



“互联网+”时代

# 大数据技术峰会

中国·深圳 | 2015.11.28-29

解码数据未来

# 社交数据在征信领域的应用探索

刘黎春     SNG, 腾讯

- 社交征信背景
- 腾讯社交网络数据
- 个体用户画像研究
- 社团圈子研究
- 模型建设及应用



# 社交征信背景



# 传统征信相关机构

## 数据的采集和提供

- 独特的数据源
- 初步的数据挖掘

数据公司

## 征信解决方案的应用

- 银行
- P2P贷款机构

征信使用方

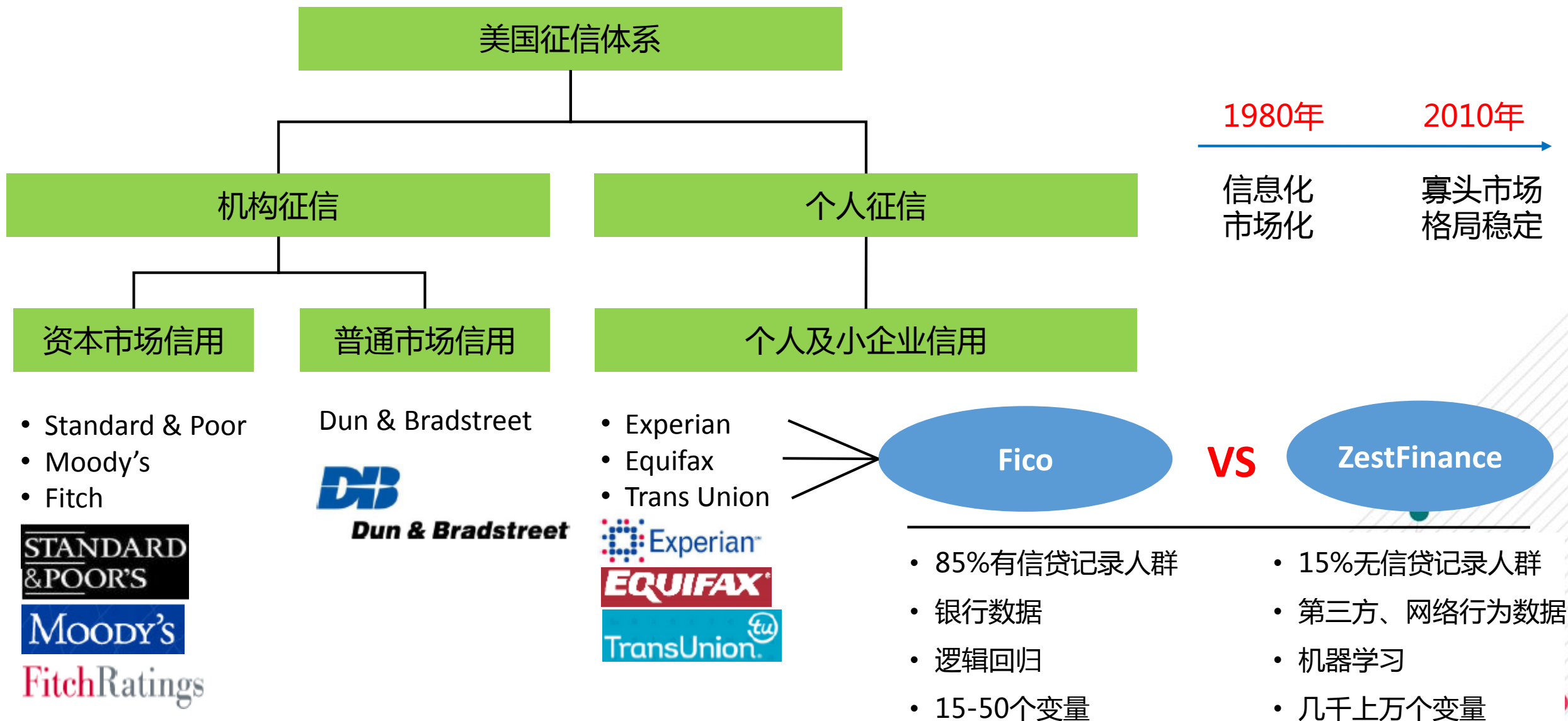
征信公司

## 数据整合和挖掘

- 多维度的数据源整合
- 征信模型建设、服务提供
- 征信解决方案提供



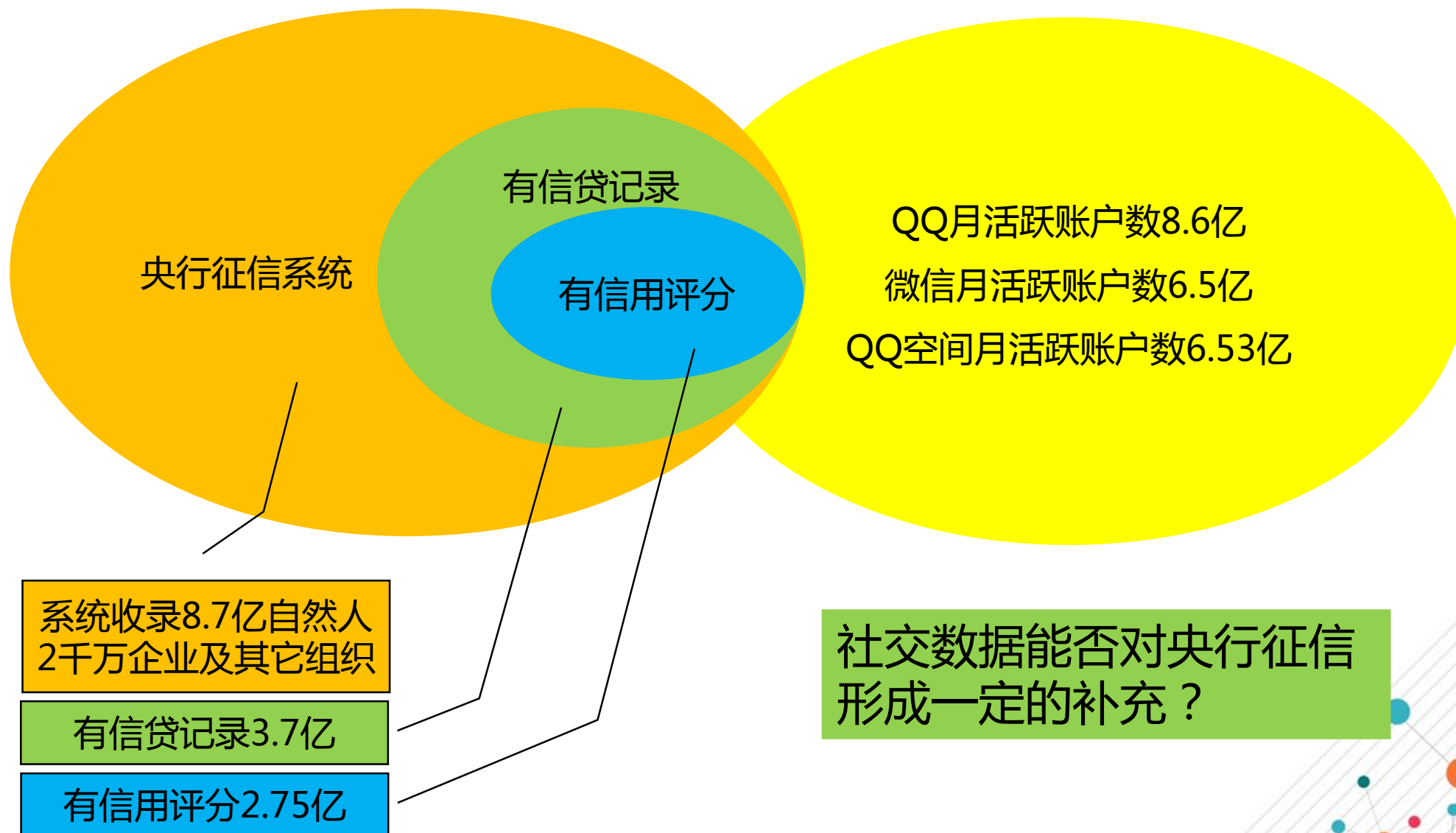
# 美国著名征信公司





# 国内征信发展历程







# 社交数据能不能用来做征信？



社交数据与信用评级有关系吗？

交易数据天然具备金融属性，社交数据有吗？

社交数据非结构化程度高，怎么挖掘并有效使用？



# 腾讯社交网络数据



# 传统征信分析维度



# QQ社交数据现状



- **QQ**
  - 月活跃**8.6亿**
  - 最高同时在线**2.39亿**
- **QQ空间**
  - 月活跃**6.53亿**



# 社交征信SWOT分析

- 覆盖人群广，数据维度多
- 大规模实时数据处理能力
