

Smartbi V10.5.12-数据挖掘

注意：（新特性列表中：+表示新增；^表示增强）

- +【挖掘】通过产品帮助功能创建节点使用示例
- +【挖掘】新增移动平均算法节点
- +【挖掘】新增指数平滑算法节点
- +【挖掘】增加正则表达式处理节点
- +【挖掘】节点能够自动连接
- +【挖掘】增加在线节点开发功能
- +【挖掘】增加日期时间节点
- +【挖掘】监控建模增加停止功能
- +【挖掘】自定义帮助指引系统
- +【挖掘】支持对源和目标节点运行结果进行数据透视
- +【挖掘】节点支持自动布局
- ^【挖掘】全表统计节点支持输出结果
- ^【挖掘】评估节点支持输出结果
- ^【挖掘】模型的预测中把预测概率输出成字段
- ^【挖掘】数据清理节点合并
- ^【挖掘】过滤和行选择节点合并
- ^【挖掘】派生列增强函数帮助提示
- ^【挖掘】SQL脚本输入表字段显示优化
- ^【挖掘】提供节点iframe扩展配置组件

具体改进点如下：

新增	增强
+【挖掘】通过产品帮助功能创建节点使用示例	^【挖掘】全表统计节点支持输出结果
+【挖掘】新增移动平均算法节点	^【挖掘】评估节点支持输出结果
+【挖掘】新增指数平滑算法节点	^【挖掘】模型的预测中把预测概率输出成字段
+【挖掘】增加正则表达式处理节点	^【挖掘】数据清理节点合并
+【挖掘】节点能够自动连接	^【挖掘】过滤和行选择节点合并
+【挖掘】增加在线节点开发功能	^【挖掘】派生列增强函数帮助提示
+【挖掘】增加日期时间节点	^【挖掘】SQL脚本输入表字段显示优化
+【挖掘】监控建模增加停止功能	^【挖掘】提供节点iframe扩展配置组件
+【挖掘】自定义帮助指引系统	
+【挖掘】支持对源和目标节点运行结果进行数据透视	
+【挖掘】节点支持自动布局	

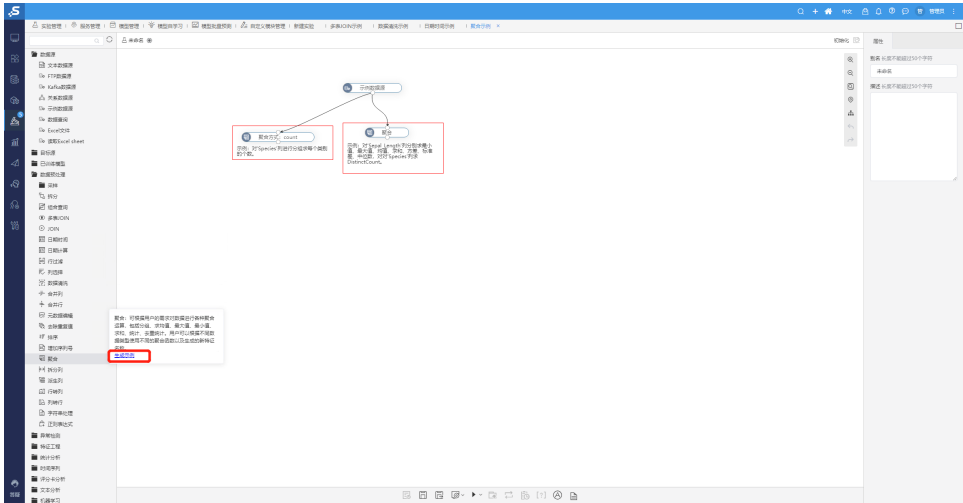
+【挖掘】通过产品帮助功能创建节点使用示例

背景介绍

为了提高产品的易用性，帮助用户快速了解某一节点的功能。新版本在节点的帮助信息后面添加创建节点使用示例的入口，针对常用的数据预处理节点和统计分析节点支持一键生成示例，帮助用户快速了解节点的使用方法。

