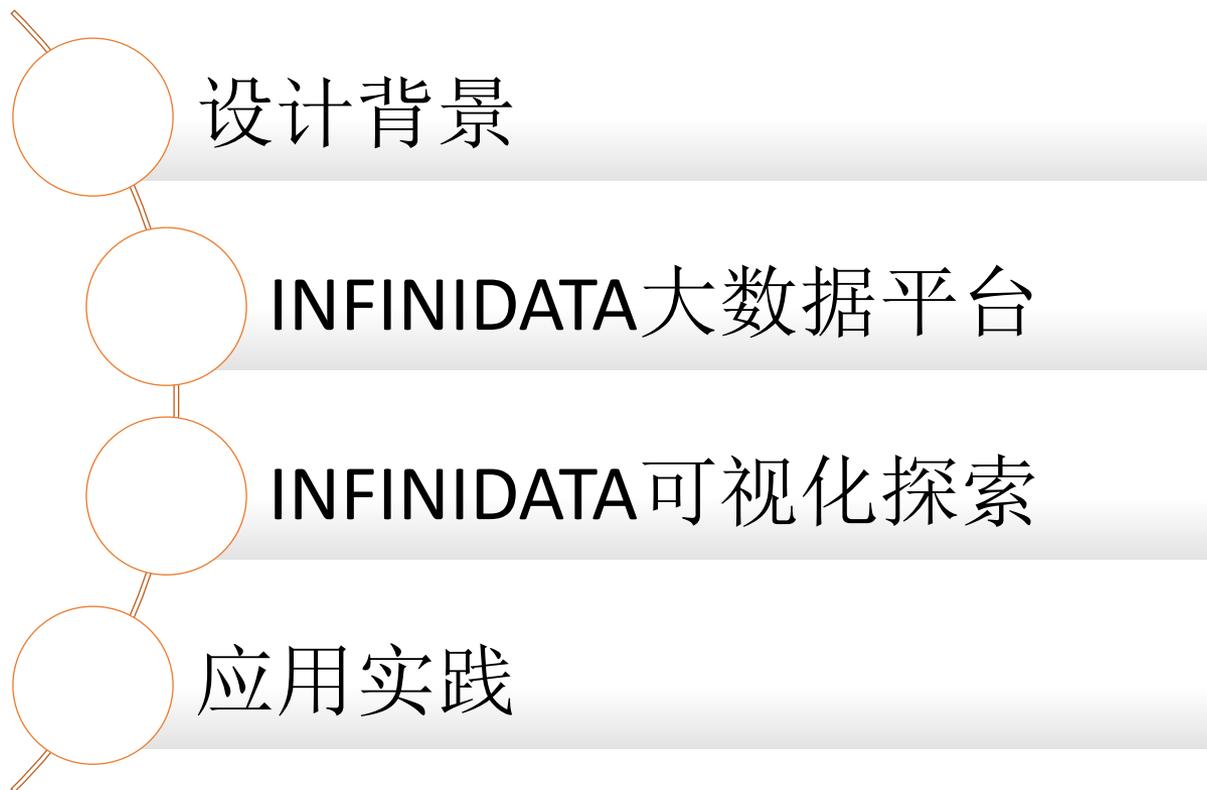


# INFINIDATA:基于Spark的统一数据管理与探索平台

熊永平

北京邮电大学网络技术研究院



## 数据应用5阶段演进模型

第五阶段

- 查询复杂度增加
- 负载混合度增加
- 数据量规模增加
- 数据模型复杂度增加
- 数据历史深度增加
- 用户数量增加
- 系统期望值增加

第四阶段

第三阶段

第二阶段

第一阶段

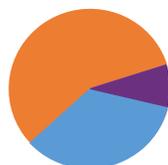
工作负载复杂度

**报表**  
发生了什么情况?



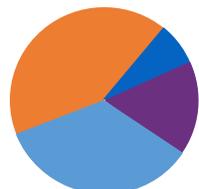
主要是批处理和预定义的查询

**分析**  
为何发生了这种情况?



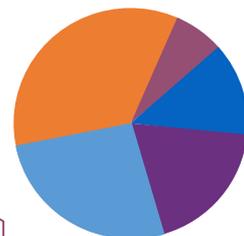
即席查询和并发查询

**预测**  
将要发生什么情况?



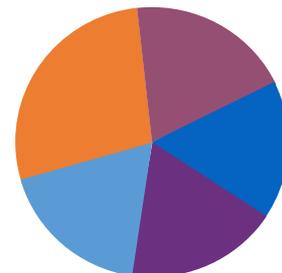
分析建模

**一线运营支撑**  
正在发生什么情况?



连续更新和流程互动

**主动事件**  
我希望发生什么情况!



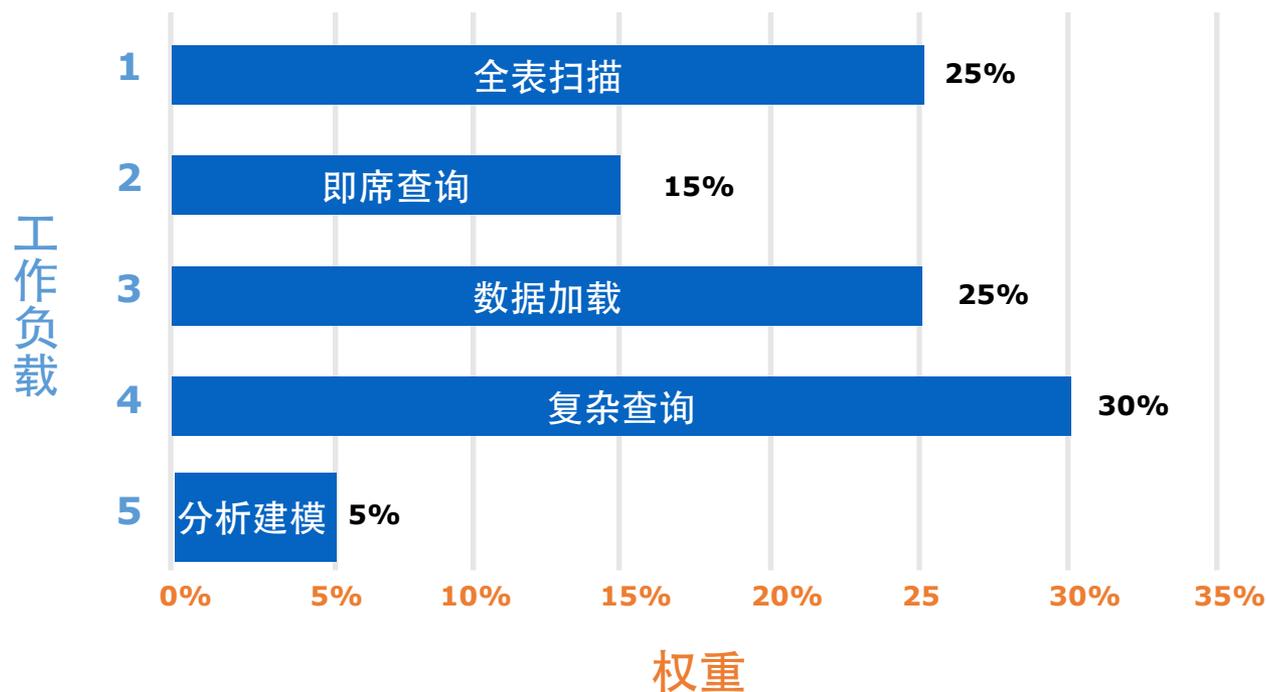
基于事件的触发

- 批处理
- 即席查询
- 分析
- 持续的更新/简短的战术性查询
- 主动触发

大部分的企业在前两个阶段

## • 典型负载

- 即席查询SQL：报表、简单查询、汇总
- 复杂检索：多字段检索、模糊检索、全文检索
- 全表扫描：离线DAG计算任务、ETL处理流程、预测等
