



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



携程基于大数据分析的实时风控体系 介绍

携程 刘江

全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



全球互联网技术大会
GLOBAL INTERNET TECHNOLOGY CONFERENCE



Contents

1

我们的挑战

2

Aegis系统架构

3

核心模块介绍

4

风控模型和策略

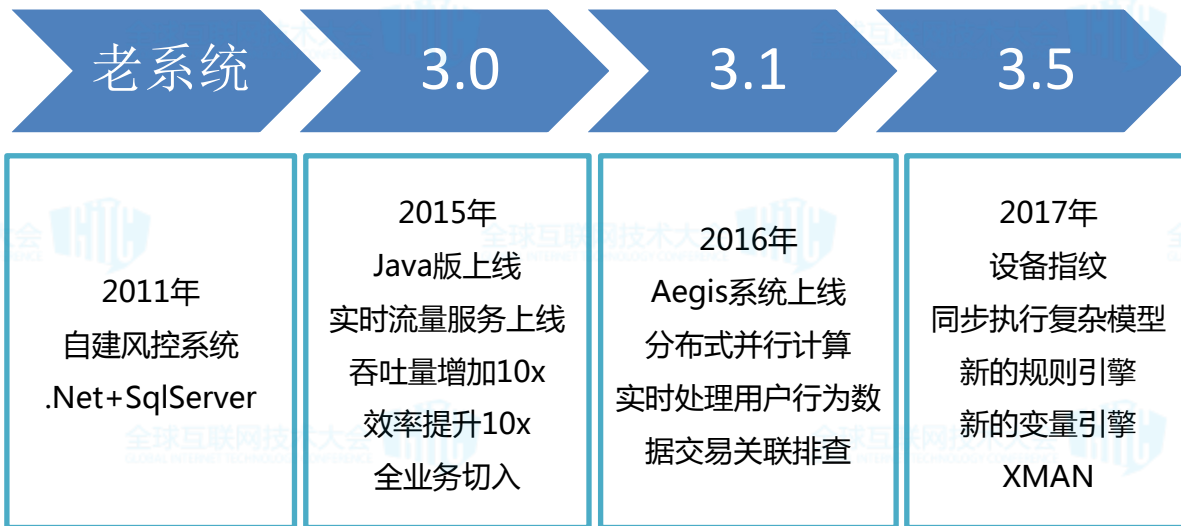
■ 携程文化：让旅游变得更幸福

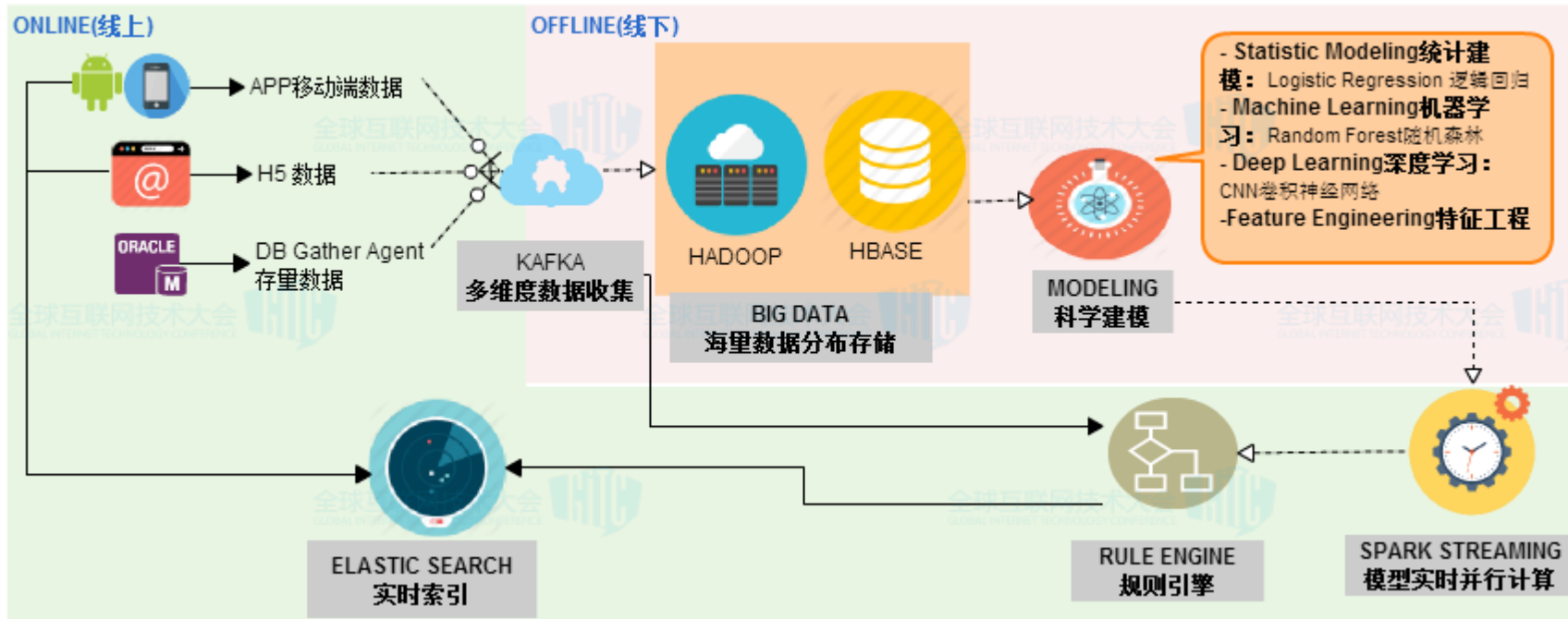
■ 风险管理文化：

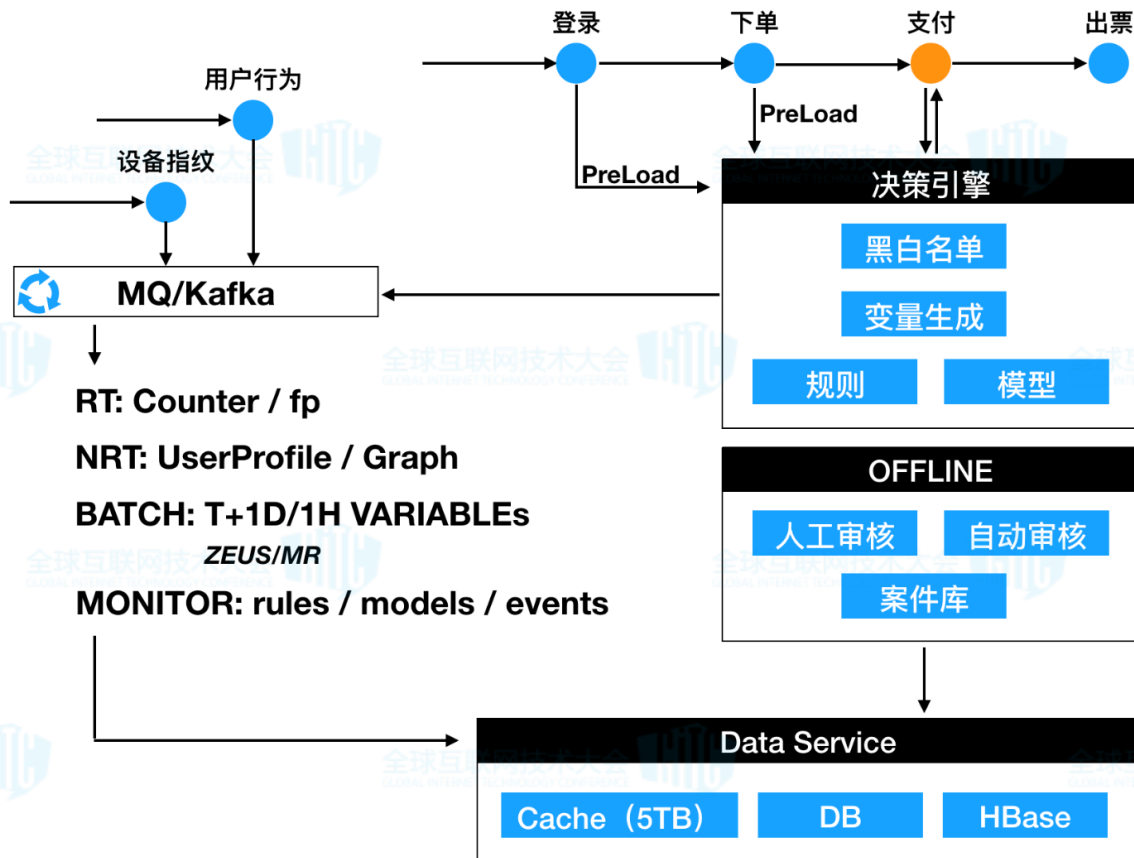
“Make the Travel More Freely
and Securely”

