



GO^{ME}

机器学习排序实践

邱宝军 杨骥 李冰青

2016年7月1日

国美大数据研究院

国美在线大数据中心

目录

➤ 01 机器学习排序

➤ 02 推荐排序实践



GOME
.COM.CN 国美在线

排序

带你嗨欢一起玩转高峰期



推荐
vs. 广告
vs. 搜索



排序→机器学习排序

搜索GITC

GITC全球互联网技术大会

What is GITC ? GITC全球互联网技术大会以汇集行业精英、促进技术交流、加深商务合作以及推动行业发展为大会宗旨。通过针对互联网技术领域专门设计的先进开放会议会...

www.thegitc.com ▾

Perfect

全球互联网技术大会_百度百科

GITC (全球互联网技术大会)项目于2013年5月正式发起，通过针对互联网技术领域专门设计的先进开放会议会展内容及形式，联合iTech Club (中国互联网技术...
baike.baidu.com/view/10900771.htm 2016-5-31

Excellent

GITC 2014新闻发布会暨开票仪式在京召开|GITC_滚动新闻 ...

2014年8月6日GITC 2014全球互联网技术大会新闻发布会暨开票仪式于北京五洲皇冠国际酒店隆重召开。作为GITC 2014全球互联网技术大会筹备召开的正式启动
tech.sina.com.cn 新浪科技 | 滚动新闻 | 2014GITC ... ▾ 2014-9-30

Good

GITC (团体儿童智力测验)_百度百科

中文名 团体儿童智力测验 外文名 GITC 优点 客观、精练、易操作、经济快速等 主要特点适用于团体快速施行,非语言量表
baike.baidu.com/subview/11415563/14249887.htm ▾ 2016-1-17

Fair

Greater Illinois Title Company 翻译此页

Bad



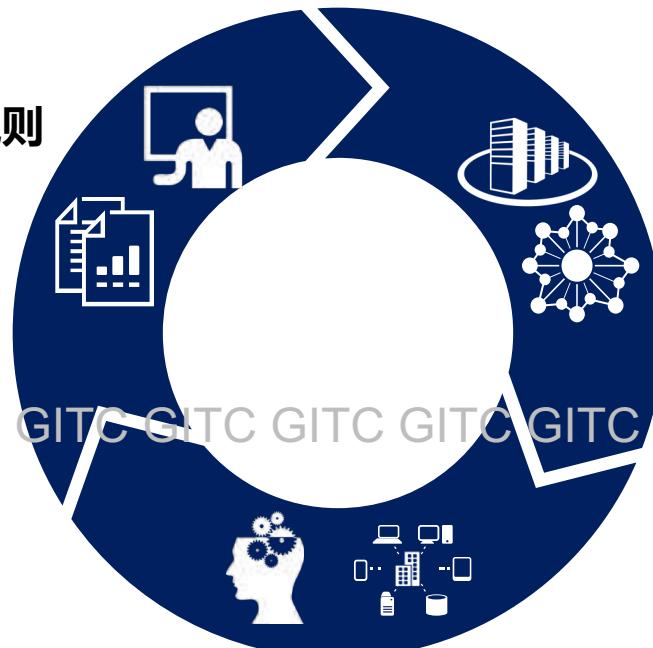
As the largest title agency in Illinois, Greater Illinois Title Company provides a single point of contact for all title and closing...

搜索词、目标文档(网站)、搜索人、时间、环境...

Query	URL	Label	Page-Rank	BM25	...	Feature_N
GITC	Thegitc.com	Perfect	0.6	0.9		100
GITC	...	Excellent	0.9	0.9		66
GITC	...	Good	0.8	0.7		75
GITC	...	Fair	0.9	0.6		32
GITC	...	Bad	0.2	0.65		120
国美	...	Perfect	0.6	0.8		150
国美	...	Perfect	0.5	0.7		50
...						