- 交互式探索：自助交互式建模



## MPP数据仓库

- TeraData、Greenplum、SAP HANA等
- BI生态成熟、非SQL任务很难支持
- 巨贵

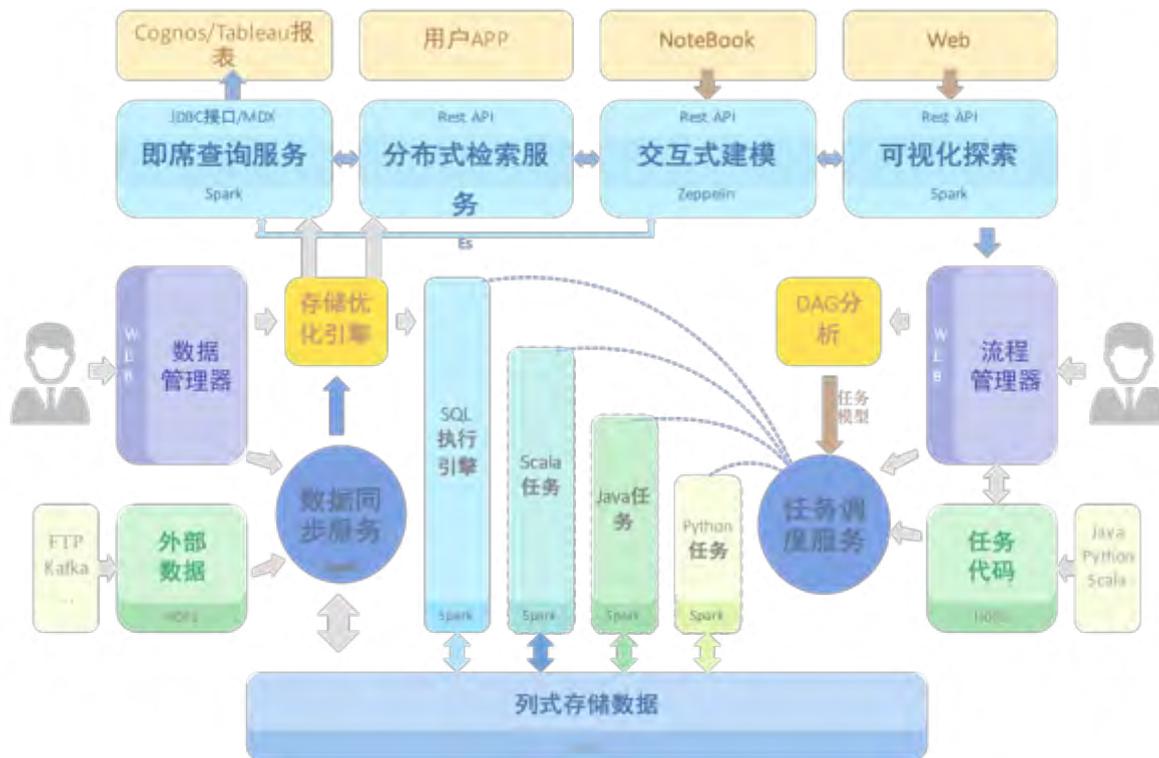
## 大数据平台

- HDP、CDH、星环等
- 技术先进、开源开放、坑多
- 暴露底层组件太多，**运维和使用技术曲线陡峭**

## 用户对大数据平台的期望

- 最好看起来像数据库一样，管理方便，使用简单
- 利用最新的大数据计算技术获得高性能和扩展能力
- 不需要掌握各种底层组件
- 兼容运行已有的数据库存储过程
- 统一管理各种数据处理任务
- 稳定可靠

## 2、INFINIDATA平台



- 特性

- 全量数据和表结构Schema自动导入
- 增删改等增量数据的智能同步
- 同时支持原始表和衍生表
- 支持对表数据和表结构Schema的变化轨迹溯源
- 自主选择存储引擎和分区分桶优化
- 数据变化自动触发相关的计算任务

- 特性

- 借鉴关系代数思想，计算流程等价于表的函数变换
- 计算流程统一管理，计算任务历史可追踪
- 计算逻辑和中间结果可共享
- 自动分析计算任务的依赖关系图进行全局调度优化
- 支持PLSQL存储过程和非SQL（Scala或Python）的复杂计算任务

- 特性

- 提供标准的JDBC访问接口，对常用的Cognos等报表服务提供driver
- 提供MDX语言的建模和OLAP分析引擎服务
- 提供标准Es接口提供数据检索服务
- 动态分层功能，将访问最频繁的数据保存在内存中，同时将很少访问的数据移至磁盘

