  
=>点击率预估 $p(\text{click} | \text{ad}, \text{user})$



缀好



缀好



主宝贝

排序分数:

$f(\text{转化率}, \text{商品价格})$

$\Rightarrow$  转化率预估

$p(\text{buy} | \text{item}, \text{user}, \text{curltem})$

相关推荐

搭配推荐



包邮Apple/苹果 iPhone 5 联通电信苹果5 iPhone5 限时抢购  
¥4699.00



Apple/苹果 iPhone 4 (8G) 苹果4代 8G 行货正品全国联保  
¥2899.00



MIUI/小米 M2小米二代小米2四核顶配 正品行货手机16G 现货  
¥2099.00



顺丰包邮Samsung/三星 GALAXY Note II N7100 Note2 牛2 正品现  
¥3894.05



洛克rock 苹果5 iPhone 5手机壳 TPU软壳 保护壳保护套 畅玩系列  
¥48.00



手机剪卡器 搭配手机一起买价格15元  
¥30.00



Zupool/触宝 苹果4 iPhone4/4S iPhone 4s保护膜屏幕贴膜 晶透系列  
¥25.00

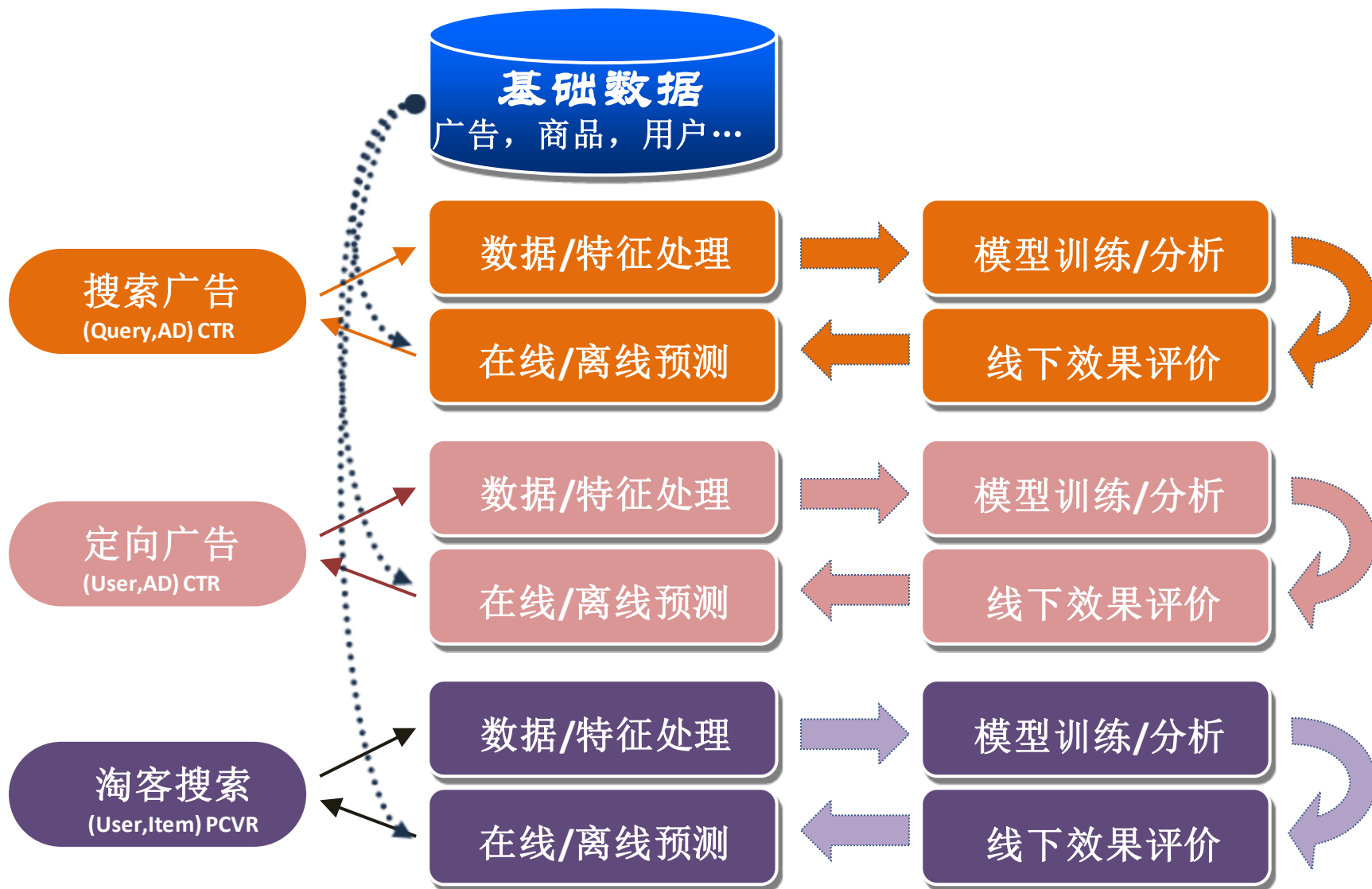


SanDisk/闪迪TF 16G TF卡 MicroSD手机内存卡存储卡70元  
¥70.00

- 点击率预估
  - 搜索广告:  $p(\text{click}|\text{ad},\text{query})$
  - 定向广告:  $p(\text{click}|\text{ad},\text{user})$
- 转化率预估
  - 淘客:  $p(\text{buy}|\text{item},\text{query})$ ,  $p(\text{buy}|\text{item},\text{user})$
  - 推荐:  $p(\text{buy}|\text{item},\text{user},\text{curitem})$
- 相同的流程
  - 基础数据+日志数据, 特征抽取和处理, 训练数据收集, 模型训练, 模型线下评估

