



促进云计算创新发展 培育信息产业新业态

第七届中国云计算大会

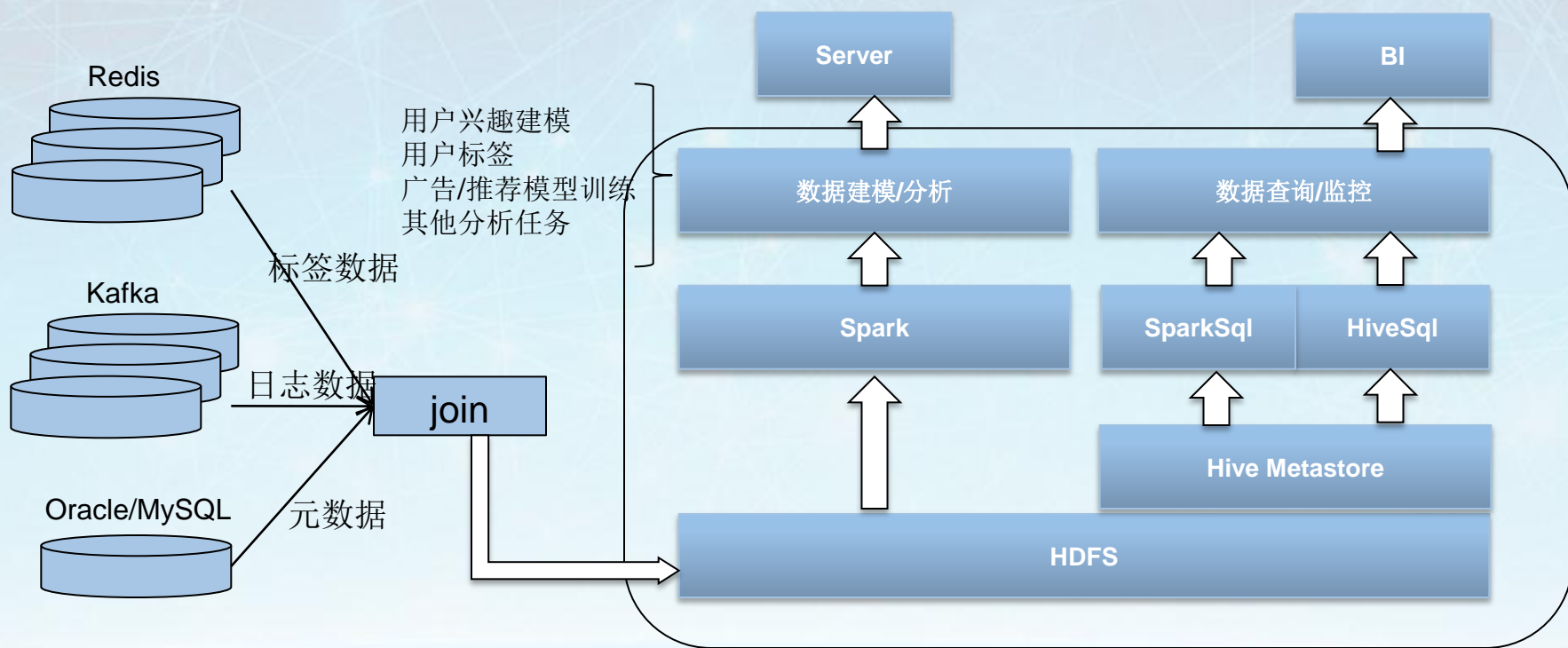
推荐算法和Spark实现

搜狐大数据中心 – 李滔

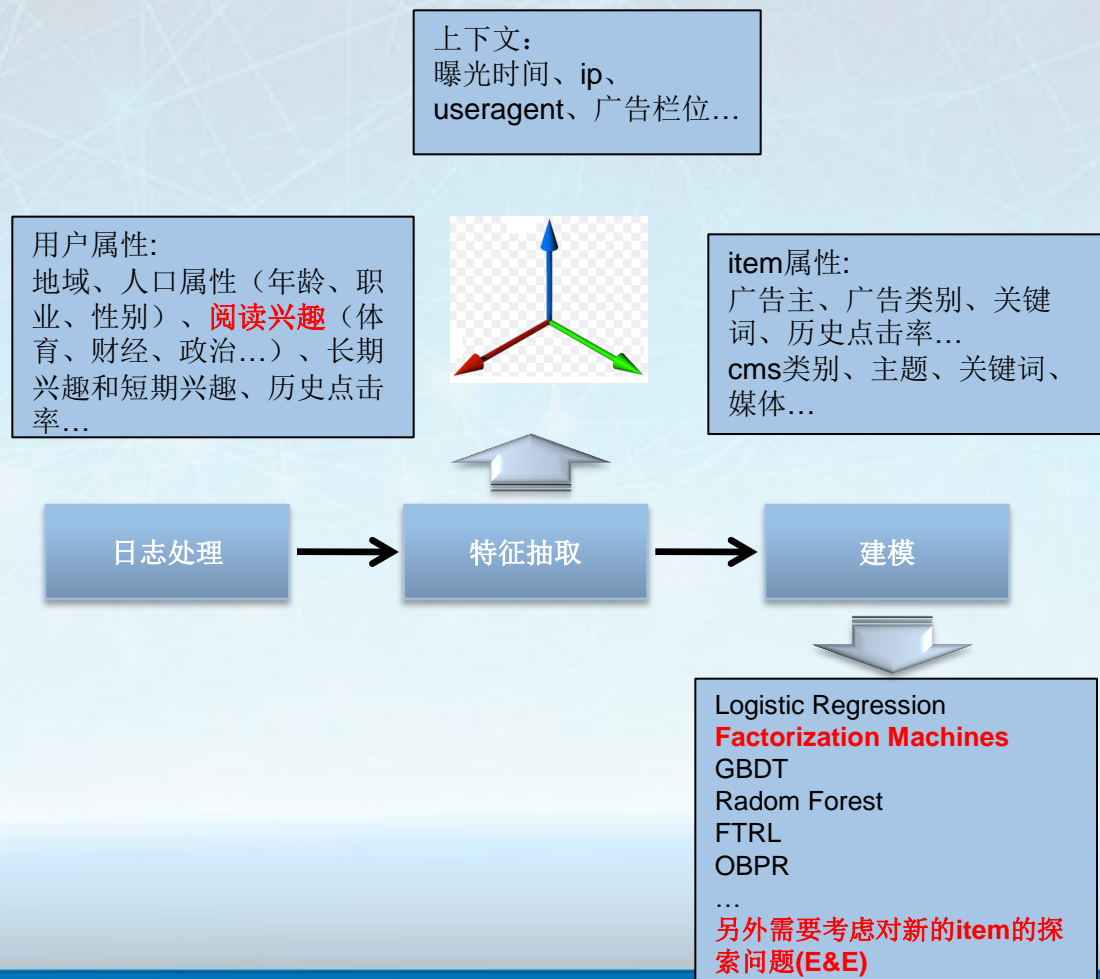
广告 vs. 新闻推荐

- 相同点
 - 都可看做一个点击率(ctr)估计的任务
 - 特征描述: 用户、商品、上下文三个维度
 - 点击率是动态变化的
- 不同点
 - 点击率: 推荐点击率是广告的10~100倍
 - 特征描述: 广告细粒度特征 vs. 推荐相对粗粒度