人工 vs 机器学习排序

人工规则



人工规则 vs. 机器学习
(广告 vs. 搜索 vs. 推荐)

机器学习排序(Learning to Rank)

- ✓ Pointwise

排序问题转化为多类分类问题或者回归问题，相关度相同的为一类

- ✓ Pairwise

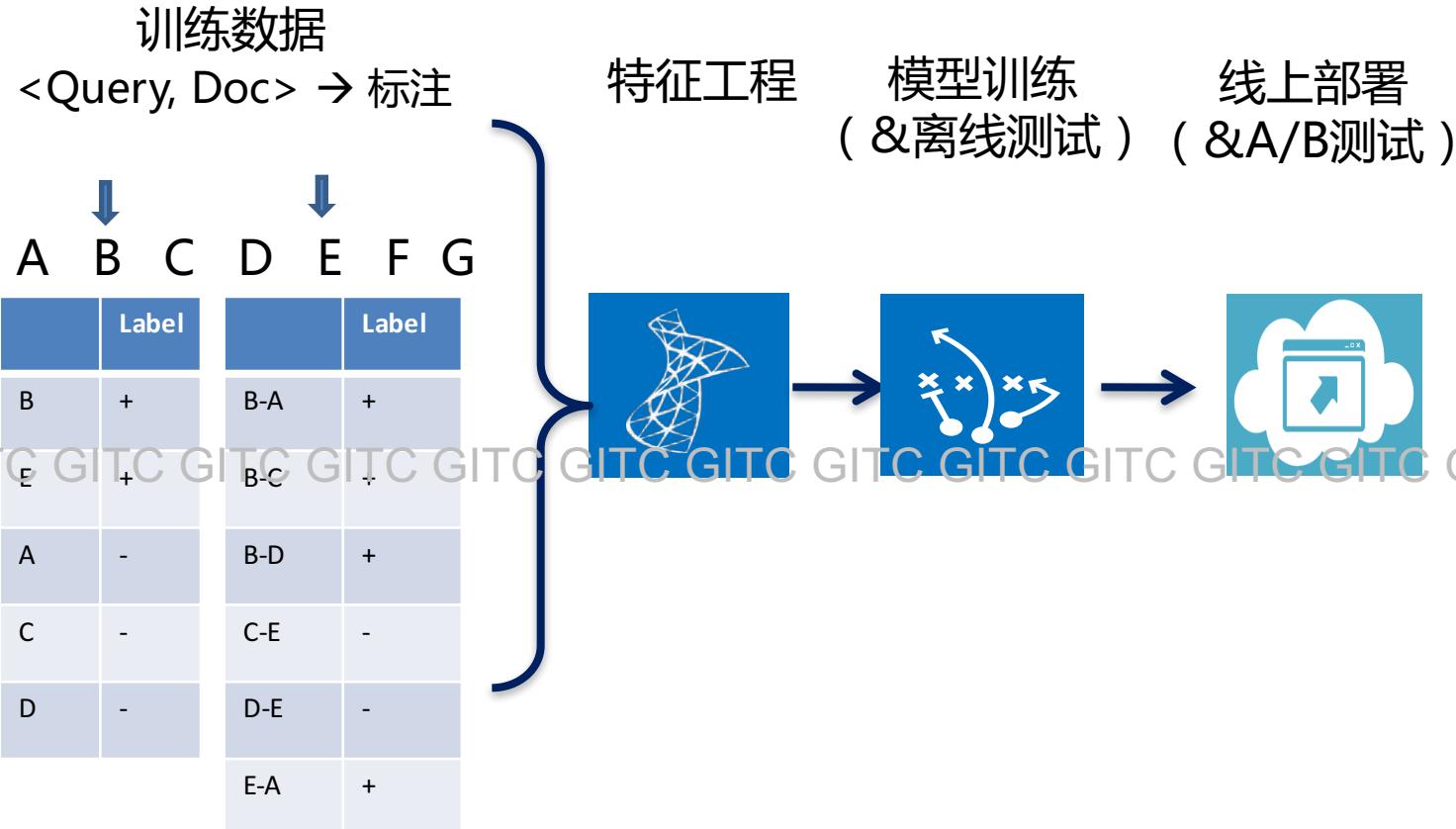
排序问题转化为二分类问题，由两两之间的偏序关系得到全局的排序

- ✓ Listwise

直接对排序结果进行优化



机器学习排序-流程



特征工程-从实体到特征

用户属性、
购买力、
品类倾向



用户

查询量、
转化率、
品类分



查询

UV、PV、
订单量、
转化率



商品

店铺评分、
店铺流量、
店铺评价



商家

当地温度、
污染指数、
人口数量



环境

.....