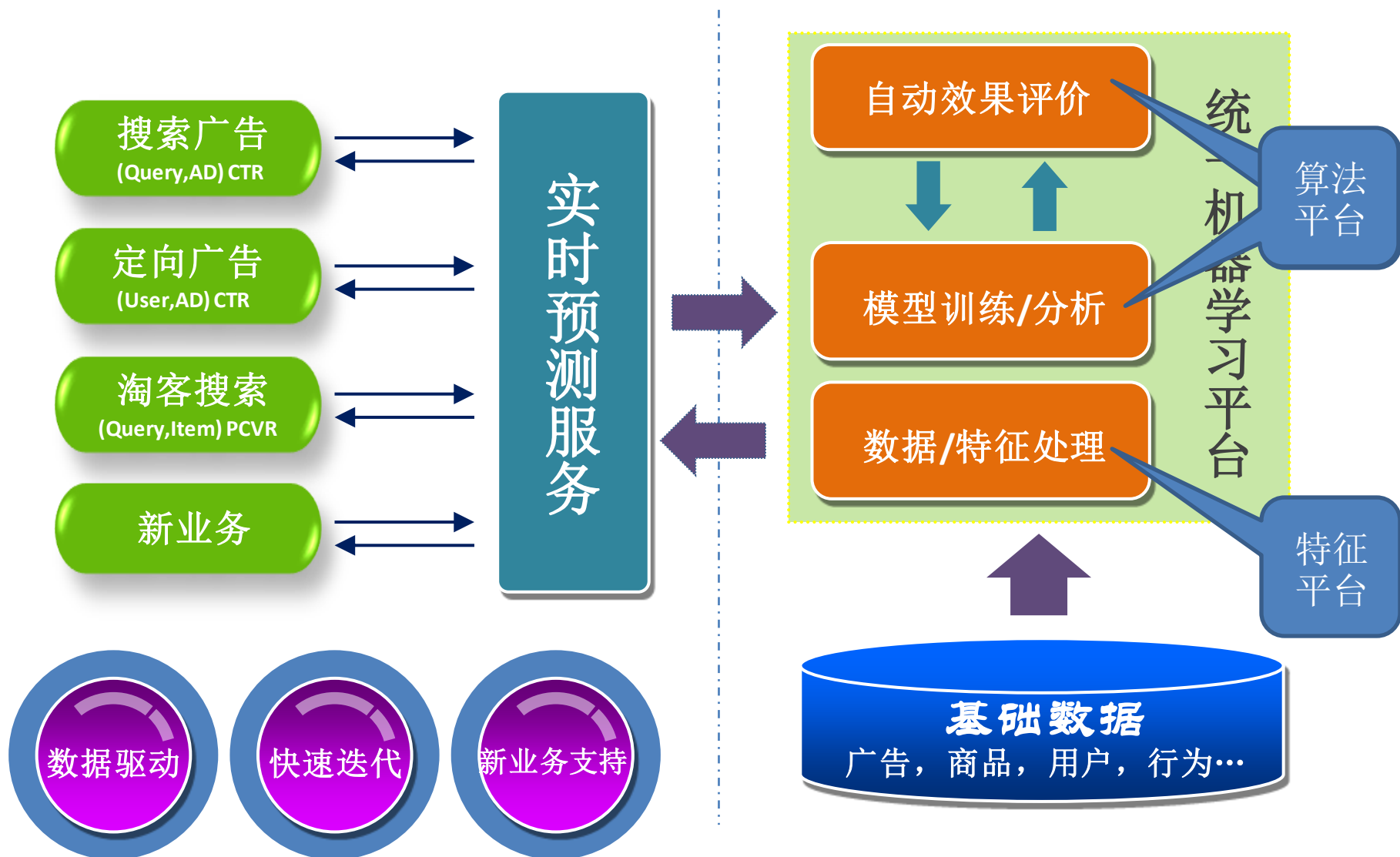


# 重复的数据流





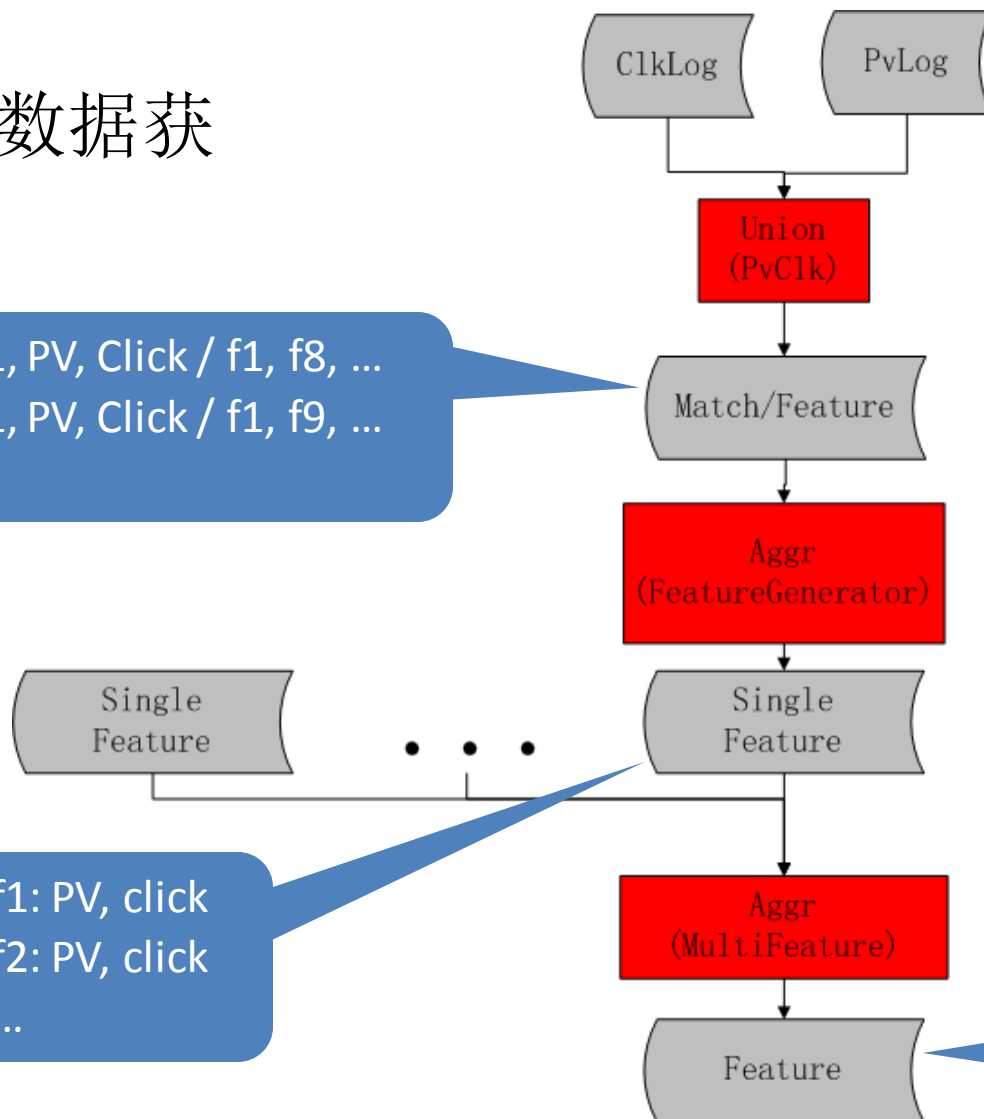
# 统一机器学习平台



- 特征处理
  - 特征抽取
    - 公用数据(Query, Ad, Item, Shop, User)
  - 特征处理
    - 变换, 规范化, 交叉组合, 离散化
- 特征分析
