

蚂蚁平台用户指南

机器学习平台

内容

目录

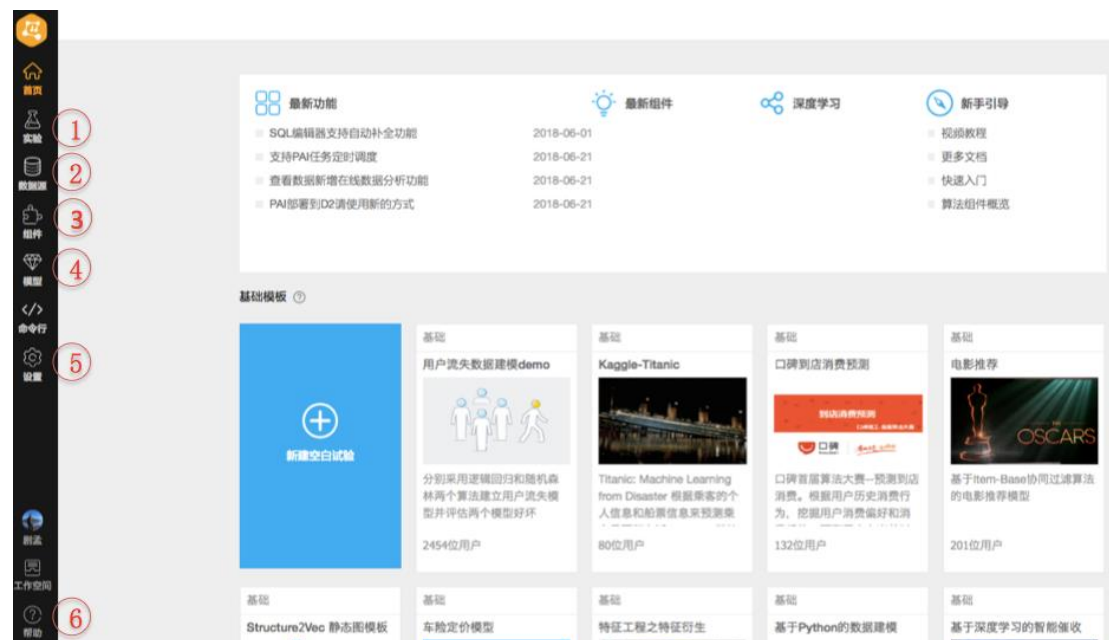
1. 蚂蚁机器学习平台介绍.....	3
1.1. 蚂蚁机器学习平台首页.....	3
2. 机器学习的一般过程.....	4
2.1. 新建实验	4
2.2. 读取数据表	7
2.3. 数据探索	11
2.3.1. 全表统计	11
2.4. 数据预处理和统计分析.....	15
2.4.1. 缺失值填充	15
2.4.2. 数据拆分	17
2.4.3. 离散值特征分析.....	20
2.5. 算法建模	23
2.5.1. 逻辑回归二分类.....	23
2.5.2. 随机森林.....	28
2.6. 模型预测及评估.....	32
2.6.1. 查看模型.....	32
2.6.2. 预测	33
2.6.3. 二分类评估	37
2.7. 结果存储	40
2.7.1. 保存模型.....	40

1. 蚂蚁机器学习平台介绍

蚂蚁金服机器学习平台-金融智能平台是可视化、分布式的机器学习平台，提供了丰富的算法能力，涵盖金融领域的所有建模场景，提供完整的从数据分析到模型训练、部署监控的一站式功能。

通常情况下，一个数据智能应用的完整链路很长，从数据的 ETL 到预处理、特征工程、模型训练、评估和服务。以前经常不得不在多个工具中来回穿梭才能完成任务，而蚂蚁机器学习平台金融智能平台致力于打造一站式的开发环境，提供了从元数据到模型部署整套流程。通过基本组件，可以搭建各个垂直场景下的解决方案，节省了大量切换环境的损耗。