功能简介

用户在鼠标移动到某一个节点之后，出现悬浮提示后会出现“生成示例”链接“，点击之后会在新窗口中生成对应的示例ETL。

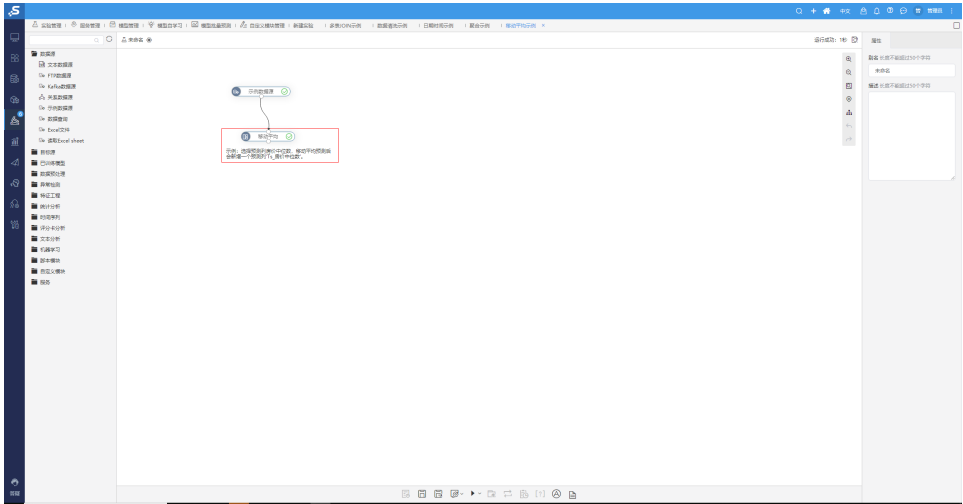


+【挖掘】新增移动平均算法节点

背景介绍

移动平均法作为经典的基础时间序列算法，可以消除部分因素的影响，显示总体走势。

功能简介



参考文档

详情请参考 [数据挖掘 - 移动平均](#)。

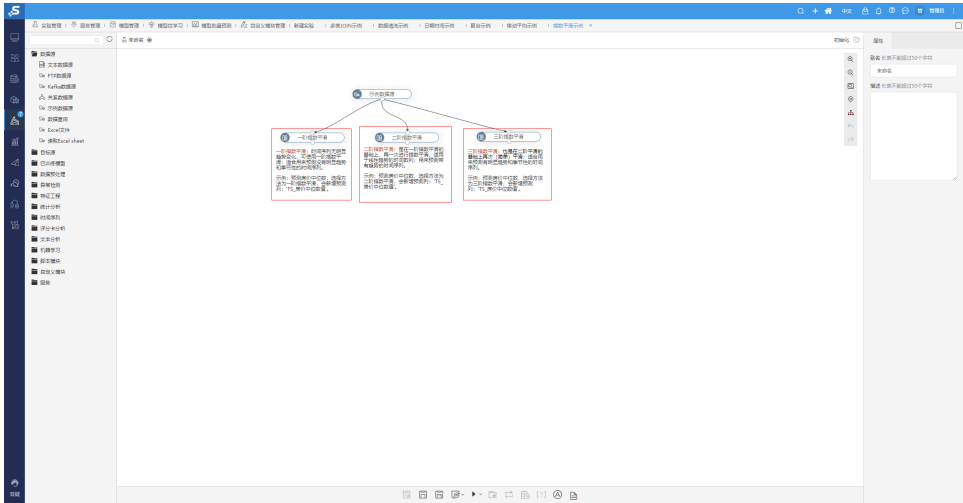
+【挖掘】新增指数平滑算法节点

背景介绍

时间序列算法中常用的算法有指数平滑法，其中指数平滑常用的几种形式有（Brown）一次指数平滑、二次指数平滑、三次指数平滑、以及（Holt）双参数、（Holt-Winters）季节性指数平滑法。