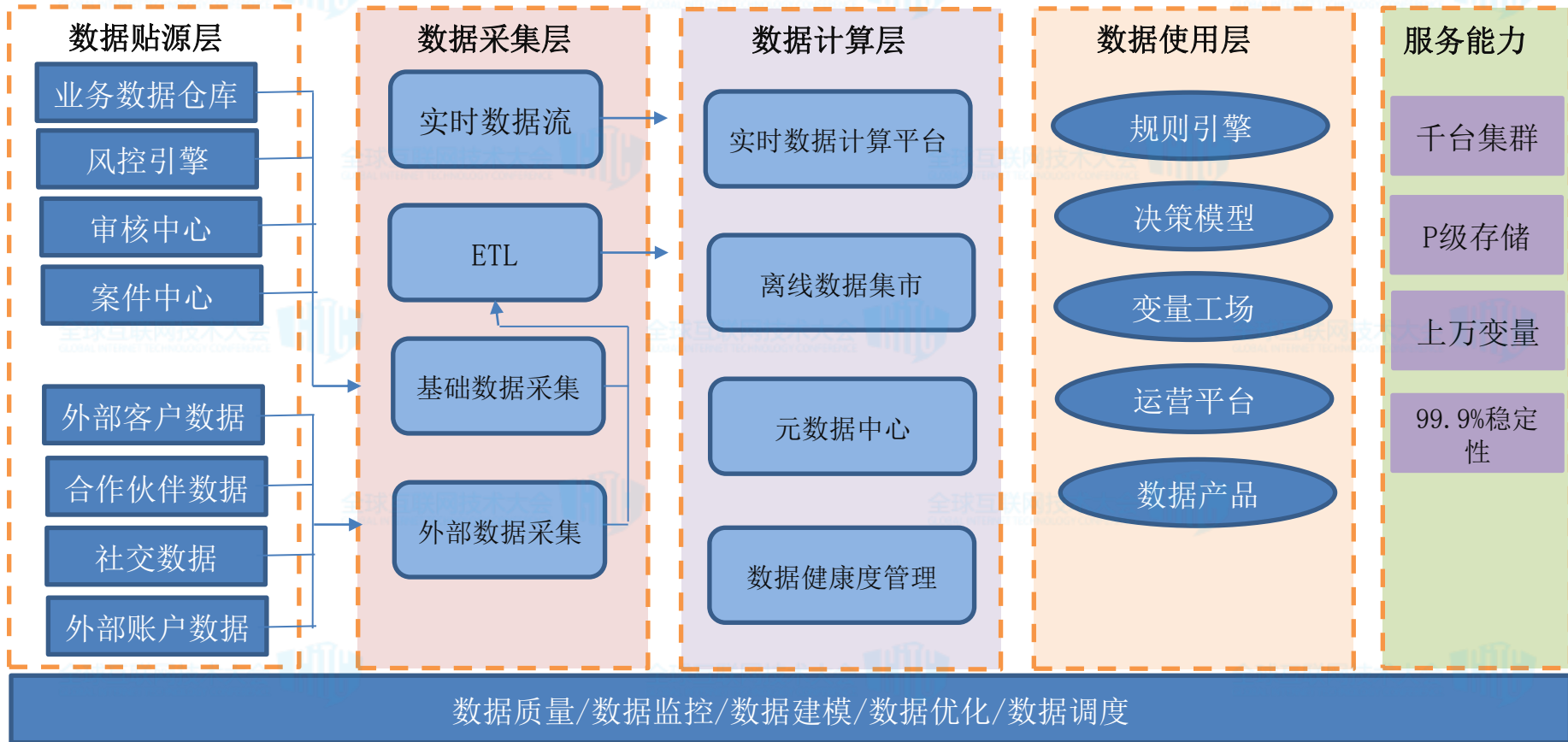


- 业务类型和数据量增长
- 需要更加自动化和智能化
- 用好设备和行为数据
- 跨海外网络的数据延迟









实时

- 黑白名单
- 数据预处理 / 变量衍生 ~ 1000-2000个
- 执行规则 ~ 400条+
- 执行模型 ~ 5-10个
- 结果计算及后处理

异步

■

- 日**亿级**交易处理能力
- 支付风控平均处理时长**小于150ms**，99.9%线600ms
- 支持DR灾备，数据分级存储，7×24H监控&预警
- 通用性强：
 - 支持支付风控、业务风控、外部合作伙伴风控支持

规则引擎
模型执行器
变量服务

实时流量

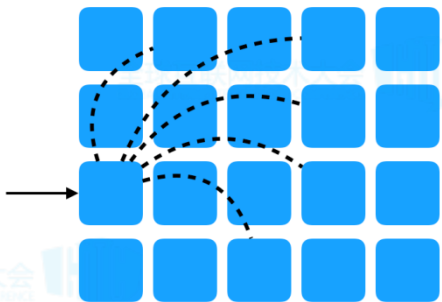
设备指纹

行为分析

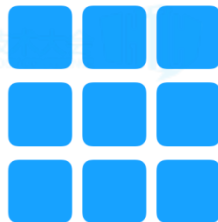
用户画像



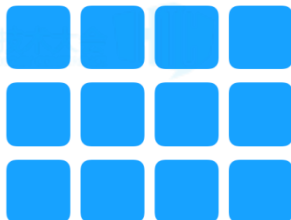
支付



机票



营销活动



Rule Engine

规则分布式并行执行

支持按业务分组

支持动态调整分组和扩容

基于Java，高吞吐量、低延迟

使用gRPC互联

	Python	JPMML	自主研发
特点	标准、开源，兼容性好	标准、开源，兼容性好。	使用Java解析并执行.dot模型文件，支持随即森林和逻辑回归算法，算法可扩展
性能	10-100ms，因需要独立部署，有网络开销	性能和Python执行.dot接近，只是可以嵌入式运行，所以稳定性比Python高	0-10ms，嵌入式执行，性能高，稳定性高

