

互联网金融的大数据应用

面包君 2017/6/20





互联网金融的发展历程

大数据在互联网金融的应用

征信体系介绍

风控反作弊欺诈模型运用

互金公司贷款授信

保险定价策略分析

量化投资应用



互联网金融在美国也叫做“FINTECH”。互联网金融的定义有广义和狭义之分。

从广义来说，互联网金融是所有使用技术来提供金融服务的代名词。无论是提供全球数字汇款服务的 SWIFT 系统还是股票交易所或清算所的门户入口，都是互联网金融的一部分。

从狭义角度来看，互联网金融多指近年来由新兴公司和金融服务业等提供的创新的具有破坏性的新型金融服务和产品，主要指的是运用新型技术提供服务的银行业务、公司金融业务、资本市场业务、金融数据分析、支付、个人财富管理。



上世纪90
年代前后

- 传统金融机构和业务信息化

上世纪90
年代中后期

- 传统业务和互联网融合：
SFNB/PayPal
货币基金

本世纪以来

- 新型互联网金融业务发展：
P2P/移动支付
/虚拟货币

DT时代

美国互联网金融生态



THE FINTECH ECOSYSTEM

Payments & Transfers



Lending & Financing



Retail Banking



Financial Management



Insurance



Markets & Exchanges

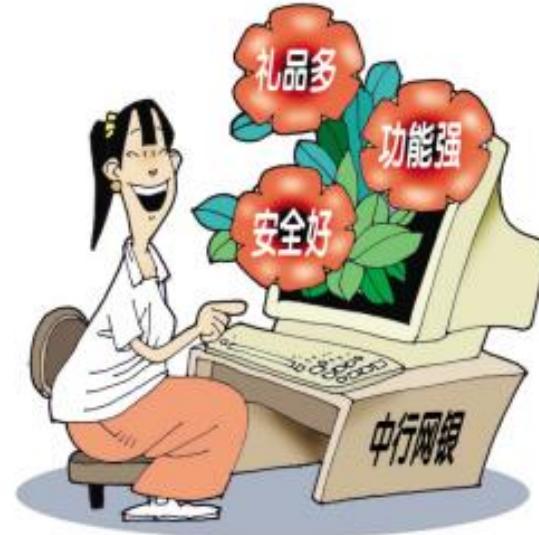


DT时代

我国互联网金融的变革



传统金融的互联网化