关键词匹配、价格匹配等

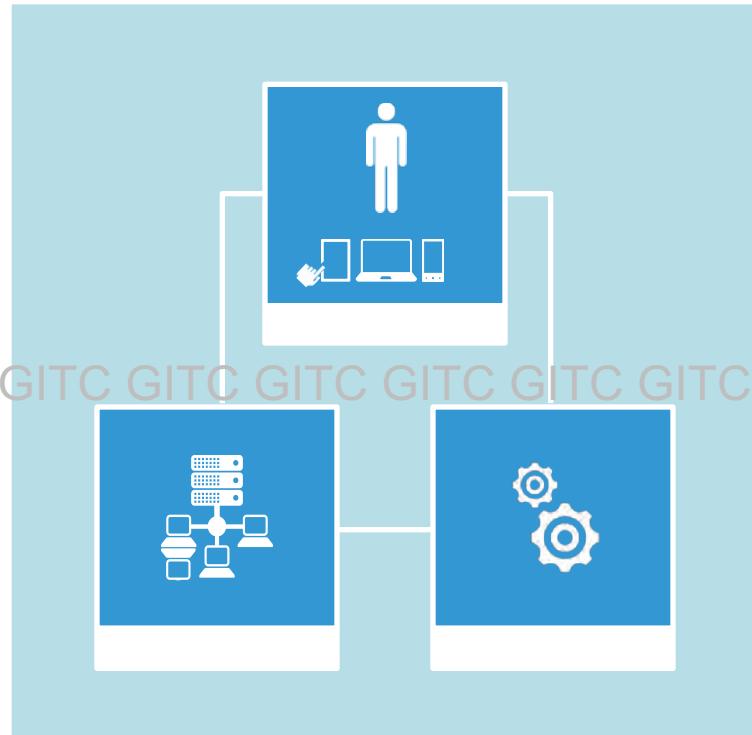
区域销量、商品天气匹配等

品类匹配、购买力匹配等

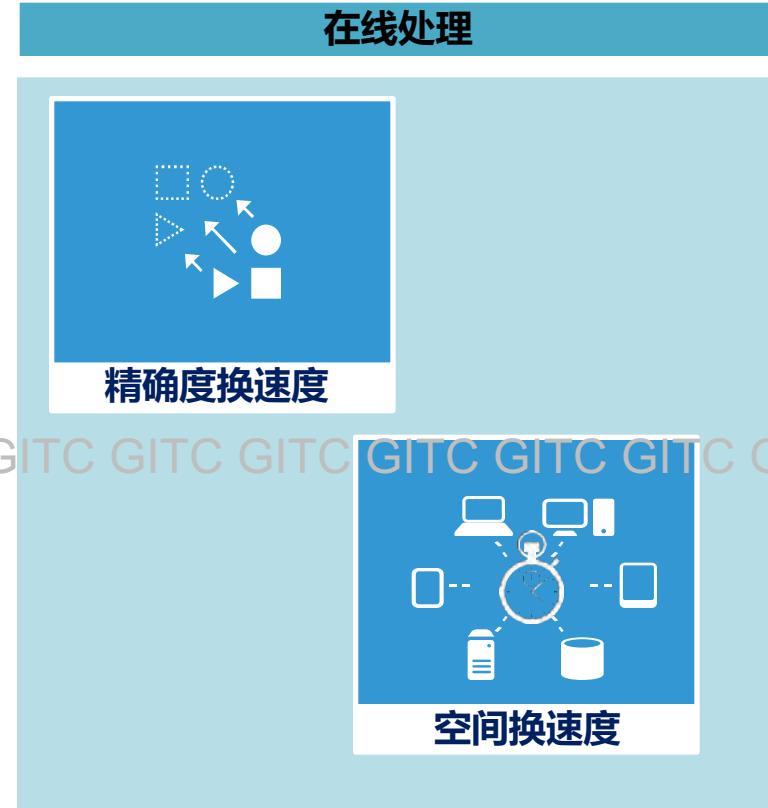
- 协同过滤
 - Deep Learning (RBM, CNN, RNN, ...)
- 回归模型 (LR, SVR...)
- 矩阵分解 (SVD, LFM, SVD++, ...)
- 马尔科夫链
 - Association Rules
 - GBDT/RF
- 聚类 (K-means, ...)
- ...