特性：

- 使用Java完全自主实现的dot模型执行器，执行耗时只有Python版本的10%
- 拥有完善的模型运行监控和熔断机制

处理提示

全部隐藏

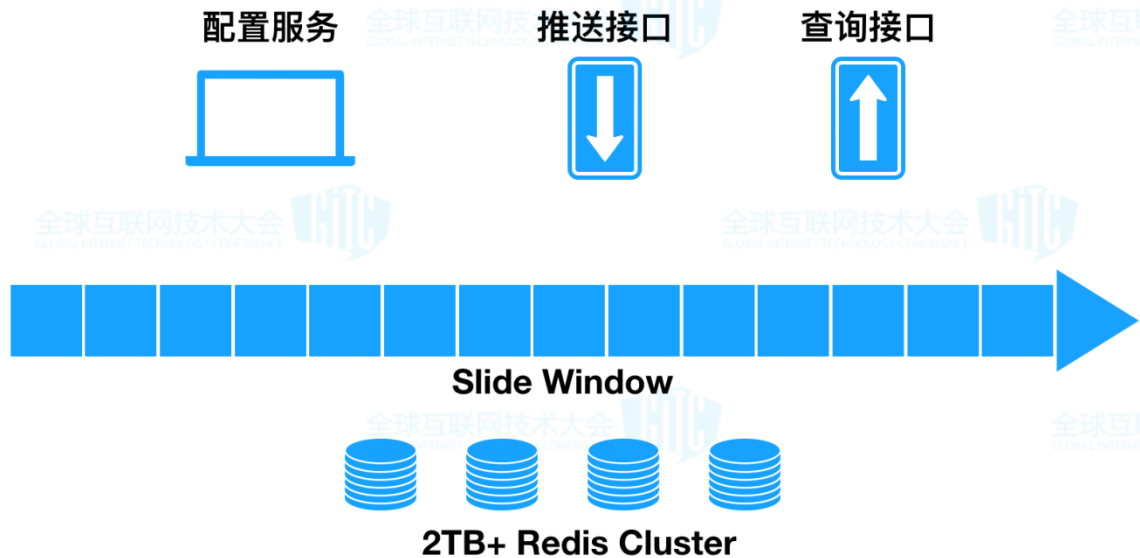
- 负面因素：
1. 证件号可能的国籍：尼日利亚|哈萨克斯坦|马来西亚|韩国|索马里|中国香港|墨西哥|土耳其|南非|荷兰|,
 - 2.
 3. 电话审核分值=0.133,注：[cutoff<=0.05]:
 4. 机票模型分值=[BOTTOM 10%]0.2655,注：[cutoff:<0.5]

处理提示

全部隐藏

- 正面因素：
1. 机票航段非单程机票类型
 2. 登机人包含持卡人姓氏或者全名

- 负面因素：
1. 证件号可能的国籍：摩尔多瓦|越南|蒙古|肯尼亚|波斯尼亚|老挝|瑞士|贝宁|印度尼西亚|黑塞哥维那|奥地利|阿富汗|斯里兰卡|,
 - 2.
 3. 电话审核分值=0.547,注：[cutoff<=0.05]:
 4. 机票模型分值=0.3374,注：[cutoff:<0.5]
 5. 行程涉及高危国家[阿联酋]