一淘
www.etao.com

淘宝网 阿里巴巴 youku 优酷

letv 乐视网
LeTV.com

拉手网 360buy 京东商城 f t

sina 新浪网
sina.com.cn

人人网 ganjia 赶集 聚美优品 逛



互联网的金融服务

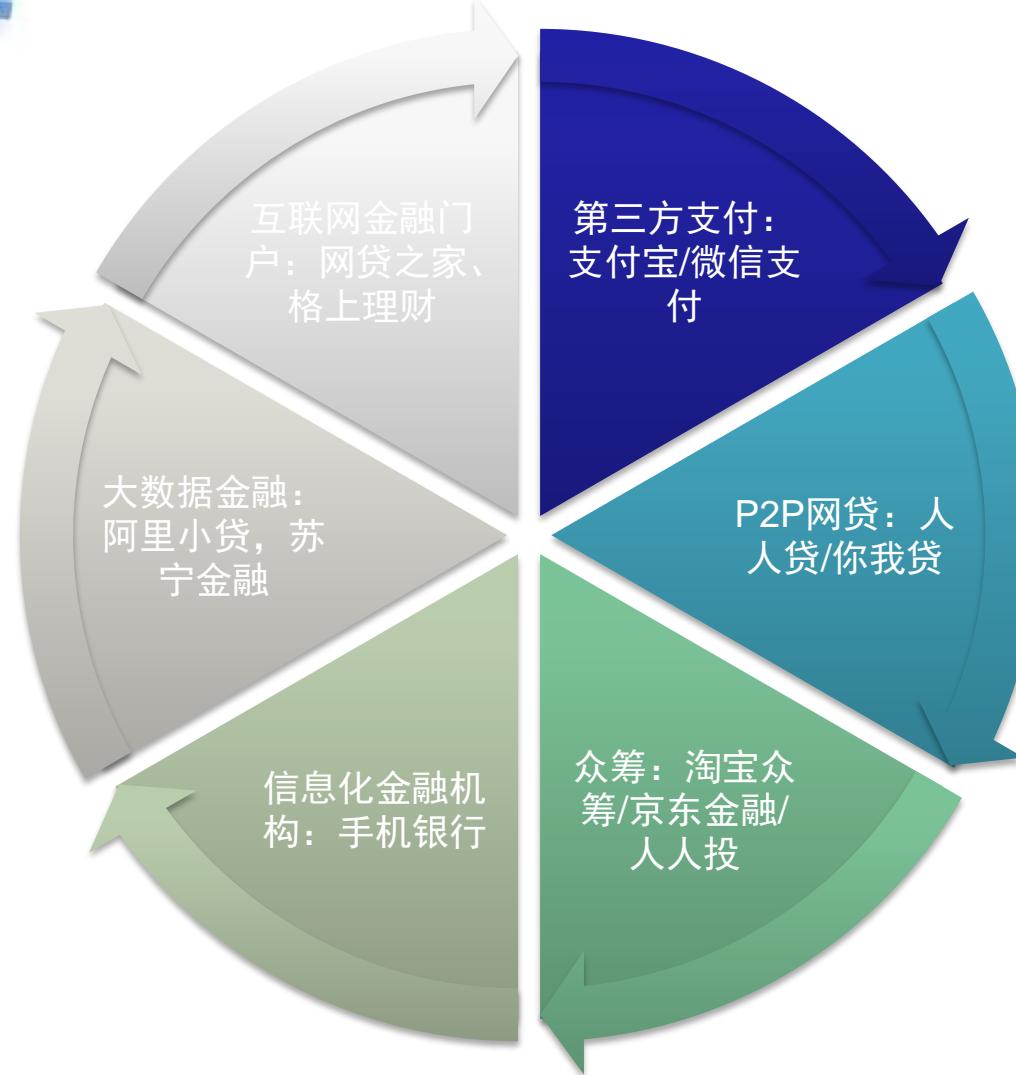
91金融超市
www.91jinrong.com

金斧子
融360 RONG360.COM

铜板街
网贷之家 您身边的网贷资讯专家
好贷 hoodai.com
银率网
存折

我国互联网金融生态

DT时代

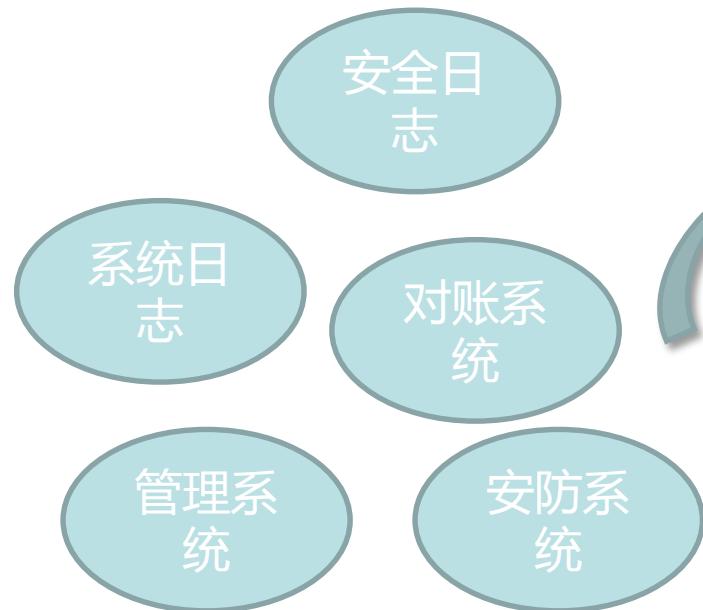


DT时代

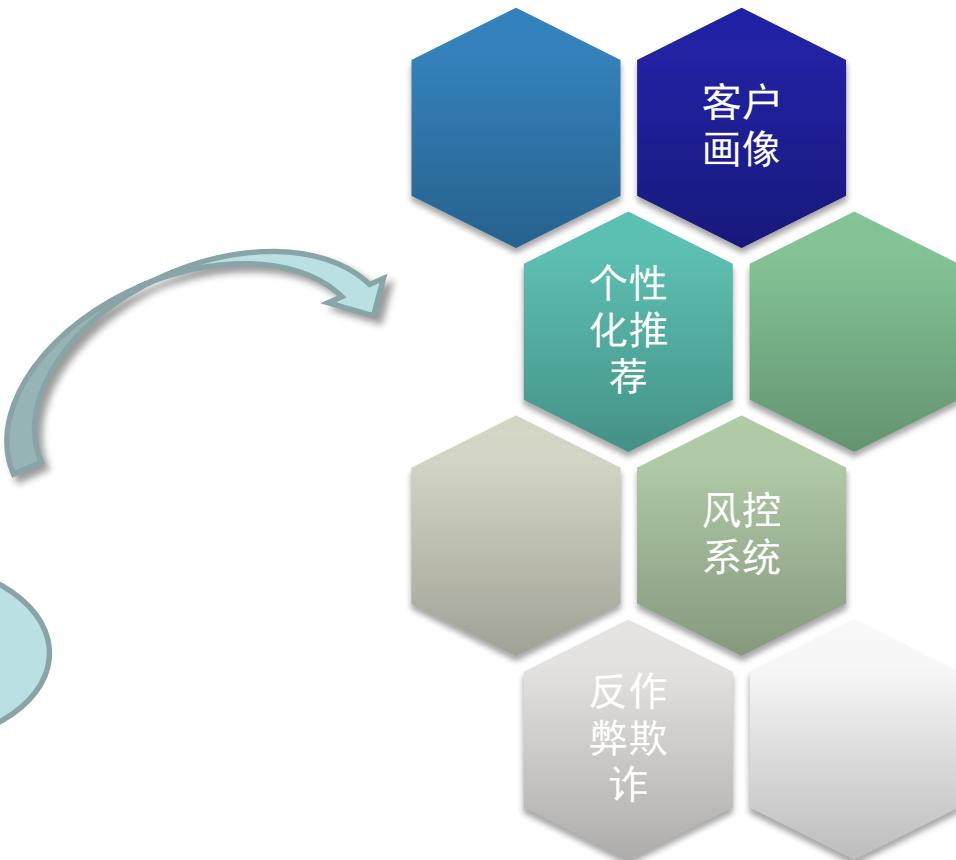
互联网金融的数据应用



IT时代



DT时代



DT时代

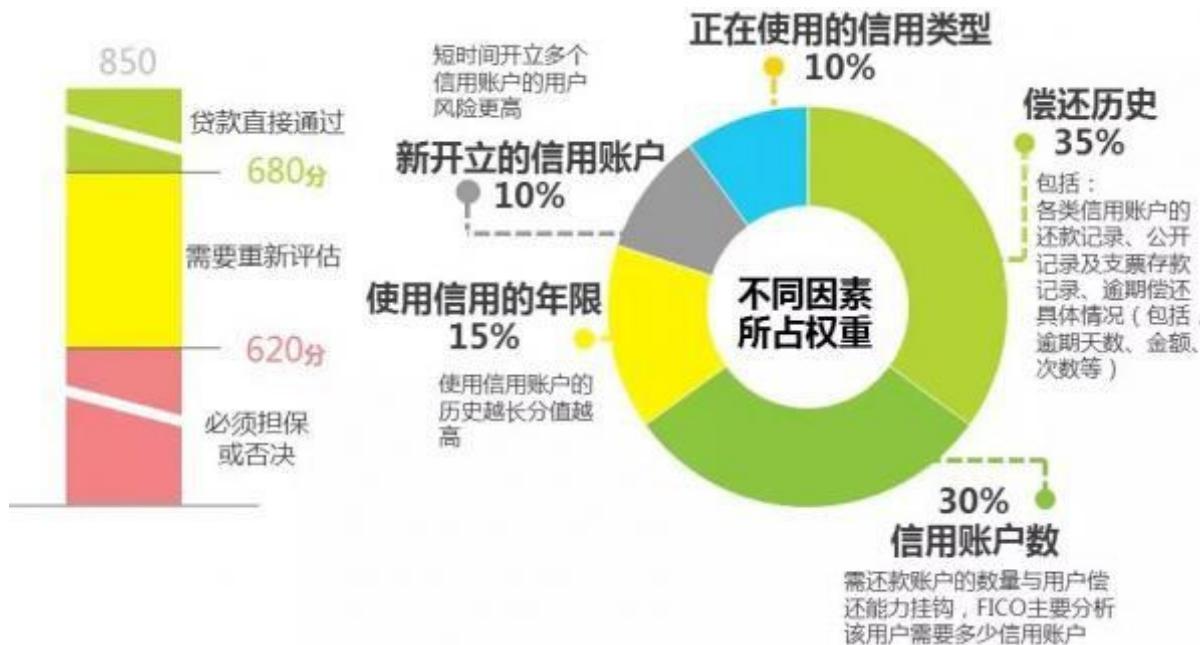
大数据在个人信用体系中的应用





FICO Score

信用分数是利用数学模型依据个人的信用报告评估银行风险大小的一个数值，一般来说数值越高风险越小。信用分数的数学模型有许多种，在银行届运用最广泛的就是FICO分数（300~850分之间），这是由FICO公司拥有的一种计算模型，其具体细节并没有披露，但是其信用分数的组成已经被总结出来了，就是下面那张图。从下图可以看出，以下几点综合，决定了您的信用分数。





计算逻辑：

1. Payment history (按时还款) :