  - 覆盖率, 均值, 方差
  - 相关系数, Chi-square, 互信息

## • 反馈特征数据获取流程

- Qid1\_Ad1, PV, Click / f1, f8, ...
- Qid2\_Ad1, PV, Click / f1, f9, ...
- ...



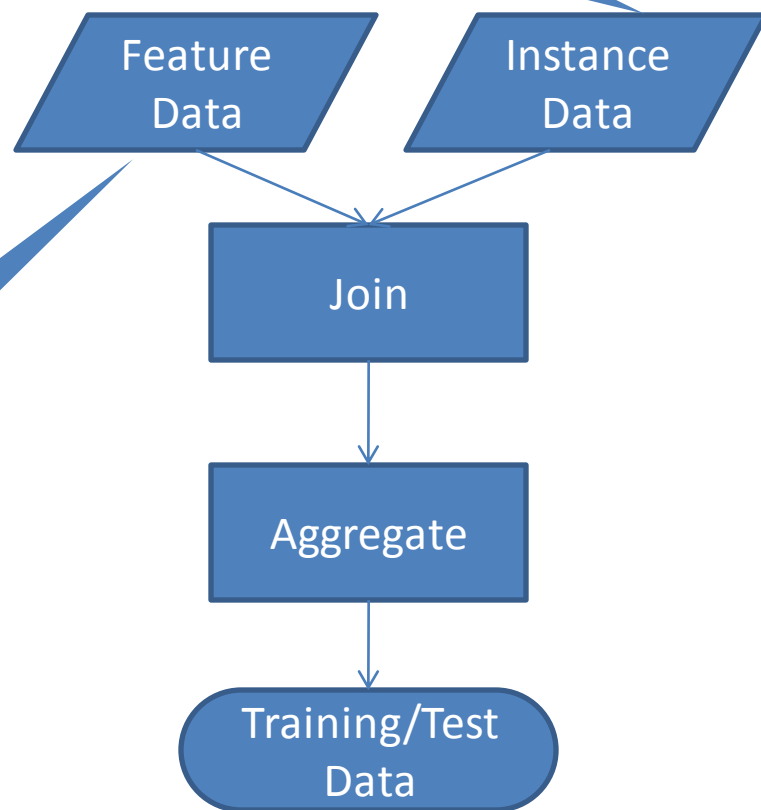
- f1: PV, click
- f2: PV, click
- ...

- f1: {values}
- f2: {values}
- f3: {values}
- ...

- 训练/测试数据生成
  - 目标和特征数据连接
  - 训练数据选择
  - 多天数据集成
- 大数据操作：Union，Aggregation，Join

- Qid1\_Ad1, PV, Click / f1, f8, ...
- Qid2\_Ad1, PV, Click / f1, f9, ...
- ...