1.1. 蚂蚁机器学习平台首页



蚂蚁机器学习平台的主要概念介绍：

1) 实验（Experiment）

实验是建模过程的一系列操作：包括读取数据、数据清洗、特征工程、训练、评估等，是可视化的操作流。

2) 数据集（Dataset）

一个实验从引入数据集开始，建模的第一步是读取数据。

3) 组件

操作流上的每一个节点，可以是一个数据集，也可以是一个数据处理操作，也可以是模型。

4) 模型 (Model)

跟机器学习中的模型是一个概念，用户通过选择数据集、选择算法、配置参数就可以自动训练一个模型。

5) 设置

设置最近使用的组件数量，设置语言。

6) 帮助

帮助文档包括蚂蚁机器学习平台的介绍，用户使用手册，常见问题 FAQ。

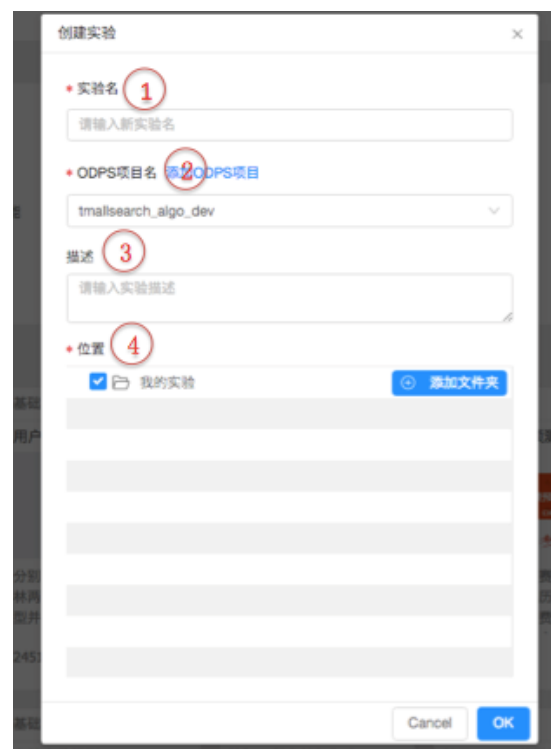
2. 机器学习的一般过程

机器学习的建模过程一般包括：读取数据、数据探查、特征工程、模型训练、模型评估，下面以实际案例来说明。

案例背景介绍：目前电信运营商面临着激烈的市场竞争，不同的电信运营商都在争夺客户。由于电信运营商前期投资成本巨大，客户数量对电信运营商来讲异常重要，有了足够多的客户才能摊低前期的投入成本，避免企业经营亏损。这个案例的目标是通过电信客户的一系列特征来预测客户是否会流失。