Counter

日查询量超100亿次

支持分钟、小时、日、月等多级精度，支持动态配置

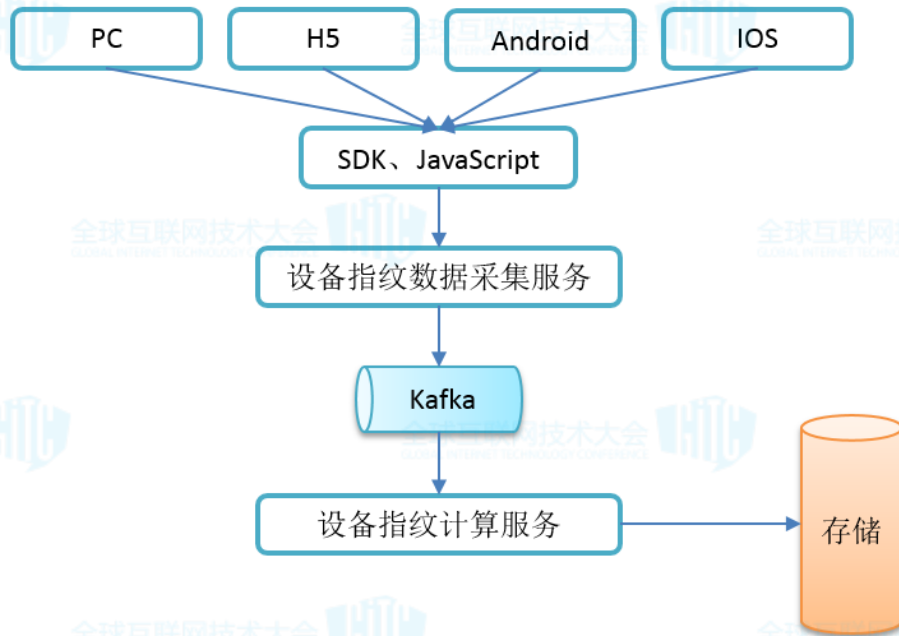
支持3个月以上的超大时间窗口

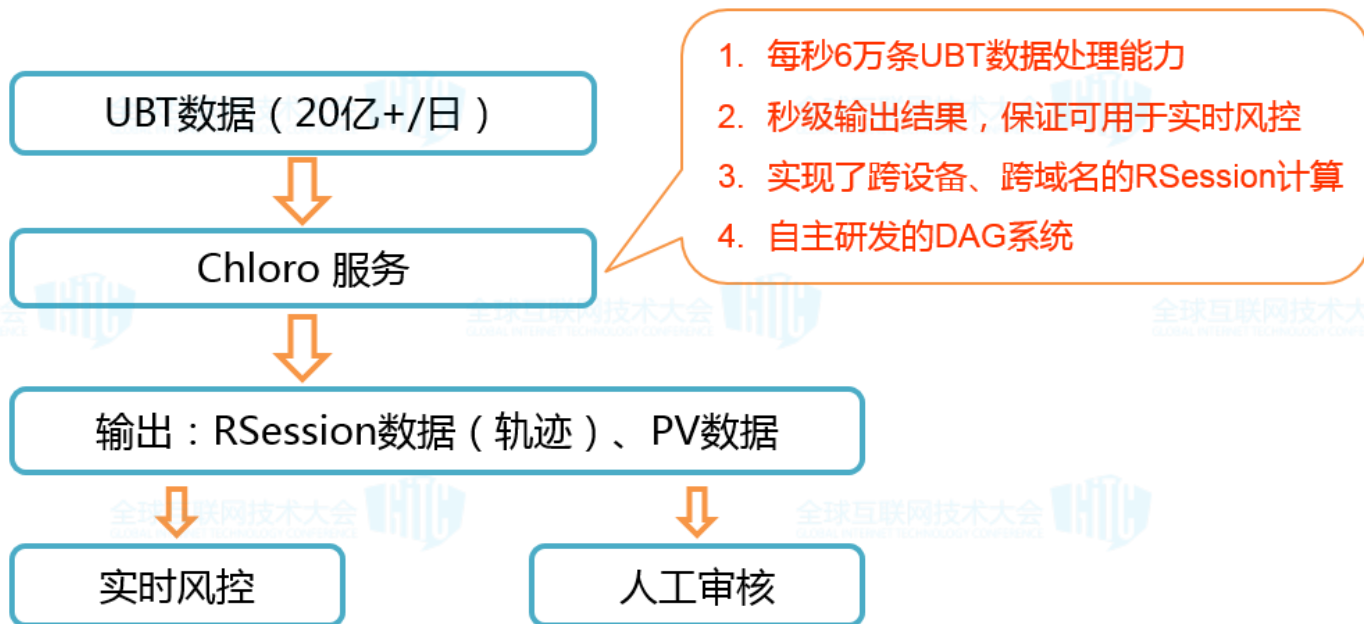
流量数据实时推送，1秒级延迟

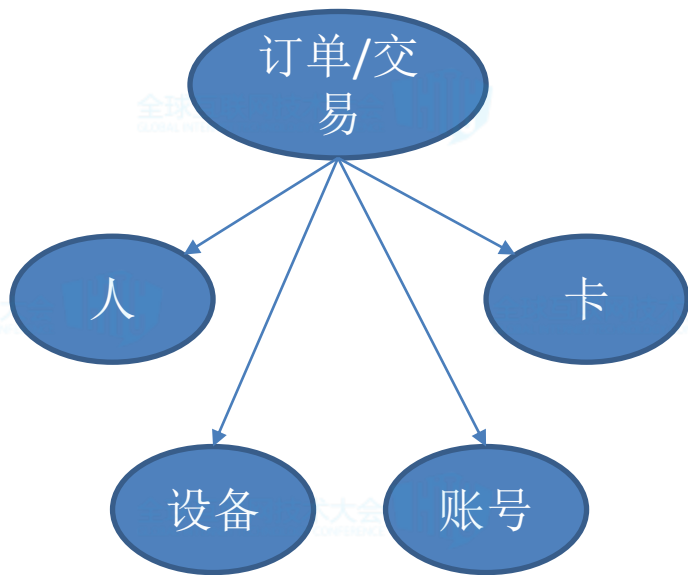
适用于限额限次、Velocity变量和Ratio变量的实时计算

特性：

- 自主研发
- 指纹准确率>99%
- 获取成功率>99%
- 全站部署







Graph

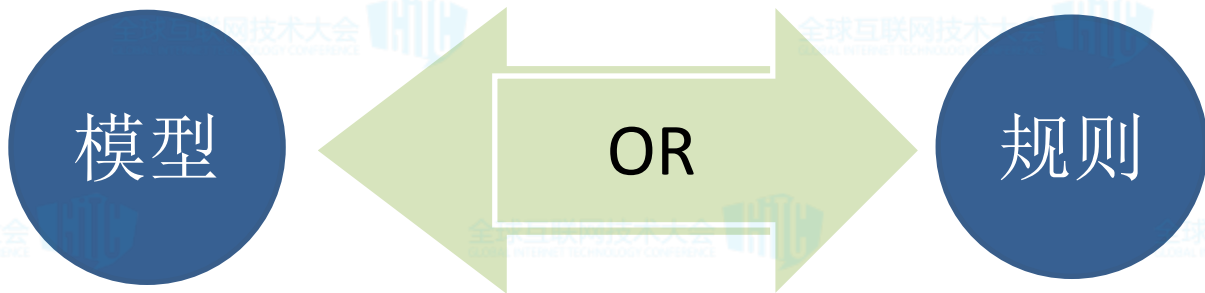
基于交易、人、设备、卡、账号等多个维度的大数据关联分析，确定关联交易。

数据用于规则、模型、和人工案件排查

基于HBase自主实现的Graph存储，50亿+交易数据，1秒级返回关联结果

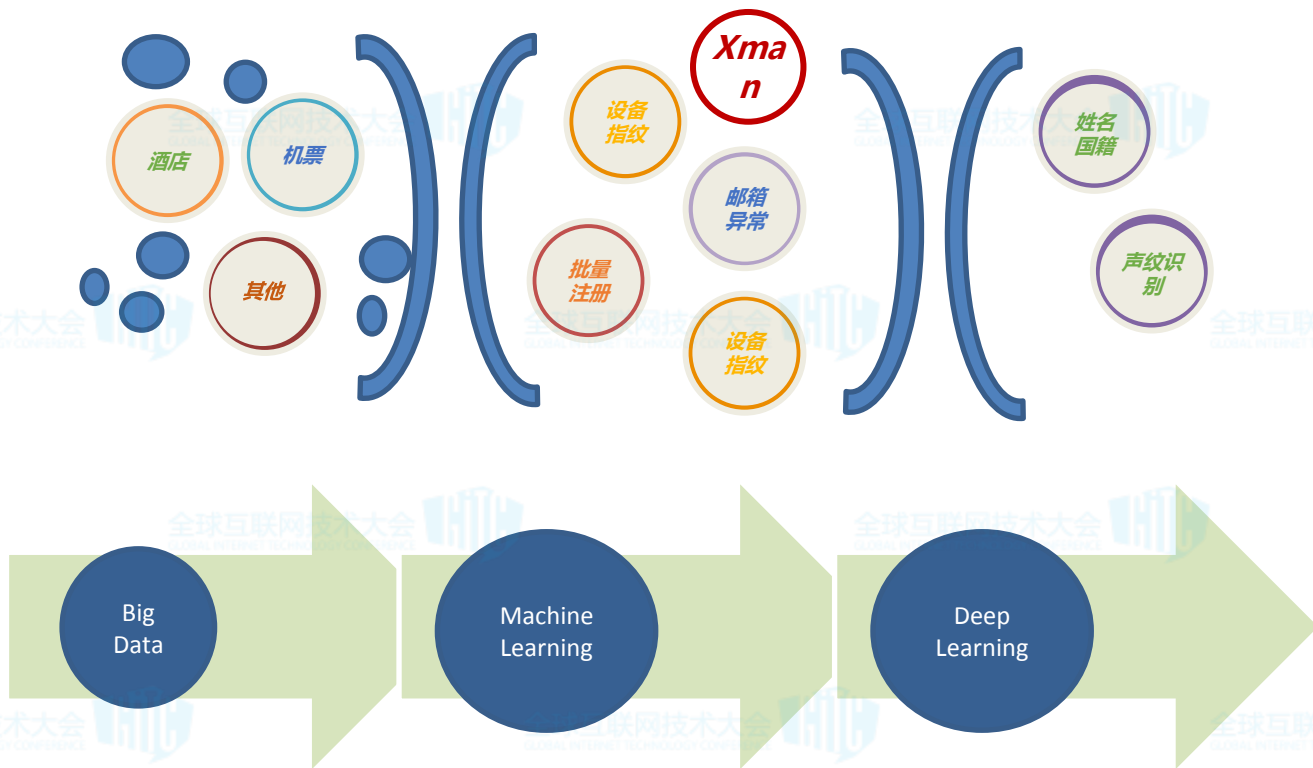
订单编号: 2529073195 命中5个: 1 命中2个: 2 命中1个: 196 全部: 199 查看所选: 0

订单号	订单日期	订单类型	录黑	UID	联系人手机	邮箱	卡号	证件号	IP	DID	ClientID	提现卡号	提现证件号	操作
2529073195	2016-07-25	高铁	F	_WeChat403168883	186			230	14.136.203.107					关联
243599878	2016-07-25	酒店	F	E14011920	135				*		1202102041			关联
2433511525	2016-07-24	酒店	F	M413128776	159				*		1200109431			关联
2469422620	2016-07-17	机票	F	*	133			230	113.0.223.16		1200108061			关联
2412163703	2016-07-15	酒店	F	M289167080	139		64007987		*		1202114021			关联
2412153645	2016-07-15	酒店	F	M289167080	139		64007987		*		1202114021			关联