基于Hadoop EcoSystem的数据分析平台



CTR预估建模过程



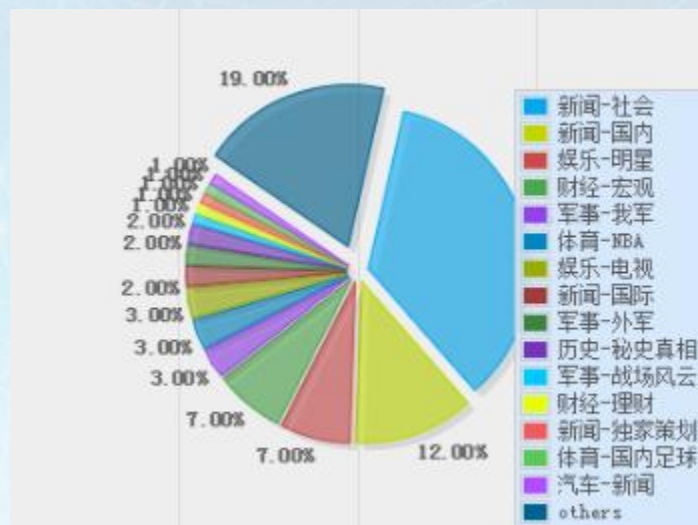
用户兴趣建模

根据用户的阅读历史对用户兴趣建模

用户的新闻阅读分布是有偏的：

- 热点新闻的巨大点击量
- 新闻的曝光是有偏的

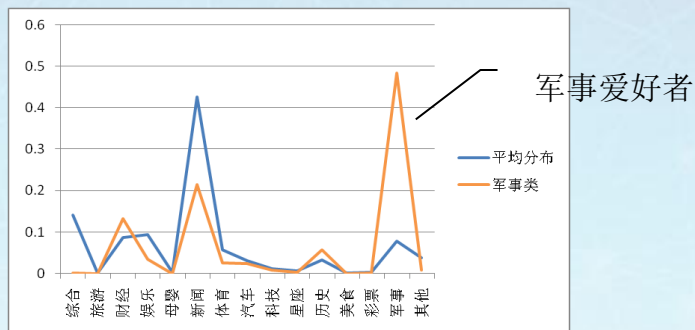
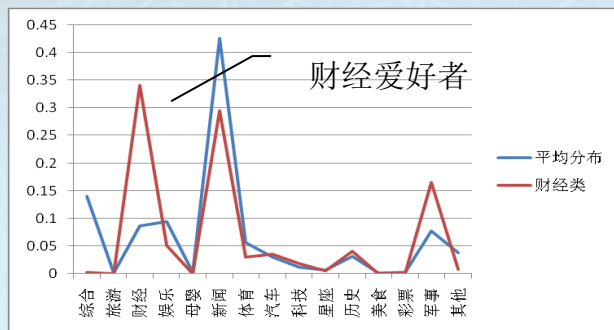
用户标签只能建立在用户相对于平均分布的偏离度上