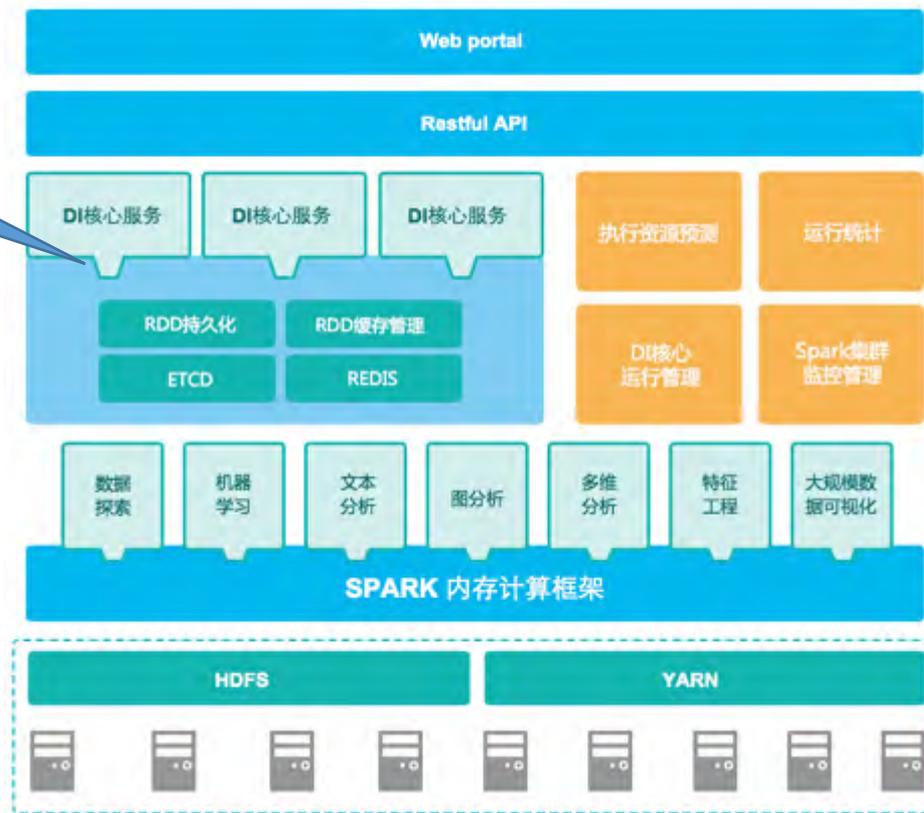
- 开源组件优化集成

- 修改hive2.3、spark2.1等相关组件的bug和源代码约80处

### 3、INFINIDATA可视化探索

## 系统框架

常驻内存服务



- 每个工程运行在一个单独的Spark环境
- Spark环境资源由YARN分配调度
- DI和Spark常驻内存，通过消息队列交互
- 利用RDD保存探索过程中的各种中间表