- 各种信用账户的还款记录
- 公开记录及支票存款记录
- 逾期偿还的具体情况，包括天数、金额、次数和时长

2. Amount Owed (信用账户) :

仍需要偿还的信用账户总数、分类账户数、余额、使用率、分期付款偿还率

3. Length of Credit History (信用年限)

4. New Credit Account (新账户数)

5. Types of credit used (信用种类)



芝麻分怎么计算的？

没有任何一个单项信息能够决定芝麻分。
芝麻分是由电脑通过多维度的因子和数据，
通过复杂模型综合计算得出



1. 身份特征 (15%)

公安实名认证

身份信息

信息稳定性

.....

2. 信用历史 (35%)

信用卡还款历史

微贷还款记录

水电煤缴费

罚单

.....

3. 履约能力 (20%)

支付账户余额

余额宝余额

车产信息

房产信息

.....

4. 人脉关系 (5%)

关系圈

朋友圈信用水平

社交影响力

.....

5. 行为偏好 (25%)

账户活跃度

消费层次

缴费层次

消费偏好

.....

数据来源：

1. 阿里的生态系统
2. 政府公共部门：
3. 合作机构：

公安、工商、税务、
移动……等；
金融
机构、同业征信等

DT时代

大数据在个人信用体系中的应用



芝麻信用，代表的是基于个人消费者及小企业互联网大数据的一系列信用评估产品，解决陌生人之间及商业交易场景中最基本的**身份可信性**问题，同时帮助**识别风险与商机**。