机器学习排序-部署效率

离线预处理



在线处理



量化与A/B测试

- 可量化目标：分类准确率、NDCG、点击率、转化率
- 假设→试验→结论



德鲁克：如果你不能衡量它，就无法增长它



GOME
COM.CN 国美在线

国美推荐A/B测试配置及量化

第一步：更改A组权重



第三步：完成分组流量配置及分组备注



第二步：更改B组配置及流量



持续观察KPI指标

ID	商品ID	商品名称	权重	CVR	在天猫的转化率	日销量	购买量	UV数	点击量	曝光量	推荐量
1	商品1	商品1	100	0.05	5%	100	50	1000	1000	1000	1000
2	商品2	商品2	20	0.05	5%	100	50	1000	1000	1000	1000
3	商品3	商品3	10	0.05	5%	100	50	1000	1000	1000	1000
4	商品4	商品4	10	0.05	5%	100	50	1000	1000	1000	1000
5	商品5	商品5	10	0.05	5%	100	50	1000	1000	1000	1000
6	商品6	商品6	10	0.05	5%	100	50	1000	1000	1000	1000
7	商品7	商品7	10	0.05	5%	100	50	1000	1000	1000	1000
8	商品8	商品8	10	0.05	5%	100	50	1000	1000	1000	1000
9	商品9	商品9	10	0.05	5%	100	50	1000	1000	1000	1000
10	商品10	商品10	10	0.05	5%	100	50	1000	1000	1000	1000
11	商品11	商品11	10	0.05	5%	100	50	1000	1000	1000	1000
12	商品12	商品12	10	0.05	5%	100	50	1000	1000	1000	1000
13	商品13	商品13	10	0.05	5%	100	50	1000	1000	1000	1000
14	商品14	商品14	10	0.05	5%	100	50	1000	1000	1000	1000
15	商品15	商品15	10	0.05	5%	100	50	1000	1000	1000	1000
16	商品16	商品16	10	0.05	5%	100	50	1000	1000	1000	1000
17	商品17	商品17	10	0.05	5%	100	50	1000	1000	1000	1000
18	商品18	商品18	10	0.05	5%	100	50	1000	1000	1000	1000
19	商品19	商品19	10	0.05	5%	100	50	1000	1000	1000	1000
20	商品20	商品20	10	0.05	5%	100	50	1000	1000	1000	1000
21	商品21	商品21	10	0.05	5%	100	50	1000	1000	1000	1000
22	商品22	商品22	10	0.05	5%	100	50	1000	1000	1000	1000
23	商品23	商品23	10	0.05	5%	100	50	1000	1000	1000	1000
24	商品24	商品24	10	0.05	5%	100	50	1000	1000	1000	1000
25	商品25	商品25	10	0.05	5%	100	50	1000	1000	1000	1000
26	商品26	商品26	10	0.05	5%	100	50	1000	1000	1000	1000
27	商品27	商品27	10	0.05	5%	100	50	1000	1000	1000	1000
28	商品28	商品28	10	0.05	5%	100	50	1000	1000	1000	1000
29	商品29	商品29	10	0.05	5%	100	50	1000	1000	1000	1000
30	商品30	商品30	10	0.05	5%	100	50	1000	1000	1000	1000
31	商品31	商品31	10	0.05	5%	100	50	1000	1000	1000	1000
32	商品32	商品32	10	0.05	5%	100	50	1000	1000	1000	1000
33	商品33	商品33	10	0.05	5%	100	50	1000	1000	1000	1000
34	商品34	商品34	10	0.05	5%	100	50	1000	1000	1000	1000
35	商品35	商品35	10	0.05	5%	100	50	1000	1000	1000	1000
36	商品36	商品36	10	0.05	5%	100	50	1000	1000	1000	1000
37	商品37	商品37	10	0.05	5%	100	50	1000	1000	1000	1000
38	商品38	商品38	10	0.05	5%	100	50	1000	1000	1000	1000
39	商品39	商品39	10	0.05	5%	100	50	1000	1000	1000	1000
40	商品40	商品40	10	0.05	5%	100	50	1000	1000	1000	1000
41	商品41	商品41	10	0.05	5%	100	50	1000	1000	1000	1000
42	商品42	商品42	10	0.05	5%	100	50	1000	1000	1000	1000
43	商品43	商品43	10	0.05	5%	100	50	1000	1000	1000	1000
44	商品44	商品44	10	0.05	5%	100	50	1000	1000	1000	1000
45	商品45	商品45	10	0.05	5%	100	50	1000	1000	1000	1000
46	商品46	商品46	10	0.05	5%	100	50	1000	1000	1000	1000
47	商品47	商品47	10	0.05	5%	100	50	1000	1000	1000	1000
48	商品48	商品48	10	0.05	5%	100	50	1000	1000	1000	1000
49	商品49	商品49	10	0.05	5%	100	50	1000	1000	1000	1000
50	商品50	商品50	10	0.05	5%	100	50	1000	1000	1000	1000