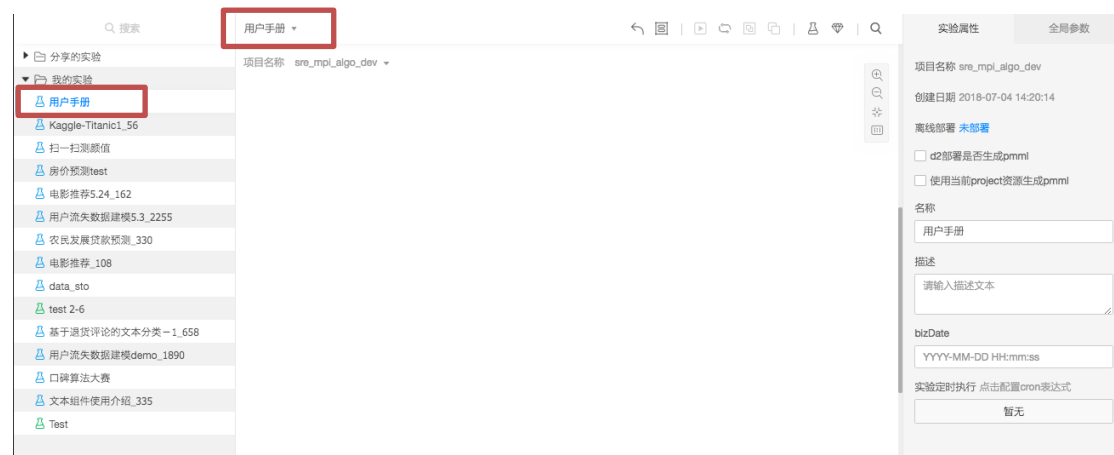
2.1. 新建实验

首先，我们需要新建一个实验，才能开始建模工作。打开 PAI 平台首页，从首页模版点击“新建空白试验”，如下图：



- 1) 实验名（必填）：对创建的实验命名。
- 2) 项目名（必填）：实验所在的项目名。
- 3) 描述（选填）：对实验内容的描述。
- 4) 位置（必填）：实验创建后放置的位置，选择一个文件夹，也可以新增文件夹。

实验创建成功后，出现如下页面：



页面右侧可以配置实验属性参数，如下图：

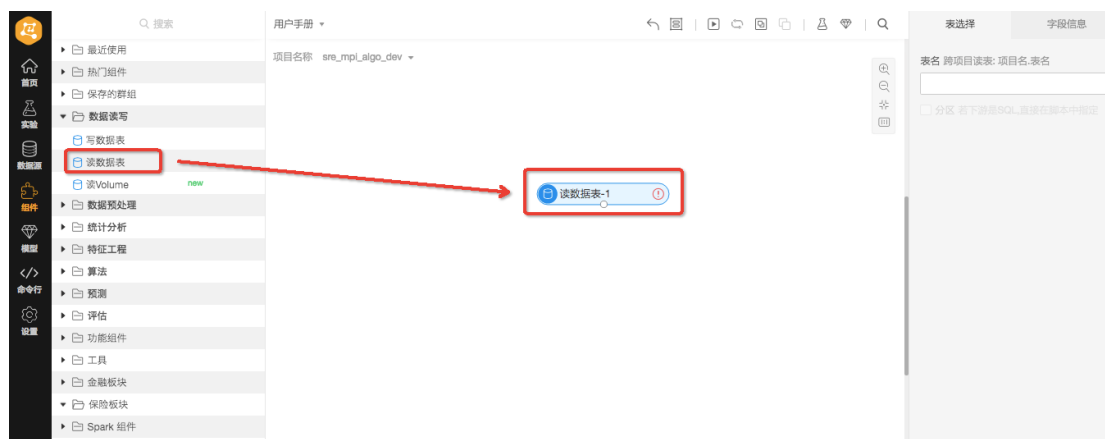


- 1) 项目名称：实验所在的项目名。
- 2) 创建日期：实验创建的日期。
- 3) 名称：创建实验的名称。
- 4) 描述：实验内容的描述。
- 5) bizDate：替换实验中的日期，按照此日期执行。
- 6) 配置实验定时执行。