- 算法模型积累较多
- 社交网络数据难以模拟，在反欺诈方面有天然优势

- 缺乏信贷相关核心数据
- 数据类型多，处理难度大
- 缺乏统一的数据标准和规范
- 数据孤岛现象严重
- 账号复杂，一人多号情况多



- 互联网金融爆发式增长，征信需求越来越大
- 大量的用户没有被央行征信系统覆盖，空间大

- 行业对社交征信的认可
- 互联网业务变化快，数据稳定性挑战较大



# 个体用户画像研究





# 用户画像主要挑战

## 1. 如何利用腾讯各种丰富的数据资源及之间的联系



社交网络



LBS日志



用户群组



多媒体数据



UGC文本



登录IP

## 2. 如何使用户画像适应各种不同的应用场景

广告  
定向



推荐  
系统



市场  
营销



信用  
评分



## 3. 如何高效的处理海量的用户数据（超过10亿的QQ用户，超过千亿级别的各类日志数据）

# 用户画像解决方案

1. 针对不同的底层数据类型设计特定的挖掘算法，挖掘用户的行为特征，形成底层标签。综合考虑不同数据来源的，形成更上层的抽象用户标签
2. 建立完善的用户画像标签体系结构，从不同维度、粒度对用户进行描述。
3. 搭建用户画像挖掘系统，基于大规模存储和机器学习计算平台，定期对全量用户数据进行计算和挖掘，并提供用户标签的使用和查询服务。



# 用户画像系统架构

## 标签应用层

TDW 离线查询

HBase 实时查询 ( 理论峰值40w/s )

