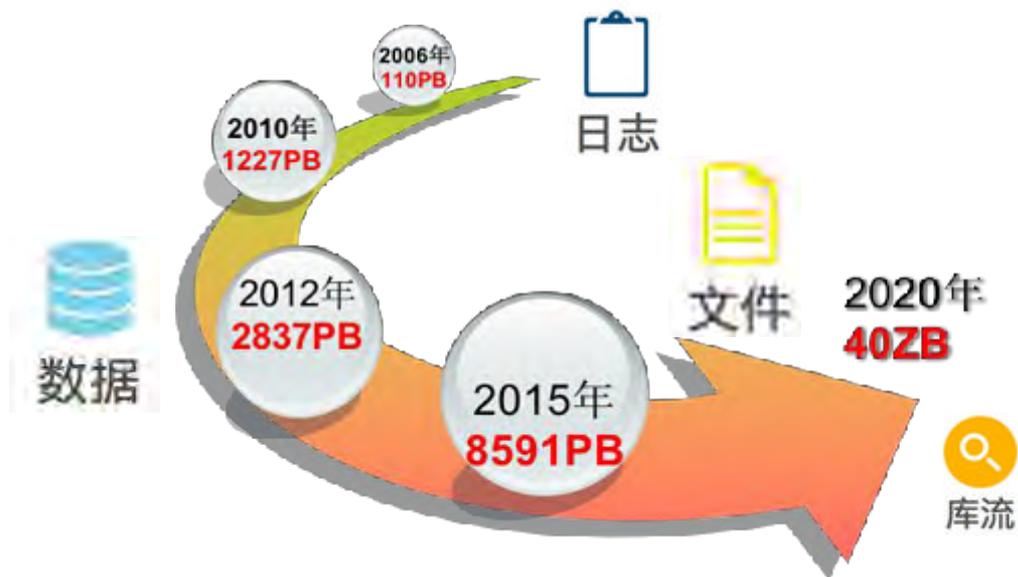


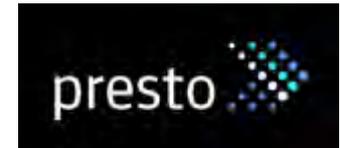
# DI: 基于SPARK的交互式数据探索与建模系统

熊永平

北京邮电大学



- 企业级应用中，BI已经满足不了业务人员的数据分析需求
- 用户越来越希望用高级数据分析方法提升业务能力
- 开发数据价值依赖于大量的懂业务的数据分析人员
- 需要强大易用的大数据分析系统



- 数据分析往往利用多种编程语言或系统管理不同的数据挖掘任务和机器学习流程
- 需要掌握从分布式系统到数据分析等众多门槛很高的工具和技能



- 针对普通的数据分析人员
- 几乎不需要编程开发分布式程序
- 提供直观易用的图形化系统界面
- 提供可扩展的数据分析手段
- 可处理大规模数据集
- 计算能力可线性增加
- 部署运维简单



性能



- 基于内存的架构极大的减少了磁盘I/O
- 通用任务上20-100x速度的提升

高效



- 精简且表达力强大的语法(如Spark2.0的Dataframe)
- 统一的编程模型
- 能用主流的编程语言—Java, Python, Scala
- 新工具减少使用的障碍 (Spark2.0支持SQL2003)

扩展性



- 通过增加机器计算能力实现整体分析能力的线性扩展
- 单个计算节点失效后自动重算
- 利用自动持久化确保整个计算过程可靠完成

利用Hadoop资产

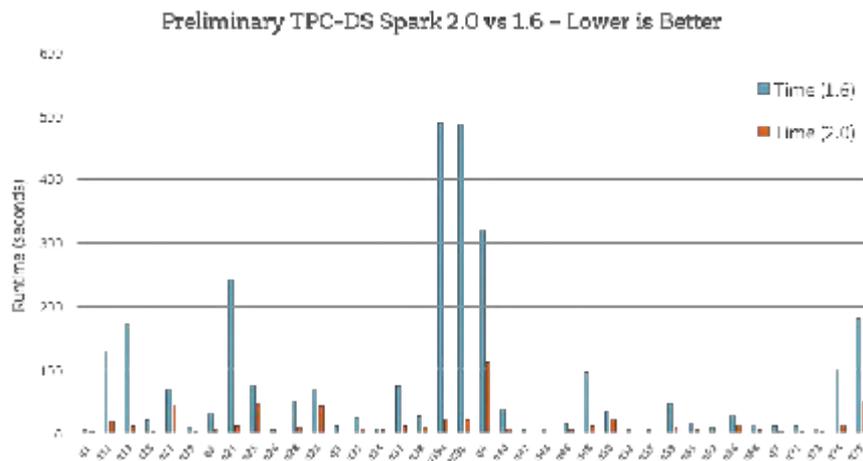
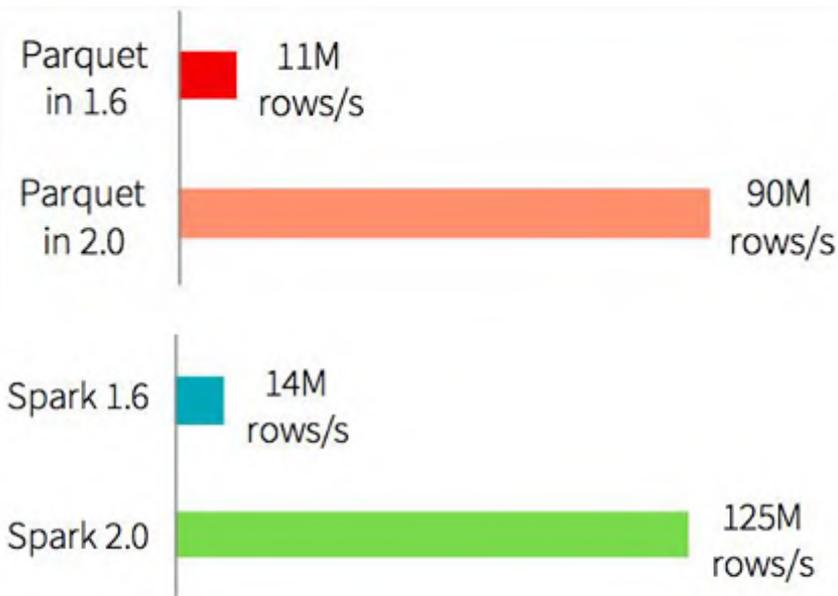
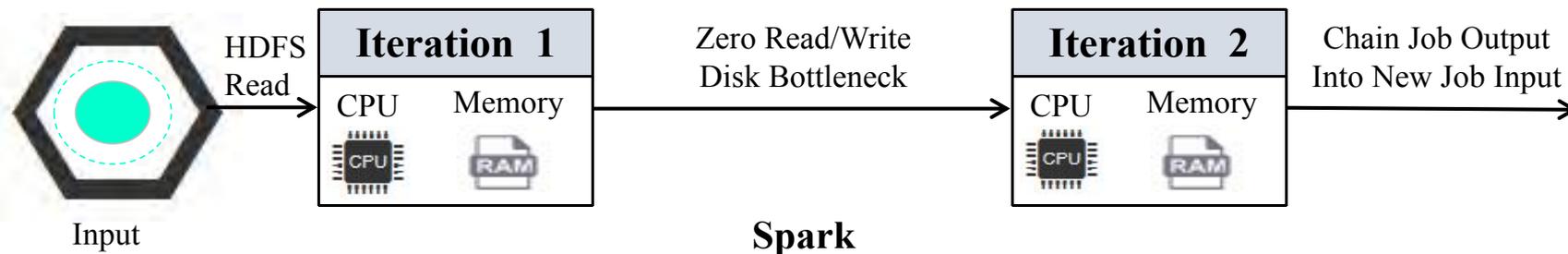