3亿实名用户，覆盖近一半的中国网民

涵盖购物、支付、投资、生活、公益等**上百种**场景数据

每天**PB级**数据，相当于5000个国家图书馆信息量

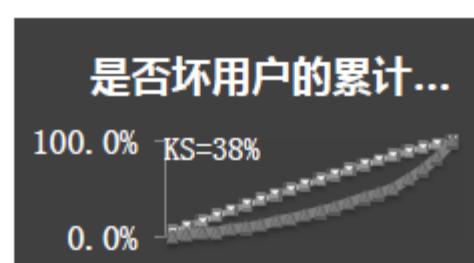


场景/结果

账户匹配度高

互联网用户识别覆盖度60%

评分区分能力好

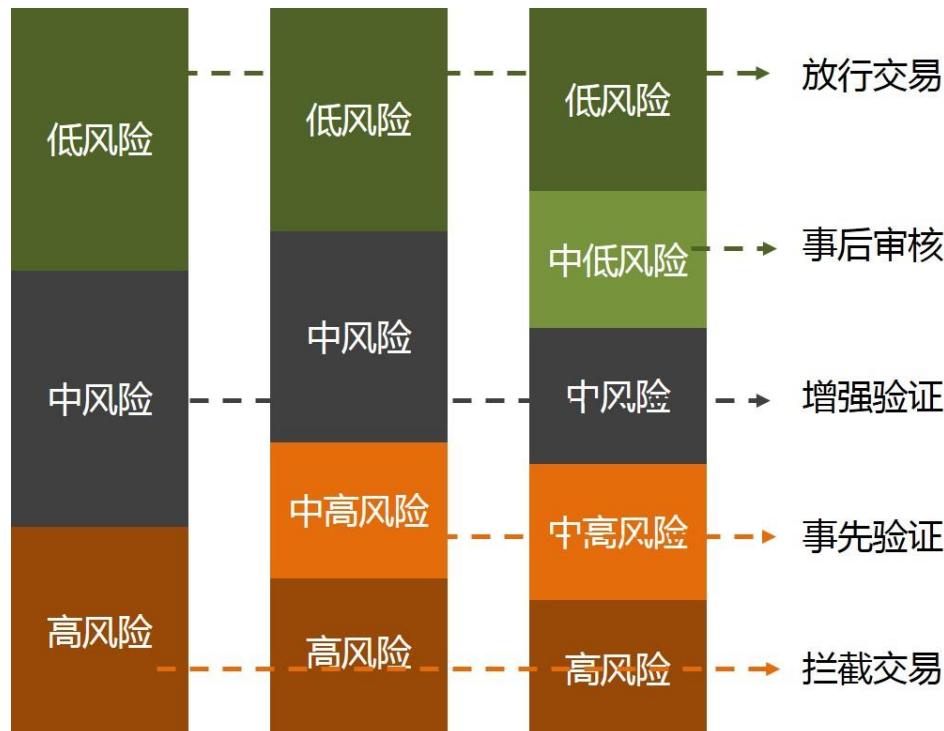


评分对于信用区分能力强
图为某小额贷款的场景下的评分效果校验



风控概念：

- 顾名思义，风控就是风险控制，最大程度地控制作弊和欺诈的发生，保障网站的正常运营和用户体验。支付风控涉及到多方面的内容，包括反洗钱、反欺诈、客户风险等级分类管理等。其中最核心的功能在于对实时交易进行风险评估，或者说是欺诈检测。如果这个交易的风险太高，则会执行拦截。由于反欺诈检测是在交易时实时进行的，在要求不能误拦截的同时，还有用户体验上的要求，即不能占用太多时间，一般要求风控操作必须控制在100ms以内，对于交易量大的业务，10ms甚至更低的性能要求都是必须的。这就需要对风控模型进行合理的设计。一般来说，要提升风控的拦截效率，就需要考虑更多的维度，但这也会带来计算性能的下降。在效率和性能之间需要进行平衡。
- 风险和作弊行为的发现、识别和处置
- 风控和反作弊是持续的博弈过程，cat-and-mouse game，时效性强，对抗性强



目前主流的风险等级划分有三种方式，三等级、四等级、五等级。

- 三等级的风险分为低风险、中风险和高风险。大部分交易是低风险的，不需要拦截直接放行。中风险的交易是需要进行增强验证，确认是本人操作后放行。高风险的交易则直接拦截。
- 四风险等级，会增加一个中高风险等级。此类交易在用户完成增强验证后，还需要管理人员人工核实，核实没问题后，交易才能放行。
- 五风险等级，会增加一个中低风险等级。此类交易是先放行，但是管理人员需要进行事后核实。如果核实有问题，通过人工方式执行退款，或者提升该用户的风脸等级。



三板斧：

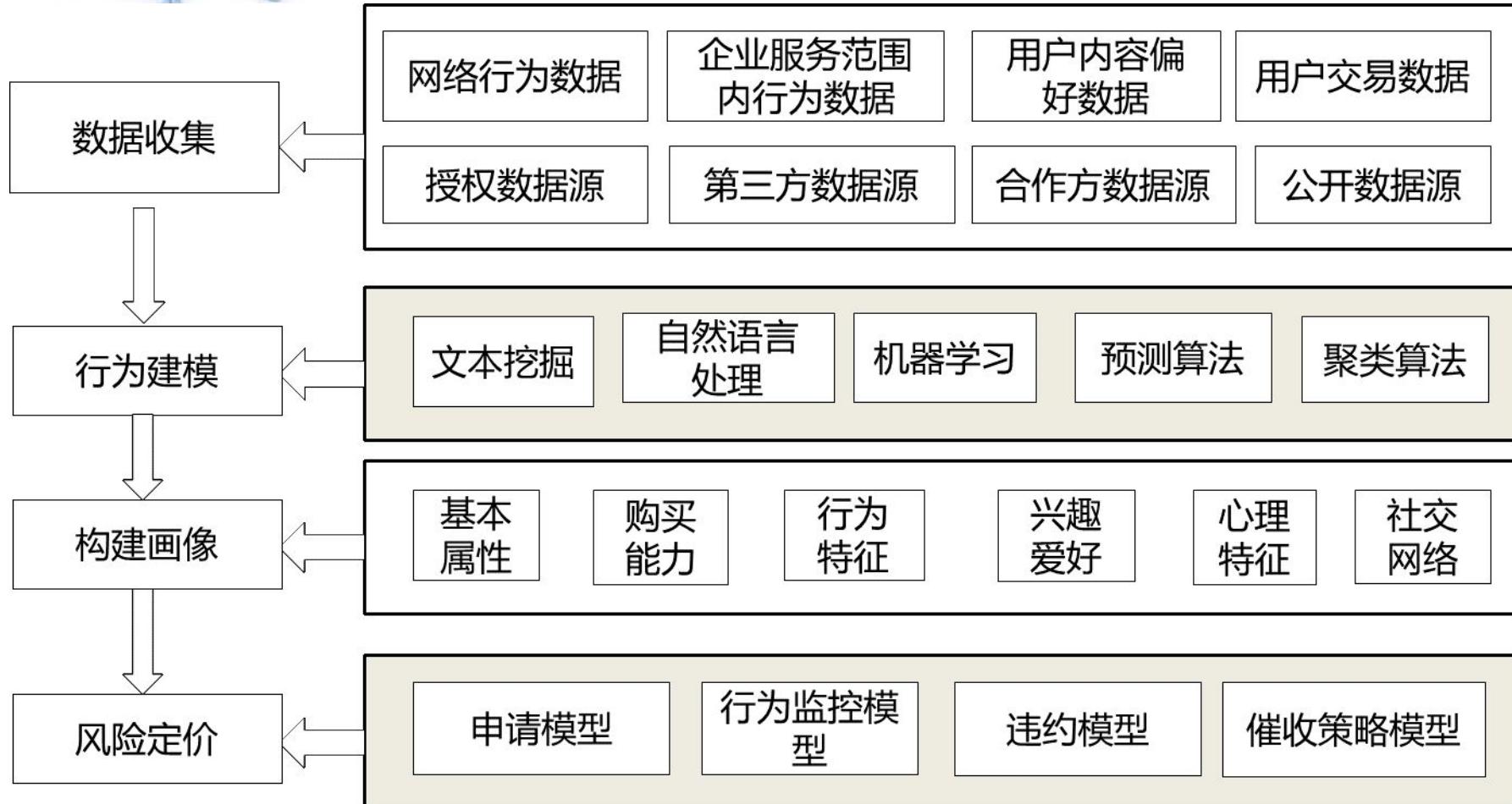
- rules；
(名单规则、操作规则、业务规则、行为规则、历史规则)
- models；
(逻辑回归、决策树、评分模型)
- strategy