2409787416	2016-07-15	酒店	F	_WeChat227024822	134				*		1202116471			关联
2409481138	2016-07-14	酒店	F	M289167080	139		64007987		*		1202114021			关联
2407634340	2016-07-13	酒店	F	_WeChat403422737	130				*		1200106311			关联
36811210	2016-07-11	消费券	F	13827797667	138	****@21cn.com			*		1200117781			关联
2398338132	2016-07-09	酒店	F	3000696130	185		63854384		*		1200101341			关联
2264859130	2016-07-09	用车	F	1102803029	136		59169486		*		0903112041			关联
2395956203	2016-07-08	酒店	F	ywsfj800720	139		1069655		*		1200109941			关联
2395603351	2016-07-08	酒店	F	E33650902	186				*		1200113941			关联
2394651946	2016-07-07	酒店	F	2056923792	186				*	82e34d2e-c	1200109621			关联
36648993	2016-07-06	消费券	F	D117876249	159	****@outlook.com			*		1202106981			关联
2391312371	2016-07-05	酒店	F	D117876249	159		46553239		*		1202106981			关联
2442197254	2016-07-04	机票	F	M268797542	*	****@163.com		230	117.40.225.183					关联
2442197240	2016-07-04	机票	F	M268797542	*	****@163.com		*	117.40.225.183	e3125fff-2				关联
2383107138	2016-07-01	酒店	F	3009781965	177				*		1202113341			关联