国美推荐A/B测试量化效果

对比A/B测试效果

开始时间: 2016-04-06 结束时间: 2016-04-30 检索查询:

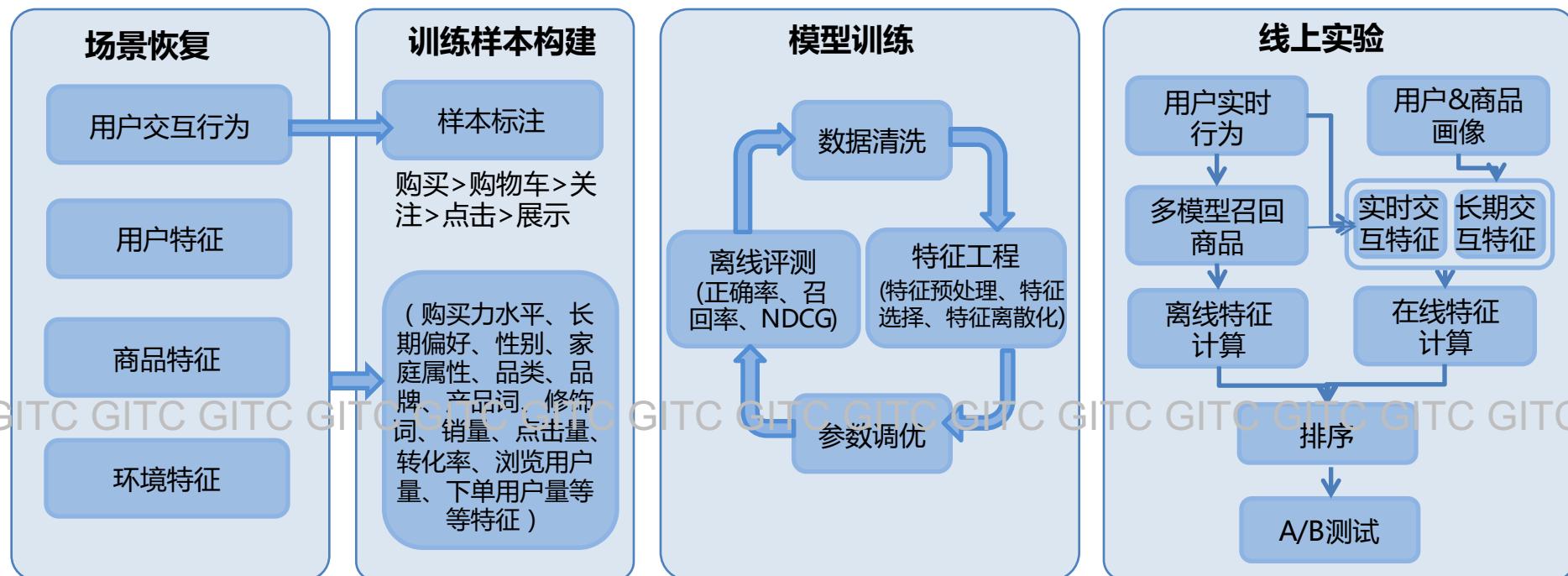
基础数据列表							
	活动编号	活动优先级	分组序号	分组权重	转化率均值	标准差	95%置信区间
<input checked="" type="checkbox"/>	1	高	03	1	0.00%	0.00%	0.00% - 0.00%
<input checked="" type="checkbox"/>	1	高	45	1	0.00%	0.00%	0.00% - 0.00%
<input type="checkbox"/>	1	高	47	1	0.00%	0.00%	0.00% - 0.00%