## 标签汇总层

不同算法、数据来源得到标签进行汇总

## 模型训练与预测层

无监督模型：  
word2vec,  
LDA, 社区发现

半监督模型：  
标签传播

监督模型：LR, Kernel  
SVM, Random Forest

基于Hadoop, Spark和GraphLab等计算平台

## 数据处理层

结构化数据统计

文本分词

LBS与POI匹配

## 原始数据层

相册说说

APP文本

群文本

操作行为

关系链

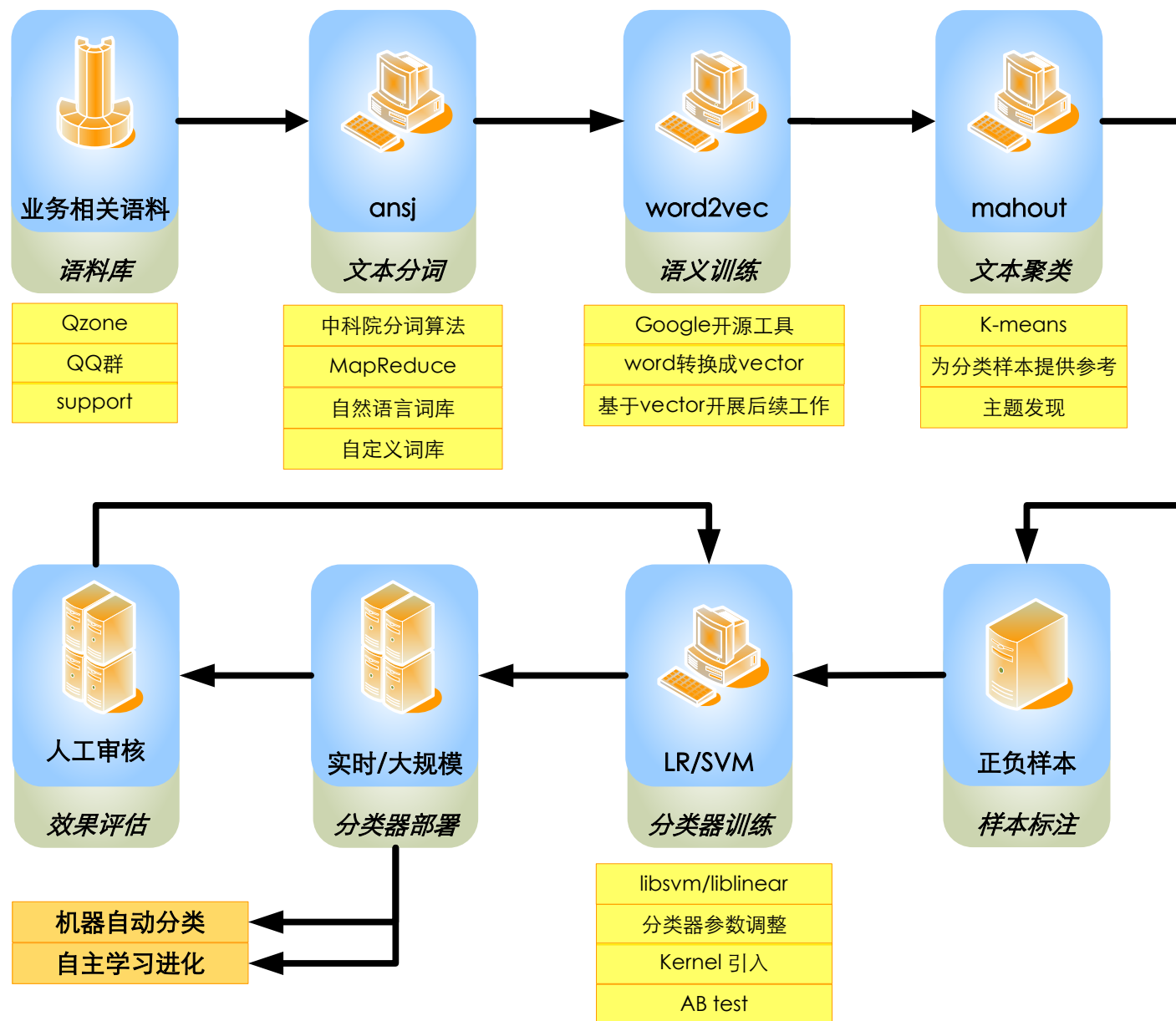
LBS数据

外部数据

TDW数据仓库



# 用户画像文本挖掘系统



每天处理**5600万**条说说, **3000万**相册标题, 讲述**4500万人**的故事

# 用户画像行业挖掘

## 思路及问题

思路1：根据用户加入的QQ群文本及其他UGC进行文本分类

存在问题：加入群只能反专业业相关兴趣，与职业并无绝对关系

思路2：判断用户工作地点，并根据工作地点推测用户行业

存在问题：同一工作地点可能存在多种不同工作行业

思路3：利用同事间好友关系网络进行行业标签传播

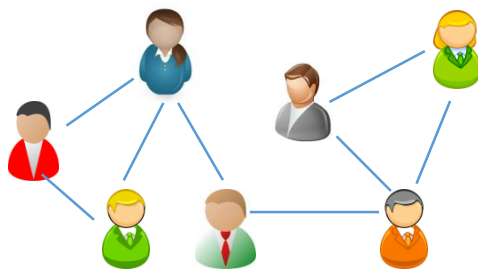
存在问题：好友关系类型比较复杂，无法确定是否为同事

## 解决方案



工作地点

LBS数据挖掘

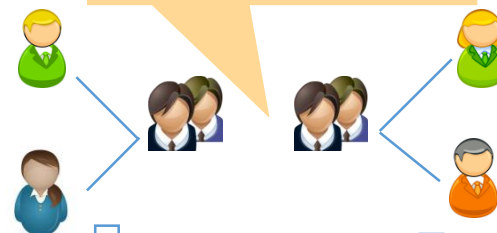


该地点工作的用户及社交网络

根据工作社团的特殊性，将部分用户的行业标签扩散给全体社团成员

Community Detection  
(FastGreedy算法)