- | | | | | |
|----------------------------|--|------------------|--------------|--------|
| • 异常监控
• 用户反馈
• 情报收集 | • 人工规则
• 机器学习模型
• 业务策略
• 网站规则 | • 离线评估
• 在线评估 | • 规则
• 模型 | • 用户反馈 |
|----------------------------|--|------------------|--------------|--------|

DT时代

大数据在风控反作弊的应用





数据：

一般业务数据：用户、商品、交易、点击、浏览、搜索、评价、服务、处罚等

安全业务数据：设备数据（UA、cookie、MAC、Umid、IMEI、IMSI）、位置数据（IP/LBS/GPS）、行为信息、生物信息、其他

算法：

机器学习：分类、聚类、graph算法异常检测

图像算法：人脸识别、OCR、图像搜索

绝大多数场景使用RF/GBDT+LR/C5.0

注意点：

- 部分高风险业务，可以投入人力审核，追求更高的准确率/召回率
- 风险(异常)占比少，属于非平衡数据集
- 对抗意识强，模型衰减快，需要结合处置手段
- 风控的成本与回报意识，平衡人力和风险
- 能够采用更复杂的算法，但需要平衡用户体验和可解释性

DT时代

大数据在风控反作弊的应用



决策树模型





评分模型



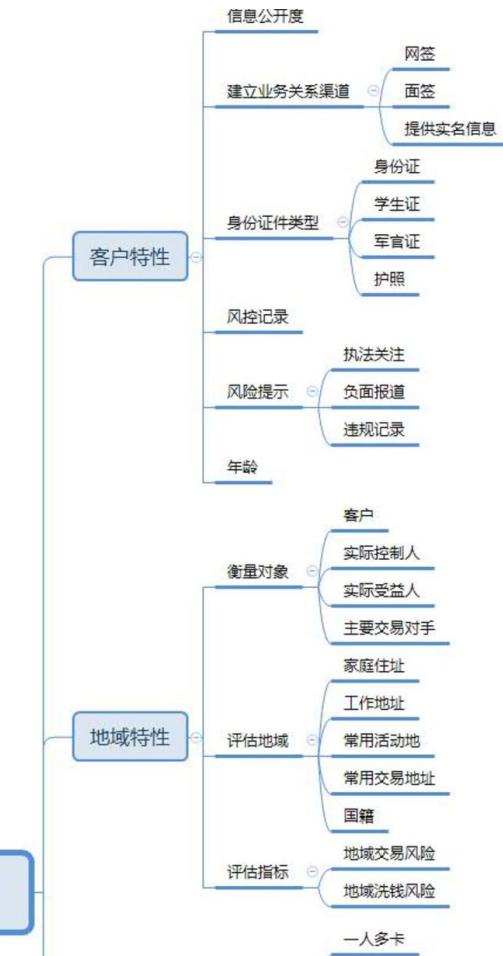
评分模型的优势在于：

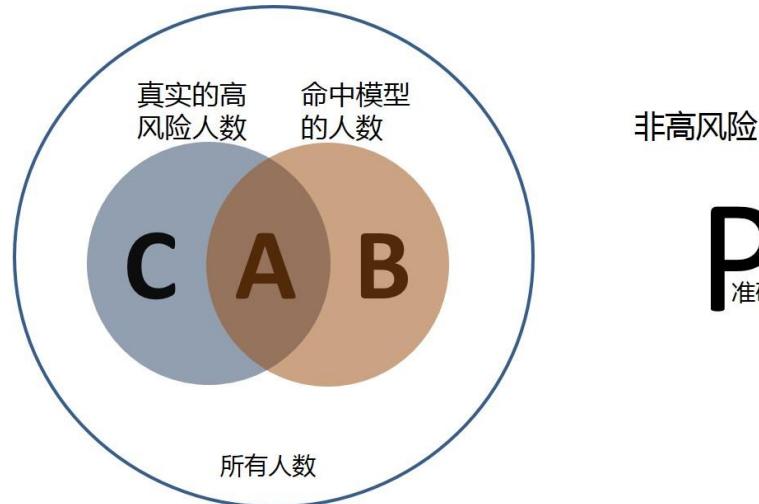
1. 性能比较高，针对交易进行指标计算，按照区间来确定风险。
2. 相对于规则，如果指标设置合理，其覆盖度高，不容易被嗅探到漏洞。
3. 理解和分析也比较容易。如果交易被拦截了，可以根据其各项打分评估其被拦截的原因。