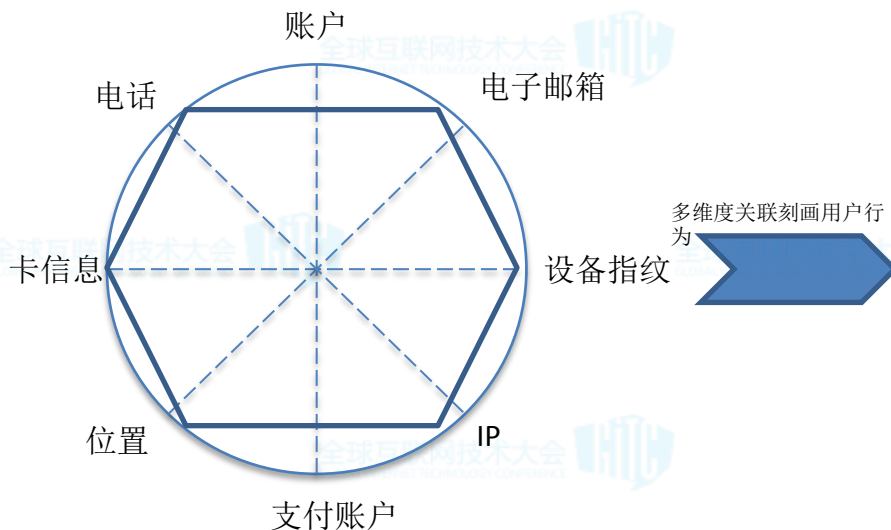
- 模型规则化
- 规则模型化

风控模型和策略



风控模型和策略

海量交易数据信号衍生



衍生方法	刻画pattern
基础衍生	高金额、快速起飞/入住等
冲突变量	信息不一致，例如发卡国和Ip国
Recency	账户年龄，最近一次交易
velocity (单、双主体)	频繁交易\换卡等
过滤条件velocity	频繁高危行为
ratio	高危行为占比，短期交易集中
个体异常	个体行为发生变化/异常
群体异常	行为相对于同地域人群异常
躲闪行为	行为有躲避风控规则的嫌疑
跳跃行为	小额试卡的行为
risktable	历史案件信息的利用

风控模型和策略

特征工程

单人游

- 下单和起飞时间之间的天数
- 手机和ID与ADcity是否冲突
- 保费, 订单金额

家庭游

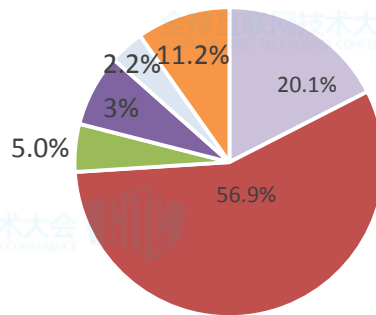
- 常旅客卡
- 订单金额
- 往返航班

好友游

- 航班类型
- 国内国外游

情侣游

- 持卡人非出行人
- 是否同省

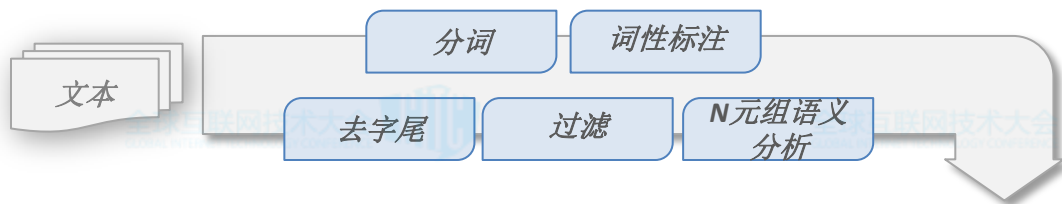


■ 国内因公 ■ 单人游 ■ 家庭游

风控模型和策略

文本信号挖掘

文本处理



NLP/文本分析技术

语音识别

- 声波数据预处理
 - 音频采样
 - 分解频带
 - 傅立叶变换创建识别码
- RNN识别音频片段字符
- 文本数据对深度学习发音预测矫正：

词汇库

- 开发词/词组库，并关联到相关主题和目标
- 考虑同义/下位词
- 使用：
 - 对特定目标识别问题焦点
 - 为模型特征工程做预备

实体识别

- 识别特殊实体类别，如人名，地点，时间，问题类别，关键名词
- 使用：
 - 识别特定种类实体
 - 对特定实体对相关词组分组

基于文本模型

- 统计分析对特定目标字词的相关性
- 根据标签的可用性应用业界最新的无监督或有监督算法
- 使用：
 - 模型特征库萃取

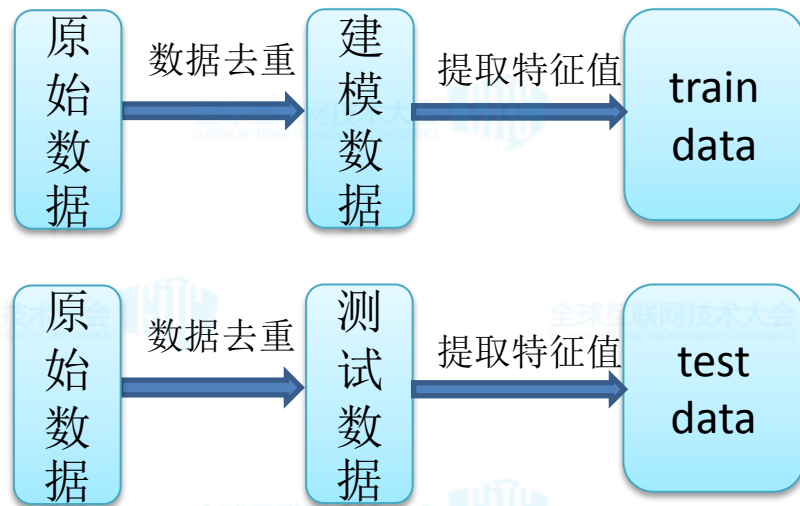
风控模型和策略

字母生成概率

变量注释

变量名

26个字母+10个数字+特殊字符频率	freq_
域名欺诈率	domain_degree
名字模式欺诈率	name_degree
正常名字模式生成可能概率	name_probability
名字复杂度	num_change
名字长度	length
生成概率（数字转移数字概率不为1）	prob_prefix_num_no_1
生成概率（数字转移数字概率为1）	prob_prefix_num_1
生成概率（只有字母）	prob_prefix_alp