功能简介

本节点实现的是布朗（Brown）的一次、二次、三次指数平滑。



参考文档

详情请参考 [数据挖掘 - 指数平滑](#)。

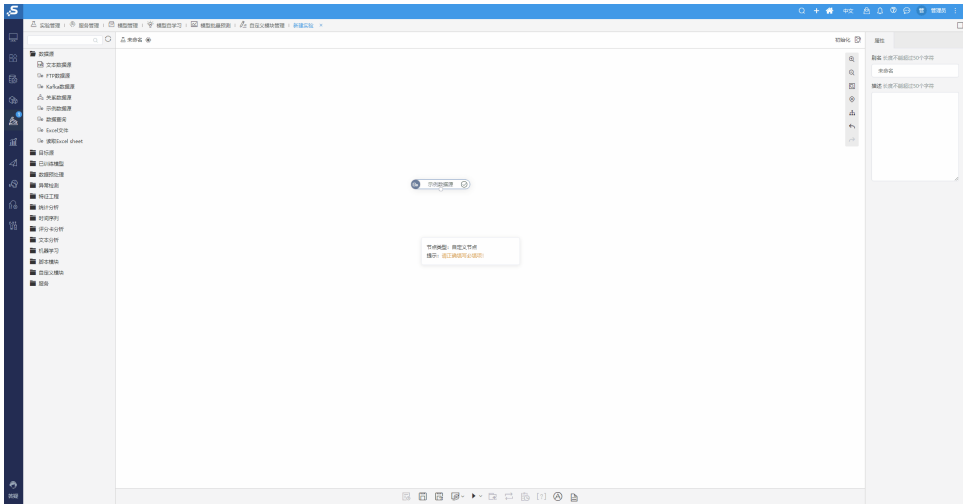
+【挖掘】增加正则表达式处理节点

背景介绍

在处理文本过程中，经常会使用正则表达式去处理，目前ETL正则处理文本功能较弱，多个ETL项目中均有实施人员提到希望增强该功能。

功能简介

用正则表达式语法的强大模式匹配能力，对字符串数据进行解析、匹配或替换。



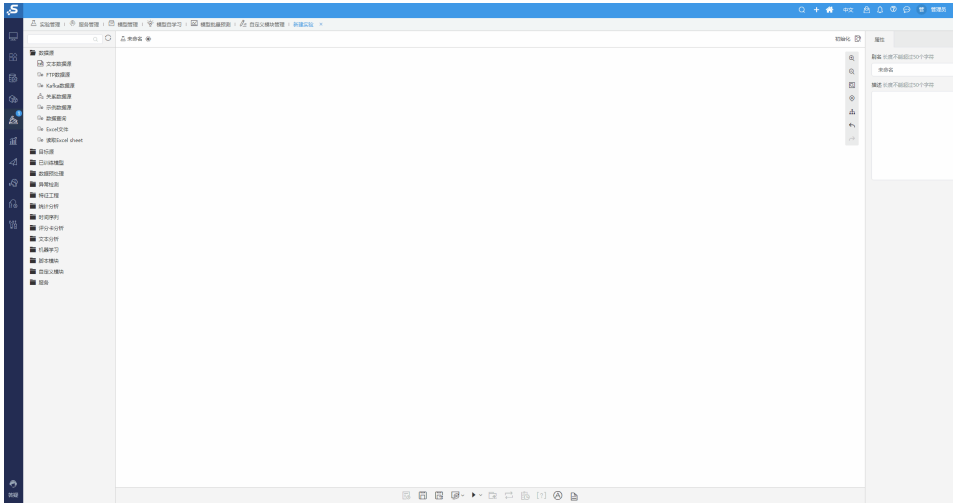
参考文档

详情请参考 [数据挖掘-正则表达式](#)。

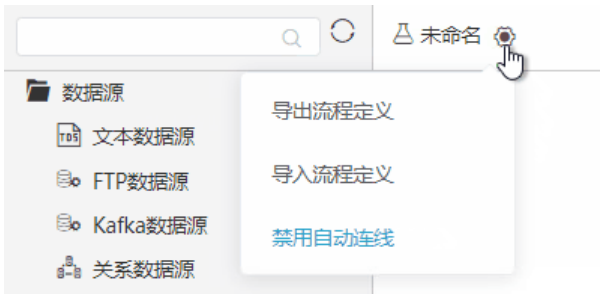
+【挖掘】节点能够自动连接

功能简介

新版本中添加了节点自动布局功能，能够让用户更加专注于建模。



当不需要自动连线功能时，可以在当前ETL中的如下入口禁用。或者可以配置系统选项：DISABLE_AUTO_CONNECTION=true



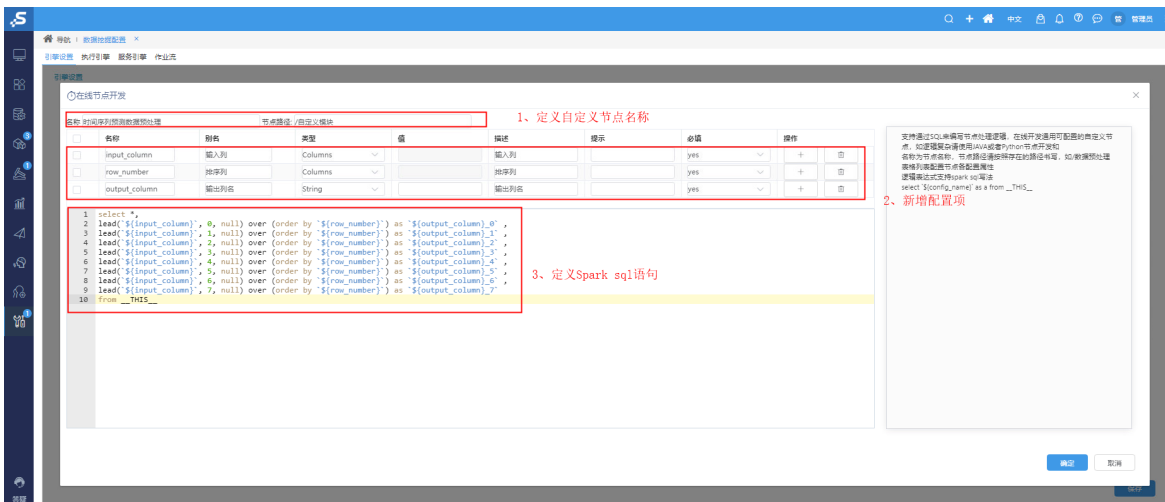
+【挖掘】增加在线节点开发功能

背景介绍

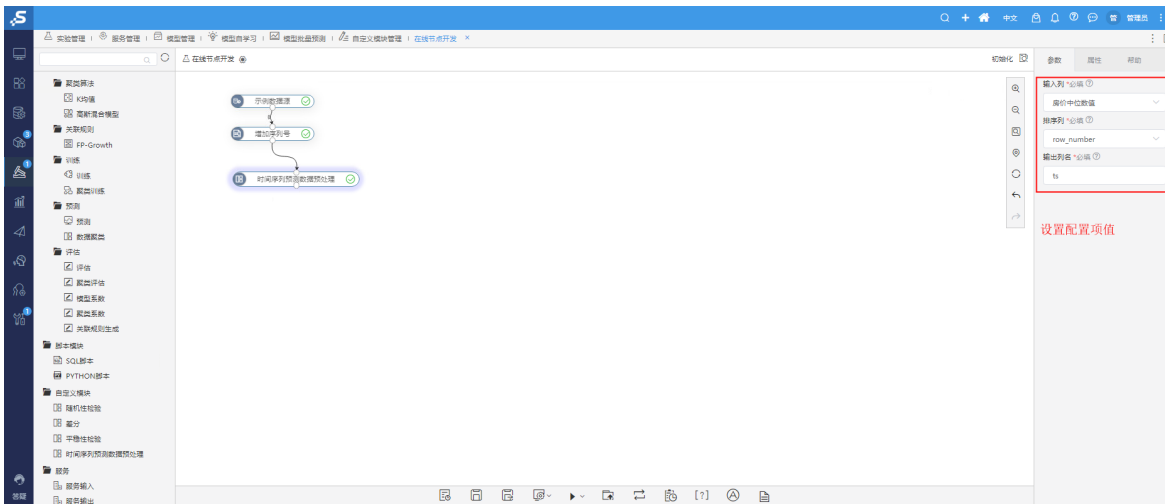
新版本中增加了在线节点开发功能，可以将部分可复用的SQL封装成预制节点进行复用，开发完成后和产品中预制节点使用方式一样。该功能入口在系统运维界面中，在/系统运维/数据挖掘配置/引擎配置/在线节点开发下，能够支持基础的配置选项，并且通过SQL来实现简单的数据处理逻辑。