存在的问题：

1. 模型真的很难建立。指标的选择是一个挑战。
2. 各个参数的调优是一个长期的过程。

风控模型





$$P_{\text{准确率}} = \frac{A}{A + B}$$

$$R_{\text{召回率}} = \frac{A}{A + C}$$

$$F_1 = \frac{2PR}{P + R}$$

以评估高风险人群的效果为例，

- Precision, 准确率，也叫查准率，指模型发现的真实的高风险人数占模型发现的所有高风险人数的比例。
- Recall, 召回率，也叫查全率，指模型发现的真实的高风险人数占全部真实的风险人数的比例。

理想情况下，我们希望这两个指标都要高。实际上，往往是互斥的，准确率高、召回率就低，召回率低、准确率高。如果两者都低，那就是模型不靠谱了。对于风控来说，需要在保证准确率的情况下，尽量提高召回率。那怎么发现实际的高风险人数呢？这就需要借助规则模型，先过滤一遍，再从中人工遴选。

DT时代

大数据在风控反作弊的应用

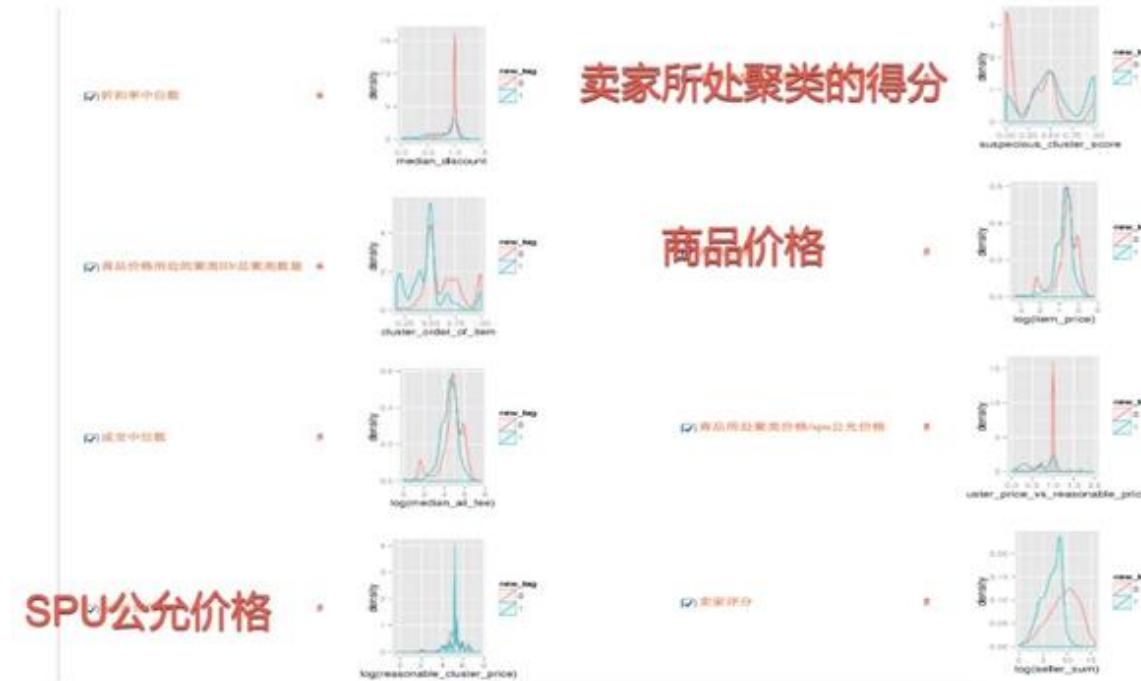


DT时代

大数据在风控反作弊的应用



Case 1：如何判断某笔交易是否虚假？

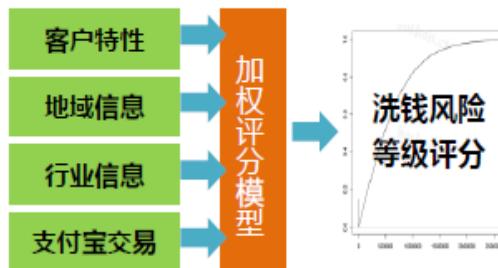




Case 2 :

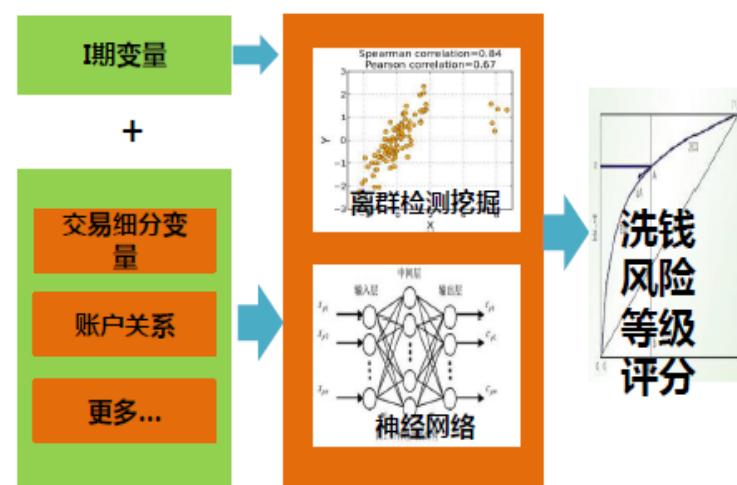
- 建立客户洗钱风险等级模型，识别高风险客户和交易行为。

模型I期：



级别	分数段	占比
低风险	[1,30)	46.4%
中低风险	[30,35)	43.7%
中风险	[35,45)	9.7%
中高风险	[40,45)	0.16%
高风险	[45,100]	0.005%
总计	[0,100]	100%