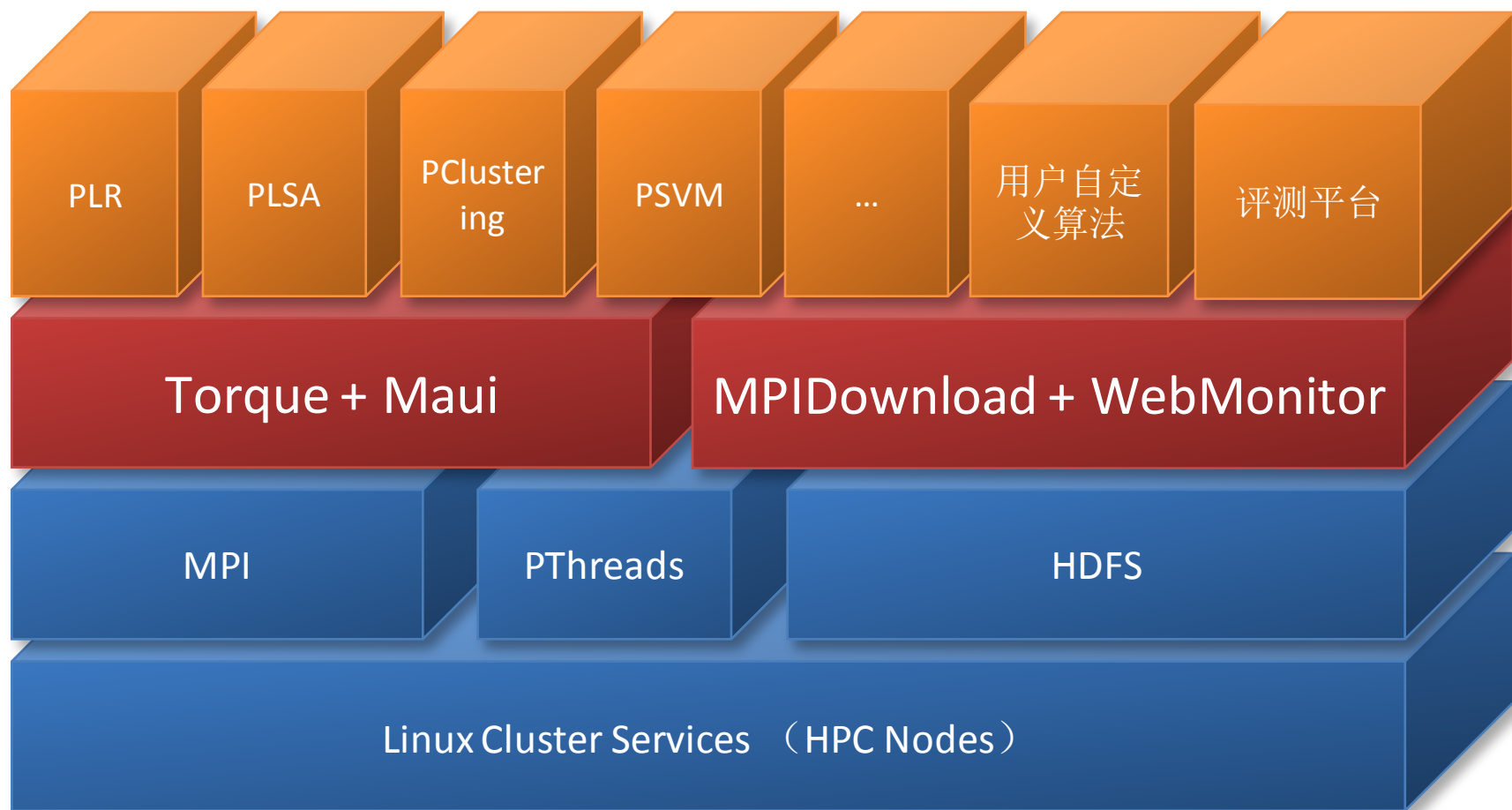
- f1: {values}
- f2: {values}
- f3: {values}
- ...