## 相关分析

- 行分析（聚类）
- 列分析（变量聚类）
- 值分析（频繁项）

## 离散矩阵分析

- 自相关分析
- 互相关分析

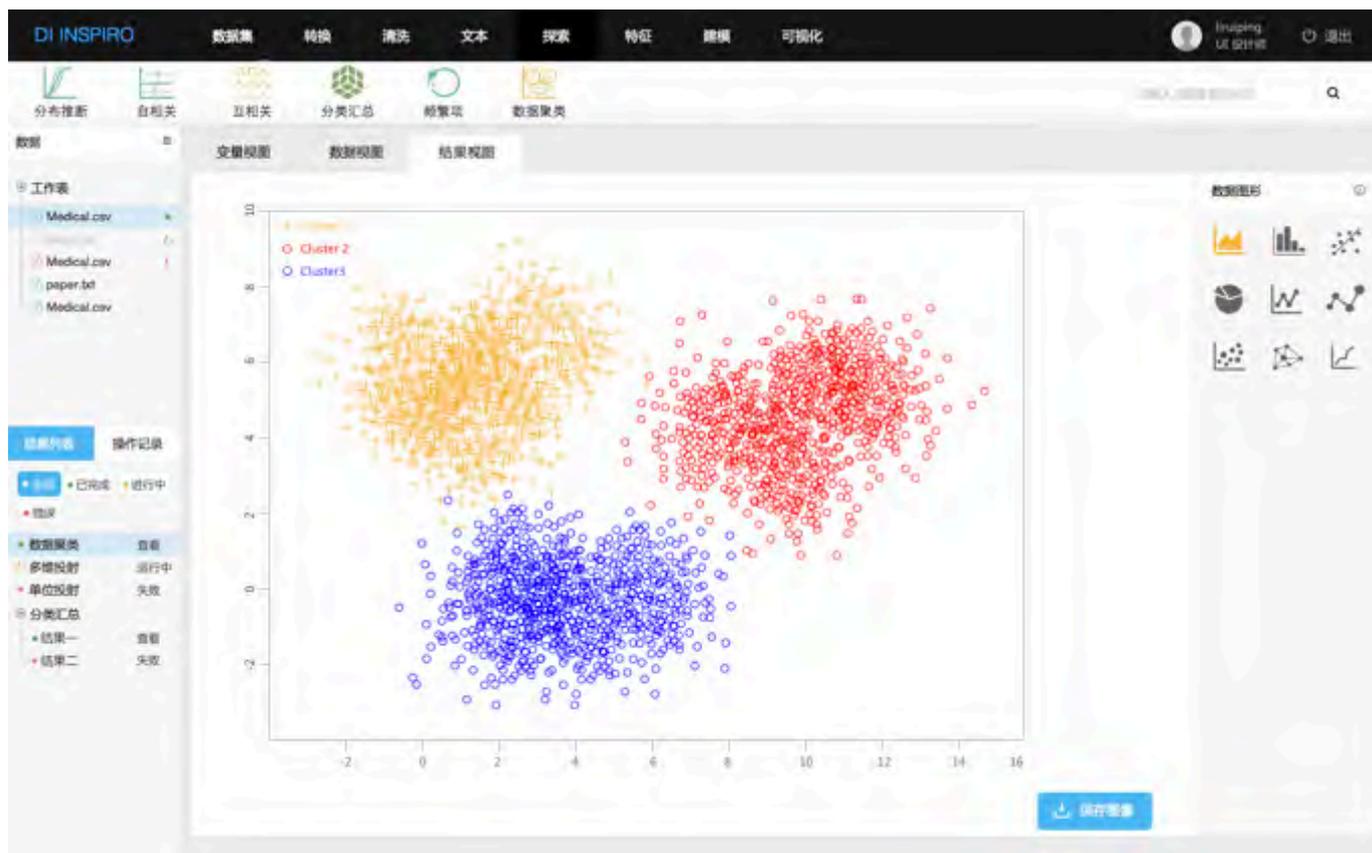
## OLAP分析

- Mondrian建模
- 多维度分析

## 可视化分析

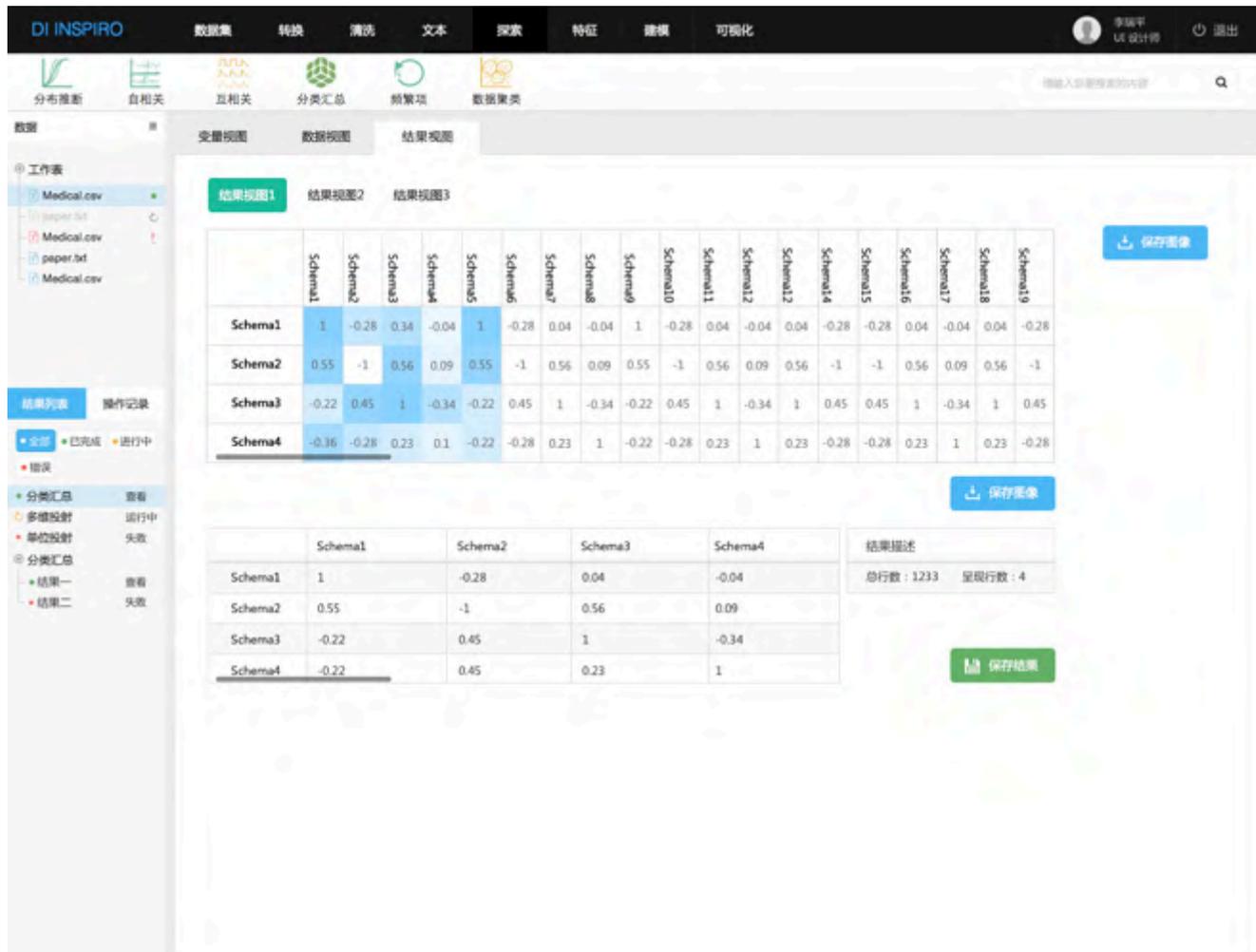
- 散点图
- 直方图
- 箱图
- 3维散点图

## K-means聚类分析



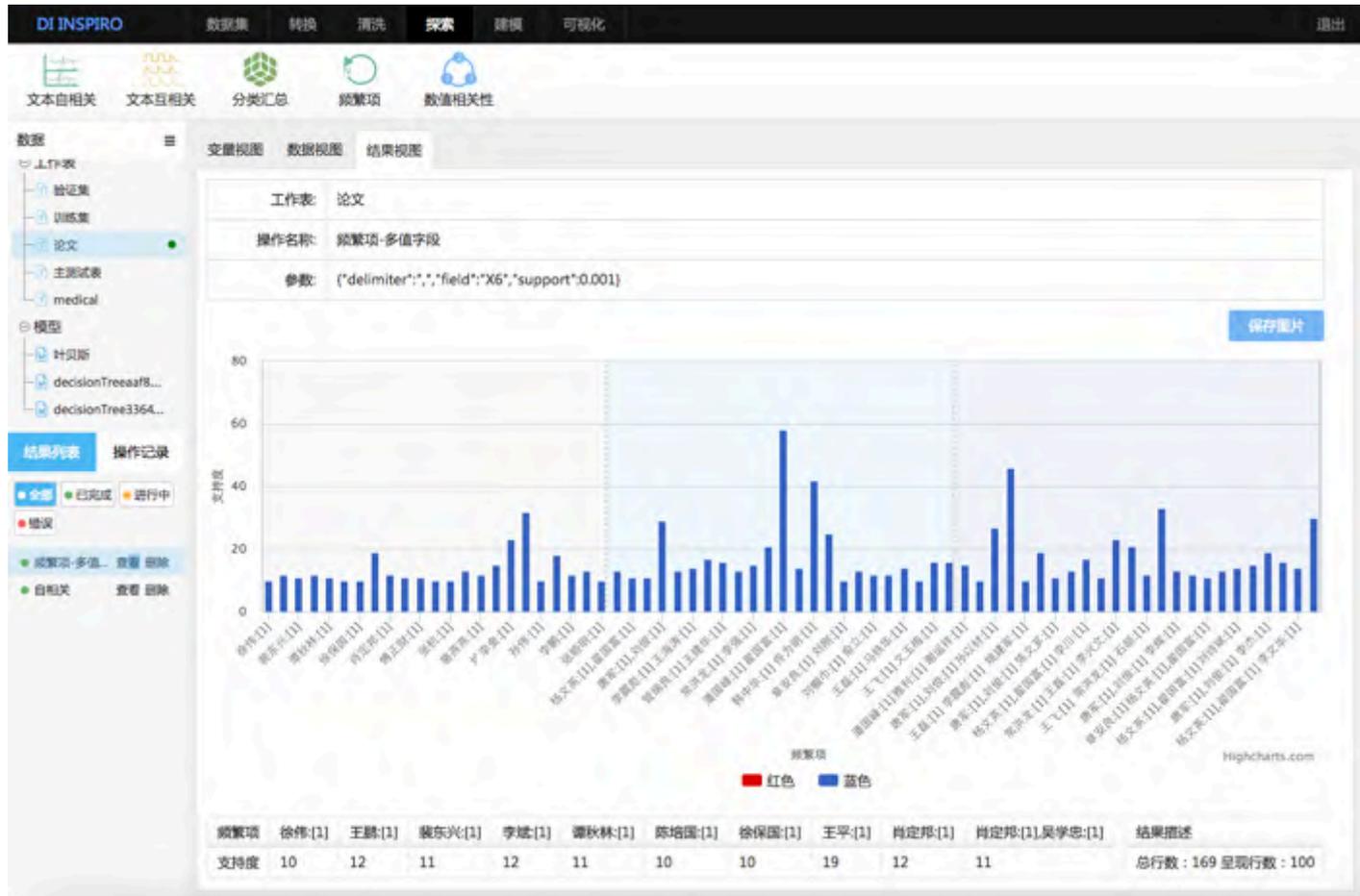
## 列相关-pearson相关性

- 热力图展现列之间相关性
- 发现基础变量和衍生变量



## 频繁项分析

属性值关联性



## 分析模型

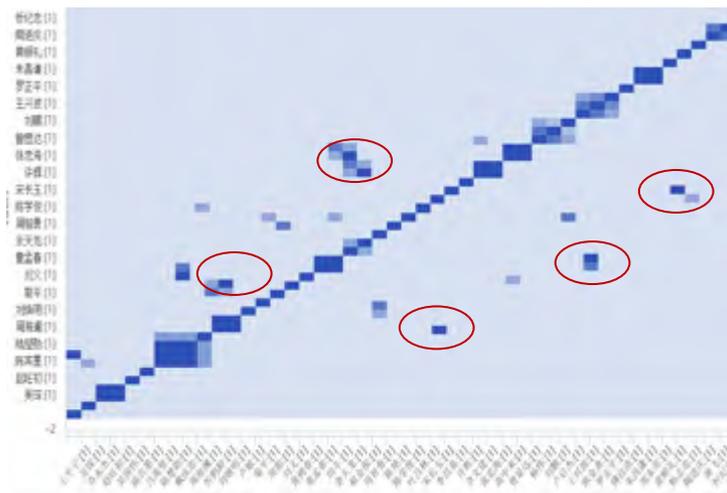
### • 图分析

- 社会网络分析
- 药品关联分析
- 公共安全

- (相同时间/机场) 乘坐相同航班的同乘分析
- (相同时间/地点) 的紧密通话客户分析

### • 科技领域

- 研发相类似技术领域的竞争对手分析
- 论文合作关系



publication_time	title	abstract	author
2000_1期	主编致辞	-	-
2000_1期	火功能汽轮机在亚临界区的试验研究	介绍通过风洞试验研究,得到在所有工况下,...	王平子 [1]
2000_1期	四杆机构机构的构型特性和传动特性	四杆机构机构的构型特性和传动特性...	葛琛 [1], 乔永志 [1]
2000_1期	循环流化床锅炉分离器的研究	着重介绍国内外分离器结构的及使用情...	彭洪和 [1]
2000_1期	300/600MW燃煤电厂输煤控制系统设计及其...	在总结上海工业自动化仪表研究所从事火电...	葛国伟 [1]
2000_1期	我国发电设备产业产品与技术发展预测	对今后10~15年内我国发电设备产业的产品...	陈其德 [1], 吕光健 [1], 魏楚敏 [1], 戴宗忠 [1]
2000_1期	论超临界汽轮机的发展	文章论述了世界上超临界火电机组的最新发...	傅海峰 [1], 孙鹤群 [1]
2000_1期	水轮机环流导叶射流网络生成技术及应用	网络生成是三维数据模拟的一个关键环节,对...	刘峰林 [1]
2000_1期	水轮发电机励磁系统问题的导向分析及应用	从同步电机励磁恒量角与功角特性出发,应用...	产鞍 [1]

## 原始数据表

	C1	C2	C3	C4	C5
R1	X <sub>1</sub>	Y <sub>1</sub>	A <sub>1</sub> ,A <sub>2</sub> ,A <sub>3</sub>	B <sub>1</sub> ,B <sub>2</sub> ,B <sub>3</sub>	Z <sub>1</sub>
R2	X <sub>2</sub>	Y <sub>2</sub>	A <sub>2</sub> ,A <sub>3</sub>	B <sub>2</sub> ,B <sub>4</sub>	Z <sub>2</sub>