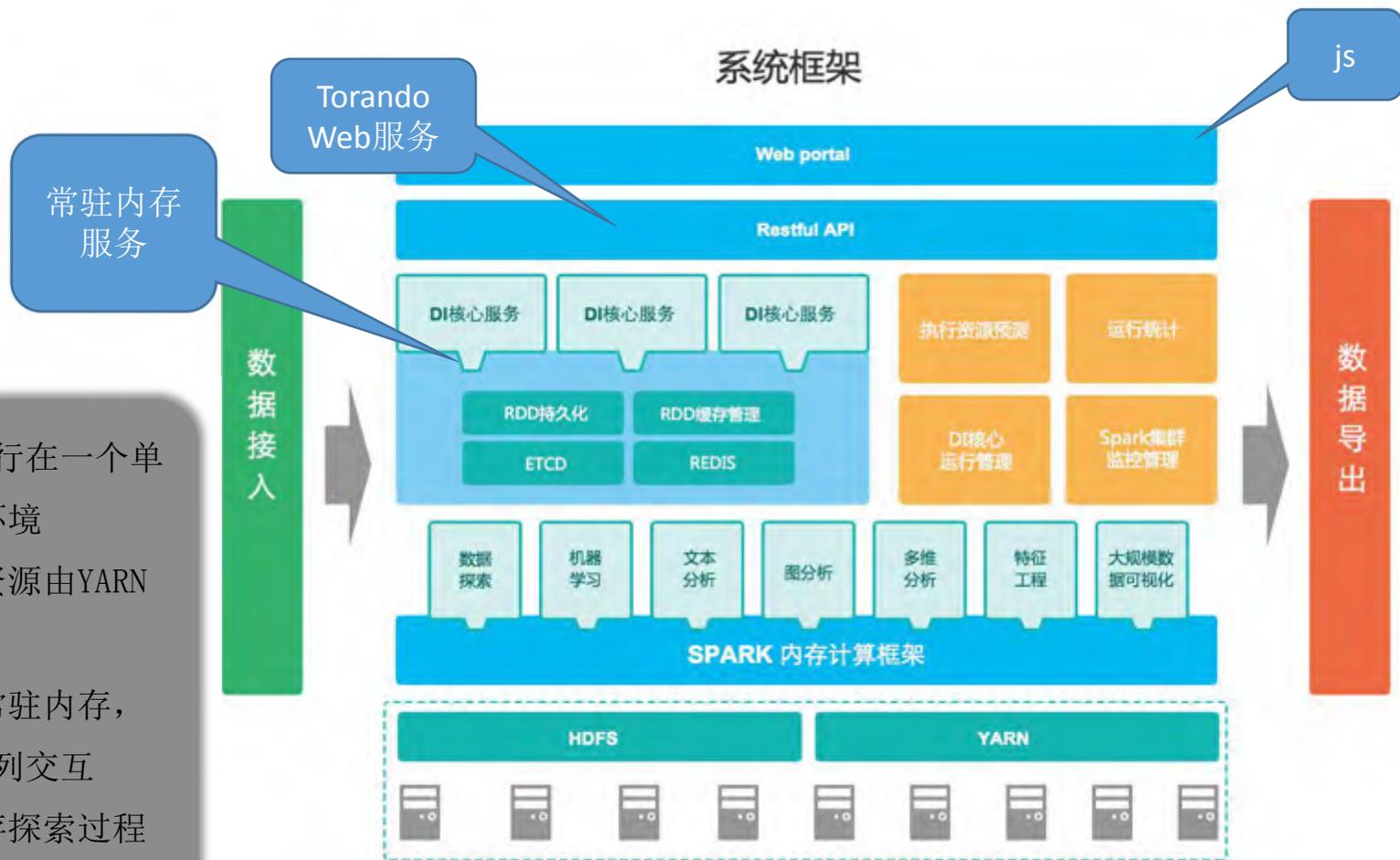
- 和已有的Hadoop生态系统能很好的工作

持续改进



- 大量的社区贡献者持续的对Spark技术栈组件进行持续的快速更新





- 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

The screenshot displays the DI INSPIRO interface. On the left, a tree view shows a project structure with folders like '数据源' and '数据表'. The main area is divided into two panes. The top pane shows a table schema with columns for ID, name, type, percentage, default value, and format. The bottom pane shows a data table with columns X1 through X7 and a '姓名' column.

#	名称	数据类型	覆盖率	默认值	分隔符	数据类型	注释
1	X1	DoubleType	100%	null	.	数值	自增
2	X2	DoubleType	100%	null	.	数值	自增
3	X3	DoubleType	100%	null	.	数值	自增
4	X4	DoubleType	100%	null	.	数值	自增
5	X5	StringType	100%	null	.	字符串	自增
6	X6	IntegerType	100%	null	.	整数	自增
7	X7	StringType	100%	null	.	字符串	自增

#	X1	X2	X3	X4	X5	X6	X7
1	4.38726985234	5.11825739016	51.9363492617	9.5055272425	China	64	博士
2	9.90873795945	5.37935188264	79.5436897973	15.2886698421	Canada	71	研究生
3	4.48357921347	5.62790548726	52.2678990074	10.0734847007	America	21	初中
4	4.0681771698	4.60285485209	50.340895849	8.67101180089	China	39	研究生
5	5.33375562027	4.5800933045	56.6687781014	9.919861915072	Japan	57	高中
6	2.08900657546	4.86072805227	44.9450328773	7.84933463473	China	15	博士
7	6.32618123902	3.22499782332	62.6309061951	8.75317906234	Canada	57	高中
8	2.54180003688	5.68094525558	42.7090001844	8.23274256246	America	26	高中
9	4.02212579774	5.94431147937	50.1106289887	9.9643727721	Canada	26	初中
10	5.97673416557	5.0070513095	59.8856708269	10.9817854349	America	52	高中
11	8.89325047995	5.72155299007	74.4962523998	14.61480343	Canada	12	大学
12	6.93782301582	4.30156951826	61.6891151791	10.6393895541	America	17	高中
13	9.5341156366	4.6758830821	77.6705798185	14.2101042723	America	68	大学
14	9.44956817723	3.20563282222	77.2478418861	12.652021895	Japan	84	小学
15	5.30634488992	5.87459071709	56.5317249496	11.180935327	England	6	小学
16	3.06532723155	4.84295708736	45.3266361577	7.50828431891	Japan	37	研究生
17	8.79142092488	3.88325912588	73.9571046244	12.6746800508	Japan	22	博士
18	1.46912134635	5.89982545846	37.3466067317	7.3691379122	England	43	大学
19	7.12430816879	5.46432062032	65.621541844	12.5886289891	America	87	大学
20	8.97402108767	4.96405896788	75.9870054384	14.1634200555	Japan	98	小学