T检验的P值: 0.003 两组差异95%的范围: 0.000 - 0.000

A/B测试转化率对比



机器学习排序-小结



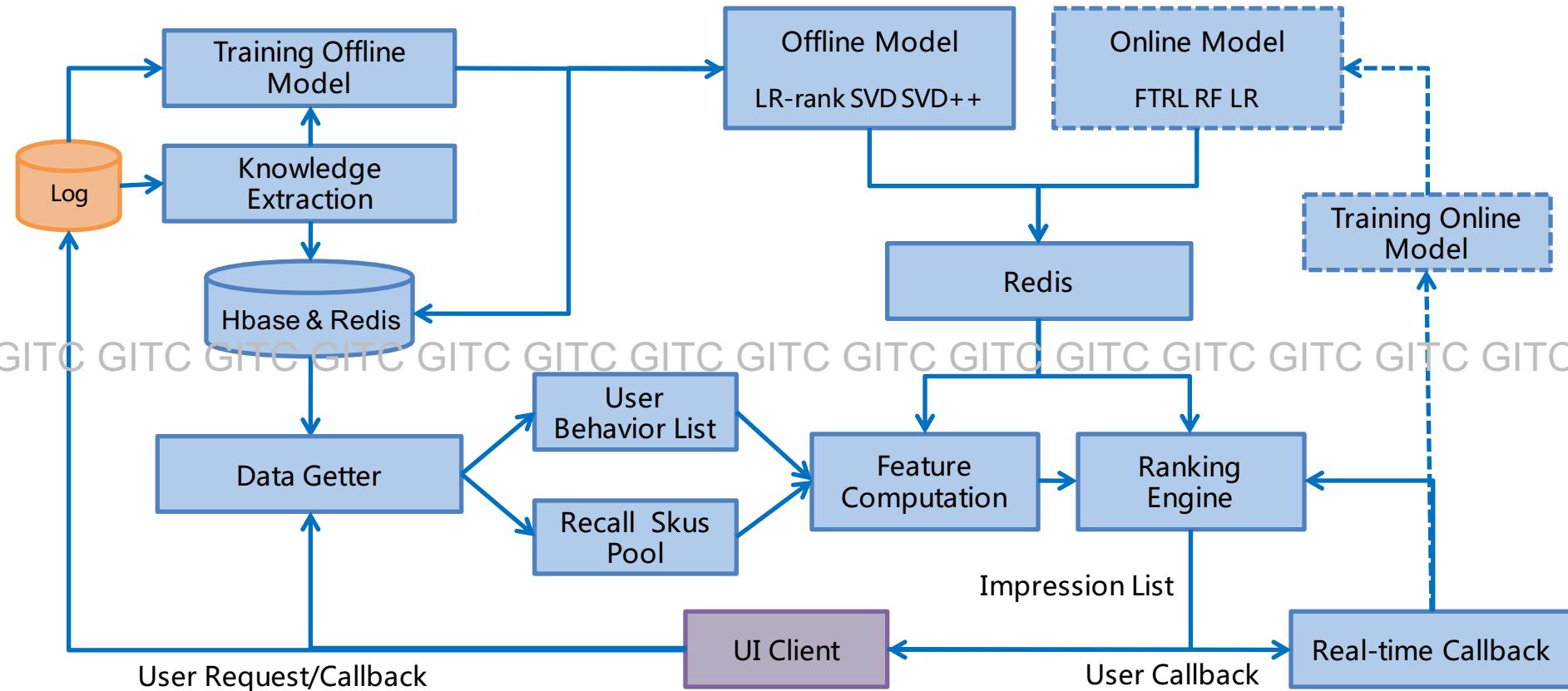
目录

➤ 01 机器学习排序

➤ 02 推荐排序实践



推荐排序实践：架构



在线实时排序

- 过滤：业务需求（无库存、成人用品、促销商品）
- 实时特征计算：人vs.商品，准实时行为特征
- 多样性：多模型、多数据源融合
- 新颖度
- CTR/CVR
- 在线A/B测试

离线数据挖掘

- 隐语义模型数据：矩阵分解
- 协同过滤：UserBasedCF、ItemBasedCF
- 用户画像：品类、品牌、性别、年龄、购买力等
- 商品画像：价格指数、适用人群和地理区域、购买周期等
- 计算机视觉特征
- 离线测试