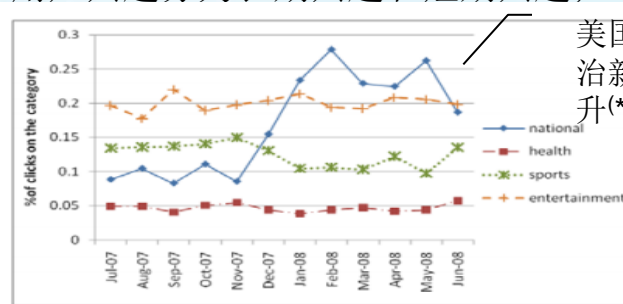
手机搜狐网一个月的新闻点击分布

用户兴趣建模

用户兴趣是相对于整体用户阅读平均分布的偏离度的度量，常用的维度包括CMS分类，topic model，关键词等：



用户兴趣分为长期兴趣和短期兴趣；

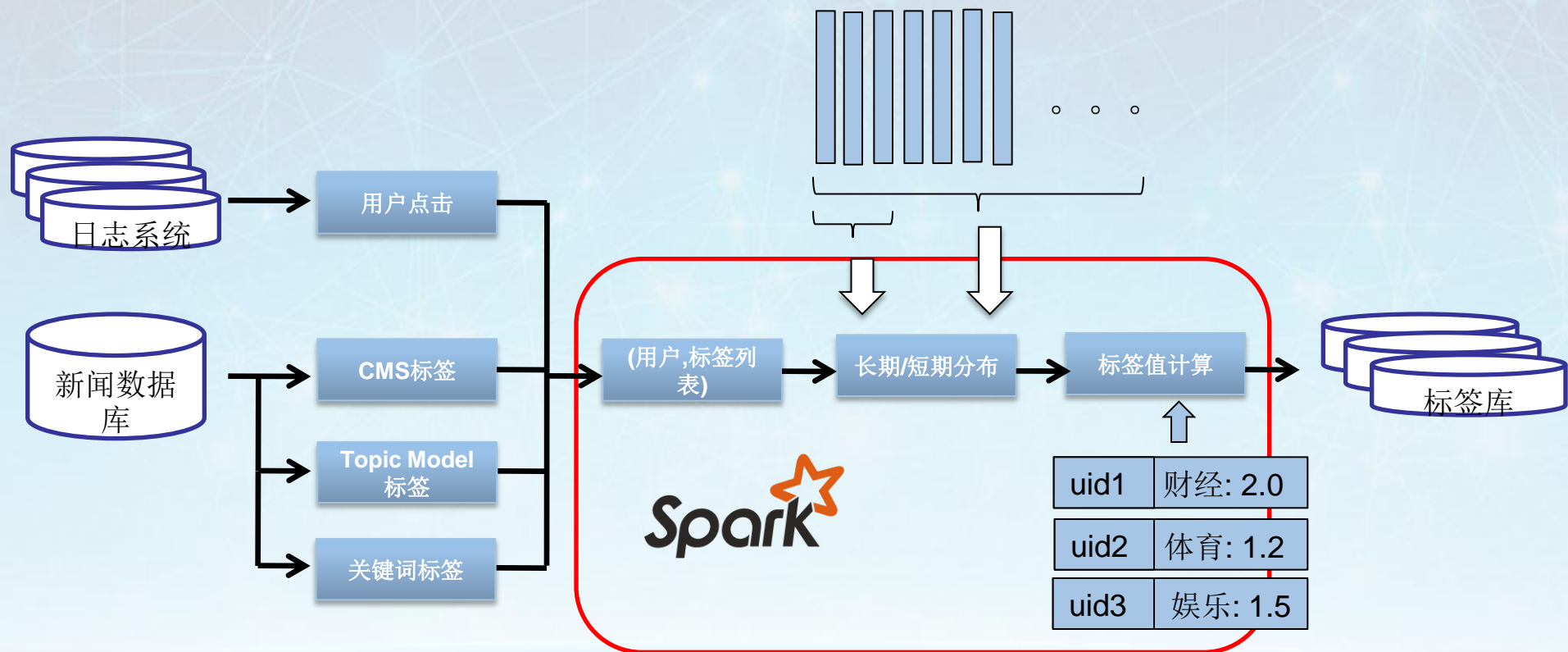


美国大选使得政治新闻的热度上升(*)

Figure 2. Interest distribution of US users over time

*Ref: J. Liu, Personalized News Recommendation Based on Click Behavior, 2010 International Conference on Intelligent User Interfaces

用户兴趣建模流程



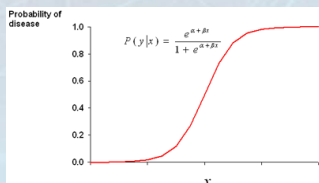
Spark带来的改善

- Spark丰富的数据处理操作
 - Map, reduce, filter, join, cogroup, ...
- 处理时间相对于Hadoop缩短近10倍
- 代码量大为减少

广告和新闻推荐建模

- 多信息源的有效利用
 - 支持多维度的特征及特征组合
 - 避免过拟合 (Overfitting)
- 探索(Exploration)和利用 (Exploitation)之间的平衡
 - Exploration: 获取信息; Exploitation: 根据当前信息决策
 - Bandit方法: ϵ -Greedy, UCB
 - Bayesian方法: Thompson Sampling
- 在线训练

Logistic Regression



Logit函数

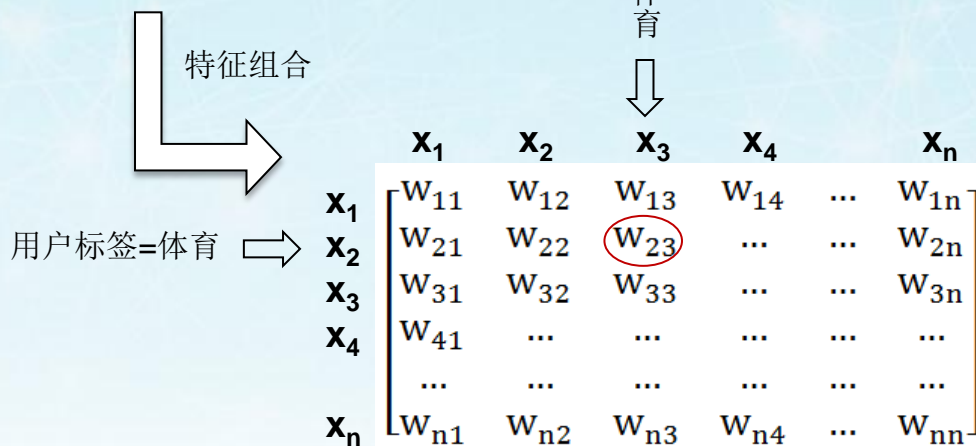
点击的概率P由特征的线性加权决定:

$$P(y|x) = \text{logit}(w_1 x_1 + w_2 x_2 + \dots + w_n x_n)$$

常用的特征组合:

用户标签和item标签组合;
曝光时间和item标签组合;