名称、简介、公告等



群文本分类

IT行业

金融行业



工作社团1

工作社团2



# 用户画像挖掘结果



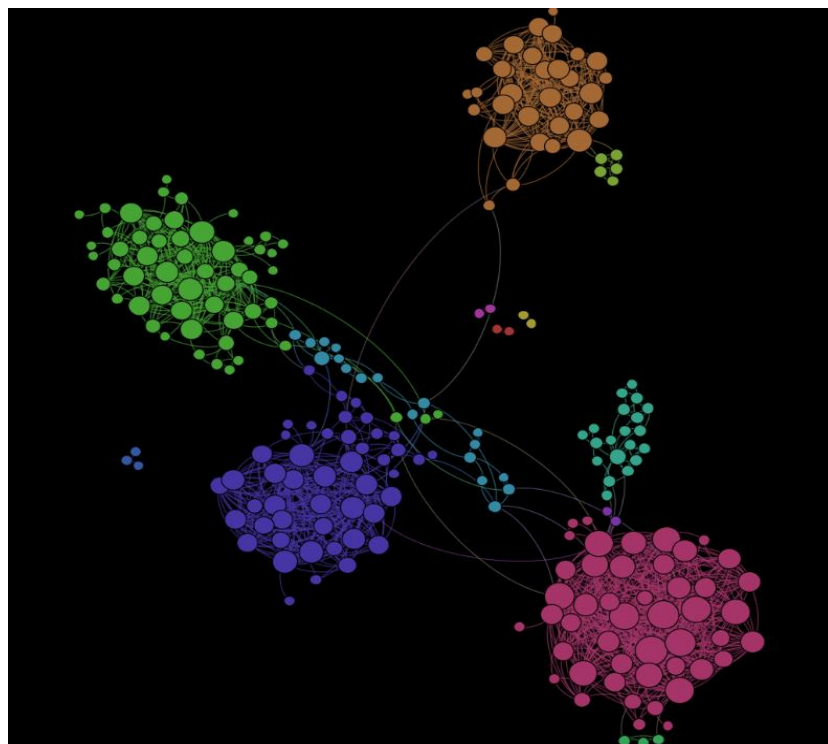


# 社团圈子研究

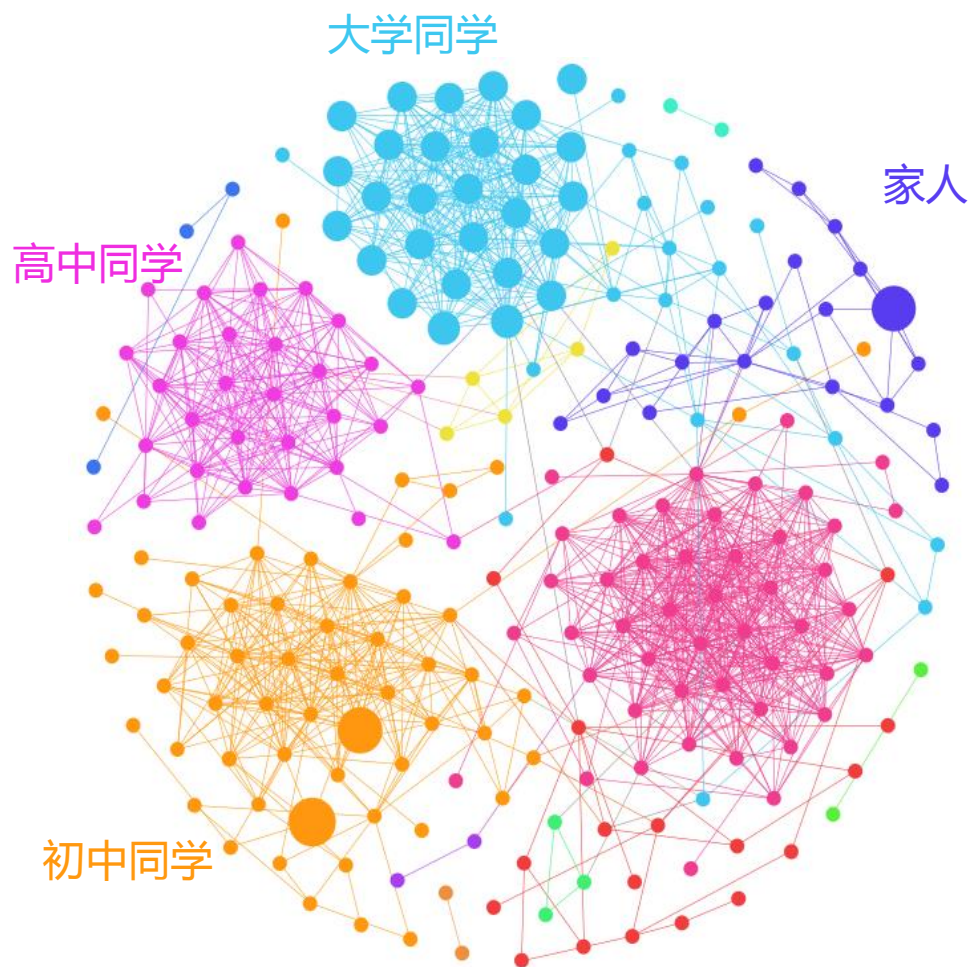