- 特征平台特点
  - 海量数据
    - 原始日志，训练数据：几十T
  - 计算逻辑简单，迭代少
  - 稳定，容错
- 基于Hadoop开发
  - 大规模MapReduce+HDFS集群
    - 阿里云梯：4000+节点，~80PB存储

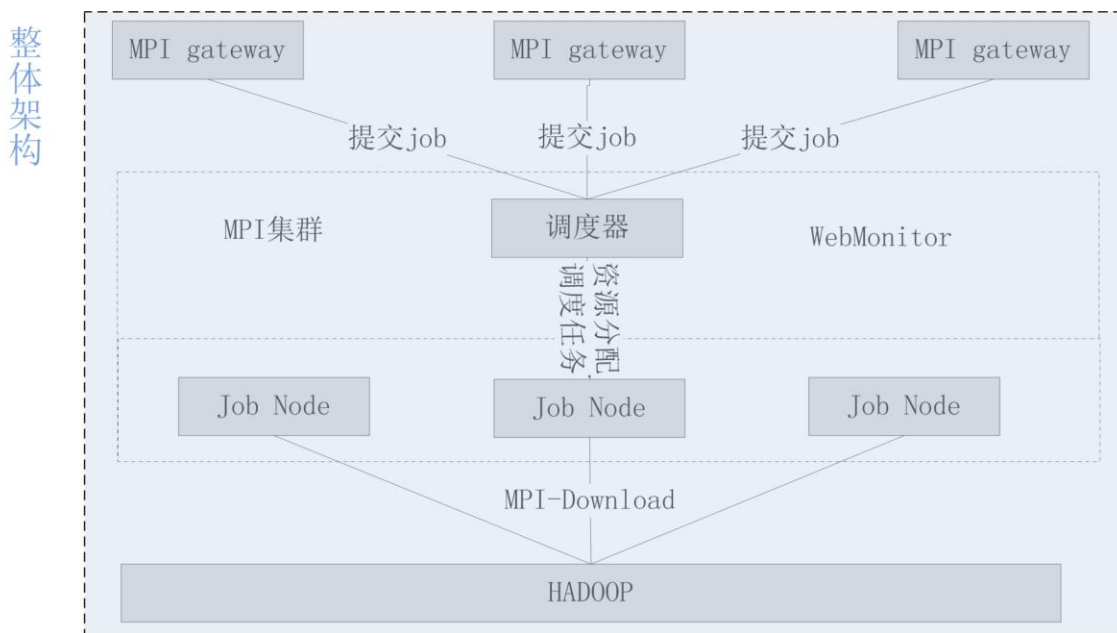
- 算法平台：模型训练和评测
  - 海量数据，大规模特征
  - 计算逻辑复杂，迭代过程很多
  - 速度要求很高
- 基于MPI框架
  - 数据常驻内存
  - 单机高性能，大内存
  - 允许自定义复杂的数据通信



- 集群概况
  - ~500台服务器
    - 单机12核（超线程24核），100G内存，千/万兆网卡
  - 生产、开发、测试集群
  - 部署MPICH2框架
  - 共享云梯存储
    - MPI-Download：统一管理与Hadoop的数据交换
      - 分布式下载，自动解压合并，自动流量控制



- 使用Torque+Maui为基础的调度器
  - Job模块自动分发，计算资源均衡
  - 调度队列区分优先级
- Cgroups确保资源隔离
  - 动态资源调整

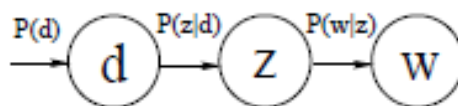


- 并行算法库
  - LR, MLR, LDA, PLSA, GBRT, SVM, MAXENT, CF, Spectral Clustering, ...