- 操作对象抽象为表
- 函数式编程思想
- 所有算子不改变原表数据
- 增加列或生成新表
- 新表单独保存

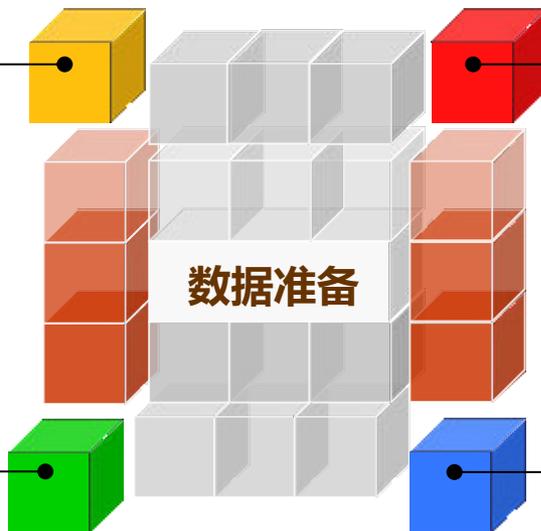
- 算子操作异步运行
- 记录算子流程DAG
- 流程持久化，实现自动批处理

## ▪ 抽样过滤

- 均匀/随机/分层抽样
- 多条件组合过滤
- 根据距离/密度/局部离群因子/类离群因子/离群点检测的聚类过滤离群点

## ▪ 去重合并

- 左右连接Join
- 取交集/并集
- 合并Merge
- 多列联合主键去重
- 识别并移除冗余



## ▪ 变量构造

- 根据已有变量拆分/组合/抽取/运算生成新变量
- 支持数学/统计/文本/日期等100多种函数运算及其逻辑组合
- 根据离散或连续分布函数/等差/等比/日期分布等生成基础列数据

## ▪ 归一化分箱

- 按根据字典数据替换
- 不同变量进行取值区间归一化
- 连续变量离散化等频/等距/Bootstrap/聚类等分箱
- 利用变量统计值填充空值



等距分箱



age	sex	病	date	money	equiBucket_age_1
53	男	手术后恶性肿瘤化学治疗	2013/11/28	5398.08	37.6-56.400000000000006
44	女	恶性肿瘤维持性化学治疗	2013/11/27	5234.93	37.6-56.400000000000006
60	男	恶性肿瘤维持性化学治疗	2013/11/27	4660.67	56.400000000000006-75.2
61	男	肺恶性肿瘤	2013/11/25	4237.37	56.400000000000006-75.2
58	女	肺恶性肿瘤	2013/11/29	3714.02	56.400000000000006-75.2
73	男	慢性肾脏病5期	2013/11/28	7380.61	56.400000000000006-75.2
78	女	肺炎	2013/11/16	21720.66	75.2-94.0000000001
63	女	慢性肾脏病5期	2013/12/2	1532.93	56.400000000000006-75.2
50	男	肺恶性肿瘤	2013/11/20	11407.83	37.6-56.400000000000006
43	女	肺恶性肿瘤	2013/11/29	2181.65	37.6-56.400000000000006
75	男	肺部感染	2013/11/25	11048.23	56.400000000000006-75.2
62	男	慢性肾脏病4期	2013/11/16	9315.45	56.400000000000006-75.2
74	男	肺恶性肿瘤	2013/12/2	7714.72	56.400000000000006-75.2
48	男	肾性贫血	2013/12/4	9267.5	37.6-56.400000000000006
22	女	肺炎	2013/11/28	10447.17	18.8-37.6
76	男	非典型分枝杆菌感染	2013/11/26	9020.44	75.2-94.0000000001

归一化

系统中采用了两种归一化方法：

1. min-max归一化  $x^* = \frac{x-x_{min}}{x_{max}-x_{min}}$
2. Z-scores归一化  $x^* = \frac{x-\mu}{\sigma}$



文本抽取

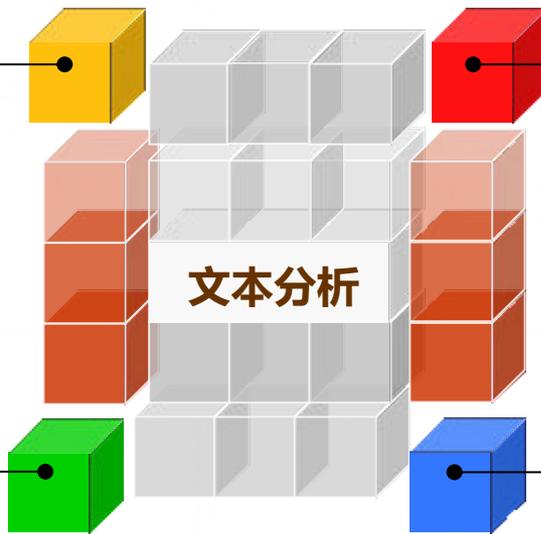
X7	extract1_X7_1
博士	-
研究生	-
初中	中
研究生	-
高中	中
博士	-
高中	中
高中	中
初中	中
初中	中
大学	-
高中	中
大学	-
小学	-
小学	-
研究生	-
博士	-

## ▪ 词句段切分

- 支持多种不同的词语切分方法和用户词典加载
- 提供段落和句子切分功能
- 提供新词发现功能，利用互信息熵、条件随机场CRF等算法

## ▪ 特征词抽取

- 提供TFIDF统计计算
- 利用TextRank算法抽取特征词
- Word2vec词向量模型



## ▪ 实体识别

- 提取标准实体, 如人名、地名、时间、日期和物理量等
- 提取特定领域的概念

## ▪ 文本挖掘

- 文本分类
- 文本聚类
- 文档矩阵主题分析

利用文档预处理、自然语言处理、主题检测等功能分析文本数据，便于数据分析人员处理非结构化文本数据。

从目标文本中按某种算法提取关键词

- TFIDF算法

$$\text{词频}(TF) = \frac{\text{某个词在文章中的出现次数}}{\text{文章的总词数}}$$

$$\text{逆文档频率}(IDF) = \log\left(\frac{\text{语料库的文档总数}}{\text{包含该词的文档数}+1}\right)$$

$$TF - IDF = TF * IDF$$

- TextRank算法

- 将待抽取关键词的文本进行分词
- 固定长度为n(通常取5)的窗口，在一个窗口中的任两个单词对应的节点之间存在一个无向无权的边
- 基于窗口分割后的边构建TextGraph图，利用PageRank计算每个节点的重要性
- 得到最重要的k个节点即提取出了k个关键词



## • 信息熵定义

□ □ □ □ □ □  $-\left(\frac{1}{2} \cdot \log\left(\frac{1}{2}\right) + \frac{1}{6} \cdot \log\left(\frac{1}{6}\right) + \frac{1}{3} \cdot \log\left(\frac{1}{3}\right)\right) = 1.5850$

□ □ □ □ □ □ □  $-\left(\frac{5}{8} \cdot \log\left(\frac{5}{8}\right) + \frac{1}{8} \cdot \log\left(\frac{1}{8}\right)\right) = 0.450561$

□ □ □ □ □ □ □ □  $-\log(1) = 0$

## • 新词识别

- 词频
- 自由度（片段所有可能左右邻的混乱程度）
  - 例句：利用公用配电**负荷**历史负载率以及中长期配电**负荷**预测结果，结合配变的容量，指出**负荷**容量不足和容量过剩的配变。
  - 片段“负荷”的所有左邻字实例为{电，出}  
熵为  $-\left(\frac{1}{2} \cdot \ln\left(\frac{1}{2}\right) + \frac{1}{2} \cdot \ln\left(\frac{1}{2}\right)\right) \approx 0.693$
  - 片段“负荷”的所有右邻字实例为{历，预，容}  
熵为  $-\left(\frac{1}{2} \cdot \ln\left(\frac{1}{2}\right) + \frac{1}{3} \cdot \ln\left(\frac{1}{3}\right) + \frac{1}{3} \cdot \ln\left(\frac{1}{3}\right)\right) \approx 1.08$
  - “负荷”的右邻字比左邻字更丰富更灵活
- 凝合度
  - 令n为文本的长度，令f(x)为字符串x在文本中出现的次数
  - 令p(x)为f(x)/n，即字符串x出现的概率
  - 定义“负载率”的可拆分为：
    - $\max\left(\frac{p(\text{负}) \cdot p(\text{载率})}{p(\text{负载率})}, \frac{p(\text{负载}) \cdot p(\text{率})}{p(\text{负载率})}\right)$