# 基于个人中心网络的关系链挖掘，QQ圈子

- 首次使用网络拓扑进行智能分圈
- 好友分圈+圈子识别+圈友识别 -> 圈子



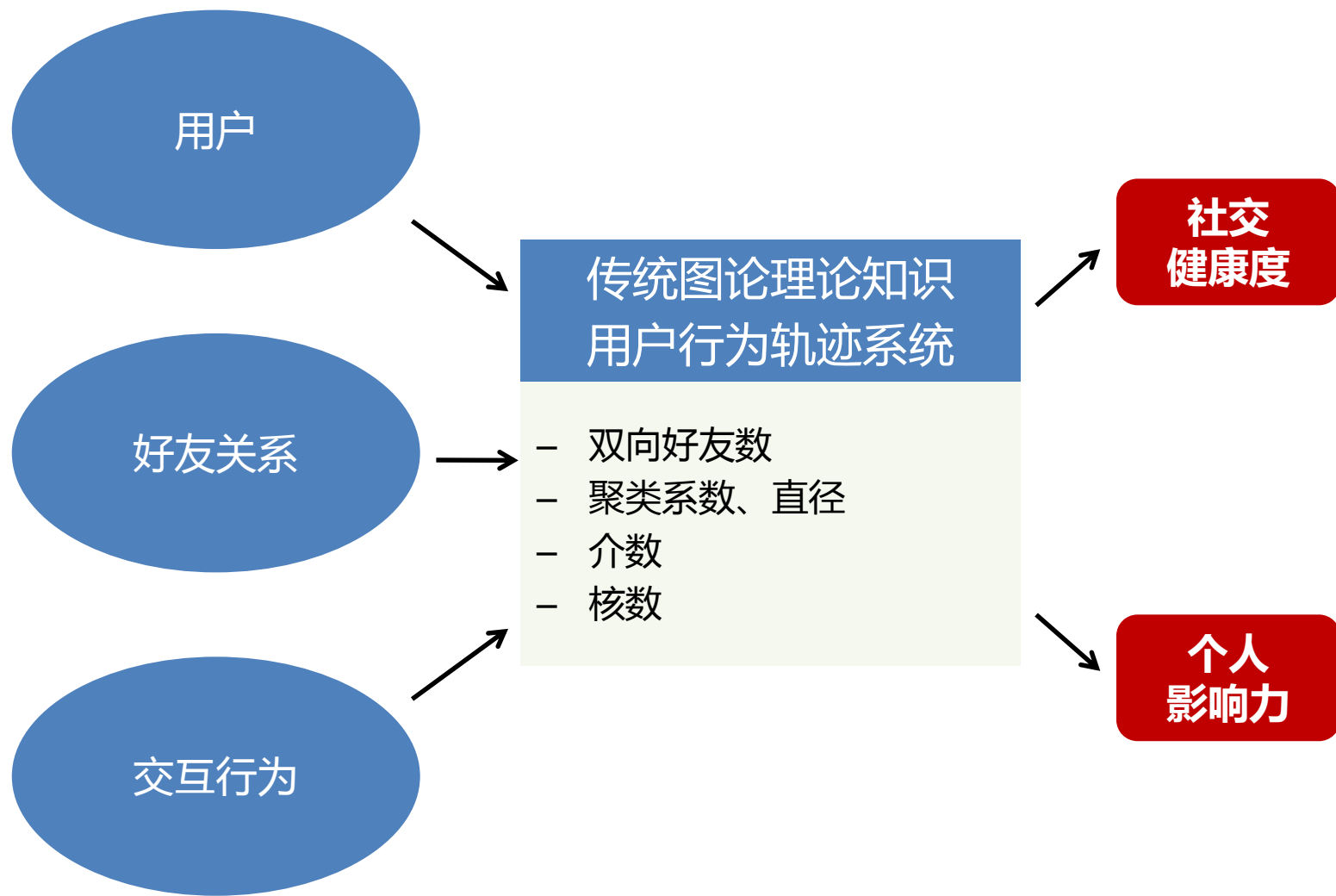
# QQ圈子的应用



- 基于用户QQ圈子的好友备注分组，挖掘学历信息
- 经验证，QQ学历数据覆盖率为74%，其中普通全日制和脱产的用户准确度均超过90%