- PLSA (Probabilistic Latent Semantic Analysis) 模型

$$p(w | d) = \sum_{z \in Z} p(w | z) p(z | d)$$

- z 代表隐藏 topics
- <d,w> 矩阵: d 先产生 z, 再由 z 产生 w



- 应用
  - 文本语义分析, 用户兴趣分析

# 基于EM算法的参数估计

– E-step

$P(z_1 | d_1, w_1) = xx$   
 $P(z_2 | d_1, w_1) = xx$   
 $\dots$   
 $P(z_k | d_1, w_1) = xx$

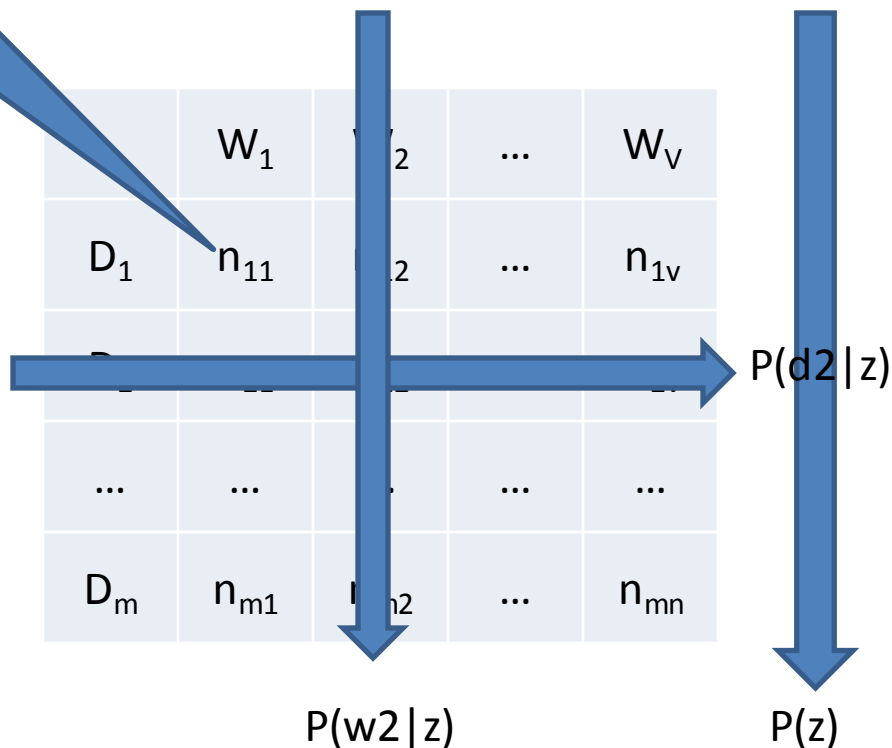
$$p(z | d, w) = \frac{p(z) p(d | z) p(w | z)}{\sum_{z' \in Z} p(z') p(d | z') p(w | z')}$$