1. 收集用户行为，包括点击、加购、关注、下单等



2. 对行为进行过滤，比如：join(白名单)、统计截断、position-bias、多次加/删购等处理



3. 制定行为评分规则，生成评分矩阵



4. 训练矩阵分解模型



5. 导入缓存

- 用户因子向量 $p_u \in \Re^f$ 和商品因子向量 $q_v \in \Re^f$
- 评分偏置(bias) $b_{uv} = \mu + b_u + b_v$ (μ 是整个网站评分的均值， b_u 和 b_v 是用户/商品较网站评分的偏差)
- 用户对某商品 v 的预测评分为 $r'_{uv} = b_{uv} + p_u^T q_v$
- SVD++ 加入隐式反馈

$$r'_{uv} = b_{uv} + q_v^T \left(|R(u)|^{-\frac{1}{2}} \sum_{j \in R(u)} (r_{uj} - b_{uj}) x_j + |N(u)|^{-\frac{1}{2}} \sum_{j \in N(u)} y_j \right)$$

■ 其中

- $q_v, y_v, x_v \in \Re^f$ 是商品的三个隐向量
- $R(u)$ 代表被用户 u 打过分的商品集合
- $N(u)$ 代表用户隐式选择的商品集合（即打分/未打分）



GOME
COM.CN 国美在线

推荐排序实践：环境特征

某地高温
橙色预警

利用天气API定位某地高温橙色预警，该地区用户进入网站，首页猜你喜欢推荐空调、风扇等商品：

猜你喜欢—你的专属推荐

换一组



海尔(Haier) KFR-23GW/12N
WA13套机 小1匹P壁挂式定频

¥1799.00

伊莱克斯(Electrolux)EAW25
FD13CA1 1匹P壁挂式定频

¥1299.00

美的 (Midea) AC120-16B
W 空调扇 (冷风机 加湿 净化)

¥328.00

美的 (Midea) AC120-16BR
W 空调扇(快速制冷 三档风)

¥429.00

青岛啤酒 纯生 330ml*24听 搭
国际啤酒企业品质听装啤酒

¥110.00

青岛啤酒 经典 500ml*12听 搭
国际啤酒企业品质听装啤酒

¥69.00

某地雾霾
红色预警

利用天气API定位某地雾霾红色预警，该地区用户进入网站，首页猜你喜欢推荐空气净化器、口罩等商品：

猜你喜欢—你的专属推荐

换一组



飞利浦 (PHILIPS) AC4076
空气净化器 (智能空气测控)

¥2499.00

格力大松 (Gree) KJFC230
A-WG 空气净化器

¥2579.00

3M9001V防雾霾口罩9002V
带呼吸阀男女骑行PM2.5粉

¥90.00

美国 3M 9021 9022 KN90防
颗粒物防尘口罩工业粉尘 防

¥110.00

养生堂 天然维生素E软胶囊2
50mg*90粒 VE ve 维生素e

¥69.00

养生堂天然维生素C咀嚼片85
0mg*80片

¥82.00



图像标注

对国美全站的商品图像进行标注：

- 最主要的工作是数据清洗，即把不能和品类对应的图片删除或者重新进行品类校准
- 统计国美全站最近一年内各品类下商品总数的分布；然后按照分布进行图片抽样



提取特征

图片的特征分为两个部分，一是通过深度学习得到的特征，二是图像色局部特征：

利用caffe训练CNN，将倒数第二层输出作为Feature Learning的结果提取出来
(Deep Learning + Transfer Learning)

利用局部特征算子(SIFT,kaze等)提取出图像的局部特征



降维

分别对前面两种特征进行降维处理：

用积量化(Quantization)的方法对CNN特征进行降维

用Fisher Vector对图像局部特征进行降维



匹配

采用最近邻搜索的方法找出每一个商品的相似商品集合



GOME
.COM.CN 国美在线

深度学习特征

手工特征：

SIFT [Lowe 99]

Spin Images [Johnson&Herbert 99]

Textons [Malik et al. 99]

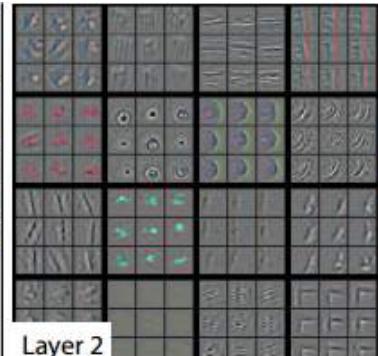
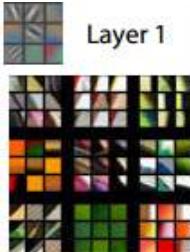
RIFT[Lazebnik 04]