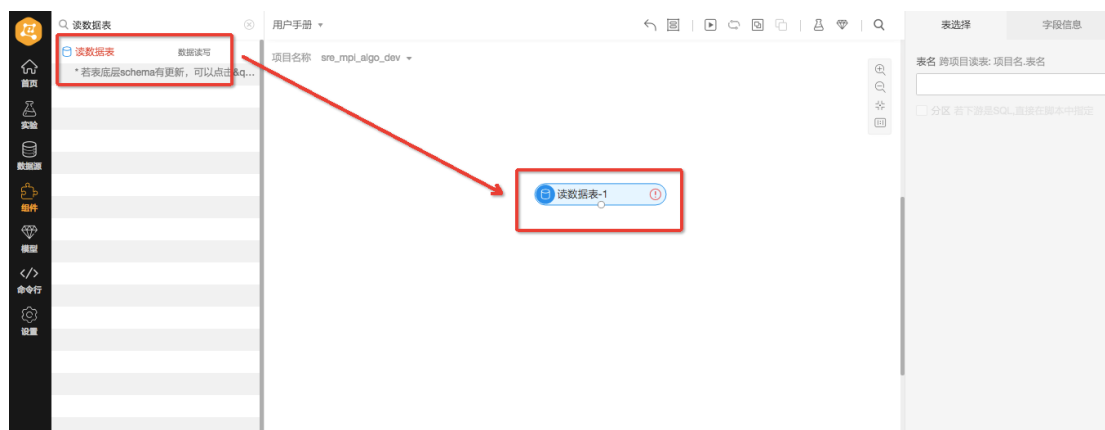
7) 设置一次全局变量，然后每个组件中都可以使用。

2.2. 读取数据表

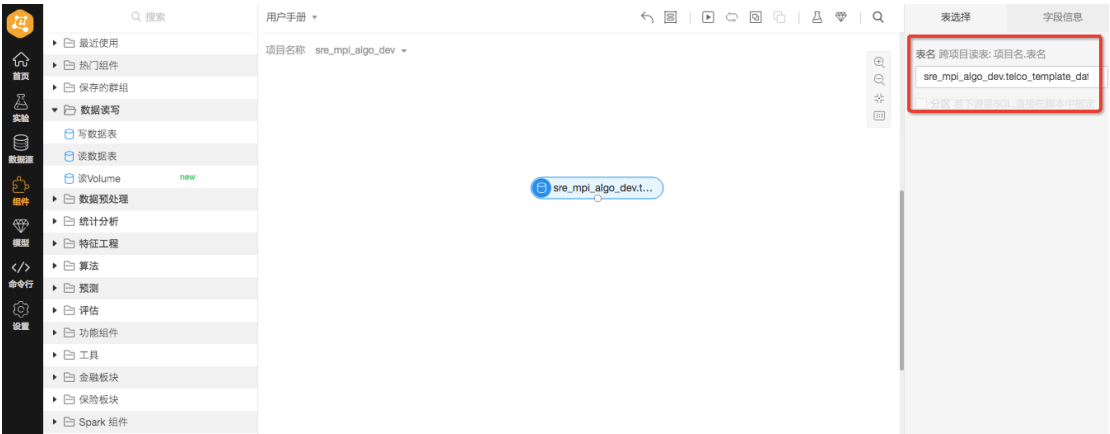
- 创建一个模型，我们第一步需要做的是引入数据集，从左侧组件栏中选择“数据读写—读数据表”或搜索“读数据表”，拖入画布中，如下图：