- 利用知网电力领域期刊全文数据库  
50G文本数据
- 抽取新词1100个

俄统国际 生物质燃料 热电联产  
能源消费 电价基 固体生物质 抽  
水蓄能 农村电气化 生物甲烷 联  
合循环 资产融资 乙醇燃料 国际电  
力 运输燃料 南卡 生物质柴油  
光伏组件 并网光伏 全球光伏 卡  
奥拉巴萨 生物质供热 并网太阳能  
能源国际 燃料混合 矿产能源部  
燃料车 税收激励 埃克森美孚 特  
许权招标 亚联邦 太阳能法令  
径流式 燃能系统

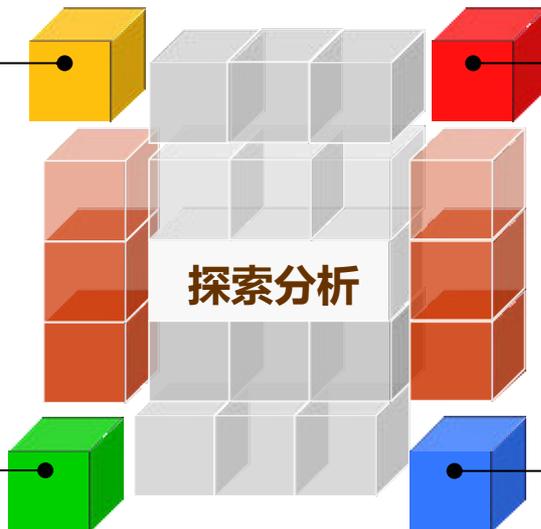


## ▪ 变量统计

- 描述性数据统计，包括常见统计量、分位数等
- 提供变量概率密度和直方图，以及分布推断
- 实现变量信息熵和信息值IV计算
- 多维分组汇总OLAP

## ▪ 行相关分析

- 关联规则分析，支持频繁项集的交互式筛选
- 对连续和离散变量的 K 均值聚类，自动估计最佳聚类数量，输出整个数据集的聚类归属和距离测量值
- 超过20种相似性距离测量方法



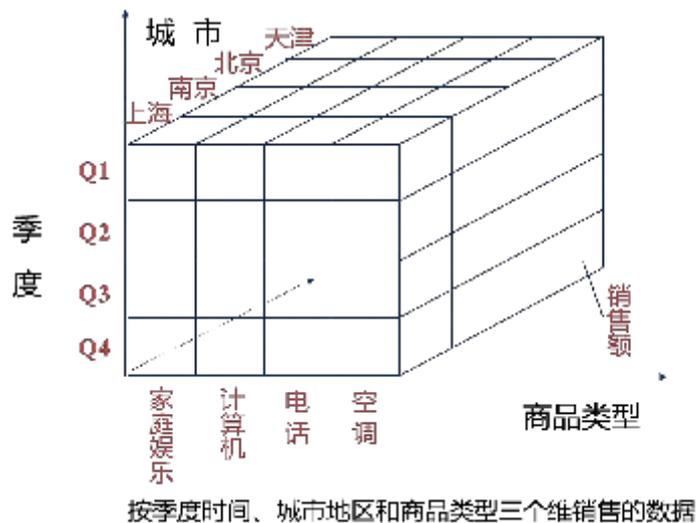
## ▪ 特征降维 (列相关分析)

- 变量聚类
- 主成分分析PCA
- 支持维度的层次划分
- 支持多种度量的计算函数

## ▪ 离散相关分析

- 相关性度量：pearson、spearman、余弦相关系数、互信息等
- 共现矩阵分析
- 自相关矩阵分析
- 互相关矩阵分析

- 利用cube进行分组汇总



Item	Color	QtySum
Table	Blue	123
Table	Red	223
Table	ALL	347
Chair	Blue	101
Chair	Red	210
Chair	ALL	311
ALL	Blue	225
ALL	Red	433
ALL	ALL	658

Item	Color	QtySum
Table	Blue	123
Table	Red	223
Chair	Blue	101
Chair	Red	210

- 提供了对数据进行多种维度、多种度量方法进行汇总展示的功能。
- 快速获取数据的宏观统计信息。

DI INSPIRO 数据表 视图 预览 预览 建模 可视化

文本自相关 文本互相关 分类汇总 维聚图 数据相关性

数据表 交互视图 数据视图 结果视图

工作表: 测试

操作名称: 分类汇总

参数: ["dimension":["X5","X6","X7"],"statistic":["sum"],["X1","avg"]]

Table: avg(X1) \* sum(X2)

sum(X2) \* X5 \* X7

X5	X6	1	2	3	4	5	6	7	8
X7									
X5									
America									
Canada	1596221347966897			1471025477038811					
China	1599036374666802			141270192712111					
England					1389782936072113	1469739644122851			
French		15427773989177297							
Japan							1495312466026359	1491546474072003	14405801712964391
Totals	2918.08781224938	1842.7778968171197	2818.728494917971	1369.7612919873118	1449.736444122851	1495.6124660426391	1491.546474072003	1440.5801712964391	14

支持维度、指标的交互式拖拽和实时计算

X5 X6 X7是维度，以交叉表的方式展现了分类汇总的统计信息

sum(X1) =

sum(X1) \* X5 \* X7

X6	1	2	4
X5			
French			
America			
Japan			
England			
French			
X7			
初中			
博士			
大学	1458.8054031759395		1140.402109314898
小学		1058.712064795208	1038.629464132992
研究生		1483.99432275955	1432.6874261400203
高中			
Totals	1458.8054031759395	3122.713687527563	3089.3158955703184
			2971.1646294348193
			1501.9210664498298
			1602

- Pearson相关系数

用来衡量两个数据集合是否在一条线上，它用来衡量变量间的线性关系。

本质上是去中心化后的余弦相似度

$$\rho_{x,y} = \frac{E(XY) - E(X)E(Y)}{\sqrt{E(X^2) - E^2(X)}\sqrt{E(Y^2) - E^2(Y)}}$$

- Spearman相关系数

用来衡量两个变量之间的关联程度与方向

$$r_s = \frac{l_{pq}}{\sqrt{l_{pp}l_{qq}}}$$

- 余弦相似度

通过计算两个向量的夹角余弦值来评价相似度

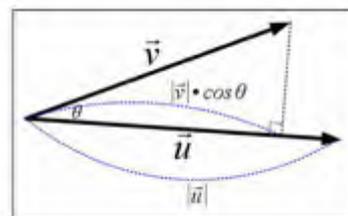
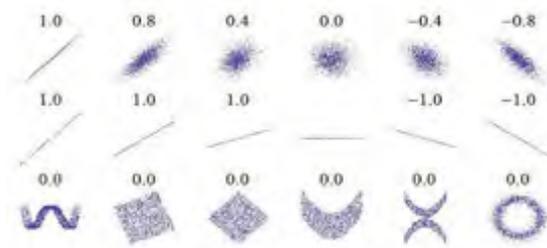
$$\cos \theta = \frac{a \cdot b}{\|a\| \cdot \|b\|}$$

- 欧氏距离

$$EUCLID = \sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

- Minkowski距离

$$MINKOWSKI(x, y) = \sqrt[p]{\sum_{i=1}^k |x_i - y_i|^p}$$



- 用途
  - 数据建模则需要快速挖掘出与目标相关度最高的特征
  - 建模前无法确定输入变量对目标变量的预测能力，把可能的基础变量和衍生变量放到模型中，但这些变量特别是衍生变量之间可能会存在相关性，会导致模型的多重共线性，从而造成模型整体预测能力的下降。
  - 例子
    - 保险业务
    - 保单号、被保险人、车牌号、**保费收入、起保日期、已赚保费**
- 解决方法
  - 主成分分析/卡方检验/岭回归等（无法解释）
  - 变量聚类
    - 采用相关系数，夹角余弦和列联系数来寻找反映元素之间亲疏关系的统计量，然后依据此把变量分为若干组。
    - 把高度相关的变量聚到一组，每一组内的变量之间信息重合度很高，互补性很弱，而组间的变量相关性则很低，信息重合度很弱，互补性很强。
    - 从每一组选择一两个最具有代表性的变量代表整个类别，参与建模。