用户地域和item标签组合;

○ ○ ○



更好的模型: Factorization Machines

$$\hat{y}(x) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle v_i, v_j \rangle x_i x_j$$



	\mathbf{x}_1	\mathbf{x}_2	\mathbf{x}_3	\mathbf{x}_4	\mathbf{x}_n	
\mathbf{x}_1	$\langle \mathbf{v}_1 \cdot \mathbf{v}_1 \rangle$	$\langle \mathbf{v}_1 \cdot \mathbf{v}_2 \rangle$	$\langle \mathbf{v}_1 \cdot \mathbf{v}_3 \rangle$	$\langle \mathbf{v}_1 \cdot \mathbf{v}_4 \rangle$...	$\langle \mathbf{v}_1 \cdot \mathbf{v}_n \rangle$
\mathbf{x}_2	$\langle \mathbf{v}_2 \cdot \mathbf{v}_1 \rangle$	$\langle \mathbf{v}_2 \cdot \mathbf{v}_2 \rangle$	$\langle \mathbf{v}_2 \cdot \mathbf{v}_3 \rangle$	$\langle \mathbf{v}_2 \cdot \mathbf{v}_n \rangle$
\mathbf{x}_3	$\langle \mathbf{v}_3 \cdot \mathbf{v}_1 \rangle$	$\langle \mathbf{v}_3 \cdot \mathbf{v}_2 \rangle$	$\langle \mathbf{v}_3 \cdot \mathbf{v}_3 \rangle$	$\langle \mathbf{v}_3 \cdot \mathbf{v}_n \rangle$
\mathbf{x}_4	$\langle \mathbf{v}_4 \cdot \mathbf{v}_1 \rangle$