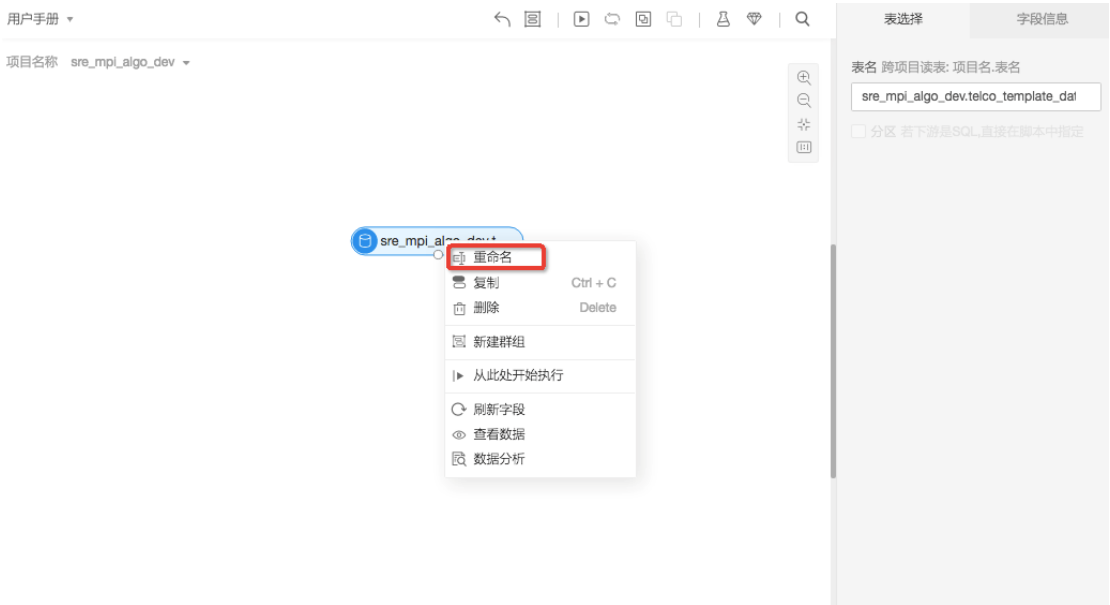
- 也可以从搜索框中搜索到这个组件，然后拖入画布，如下图：



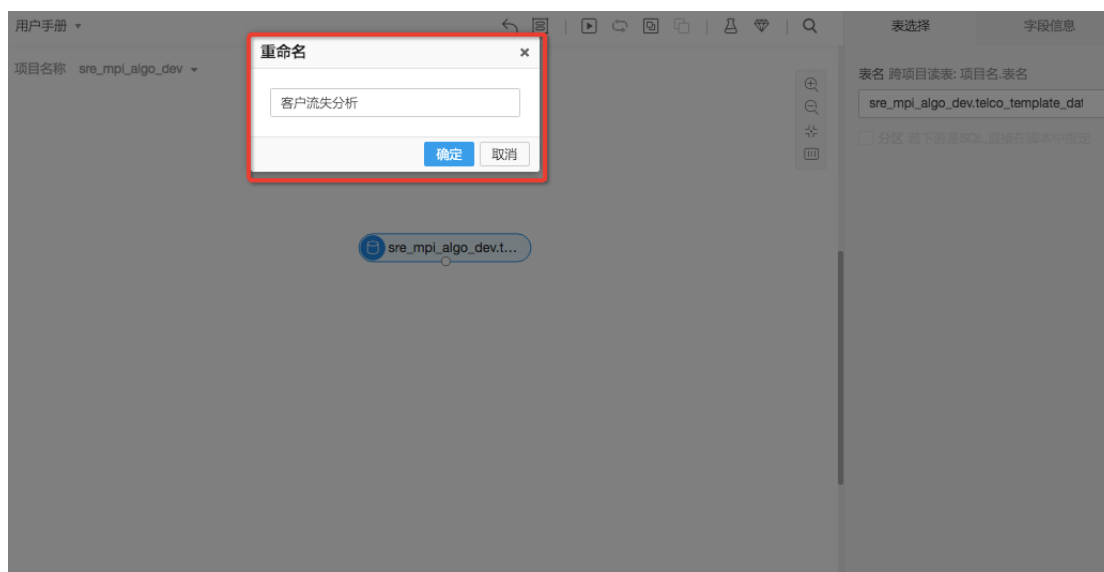
- 在“表选择—表名”输入 ODPS 数据表名，确保数据表存在于项目 Project 中，系统自动读取数据，如下图：