# 社交网络拓扑的应用





# 模型建设及应用



# 自然人计算

## ➤ 自然人计算模型的基本假设：

### 同时在线假设

- 同一设备登录的号码对可能属于同一自然人

### 互斥假设

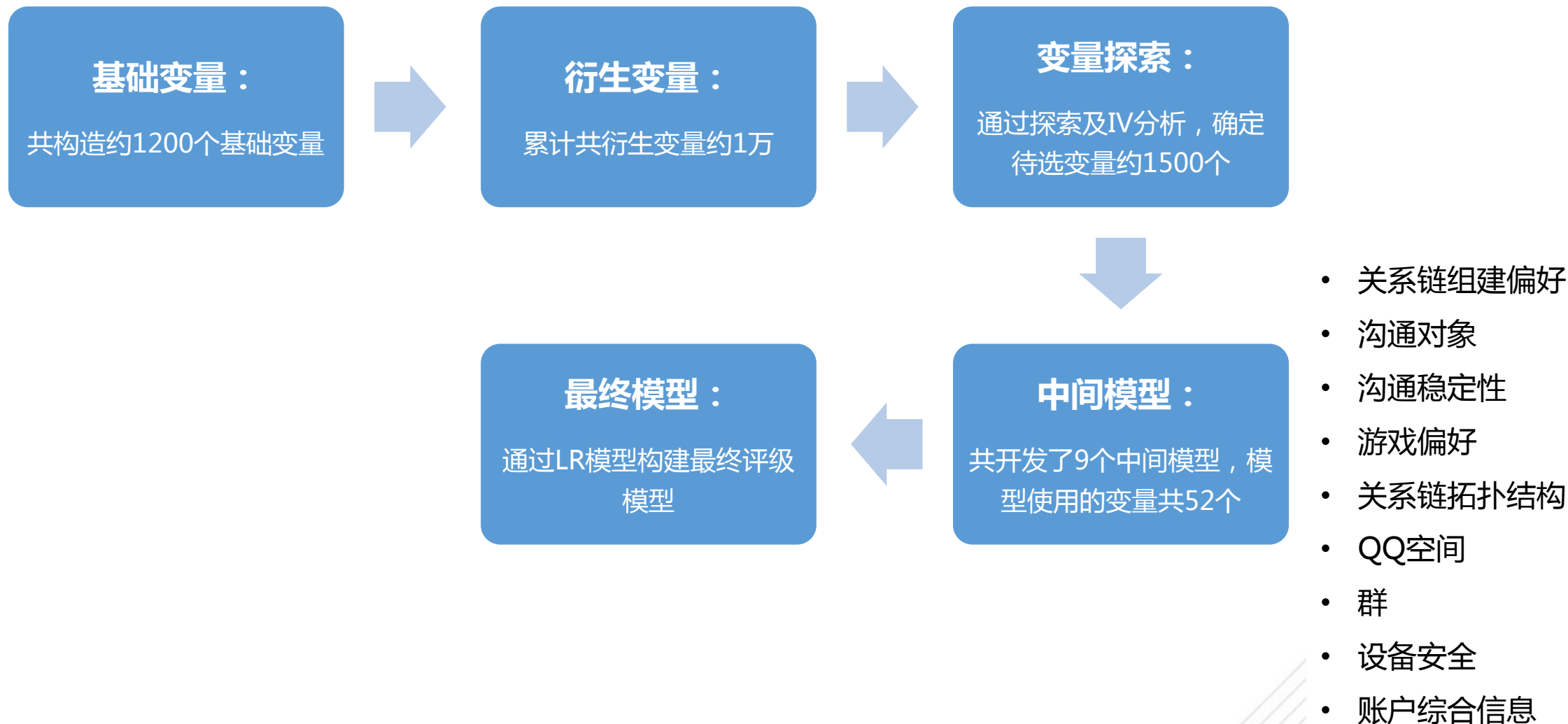
- 包括异地互斥和消息互斥，排除异地同时在线、互发消息的号码对