## 论文数据

publication_time	title	abstract	author
2000_1期	主编致辞	-	-
2000_1期	大功率汽轮机低压调节阀的试验研究	介绍通过吹风试验研究,得到在所有工况下都...	王平子:[1]
2000_1期	四杆同源机构的构成特性和传动角特性	四杆曲线同源机构与原始机构的对应杆长具...	吴琛:[1],乔永杰:[1]
2000_1期	循环流化床锅炉分离器的研究	着重介绍国内外气固分离器的结构及使用情...	赵旺初:[1]
2000_1期	300/600MW燃煤电厂输煤控制系统设计及...	在总结上海工业自动化仪表研究所从事火力...	吴国伟:[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...	陈宾墨:[1],吕兆璧:[1],陆楚勋:[1],戴庆忠:[1]
2000_1期	论超临界火电机组的发展	简要论述了世界上超临界火电机组的最新发...	周海澜:[1],忻鹤龄:[1]
2000_1期	水轮机环列导水叶栅网格生成技术及应用	网格生成是流场数值模拟的一个关键环节,网...	刘焕明:[1]
2000_1期	水轮发电机临界失步问题的导纳分析及应用	从同步电机电势相量图与功角特性出发,用导...	卢敬:[1]

自相关

工作表:

选择变量:

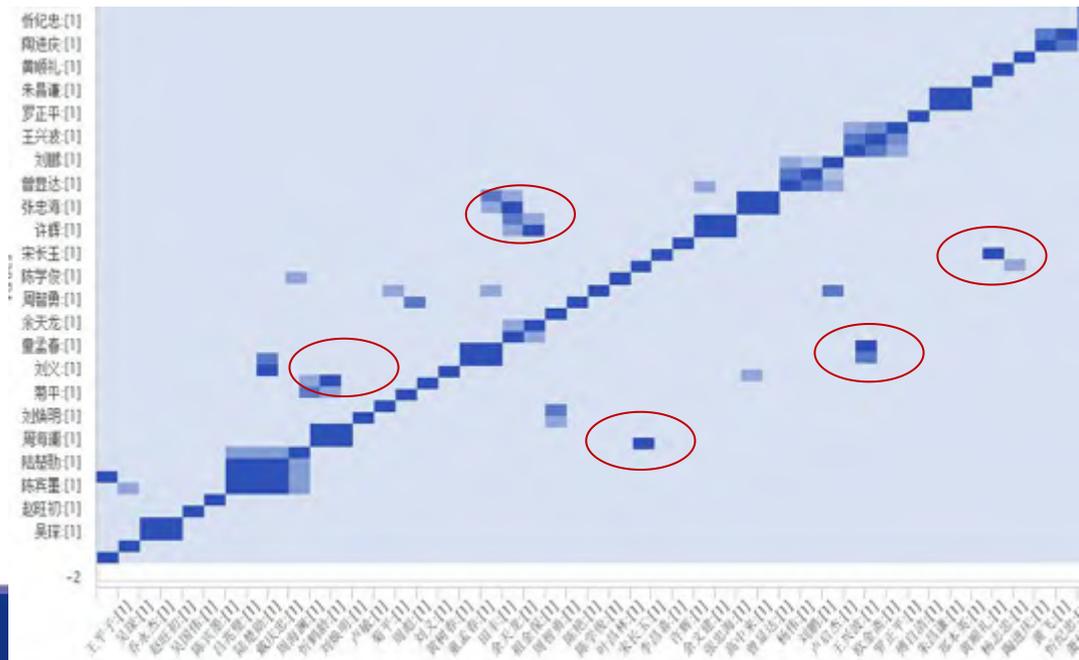
分隔符:

相关性方法:

计算方法:

取消

探索的是作者之间的合作关系  
相关性方法支持pearson、cos



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

	R1	R2	R3	R4	R5
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	

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

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

### 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

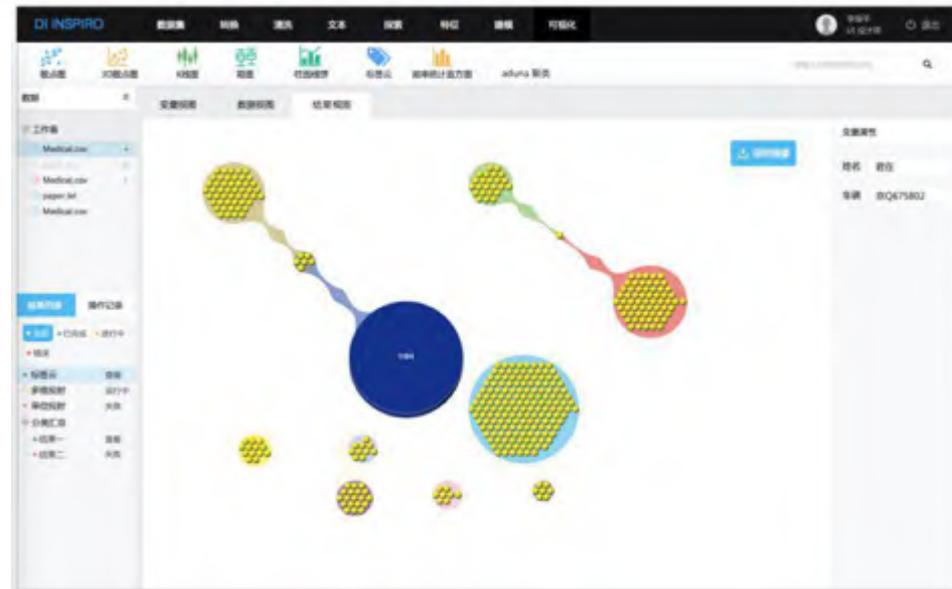
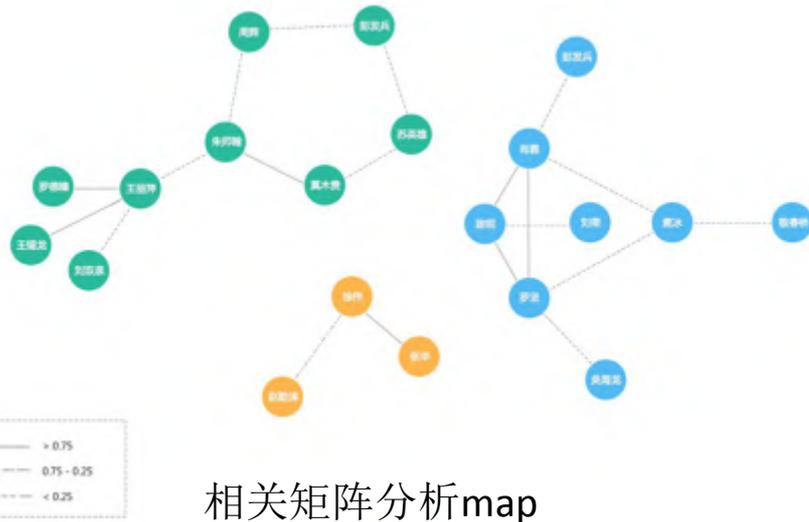
pearson相关性