模型II期：

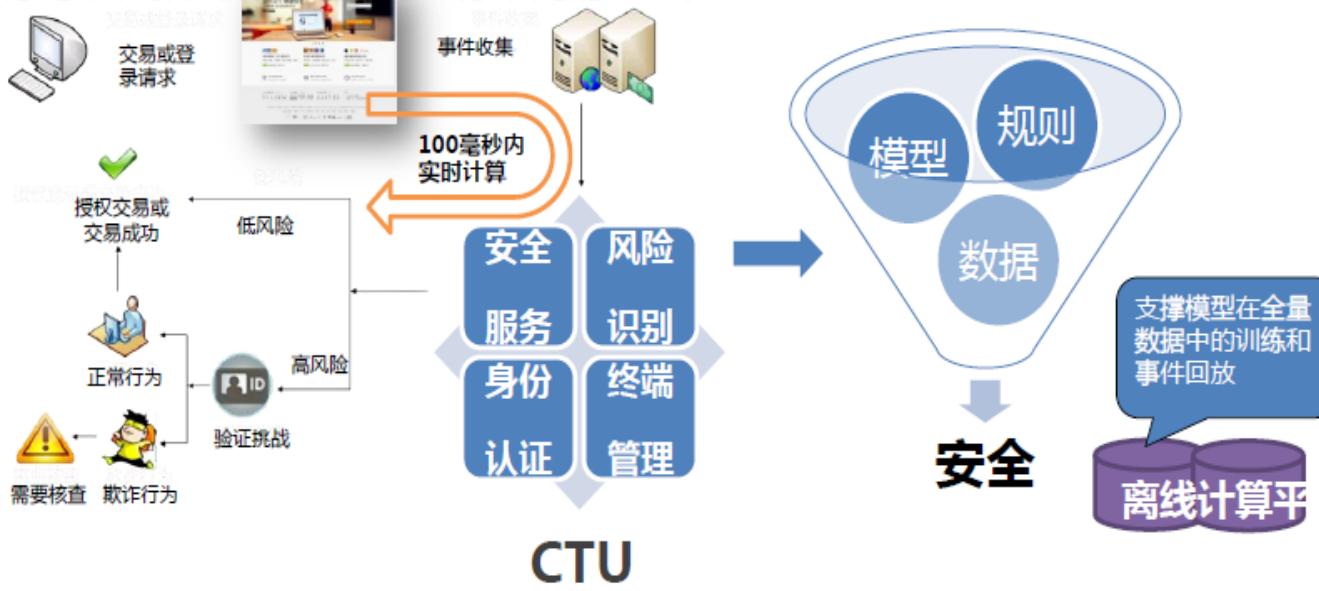


- 更多维度变量
- 更智能的方法
- 更精准的结果



Case 3 :

智能实时风险监控系统是目前国内最先进的网上支付风险实时监控系统之一，通过数据分析、数据挖掘等技术进行风险嫌疑数据的抓捕，发现异常的或有风险的操作行为，根据风险级别不同进行不同的处理。

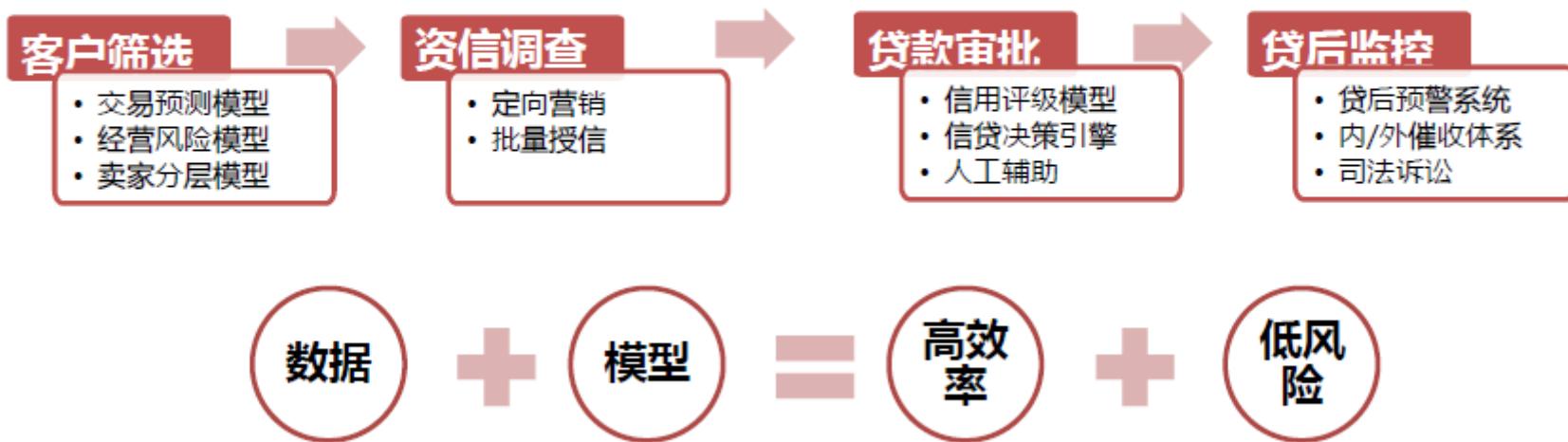


DT时代

大数据在小微贷中的应用



基于大数据的信贷模式是阿里小贷的最大优势，在这个模式下，客户立即申请立即获贷，不良率低于传统银行。



大数据在保险产品中的应用



- 2010年6月上线，解决退货纠纷
 - 5毛钱的保费，10块钱的保额，堪称“小而美”金融产品的典范
- 4年来，保单量以每年100%的速度增长
- 2013年全年保单量突破**10亿**单，双11当天保单出单量**1.2亿笔**
 - 覆盖淘宝全网所有实物商品类目，占实物订单量16%