## ➤ 自然人计算效果（准确率：84.9%；覆盖率：75.0%）

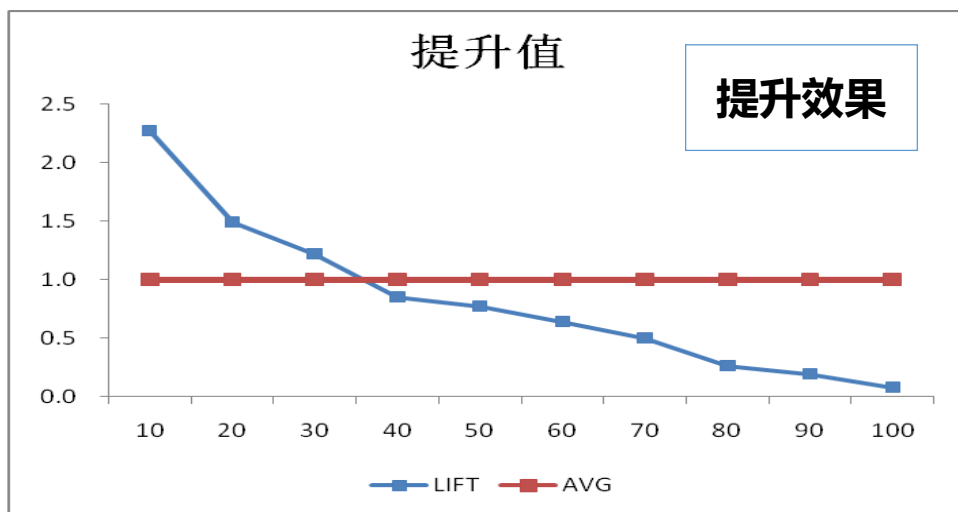
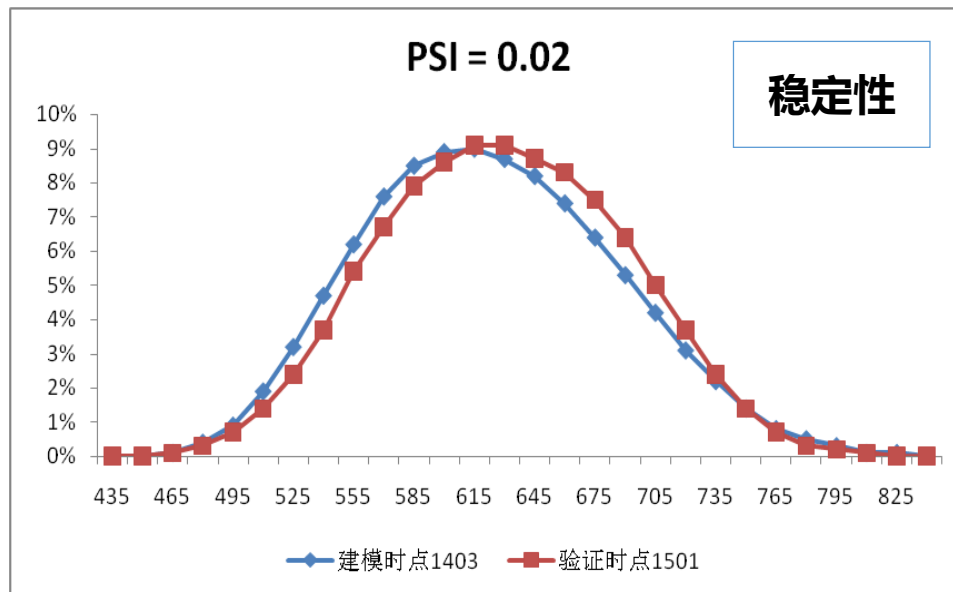
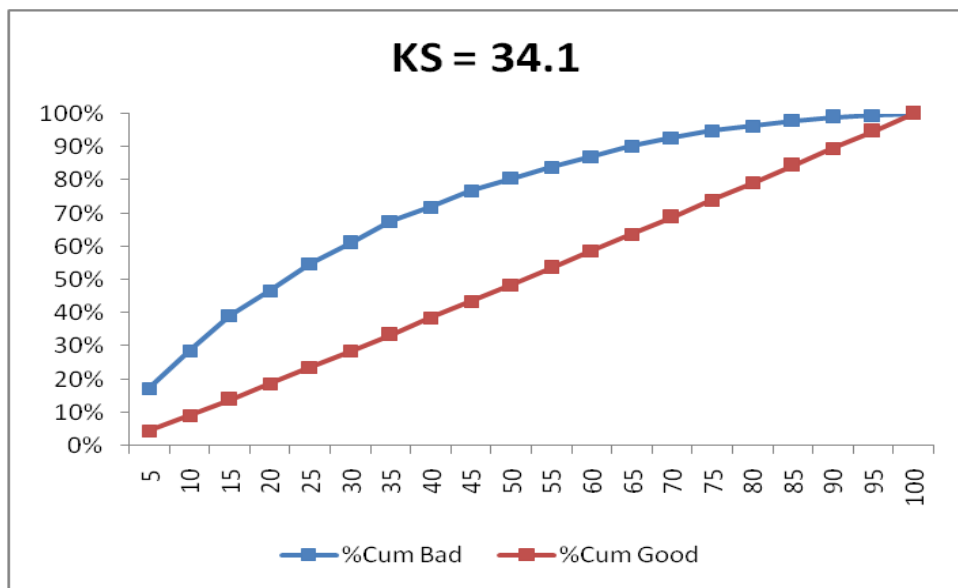