cos相关性

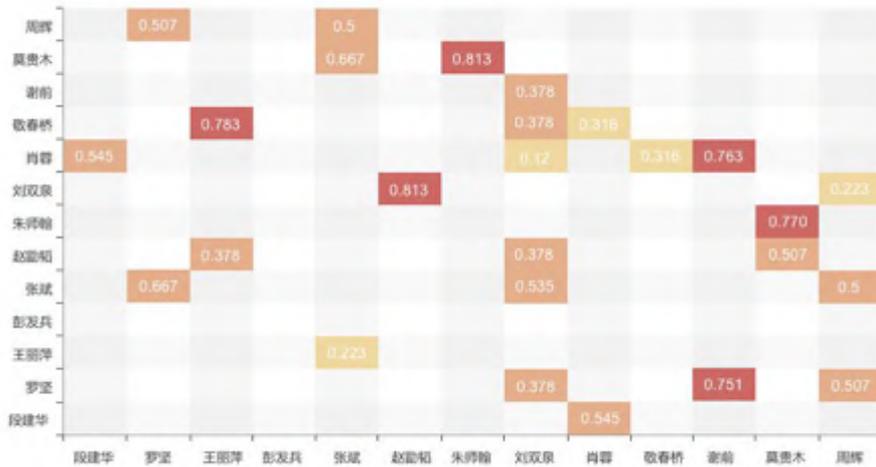
	谢前	敬春桥	罗坚	肖蓉	刘双泉
谢前	1	0	0.816	0.577	0.577
敬春桥	0	1	0	0	0
罗坚	0.816	0	1	0	0.707
肖蓉	0.577	0	0	1	0
刘双泉	0.577	0	0.707	0	1

共现相关性

	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	



Aduna图合作关系可视化



原始数据表

案件编号	人员	车牌	地点	金额
344573	周辉, 朱师翰	湘A65N90, 湘A1661K	桂花路	20000
344574	苏英雄, 莫贵木, 罗坚, 肖蓉, 谢前	湘A2KA19, 湘AT0101, 湘J7ZH85, 湘AL5S84, 湘AA0390	人民路	50000
344575	张斌, 刘荣	湘A632AK, 湘AUQ852	开元西路	100000
344576	罗德臻, 王丽萍	湘A9HT68, 湘A65N90	北斗路	70000
344577	吴海龙, 张华	湘AVB590, 湘AA1391	北斗路	10000

矩阵变换

	344573	344574	344575	344576	344577
湘ZG0069	0	0	0	0	0
湘A2KA19	0	1	0	0	0
湘AA0390	0	1	0	0	0
湘A65N90	1	0	0	1	0
湘AA1391	0	0	0	0	1

cos相关性

	湘ZG0069	湘A2KA19	湘AA0390	湘A65N90	湘AA1391
湘ZG0069	1	0	0	0	0
湘A2KA19	0	1	1	0	0
湘AA0390	0	1	1	0	0
湘A65N90	0	0	0	1	0
湘AA1391	0	0	0	0	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

	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

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
344564	彭发兵, 周辉, 苏英雄	湘A2ZB92, 湘B2HL12, 湘A2KA19	人民路	70000
344574	苏英雄, 莫贵木, 罗坚, 肖蓉, 谢前	湘A2KA19, 湘AT0101, 湘J7ZH85, 湘AL5S84, 湘AA0390	人民路	50000
344563	王丽萍, 刘双泉	湘A65N90, 湘A1661K	寿昌路	100000
344581	张华, 莫木贵, 许卫	湘AA1391, 湘AT0101, 湘ZG00069	板仓南路	100000

矩阵变换

	板仓南路	开元西路	寿昌路	人民路	湘江东路
湘ZG0069	3	0	0	0	0
湘A2KA19	0	0	0	2	0
湘AA0390	0	0	0	1	0
湘A65N90	0	0	1	0	1
湘AA1391	3	0	0	0	0

Cos互相关

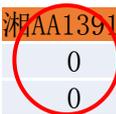
	湘ZG0069	湘A2KA19	湘AA0390	湘A65N90	湘AA1391
湘ZG0069	1	0	0	0	1
湘A2KA19	0	1	1	0	0
湘AA0390	0	1	1	0	0
湘A65N90	0	0	0	1	0
湘AA1391	1	0	0	0	1

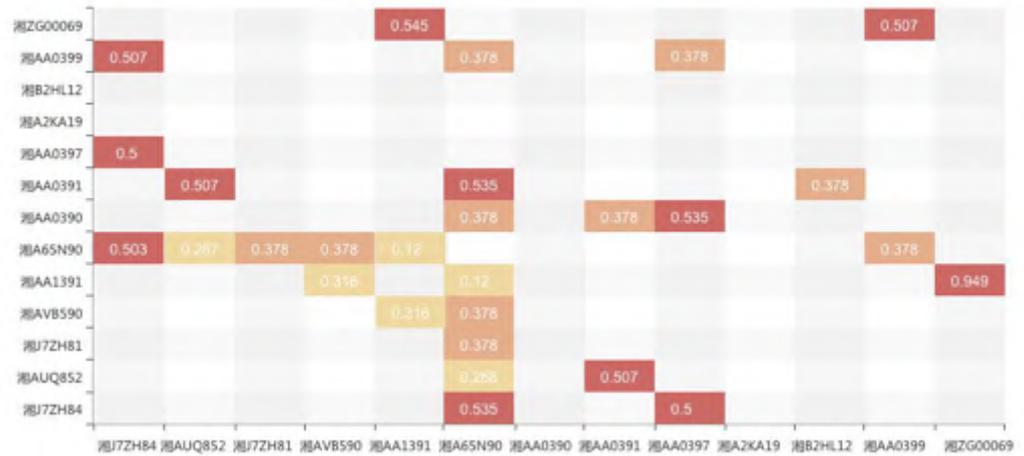
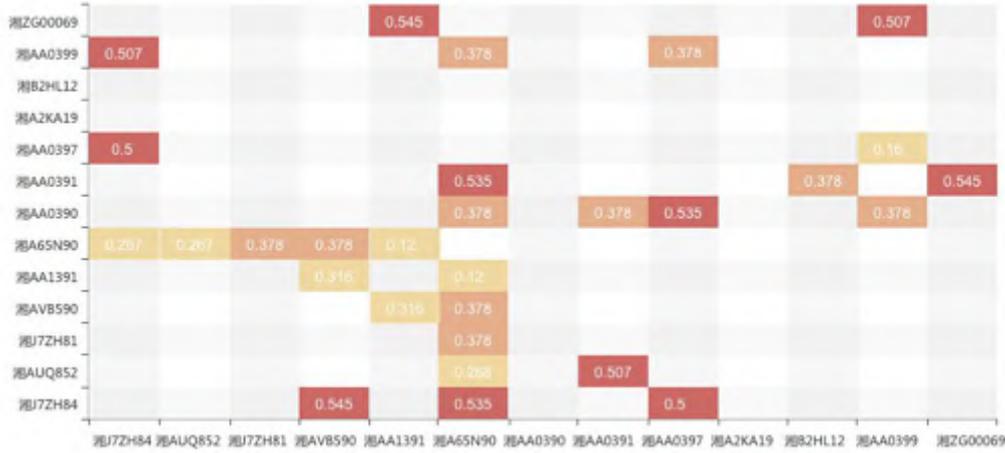
cos自相关

	湘ZG0069	湘A2KA19	湘AA0390	湘A65N90	湘AA1391
湘ZG0069	1	0	0	0	0
湘A2KA19	0	1	1	0	0
湘AA0390	0	1	1	0	0
湘A65N90	0	0	0	1	0
湘AA1391	0	0	0	0	0

	湘ZG0069	湘A2KA19	湘AA0390	湘A65N90	湘AA1391
湘ZG0069		0	0	0	3
湘A2KA19	0		1	0	0
湘AA0390	0	1		0	0
湘A65N90	0	0	0		0
湘AA1391	3	0	0	0	