GLOH[Mikolajczyk&Schmid 05]

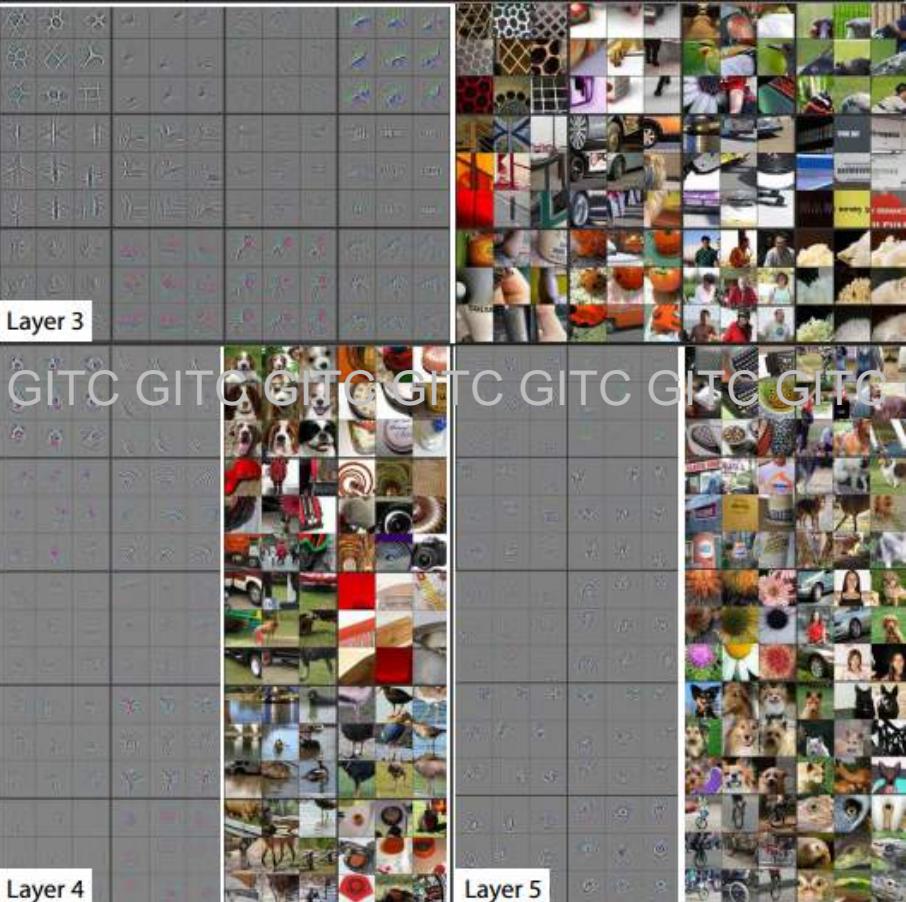
HOG[Dalal&Triggs 05]

GITC GITC GITC GITC GITC GITC GITC GITC

深度学习特征：



Zeiler&Furges, ECCV2014





推荐排序实践：提升推荐个性化及转化率

未经机器学习排序，网站首页猜你喜欢结果

猜你喜欢—你的专属推荐

换一组



乐扣乐扣(300ML)保温杯
¥89.00



乐扣乐扣保温杯拉恩马克杯水杯
¥122.00



膳魔师(THERMOS)保温保冷杯JNS-
¥172.00



THERMOS/膳魔师 高真空保温保冷
¥175.00



乐扣乐扣保温壶保温杯水杯豪华
¥239.00



富光保温杯保温壶家用商务办公
¥38.00

使用机器学习提升推荐转化率

利用机器学习排序就是从数据中
自动学习模式，在若干限定条件
下，找出全局最优或者局部最优
的近似值

利用机器学习排序后，网站首页猜你喜欢推荐商品

猜你喜欢—你的真实喜好

换一组



贝思客 生日蛋糕 带馅黑巧克力
¥65.00



贝思客 生日蛋糕 干果梅挞芝士
¥66.00



日日特价41003 西门子
¥588.00



凌志H4440行车记录仪 行车夜视
¥389.00



南枫人 冰丝第三件套被子
¥79.00



梦纳云家纺 个性潮流系列纯棉四
季被子
¥109.00

推荐排序实践：A/B测试监测分析



谢谢

Thanks!