R3	X <sub>3</sub>	Y <sub>3</sub>	A <sub>1</sub> ,A <sub>4</sub> ,A <sub>5</sub>	B <sub>2</sub> ,B <sub>3</sub> ,B <sub>6</sub>	Z <sub>3</sub>
R4	X <sub>4</sub>	Y <sub>4</sub>	A <sub>2</sub> ,A <sub>5</sub>	B <sub>1</sub> ,B <sub>4</sub>	Z <sub>4</sub>
R5	X <sub>5</sub>	Y <sub>5</sub>	A <sub>3</sub> ,A <sub>4</sub>	B <sub>1</sub> ,B <sub>5</sub>	Z <sub>5</sub>

投影



## 矩阵变换

	R <sub>1</sub>	R <sub>2</sub>	R <sub>3</sub>	R <sub>4</sub>	R <sub>5</sub>
A <sub>1</sub>	1	0	1	0	0
A <sub>2</sub>	1	1	0	1	0
A <sub>3</sub>	1	1	0	0	1
A <sub>4</sub>	0	0	1	0	1
A <sub>5</sub>	0	0	1	1	0

## 共现相关性

	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
A <sub>1</sub>		1	1	1	1
A <sub>2</sub>	1		2	0	1
A <sub>3</sub>	1	2		1	0
A <sub>4</sub>	1	0	1		1
A <sub>5</sub>	1	1	0	1	

## cos相关性

	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
A <sub>1</sub>	1	0.408	0.408	0.5	0.5
A <sub>2</sub>	0.408	1	0.667	0	0.408
A <sub>3</sub>	0.408	0.667	1	0.408	0
A <sub>4</sub>	0.5	0	0.408	1	0.5
A <sub>5</sub>	0.5	0.408	0	0.5	1

## pearson相关性

	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
A <sub>1</sub>	1	-0.167	-0.167	0.167	0.167
A <sub>2</sub>	-0.167	1	0.167	-1	-0.167
A <sub>3</sub>	-0.167	0.167	1	-0.167	-1
A <sub>4</sub>	0.167	-1	-0.167	1	0.167
A <sub>5</sub>	0.167	-0.167	-1	0.167	1