# 变量衍生与模型结果



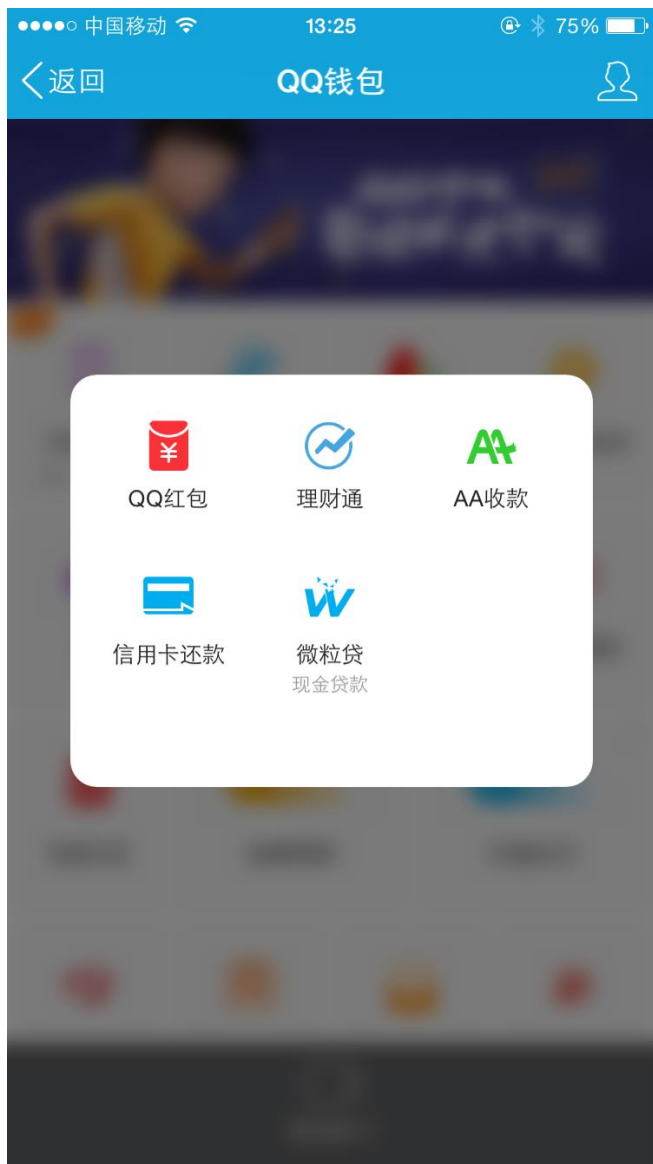
# 模型整体效果



- 模型具有较好的区分能力
- 模型在不同时点分数分布较稳定



# 微粒贷应用



Thank you