运费险定价演进

一口价时代

- 保费按5%费率统一收
- 保额10元，保费0.5元

简单粗糙，劣币驱逐良币

精算定价时代

- 以历史出险率为唯一定价因子

解释方便，定价偏差大

大数据定价时代

- 百万ID特征，实时特征
- 广告机器学习方法

预测很准，难于解释

数据定价时代

- 以30+因子统计建模，预测退货率

预测较准，解释性差





程序化交易策略的基本流程





从事量化交易后要考虑的几个问题：

- 如何找出合适的策略？
- 回测前辨别策略的优劣？
- 策略的回测？
- 怎么执行策略？
- 交易过程中扩大规模增加收入？
- 管理无法避免的亏损？



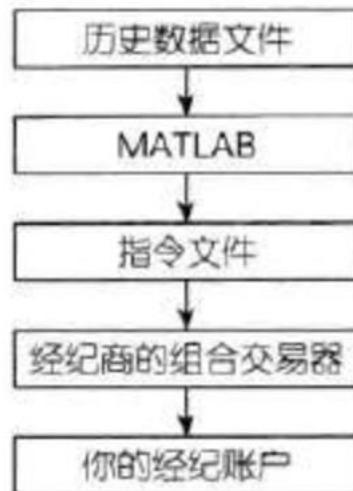
- 开发工具：
MATLAB/Python/C++
- Python环境：
pandas数据统计，numpy数值处理，matplotlib绘图，sklearn
机器学习库
- 编辑器：
Eclipse+Pydev
- 数据来源：
东方财富/wind/新浪财经/集思录/优矿等
- 推荐资料
<https://zhuanlan.zhihu.com/p/24220361>

DT时代

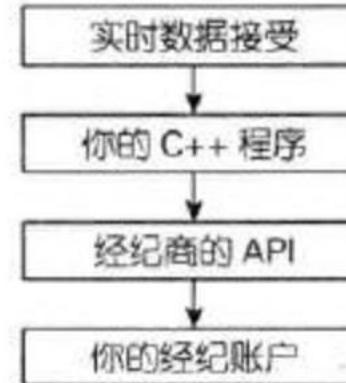
大数据在量化投资中的应用



半自动交易系统



全自动交易系统



DT时代

大数据在量化投资中的应用



TradeBlazer - [工作区 11 行]

文件(F) 视图(V) 窗口(W) 帮助(H)

面板

超级图表

行情报价

TB公式

帮助与链接

嵌入式文档

帐户管理

状态栏

UQER 优矿

首页 开始研究 我的交易 研究数据 量化社区 实盘大赛 量化学堂 帮助 面包君 保存 全部运行 重启 另存为lib

搜索文件夹/文件名称 新建

已开启的Notebook

市场强弱择时-10日均线

所有Notebook

基于市场强弱的因子择时策略0
基于彼得·林奇选股法的改进
市场强弱择时-10日均线
凤鸣朝阳 - 股价日内模式分析
Simple MACD
姚大大带你飞翔在股市
海龟+羊驼 神兽拼接可行吗?
自适应海龟系统
超简海龟策略
基于市场强弱的因子择时策略
羊驼策略
破解股市泡沫之谜——对数周期幂率 (LPPL) 模型
今天大盘熔断下跌，后市如何——based on LPPL anti-bubble model
尝试下新闻协作

43
44
45
46
47
48

```
if yesterday_tn < test_neck:  
    sell_list = account.security_position  
    for stk in sell_list:  
        order_to(stk, 0)  
  
else:  
    return
```

年化收益率 1.0% 基准年化收益率 12.6% 阿尔法 0.8% 贝塔 0.02 夏普比率 1.15 收益波动率 0.9%

信息比率 -0.99 最大回撤 0.4% 换手率 1.06

回测详情 开始交易

累计收益率

普通 对数轴 相对收:

10.00%
5.00%
0.00%

9480.0 仓位 267.5/2.90%
现手 41 速涨 0.00%
总手 171948 开盘 9067.5
持仓 538771 最高 9750.0
日增 0 最低 9052.5
外盘 比例 0.0 结算价 0.0
昨收 9212.5 中价

合约 手数 价格... CNUS 1 对手价
9462.5 买入 9480.0 卖出
多单1手
外盘自动开平
<= 471
简单 对价限 排队限
A50E09 1手多仓 *50
止损开仓 滑线下单 设条件单

持仓 委托 成交 条件单 止损单 资金 合约
品种 / 合约号 多空 手数 可用 开仓... 逐笔浮盈 赢利 t
A50 CNUS 多 1 1 9425 50USD 9

平33% 平50% 平100% 反手 止损
时间 / 合约 状态 买卖 开平 委托价 委托量 可撤 已成交 报价
下单 外盘 银票 黄金 11:03:01 Local

webStock 需要自设盈利功能的用户, 请下载 v3.4

感谢您的关注

知乎专栏：数据分析侠