功能简介

节点在线开发编辑界面如下：



使用在线自定义节点：



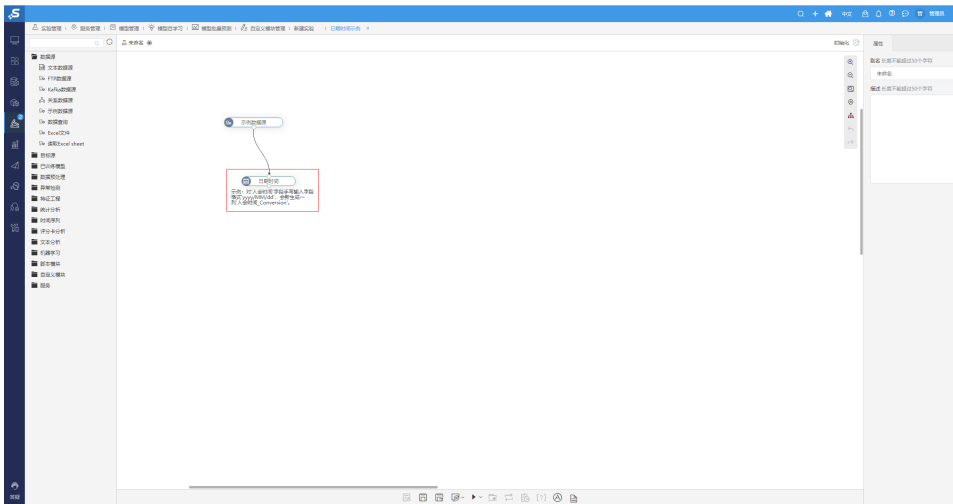
+【挖掘】增加日期时间节点

背景介绍

用户在录入日期类型数据时，经常会遇到日期格式不统一的问题。

功能简介

新版本为了方便用户统一日期时间格式，提供了日期类型字符转换功能。



详情请参考 [数据挖掘-日期时间](#)。

+【挖掘】监控建模增加停止功能

背景介绍

异常情况下，引擎中运行的实验节点会存在无法停止运行的现象。因此，新版本在实验监控界面中提供了停止运行的功能。

功能简介

可以在系统监控下的实验监控下看到具体的引擎中运行的实验的记录，根据需要将其停止。



+【挖掘】自定义帮助指引系统

背景介绍

产品的节点通常需要，帮助用户快速了解某一节点的功能。新版本在节点的帮助信息后面添加创建节点使用示例的入口，针对常用的数据预处理节点和统计分析节点支持一键生成示例，帮助用户快速了解节点的使用方法。

功能简介

用户在鼠标移动到某一个节点之后，出现悬浮提示后会出现对应的自定义帮助信息。

功能简介



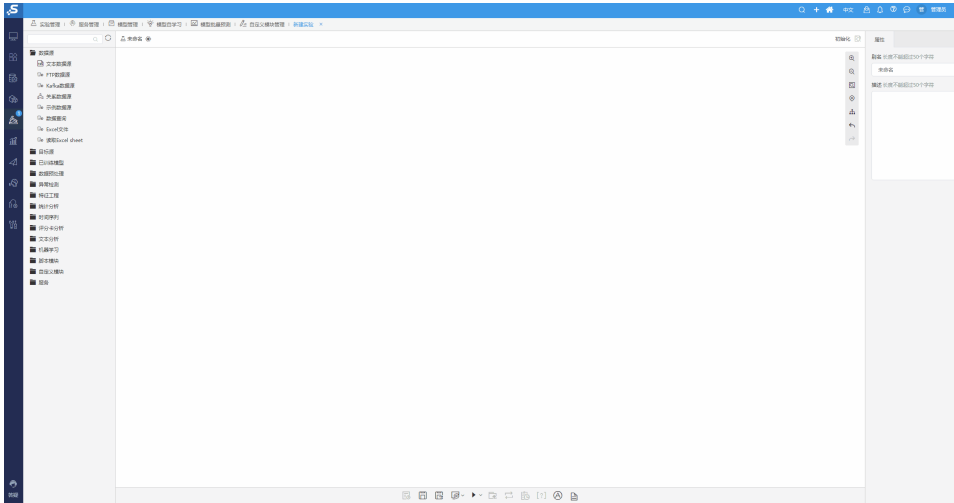
+【挖掘】支持对源和目标节点运行结果进行数据透视

背景介绍

旧版本的数据挖掘实验的可视化实现方式不足，做数据探索需要导出数据到BI，或者用Python去做。所以增加数据透视功能，右键选择运行好的节点后能点击数据透视，进入一个透视分析界面，对该节点的运行结果进行透视分析。透视分析界面沿用BI分析展现模块的透视分析功能。