\mathbf{x}_n	$\langle \mathbf{v}_n \cdot \mathbf{v}_1 \rangle$	$\langle \mathbf{v}_n \cdot \mathbf{v}_2 \rangle$	$\langle \mathbf{v}_n \cdot \mathbf{v}_3 \rangle$	$\langle \mathbf{v}_n \cdot \mathbf{v}_4 \rangle$...	$\langle \mathbf{v}_n \cdot \mathbf{v}_n \rangle$

每维特征对应一个p维权重向量, 参数数量为 $n \cdot (p+1) + 1$

Recommender Data	Feature vector x												Target y
(A, TI, H, {C}, 5)													5 $y^{(1)}$
(A, NH, S, {}, 3)													3 $y^{(2)}$
(A, SW, N, {B, C}, 1)													1 $y^{(2)}$
(B, SW, N, {A, C}, 4)													4 $y^{(3)}$
(B, ST, H, {}, 5)													5 $y^{(4)}$
(C, TI, S, {A}, 1)													1 $y^{(5)}$
(C, SW, H, {A, B}, 5)													5 $y^{(6)}$
	A	B	C	TI	NH	SW	ST	S	N	H	A	B	C
	User			Movie				Mood			Watched with		
$x^{(1)}$	1	0	0	1	0	0	0	0	0	1	0	0	1
$x^{(2)}$	1	0	0	0	1	0	0	1	0	0	0	0	0
$x^{(3)}$	1	0	0	0	0	1	0	0	1	0	0	0.5	0.5
$x^{(4)}$	0	1	0	0	0	1	0	0	1	0	0.5	0	0.5
$x^{(5)}$	0	1	0	0	0	0	1	0	0	1	0	0	0
$x^{(6)}$	0	0	1	1	0	0	0	1	0	0	1	0	0
$x^{(7)}$	0	0	1	0	0	1	0	0	0	1	0.5	0.5	0

优点:

General purpose模型, 覆盖率像MF, SVD++等多种传统模型;

方便加入多种特征;