共现互相关





- 实质
  - 二分图
- 很多应用领域
  - 医疗领域
    - 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

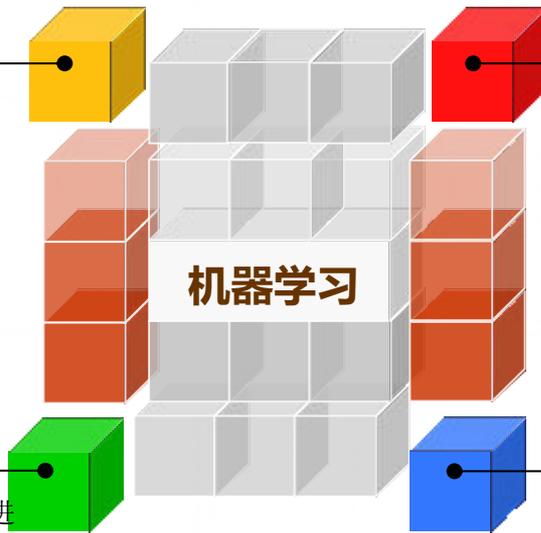


## ▪ 线性模型

- 支持线性回归和逻辑回归模型，支持任意程度嵌套效果、多项式和曲线效果
- 支持自动模型选择，提供丰富的模型诊断结果和自动模型评估。

## ▪ 模型管理

- 利用分区数据集对预测模型结果进行测试比对，快速确定最佳模型
- 利用提升表、ROC 图 表、协调统计和错误分类表，通过可视化评估和验证指标来验证结果
- 支持模型的持久化存储



## ▪ 决策树和随机森林

- 支持包含分类和连续特征的分类树和回归树
- 提供成本复杂性、C4.5 和减少误差的自动修剪并基于保留最优树
- 支持二分变量、名义变量和连续变量的随机森林、自动组合多个决策树预测单个目标
- 自动分配独立模型训练任务，自动智能调整参数设置确定最佳模型

## ▪ 神经网络和支持向量机

- 支持二分变量、名义变量和连续变量的神经网络
- 提供智能默认的大部分神经网络参数, 如激活和误差函数, 定制神经网络结构和加权
- 支持二分变量的支持向量机模型, 线性和多项式内核模型训练

模型建立—无需编写代码，在web页面可配置模型参数进行学习、评价、调优、存储

模型评价—在界面上即可对生成的模型进行评价。

模型预测—在界面上选取保持的模型及要进行预测的数据集即可生成预测数据

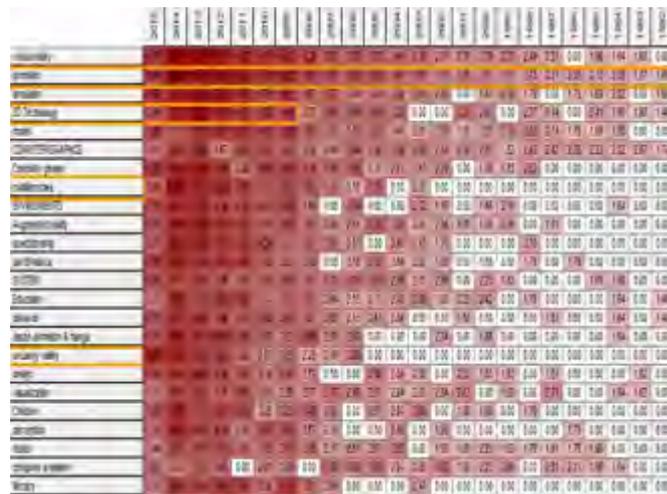
The screenshot displays the DI INSPIRO software interface for building a decision tree model. The main workspace shows a data table for a 'car' dataset with the following columns: 价格 (Price), 维修费用 (Repair Cost), 车门数量 (Number of Doors), 载人数 (Number of Passengers), 后备箱大小 (Trunk Size), 安全指数 (Safety Index), and 评价 (Evaluation). The table contains 10 rows of data, all with 'vhigh' for price and repair cost, and 'unacc' for evaluation.

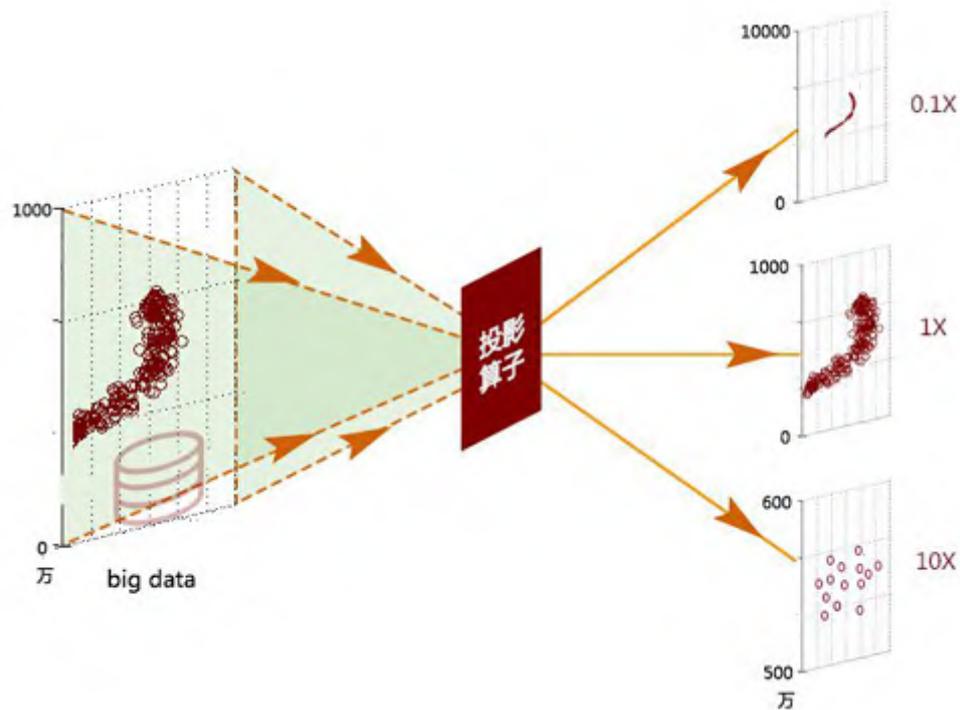
Overlaid on the table is a decision tree visualization. The root node is '价格' (Price), which splits into 'mid', 'low', and 'high' branches. The 'low' branch leads to a leaf node 'good'. The 'high' branch leads to a leaf node 'acc'. The 'mid' branch leads to a node '维修费用' (Repair Cost), which splits into 'mid' and 'high' branches. The 'high' branch leads to a leaf node 'acc'. The 'mid' branch leads to a node '后备箱大小' (Trunk Size), which splits into 'small' and 'big' branches. The 'small' branch leads to a leaf node 'acc'. The 'big' branch leads to a node '安全指数' (Safety Index), which splits into 'low' and 'high' branches. The 'low' branch leads to a leaf node 'acc', and the 'high' branch leads to a leaf node 'vgood'.

The right sidebar shows the model parameters for the decision tree: 操作名称: 决策树 (Operation Name: Decision Tree) and 参数: ("delimiter": ",", "field": "X5", "support": 0.001). Below the parameters is a '保存截图' (Save Screenshot) button.