功能简介

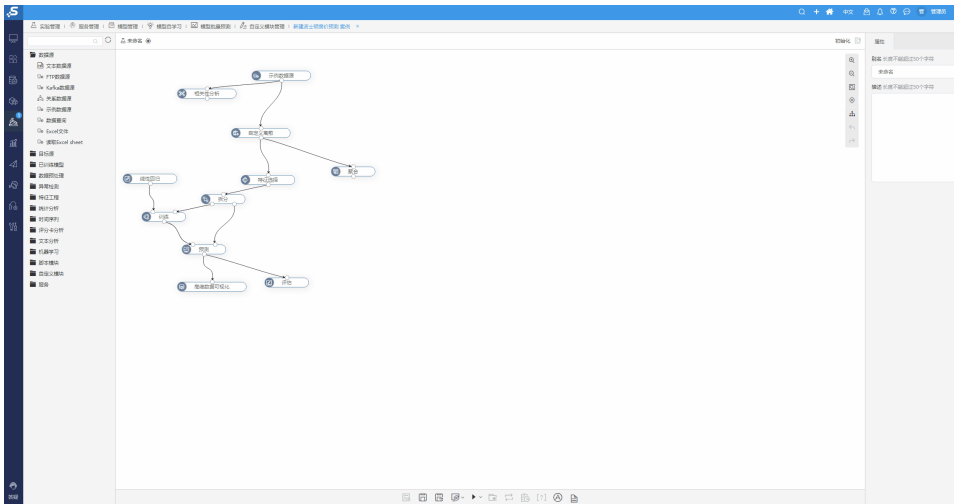
用户可以在关系数据源和关系目标源节点上右键，点击数据透视，打开一个透视分析界面，直接对数据源中的数据进行分析。



+【挖掘】节点支持自动布局

功能简介

新版本中添加了节点自动布局功能，能够让用户更加专注于建模。



^【挖掘】全表统计节点支持输出结果

功能简介

旧版本的全表统计节点不支持输出统计结果，新版本中对该节点进行了改进，将各统计结果转换成Dataset输出，方便后续监控和图表展现分析等。

Col_Name	Row_Count	min	max	mean	Lower_Quartile	media	Upper_Quartile	Distinct_Count
编号	10000	1	10000	5000.5	2500	5000	7500	10000
客户ID	10000	15565701	15815690	1.56909405694E7	15628420	15690731	15753215	10000
信用分	10000	350	850	650.5288	584	652	717	460
年龄	10000	18	92	38.9218	32	37	44	70
年限	10000	0	10	5.0128	3	5	7	11
账户余额	10000	0.0	250898.09	76485.88928799996	0.0	97157.96	127638.35	6382
月均AUM	10000	3323.6	286391.9	90579.54530999999	31283.0	105084.0	134986.0	9954
购买理财产品数量	10000	1	4	1.5302	1	1	2	4
是否持有信用卡	10000	0	1	0.7055	0	1	1	2
是否活跃	10000	0	1	0.5151	0	1	1	2
年收入	10000	11.58	199992.48	100090.23988100006	50972.6	100183.05	149381.32	9999
是否流失	10000	0	1	0.2037	0	0	0	2

【挖掘】评估节点支持输出结果

功能简介

在挖掘实验的建模流程中，模型训练以后进行验证时通常需要评估节点对模型进行评估。目前的评估节点只支持查看评估指标，不支持输出评估结果，在有大量模型训练的情况下难以把评估结果落地进行分析。所以需要在评估节点增加输出端口，输出评估指标结果。

分类算法评估指标：

key	value
algorithm.initMode	k-means==初始化模式
algorithm.seed	2==随机种子
algorithm.sol	1.0E-4==收敛阈值
algorithm.maxiter	50==迭代次数
features.name	主营业务净利润Normalized,资产负债率Normalized,销售毛利率Normalized,资产净利润Normalized,销售净利润Normalized
cluster	("0":["values":["-0.2062488961224452,0.1423015661752583,-0.31290757237235517,-0.12036746257206118,-0.210763714...
classes	("["schema":"d","hashCode":0,"fields":{"name":"prediction","dataType":"numeric","integral":0,"ordering":0,"exactNumeric...
silhouette	0.63977
SSE	11926.57105
CHI	1.90856
algorithm.normalization	("method":"StandardScaler","pnorm":2,"std":false,"mean":false)==归一化
algorithm.k	4==K值