风控模型和策略

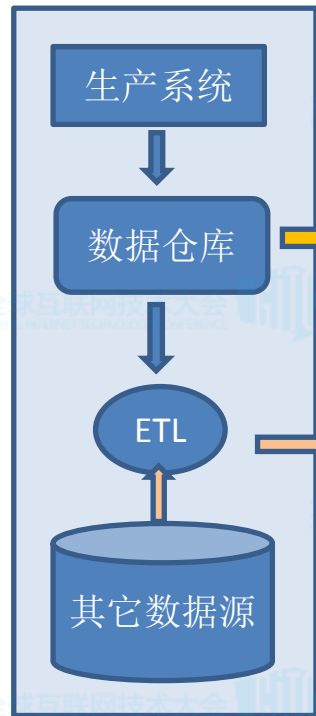
模型工厂



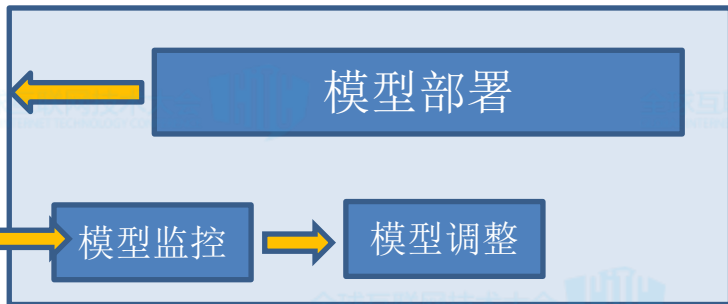
风控模型和策略

模型生命周期

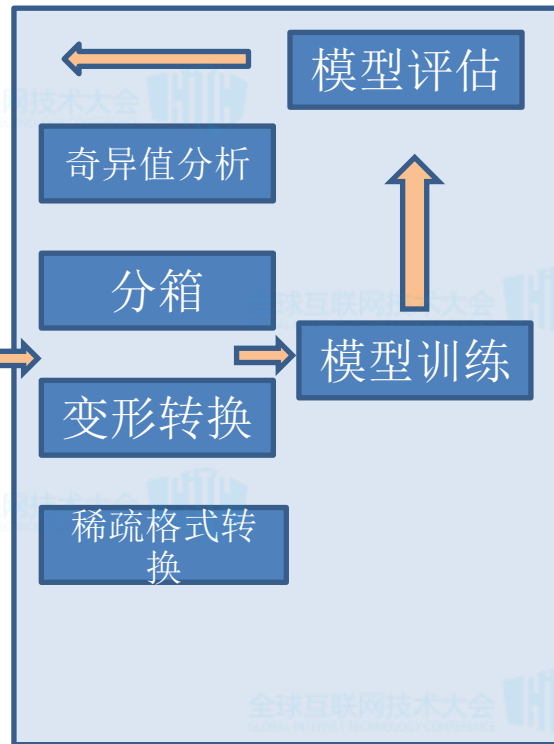
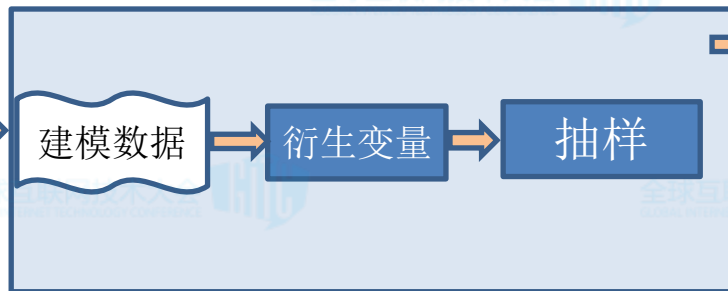
数据集成



模型优化



数据准备



模型开发

风控模型和策略

基于统计分析&机器学习的欺诈策略

欺诈交易识别分二步: 1. 欺诈特征生成; 2. 机器学习模式识别

欺诈特征

- **Velocity** – E.g. 2 连续交易发生在很多时间内



- **Distance to home** – E.g. 持卡人地理位置和常用地址距离很大



- **Transaction time** – E.g. 发生在凌晨的交易欺诈率高



- Etc

模型& 评分

- 利用传统模型方法如逻辑回归, 神经网络, 矩阵因子分解, K最近邻法等

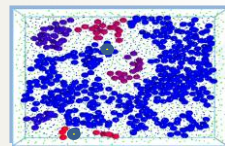
- Matrix Factorization



- Random Forest



- KNN





技术化 not 简单化

自动化 not 排他化

智能化 not 掉包化

多样化 not 单一化

战略化 not 短期化



提供SAAS服务

提供风控服务和设备指纹服务



欧洲数据中心

更好的服务于海外合作伙伴



技术迭代升级

谢谢