### 理赔数据表

案件编号	人员	车牌	地点	金额
344561	段建华, 张华, 许卫	湘A2BA32, 湘AA1391, 湘ZG00069	板仓南路	20000
344562	罗坚, 肖蓉	湘J7ZH83, 湘AL5S85	开元西路	50000
344563	王丽萍, 刘双泉	湘A65N90, 湘A1661K	寿昌路	100000
344564	彭发兵, 周辉, 苏英雄	湘A2ZB92, 湘B2HL12, 湘A2KA19	人民路	70000
344565	张斌, 王丽萍	湘AT8137, 湘A65N90	湘江东路	10000

### 矩阵变换

	344567	344568	344569	344570	344571
谢前	0	1	1	0	1
敬春桥	0	0	0	1	0
罗坚	0	1	0	0	1
肖蓉	0	0	1	0	0
刘双泉	0	1	0	0	0

	谢前	敬春桥	罗坚	肖蓉	刘双泉
谢前	1	-0.612	0.667	0.408	0.408
敬春桥	-0.612	1	-0.408	-0.25	-0.25
罗坚	0.667	-0.408	1	-0.408	0.612
肖蓉	0.408	-0.25	-0.408	1	-0.25
刘双泉	0.408	-0.25	-1	-0.25	1

### cos相关性

	谢前	敬春桥	罗坚	肖蓉	刘双泉
谢前	1	0	0.816	0.577	0.577
敬春桥	0	1	0	0	0
罗坚	0.816	0	1	0	0.707
肖蓉	0.577	0	0	1	0
刘双泉	0.577	0	0.707	0	1

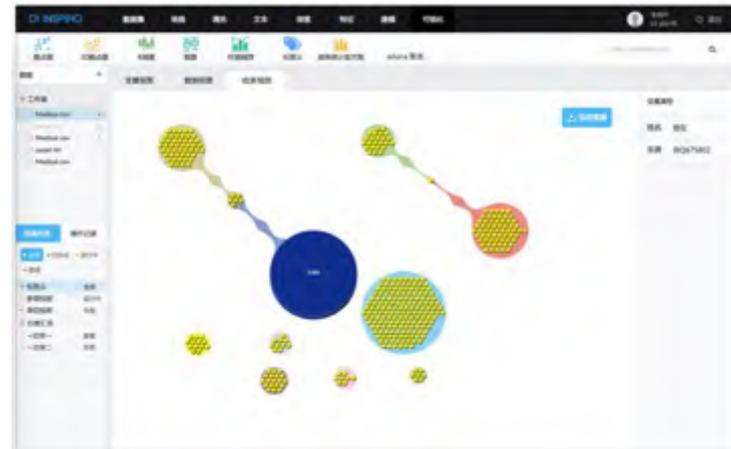
### pearson相关性

	谢前	敬春桥	罗坚	肖蓉	刘双泉
谢前		0	2	1	1
敬春桥	0		0	0	0
罗坚	2	0		0	1
肖蓉	1	0	0		0
刘双泉	1	0	1	0	

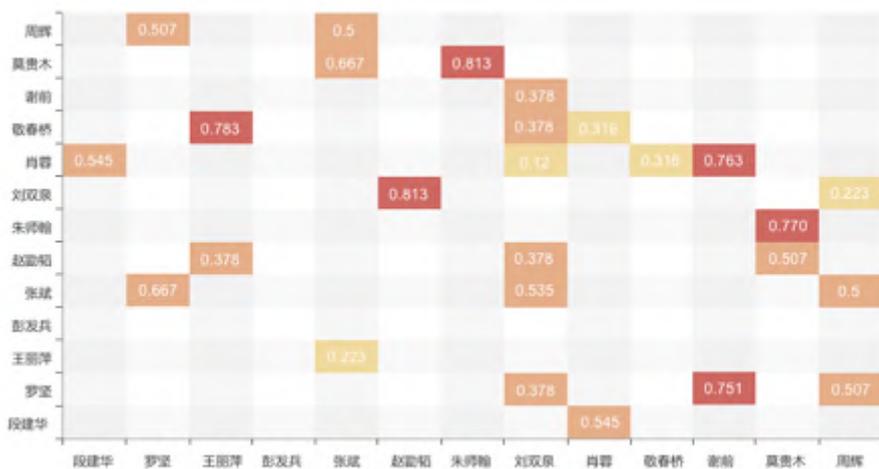
### 共现相关性



相关矩阵分析map



Aduna图合作关系可视化



### 原始数据表

	C1	C2	C3	C4	C5
R1	X <sub>1</sub>	Y <sub>1</sub>	A <sub>1</sub> ,A <sub>2</sub> ,A <sub>3</sub>	B <sub>1</sub> ,B <sub>2</sub> ,B <sub>3</sub>	Z <sub>1</sub>
R2	X <sub>2</sub>	Y <sub>2</sub>	A <sub>2</sub> ,A <sub>3</sub>	B <sub>2</sub> ,B <sub>4</sub>	Z <sub>2</sub>
R3	X <sub>3</sub>	Y <sub>3</sub>	A <sub>1</sub> ,A <sub>4</sub> ,A <sub>5</sub>	B <sub>2</sub> ,B <sub>3</sub> ,B <sub>6</sub>	Z <sub>3</sub>
R4	X <sub>4</sub>	Y <sub>4</sub>	A <sub>2</sub> ,A <sub>5</sub>	B <sub>1</sub> ,B <sub>4</sub>	Z <sub>4</sub>
R5	X <sub>5</sub>	Y <sub>5</sub>	A <sub>3</sub> ,A <sub>4</sub>	B <sub>1</sub> ,B <sub>5</sub>	Z <sub>5</sub>

投影

### 矩阵变换

	B <sub>1</sub>	B <sub>2</sub>	B <sub>3</sub>	B <sub>4</sub>	B <sub>5</sub>
A <sub>1</sub>	1	2	2	0	0
A <sub>2</sub>	1	2	1	2	0
A <sub>3</sub>	1	2	1	1	1
A <sub>4</sub>	1	1	1	0	1
A <sub>5</sub>	1	1	1	1	0

### 共现互相关性

	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
A <sub>1</sub>		3	3	3	3
A <sub>2</sub>	3		5	0	0
A <sub>3</sub>	3	5		2	0
A <sub>4</sub>	3	0	2		3
A <sub>5</sub>	3	0	0	3	