回归算法评价指标：

key	value
rmse	0.50911
r2	0.8111
adjusted r2	0.81091
mae	0.37373
mse	0.2592

聚类算法评价指标：

key	value
algorithm.initMode	k-means==初始化模式
algorithm.seed	2==随机种子
algorithm.sol	1.0E-4==收敛阈值
algorithm.maxIter	50==迭代次数
features.name	主营业务净利润Normalized, 资产负债率Normalized, 销售毛利率Normalized, 资产净利润Normalized, 销售净利润Normalized
cluster	["0":["values":["-0.2062488961224452,0.1423015661752583,-0.31290757237235517,-0.12036746257206118,-0.210763714...
classes	[[{"schema":"_hashCode":0,"fields":{"name":"prediction","dataType":"numeric","integral":0,"ordering":0,"exactNumeric...
silhouette	0.63977
SSE	11926.57105
CH	1.90856
algorithm.normalization	("method":"StandardScaler","pnorm":"2","std":"false","mean":"false")==归一化
algorithm.k	4==K值

【挖掘】模型的预测中把预测概率输出成字段

功能简介

数据挖掘中的分类模型完成预测时，预测输出结果中包含了预测的概率，但目前这个概率字段格式不好，无法直接用于后续节点中进行处理，导致在项目中需要另写脚本提取相关预测概率的数据。所以需要支持把概率输出成如数值型的字段，更方便对预测结果进行分析处理。

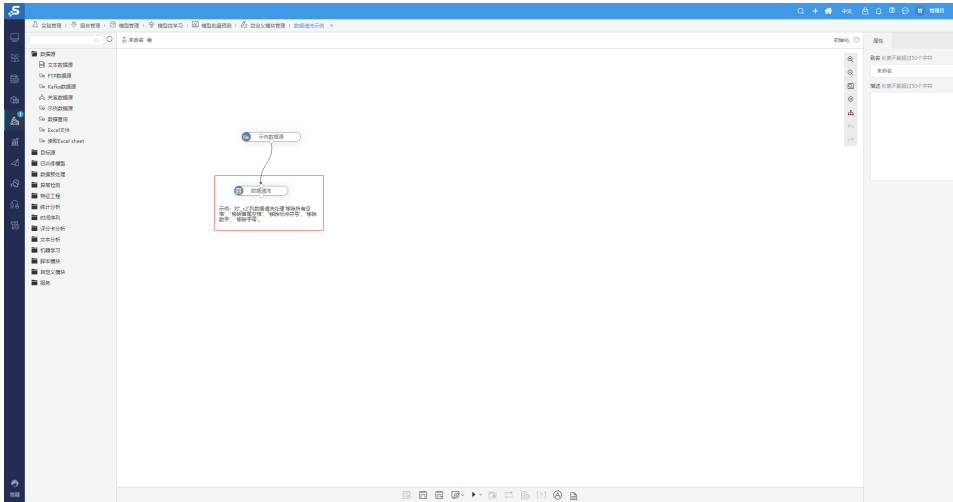
probability	probability_0	probability_1	probability_2	probability_3	probability_4	probability_5	
[3.868546872736931E-9,0.9999882732818455,...	3.868546872736931E-9	0.9999882732818455	3.155956714254193E-9	8.973944220181147E-6	3.0533686530177544E-10	2.6524570584346917E-6	9.
[3.868546872736931E-9,0.9999882732818455,...	3.868546872736931E-9	0.9999882732818455	3.155956714254193E-9	8.973944220181147E-6	3.0533686530177544E-10	2.6524570584346917E-6	9.
[0.0012767866651881732,4.136273384540328...	0.0012767866651881732	4.136273384540328E-4	0.011620101188780036	0.002978593473710343	0.11787336190884666	0.010406000238996363	0.0
[4.457082750190825E-5,8.042013605604958E...	4.457082750190825E-5	8.042013605604958E-6	3.277805309685768E-5	4.9319315165899116E-5	9.85231160754463E-4	0.015069393997299797	0.0
[1.5818079747150615E-5,5.2861375387163397...	1.5818079747150615E-5	2.8613753871633973E-6	4.47847200960279E-6	2.0936887986554473E-5	9.9837009002287E-5	0.011443113475670934	0.0
[1.3808339616809233E-7,6.590283410759679...	1.3808339616809233E-7	6.590283410759679E-7	0.9998825278746927	9.962968865749604E-5	1.419880584234422E-7	2.2220018041724596E-6	1.4
[4.206255237211825E-6,6.82572314159652...	4.206255237211825E-6	6.82572314159652E-6	1.9411112442117786E-6	1.52048984299023E-5	8.255228151329331E-8	0.0025413836855004753	0.0
[6.893801173346397E-6,1.9602604903233702...	6.893801173346397E-6	1.9602604903233702E-6	6.814469066153104E-5	1.3744064545571095E-5	6.038901653393605E-8	0.0017846489729278802	0
[1.3310193722239072E-8,0.999894515566592...	1.3310193722239072E-8	0.9998945155665924	3.04776886251785E-9	1.095877698671578E-6	6.985056992789256E-10	1.0373615434135084E-4	6.
[1.325449874643394E-9,0.9999968502292432...	1.325449874643394E-9	0.9999968502292432	9.550048507061064E-10	2.4717522352494705E-6	1.3489384223906163E-10	6.469798140402787E-7	2.8
[3.868546872736931E-9,0.9999882732818455...	3.868546872736931E-9	0.9999882732818455	3.155956714254193E-9	8.973944220181147E-6	3.0533686530177544E-10	2.6524570584346917E-6	9.
[2.8568448632511486E-8,3.389269670434608...	2.8568448632511486E-8	3.389269670434608E-8	3.6769638753371366E-9	6.338483246858044E-7	1.4456015662217E-11	1.6218316044219465E-6	0
[4.332408782865055E-7,5.982660931493343...	4.332408782865055E-7	5.982660931493343E-6	8.135316317479123E-7	4.6532778383836627E-7	1.5357650840139042E-9	3.9303737817294145E-5	0
[6.893801173346397E-6,1.9602604903233702...	6.893801173346397E-6	1.9602604903233702E-6	6.814469066153104E-5	1.3744064545571095E-5	6.038901653393605E-8	0.0017846489729278802	0
[3.868546872736931E-9,0.9999882732818455...	3.868546872736931E-9	0.9999882732818455	3.155956714254193E-9	8.973944220181147E-6	3.0533686530177544E-10	2.6524570584346917E-6	9.