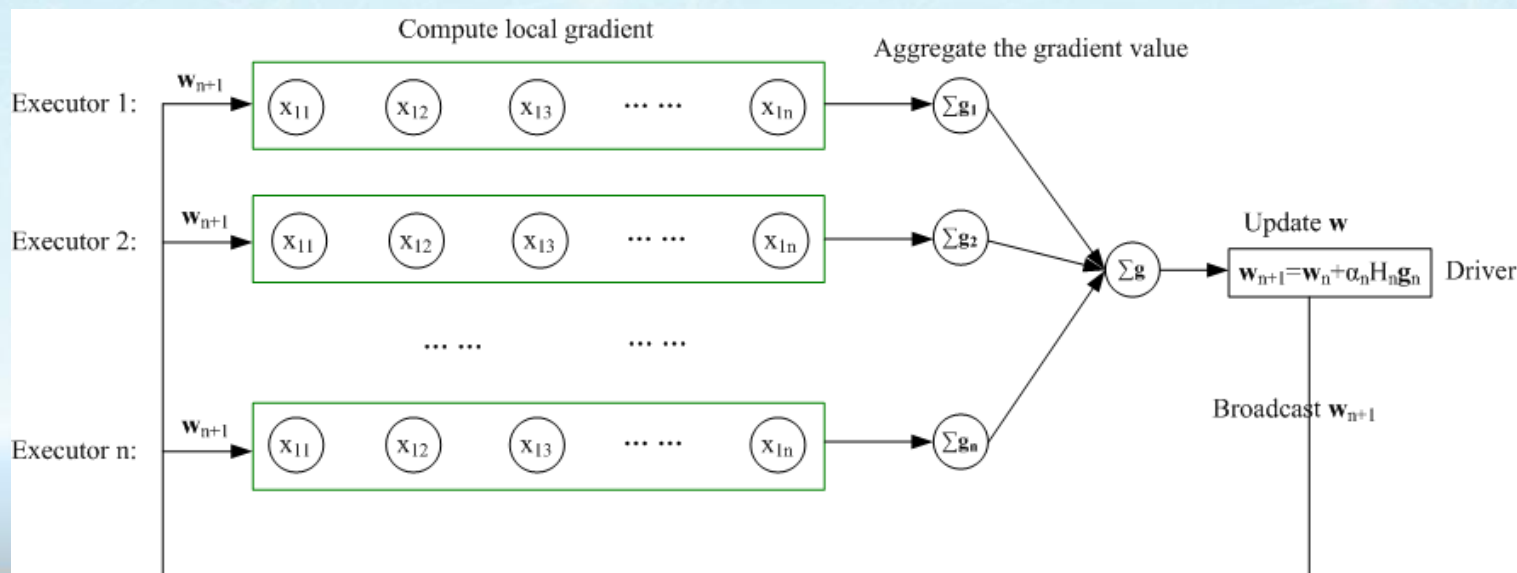
优化过程简单

Factorization Machines – 基于Spark的实现

实现了一个基于OWLQN的算法，以期得到稀疏解：

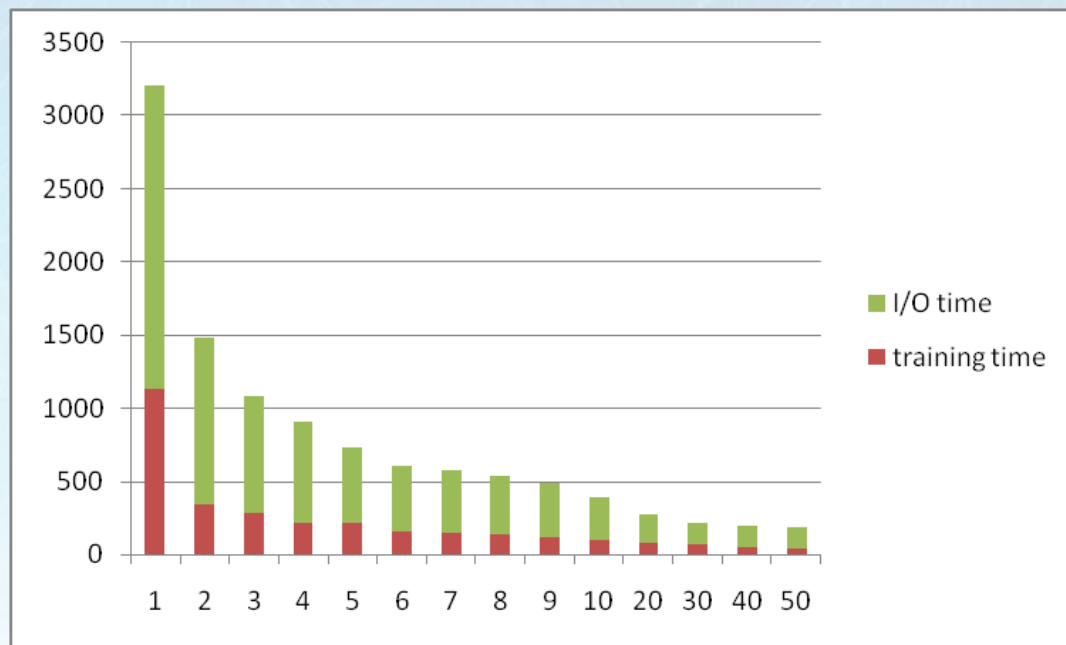
$$\hat{y}(x) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle v_i, v_j \rangle x_i x_j + \lambda \cdot (\sum_{i=1}^n |w_i| + \sum_{i=1}^n |v_i|)$$

Batch Gradient的算法非常适合在Spark上实现



Factorization Machines – 基于Spark的实现

对一个8GB的数据集（约40M样本），比较并行度和处理时间的关系。（设置样本RDD的partition数量为executor数量的3倍）



Executor 数量和处理时间(s)的关系

Bayesian Factorization Machines

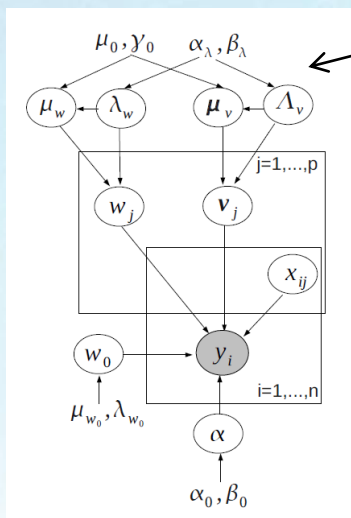
- 为何要Bayes化:
 - 通过Bayes模型我们可以得到一个ctr的分布而非仅仅一个估计值
 - 有利于探索和利用问题:
 - 有利于控制投放速率