### pearson互相关性

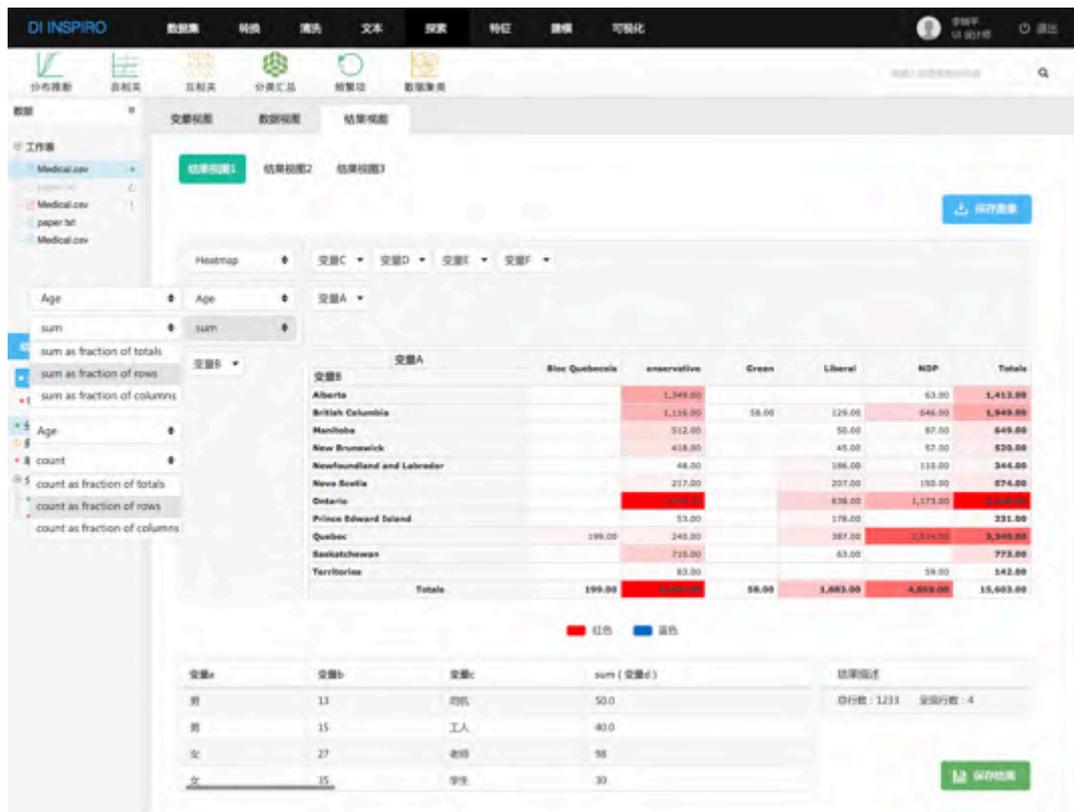
	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
A <sub>1</sub>	1	0.25	0.354	0.548	0.548
A <sub>2</sub>	0.25	1	0.707	-0.548	0.548
A <sub>3</sub>	0.354	0.707	1	0	0
A <sub>4</sub>	0.548	-0.548	0	1	-0.2
A <sub>5</sub>	0.548	0.548	0	-0.2	1

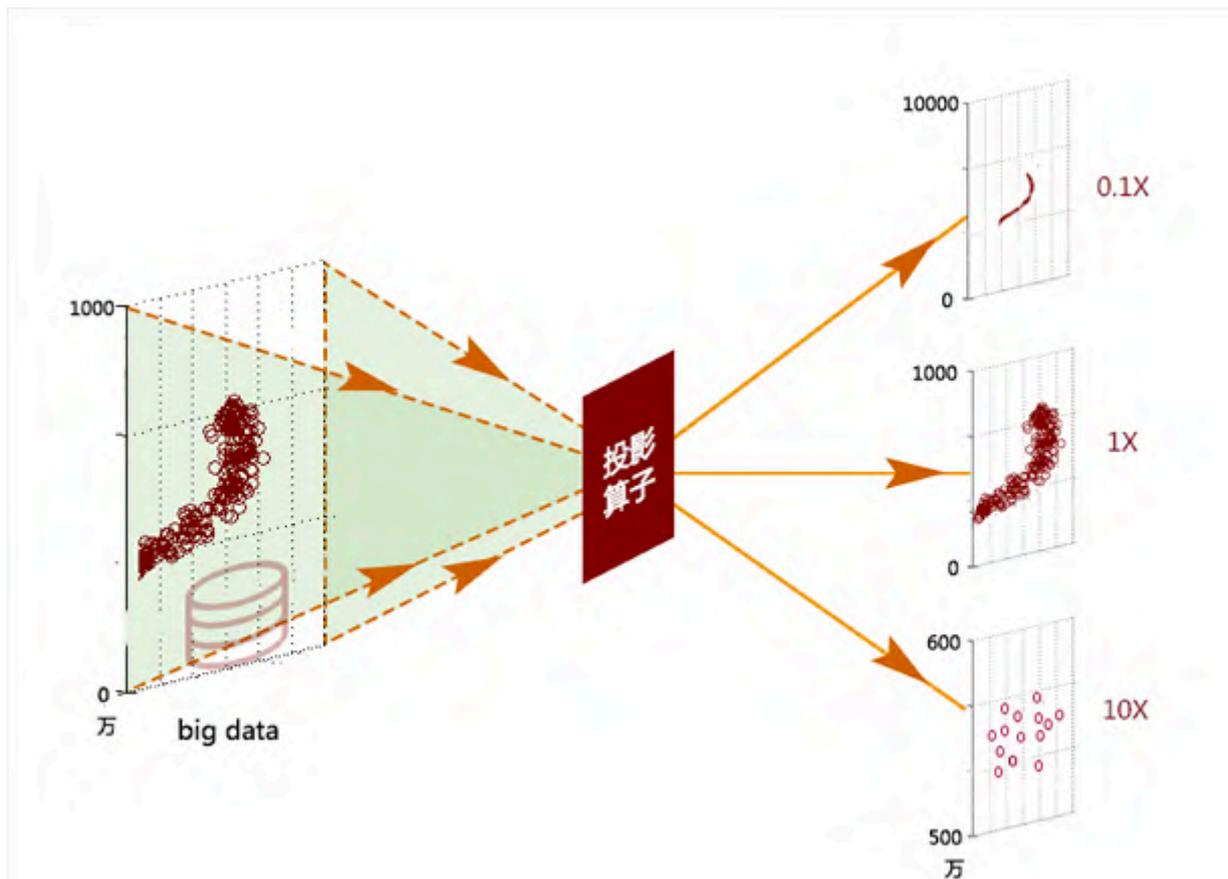
### cos互相关性

	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>	A <sub>4</sub>	A <sub>5</sub>
A <sub>1</sub>	1	0.7	0.783	0.849	0.849
A <sub>2</sub>	0.7	1	0.894	0.566	0.849
A <sub>3</sub>	0.783	0.894	1	0.791	0.791
A <sub>4</sub>	0.849	0.566	0.791	1	0.8
A <sub>5</sub>	0.849	0.849	0.791	0.8	1