在真实汽车数据集上利用决策树发现用户对车的评价

- 散点图
  - 2D、3D散点图
    - 快速的对大规模数据生成非失真的可视化图形，并支持无级缩放
- 聚类图
  - 直观查看数据的聚类结果
- 热力图
  - 相关性分析结果展示
- 基本图形
  - 点线柱图
  - 箱图、直方图—方便了解数据的总体分布



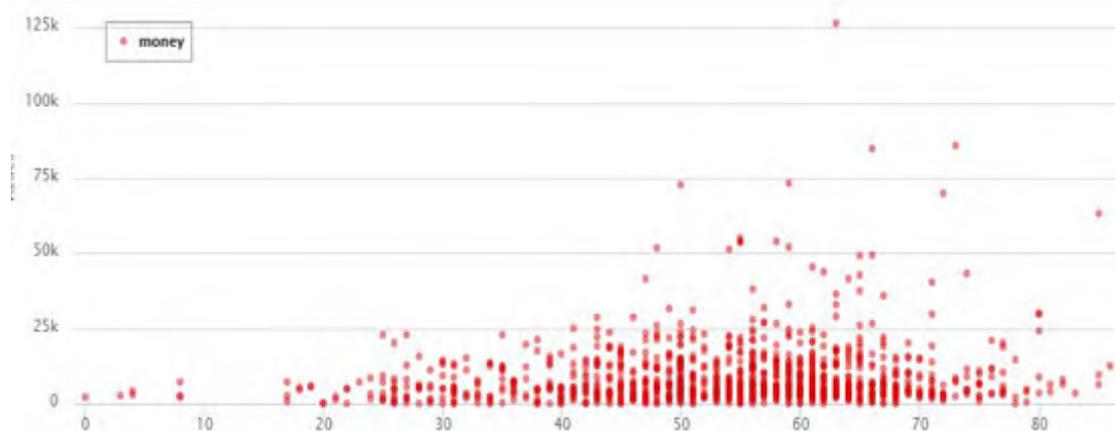


- 基于保持数据概率分布不变的思想
- 将原始数据根据缩放级别和距离远近将原始数据映射成特定显示区域的矩阵

## 医疗数据（真实数据）

年龄	性别	疾病种类	就诊日期	医疗费用
53	男	手术后恶性肿瘤化学治疗	2013/11/28	5398.08
44	女	恶性肿瘤维持性化学治疗	2013/11/27	5234.93
60	男	恶性肿瘤维持性化学治疗	2013/11/27	4680.67
61	男	肺恶性肿瘤	2013/11/25	4237.37
58	女	肺恶性肿瘤	2013/11/29	3714.02
73	男	慢性肾脏病5期	2013/11/28	7380.61
78	女	肺炎	2013/11/16	21720.56

患者年龄和医疗费用的相对分布？

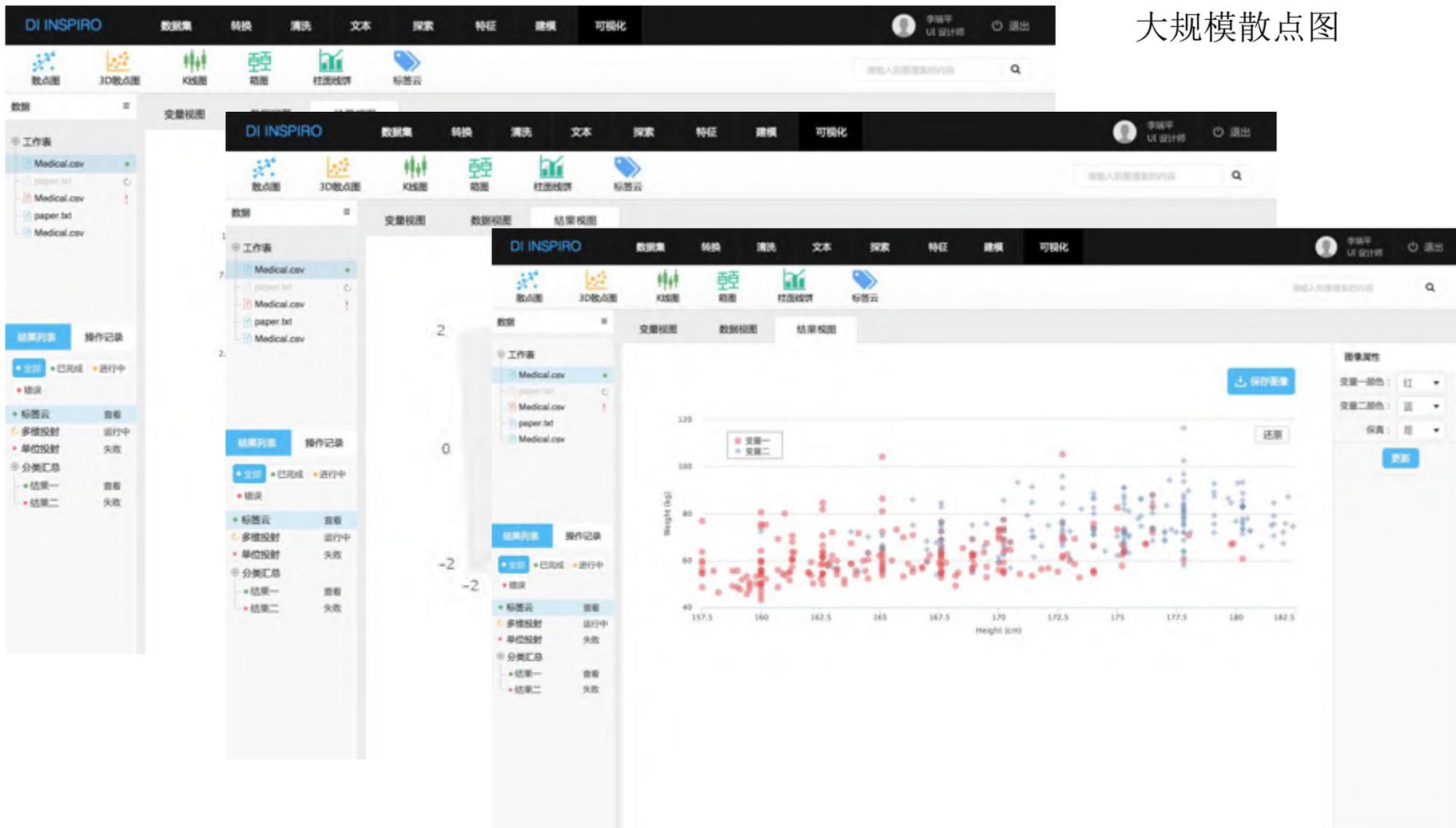


X轴：年龄，  
 Y轴：医疗费用  
 是否保真：系统会自动选用合适的粒度展示两者之间的关系。

上图反映了：

1. 患者主要集中在40岁~50岁之间
2. 医疗费用会在40~50岁之间出现大幅度波动，且总体趋势升高

## 大规模散点图



## • DI设计目标

- 旨在为普通数据分析人员提供直观易用的大数据分析处理工具
- 无需掌握复杂的分布式集群和编程，而是直接利用web操纵大规模数据集

## • 特性

- 图形化Spark大数据分析系统
  - 基于SPARK实现的图形化交互式大数据分析系统
- 插件式架构设计
  - 插件式设计，支持分析算子的动态加载和删除，具有高度扩展性
- 分析功能服务化设计
  - 数据分析功能提供Restful API接口，便于与其它系统集成
- 自适应的交互响应速度
  - 利用数据集分析历史，根据待分析文件规模对所需计算资源估计并动态分配执行资源，确保交互响应速度



- 开发状态
  - 主要开发语言：python、scala
  - 核心代码行数：3500行
- 项目数据
  - 从2016年1月提出构想到开发出第一版，耗时8个月，目前还没有完全稳定
  - 2017年初计划开放公测和免费使用
  - 等代码优化到一定程度后，计划将项目开源

谢谢！

[www.datainsight.tech](http://www.datainsight.tech)

**BDTC** 2016 中国大数据技术大会  
Big Data Technology Conference 2016