^【挖掘】数据清理节点合并

功能简介

新版本的数据清洗合并了旧版本中的空值处理、值替换、数据清理节点。通过该节点，可以实现以下几点功能：

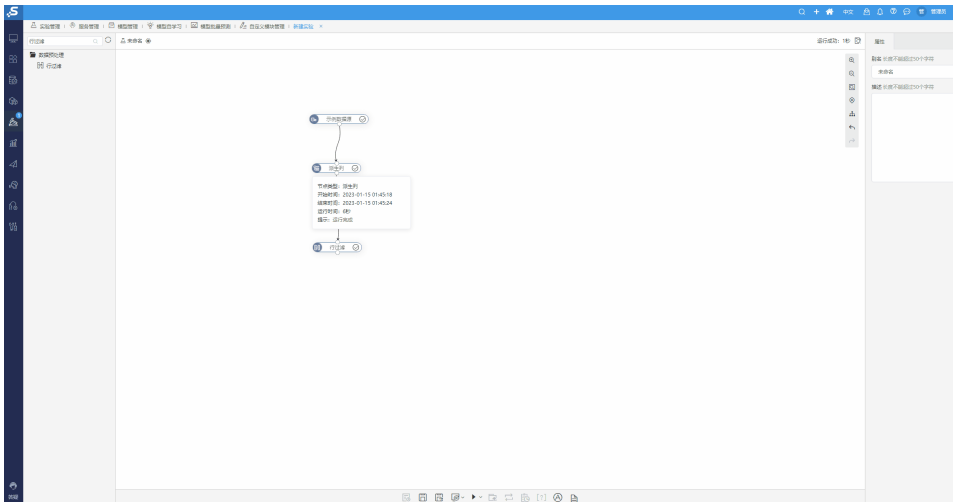
- （1）空值替换为均值、最大频数或者用户自定义的值等，实现空值的填充或者过滤；
- （2）移除字符串中空格、标点符号、字母、数字等不必要的字符，或设置大小写方式。



^【挖掘】过滤和行选择节点合并

功能简介

新版本的过滤节点，整合了旧的过滤节点和旧的行选择节点。提供了两种类型的筛选器。基本筛选器可以根据用户需求设置不同的筛选或者删除条件，选择不同数量的行；自定义筛选器通过写SQL语句(片段)，对数据按照过滤表达式进行筛选。