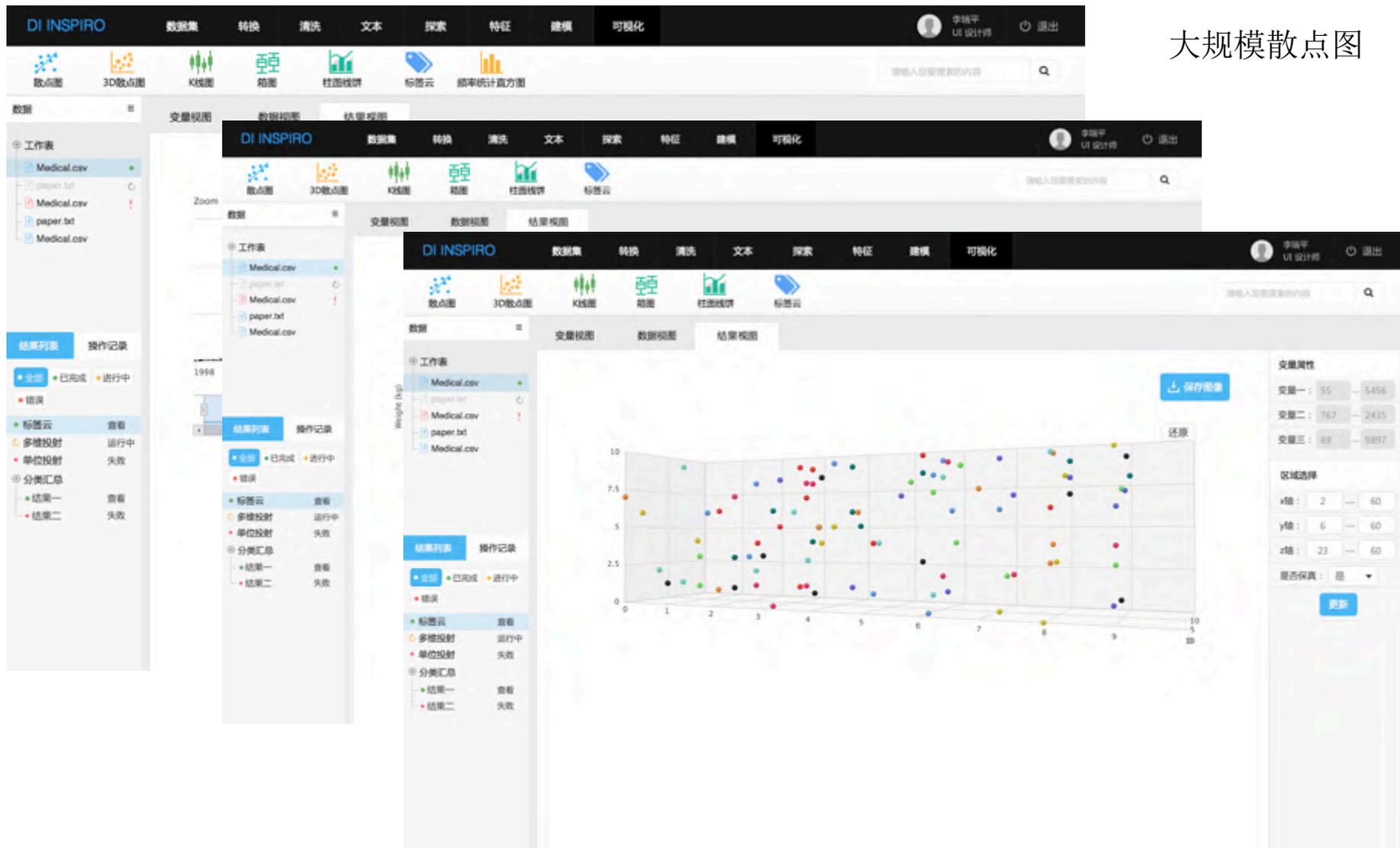
## 多维度汇总分析

- 基于MDX的数据模型
  - Mondrian引擎





- 基于保持数据概率分布不变的思想
- 将原始数据根据缩放级别和距离远近将原始数据映射成特定显示区域的矩阵



大规模散点图

## 4、应用实践



## 需求

- 100多个维度、60多个指标，
- 单表数据数据量大小为6-30G
- 修改Hive JDBC支持Cognos报表
- DAG任务流程运行生成事实表和维度表

	2016年	2017年	年份
本年投保数量	25863	1032	26895
本年未结赔付款件数	25285	1328	26613
本年未结赔款	1352067273.070000	24478550.950000	1376545874.020000
本年未结赔付款件数	24932	1289	26221
本年已结赔付款件数	186057	2939	188996
本年未结赔款	.000000	.000000	.000000
本年已结赔款	595634179.240000	16753960.080000	612388139.320000
本年非零赔款已结赔付款件数	9547	430	9977
本年非零赔款已结赔付款件数	151943	2702	154645
本年新车购置价	200439283754.810060	28548525007.000008	228987808761.810090
本年赔单赔付率	78569	2706	81275
本年有效报案数量	221602	4556	226158
本年已赔单保费	24492138.164092	465256.912296	24957395.076388
本年有效报案数量/本年已赔单保费	873	351	1224
本年标准保费	6991574091.370002	891420936.090000	7882995027.460002
本年商业险标准保费	6991552660.020002	891420936.090000	7882973596.110002
本年NCD调整保费	828165007560.100100	63019723564.250015	691204731124.350100
本年NCD续保调整保费	604700562090.704470	60221147865.475906	664921709956.180300
本年NCD续保续用调整保费	603032356009.099980	60042105420.350014	663074463429.450200

## 需求

- 保险客户信息真实性较低，无法服务于精准化营销和客服资源的精准化投放
- 每天导入来自车管所、电信公司、俱乐部、分布在承保理赔各环节的碎片化信息
- DAG流程处理各来源数据并进行交叉核验，生成用户画像

**客户业务信息**

信息采集类型	信息采集的名称	数据来源类型	客户名称
证件类型	证件号码	手机号码1	手机号码
手机号码3	手机号码4	手机号码	手机号码来源依据

**客户基本信息**

民族	性别	出生
座机号码	单位名称	家庭住址
兴趣	爱好	

**客户车辆信息**

车牌号码	发动机号	车架号	排量
车辆厂牌型号	号牌种类	VIN码	初登日期
号牌底色	车辆种类	车辆使用性质	整备质量
核定载客人数	核定载质量	准牵引总质量	

Copyright © 2008-2017 XXXXXX版权所有



子曰：人而无信，不知其可也  
---- 《论语·为政》

如果没有加这个微信，不知道还可不可以做大数据？

谢谢！  
欢迎交流！