– M-step

$$p(d | z) \propto \sum_{w \in W} n(d, w) p(z | d, w)$$

$$p(w | z) \propto \sum_{d \in D} n(d, w) p(z | d, w)$$

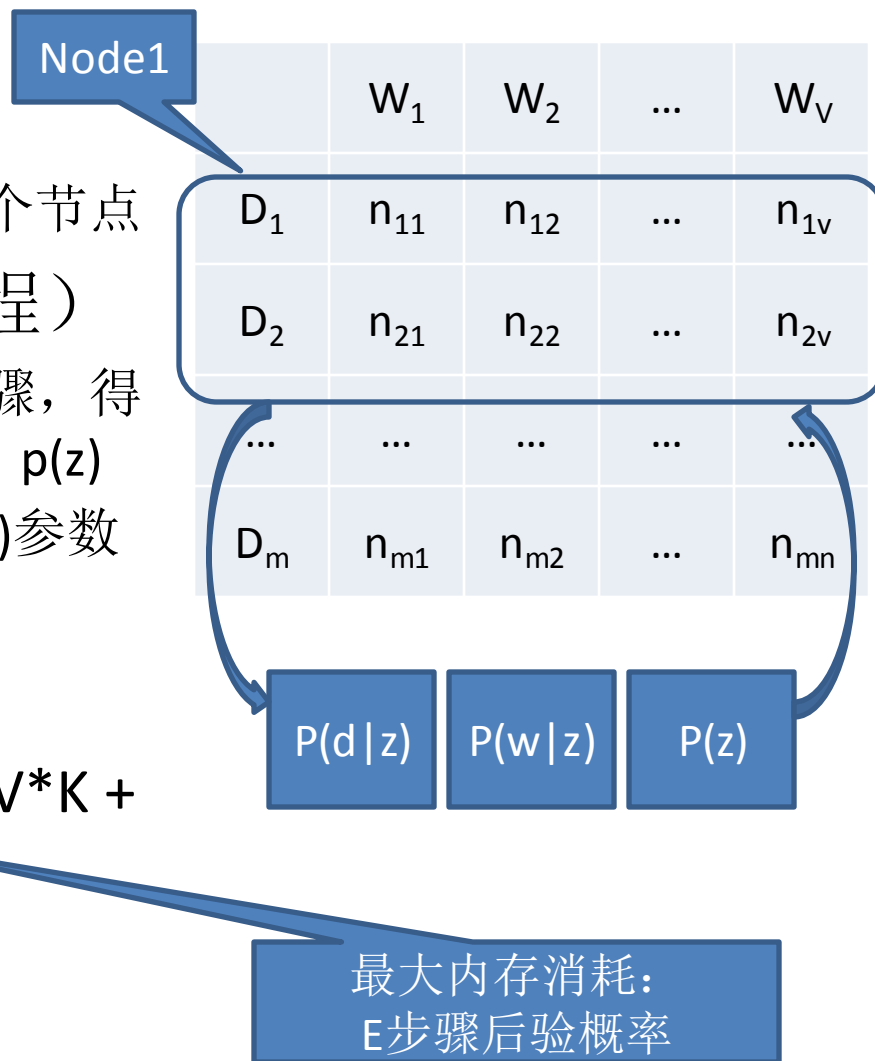
$$p(z) \propto \sum_{d \in D} \sum_{w \in W} n(d, w) p(z | d, w)$$





# 基于MPI的并行PLSA

- 数据并行化
  - $\langle d, w \rangle$  矩阵按行( $d$ )均匀划分到每个节点
- 计算并行化 ( $N$ 个节点/进程)
  - 节点同步完成单机数据的E-M步骤, 得到新一轮的 $p(d|z)$ 和局部 $p(w|z)$ ,  $p(z)$
  - 节点间通信得到全局 $p(w|z)$ ,  $p(z)$ 参数更新
- 单节点内存占用
  - $|\langle d, w \rangle|/N + |\langle d, w \rangle| * K/N + V * K + K * D/N + K$



- 内存优化
  - E步M步交叉进行，不存储后验概率 $p(z|d,w)$ 
    - 计算完每个词的 $p(z|d,w)$ ，直接更新对应的 $p(d|z)$ ， $p(w|z)$ 参数
    - 需要存储新旧两个模型 $p(d|z), p(w|z)$
  - 按行(d)扫描数据计算，省下内存 $p(d|z)$ 
    - 每行计算完直接在旧 $p(z|d)$ 上更新新值
- E-step计算优化

$$p(z|d,w) = \frac{p(z)p(d|z)p(w|z)}{\sum_{z' \in Z} p(z')p(d|z')p(w|z')}$$



$$p(z|d,w) = \frac{p(w|z)p(z|d)}{\sum_{z' \in Z} p(z'|d)p(w|z')}$$

# 基于MPI的自动评测

- 基于MPI的模型评测
  - 无需回传model到hadoop
  - 训练完直接评测，也可单独评测
  - 自动化的参数调整
- 指标种类
  - MAE, MSE, Accuracy, Precision, Recall, AUC, GAUC

- 使用方
  - 阿里妈妈事业部，搜索事业部
  - 天猫算法，交易算法团队
- 集群稳定性 **>99.9%**
- 集群资源CPU周级峰值平均利用率 **60%+**，平均利用率 **30%+**
- 周Job数 **1000+**



- 基于MPI的LR/MLR模型
  - 搜索广告ctr预估
  - 淘客搜索CVR预估
  - Tmall推荐融合排序
  - 定向广告CTR预估
- 基于MPI版本PLSA模型
  - 外投广告个性化
  - 搜索个性化

Thanks!