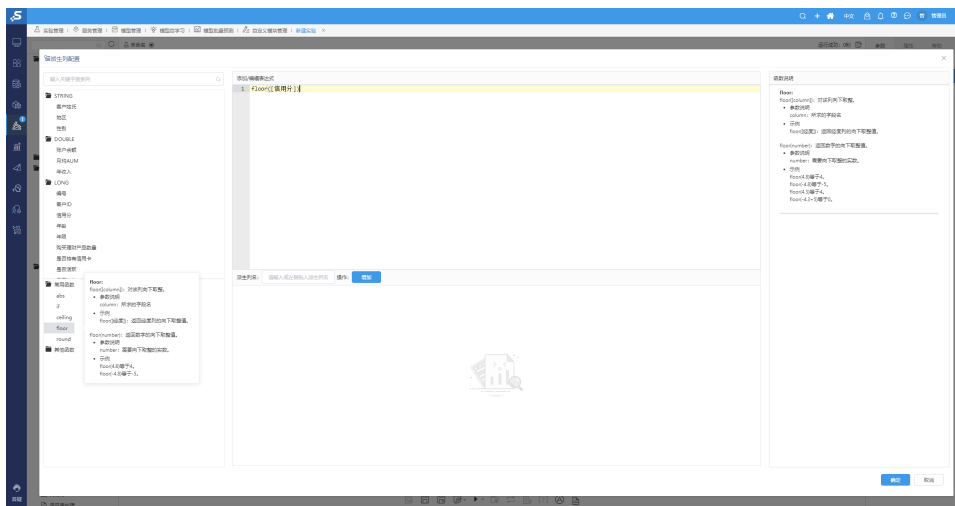
^【挖掘】派生列增强函数帮助提示

背景介绍

旧版本的派生列节点编辑界面中没有任何的指引，业务用户必须知道有这个函数，才能够写表达式，而且还没有语法规则和示例。

功能简介

新版本中对左侧函数资源树进行了优化，将用户常用的函数提前，并在鼠标悬浮时，以及左侧面板中显示对应的函数说明信息。



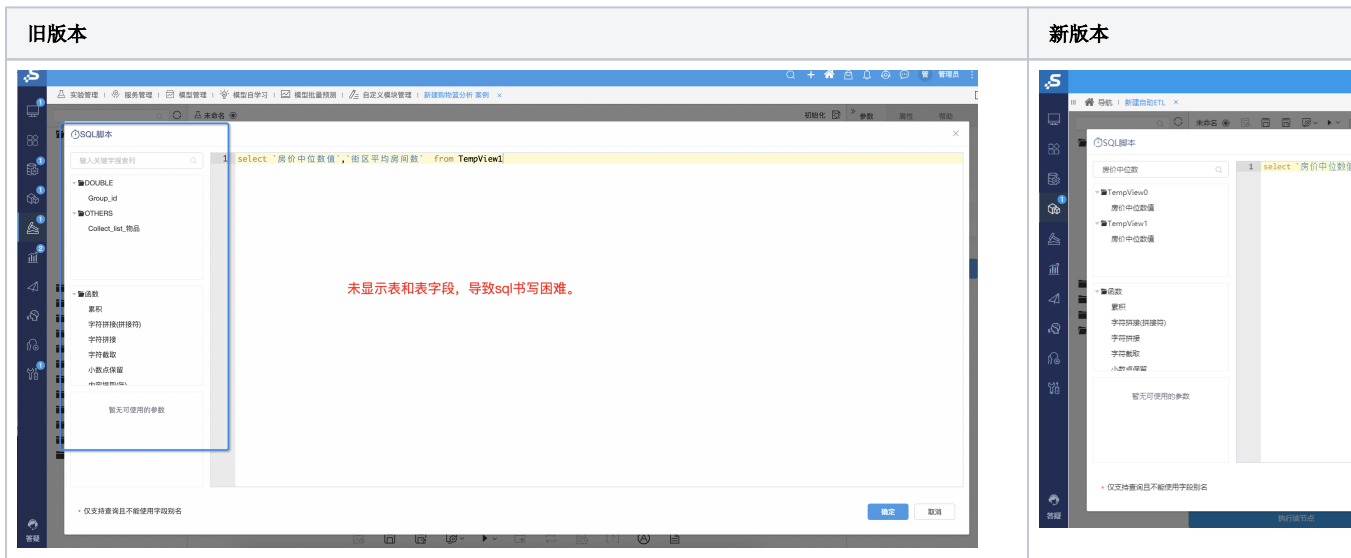
^【挖掘】SQL脚本输入表字段显示优化

背景介绍

SQL脚本节点主要通过执行Spark SQL语句，实现数据的查询。旧版本的SQL脚本节点中，没有显示节点相关的表及其字段，导致SQL书写困难。

功能简介

新版本中显示节点相关的表及其字段，用户可以直接拖拽表和字段，方便用户书写Spark SQL相关语句。



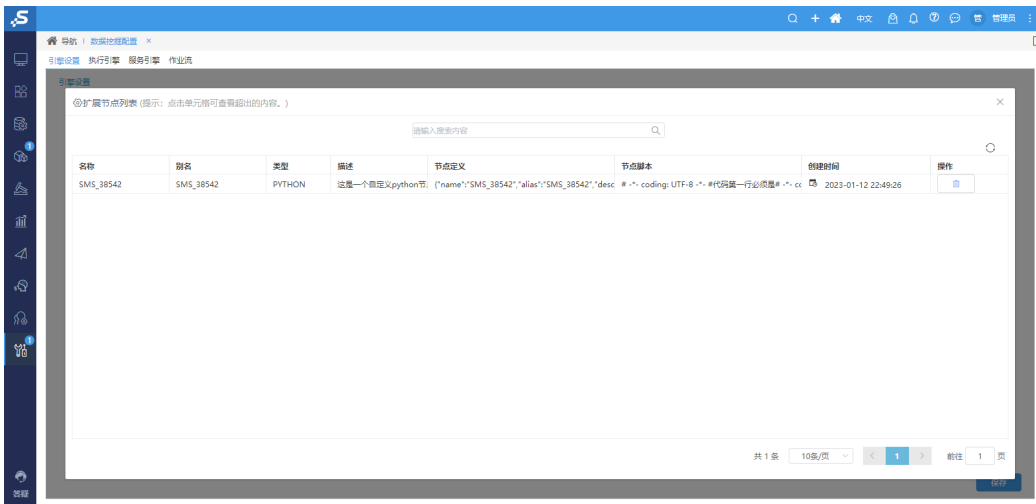
^【挖掘】提供节点iframe扩展配置组件

背景介绍

项目上有时需要一些特殊的配置项来进行自定义节点（自定义Java节点和自定义Python节点）开发，目前产品内置的配置项控件无法满足。

功能简介

因此，新版本中需要提供iframe扩展配置项，可以通过V的机制进行扩展开发，满足个性化需求。



参考文档

详情请参见：[如何自定义java节点。](#)