Bayesian Factorization Machines – 原理

$$\hat{y}(\mathbf{x}) := w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle \mathbf{v}_i, \mathbf{v}_j \rangle x_i x_j \quad (1)$$



Graph
Model



μ, λ 的先验分布为Gamma分布

假设 w, v 服从 $N(\mu, \lambda)$ 的高斯分布

预测输出:

1) Sampling w, v according to:

$$w \sim N(\mu_w, \lambda_w)$$

$$v \sim N(\mu_v, \lambda_v)$$

2) Compute y based on equation (1)

这样我们得到的 w, v 是一个分布而非固定值，由此能得到预测值 y 的分布

Bayesian Factorization Machines – 基于Gibbs采样求解

设定参数 $\theta_0 = \{\gamma_0, \mu_0, \alpha_0, \beta_0, \alpha_\lambda, \beta_\lambda\}$

初始化 $\theta \sim \mathcal{N}(0, \sigma)$

以指定迭代次数重复下列步骤:

Step1: Update μ_θ

$$p(\mu_\theta | y, X, \theta, \theta_H \setminus \{\mu_\theta\}, \theta_0) = \mathcal{N}(\mu_{\mu_\theta}, \sigma_{\mu_\theta}^2)$$

其中 $\mu_{\mu_\theta} = \sigma_{\mu_\theta}^2 \lambda_\theta \left(\sum_{j=1}^p \theta_j + \gamma_0 \mu_0 \right)$, $\sigma_{\mu_\theta}^2 = ((p + \gamma_0) \lambda_\theta)^{-1}$

Step2: Update λ_θ

$$p(\lambda_\theta | y, X, \theta, \theta_H \setminus \{\lambda_\theta\}, \theta_0) = \Gamma(\alpha_\theta, \beta_\theta)$$

其中 $\alpha_\theta = (\alpha_\lambda + p + 1) / 2$, $\beta_\theta = \left(\sum_{j=1}^p (\theta_j - \mu_\theta)^2 + \gamma_0 (\mu_\theta - \mu_0)^2 + \beta_\lambda \right) / 2$

Step3: Update α

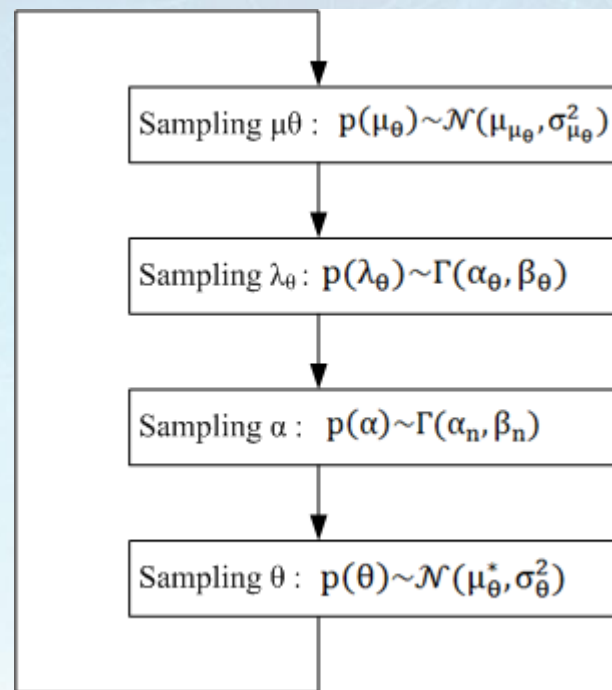
$$p(\alpha | y, X, \theta, \theta_H \setminus \{\alpha\}, \theta_0) = \Gamma(\alpha_n, \beta_n)$$

其中 $\alpha_n = (\alpha_0 + n) / 2$, $\beta_n = \left(\sum_{i=1}^n (y_i - y(x_i, \theta))^2 + \beta_0 \right) / 2$

Step4: Update θ

$$p(\theta | y, X, \theta \setminus \{\theta\}, \theta_H, \theta_0) = \mathcal{N}(\mu_\theta^*, \sigma_\theta^2)$$

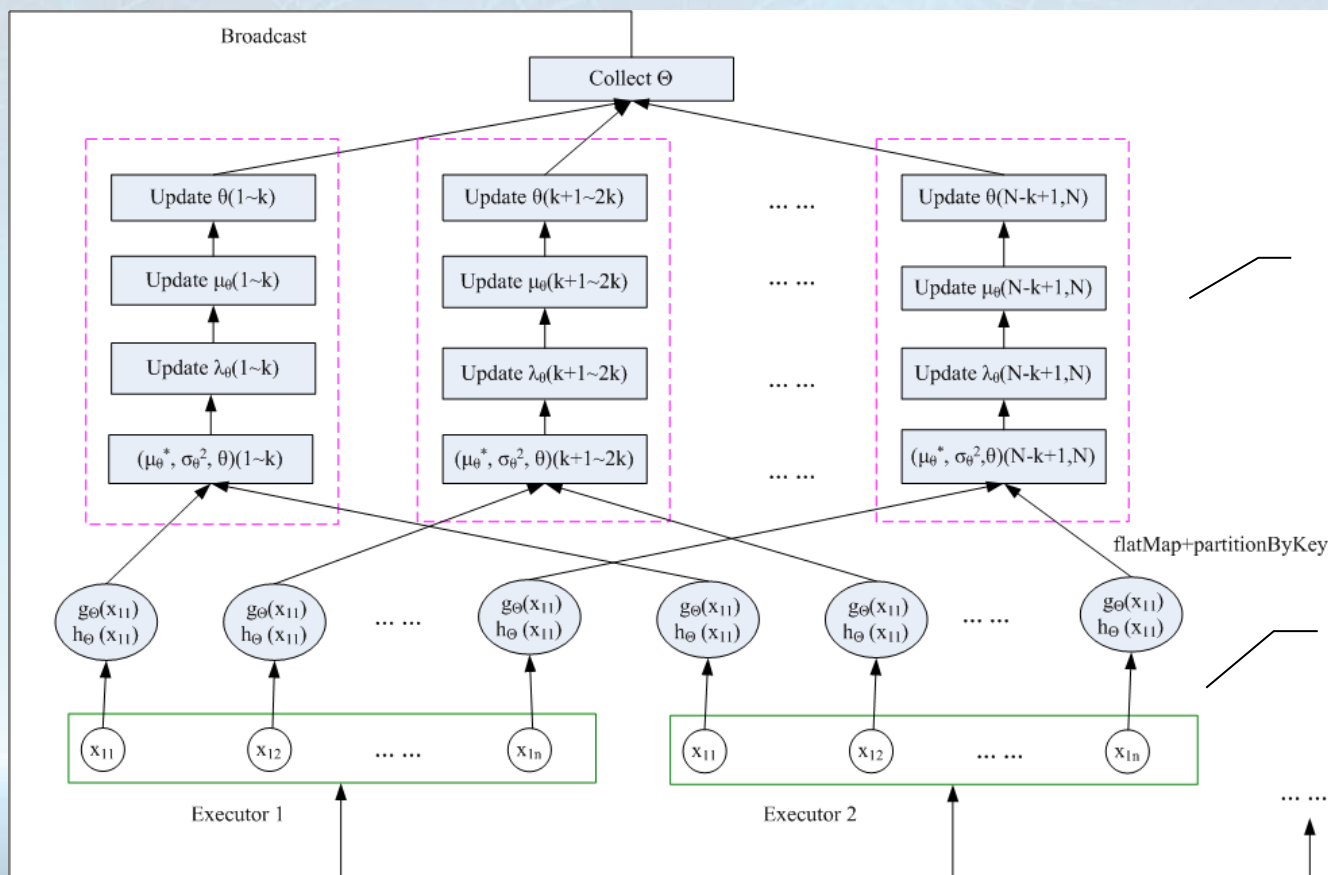
其中 $\mu_\theta^* = \sigma_\theta^2 \left(\sum_{i=1}^n \alpha(y_i - g_\theta(x_i)) h_\theta(x_i) + \mu_0 \lambda_\theta \right)$, $\sigma_\theta^2 = \left(\alpha \sum_{i=1}^n h_\theta(x_i)^2 + \lambda_\theta \right)^{-1}$



扫描整个数据集

基于Gibbs采样的算法步骤

Bayesian Factorization Machines – 基于Spark实现



每个executor更新一个参数子集

分布式计算预测值:

$$y(x|\Theta) = g_{\theta}(x) + \theta h_{\theta}(x)$$

时间复杂度: $O(kN_z)$

Summary

- 广告和推荐中的ctr估计是机器学习的典型应用；
 - 模型能有效地表示多样性的特征
 - 预测对象的动态变化需要模型有探索(Exploration)能力
- Spark可以作为ctr建模的有效的计算平台
 - 用户标签处理需要做大量日志数据的整合关联
 - ctr建模算法往往都是迭代运算



The 7th China
Cloud Computing
Conference

谢谢！