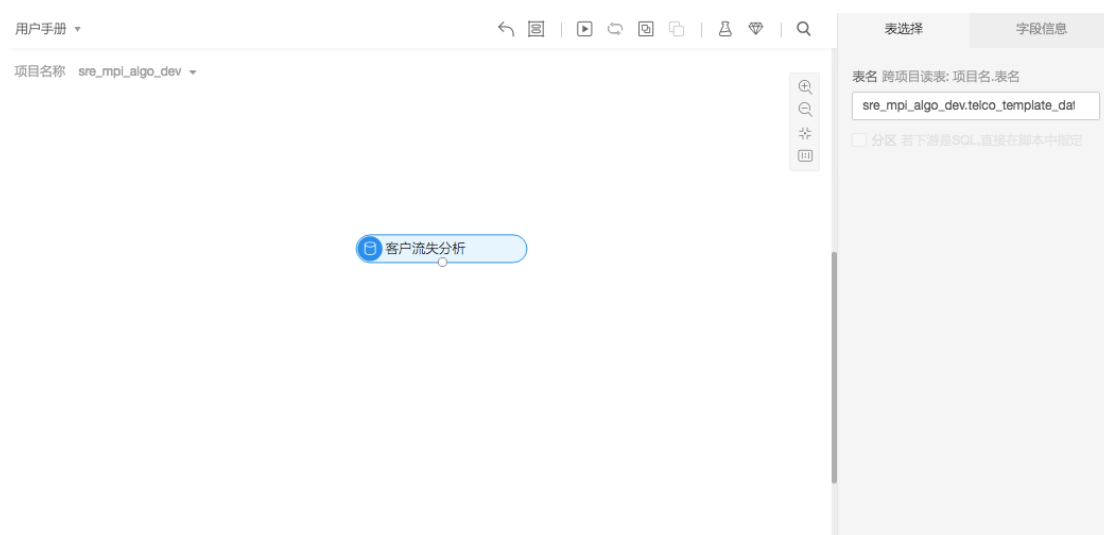
- 对数据表进行重命名，方便对数据的理解，右键点击组件，在下拉菜单中选择“重命名”，如下图：



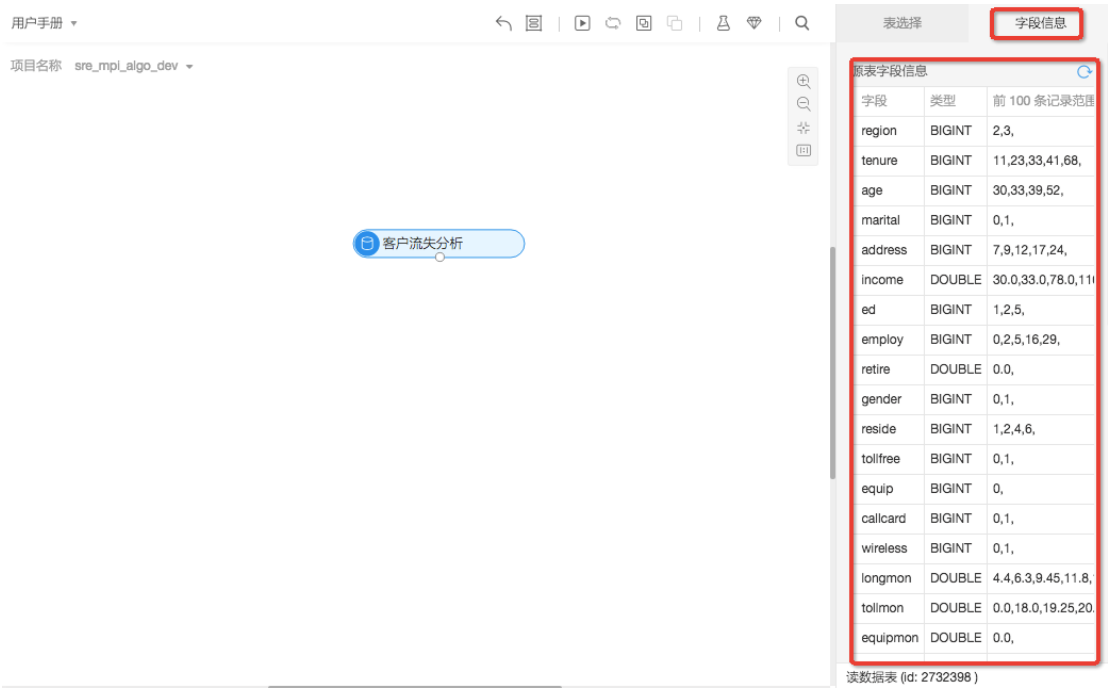
- 在弹出的提示框中输入“客户流失分析”，如下图：



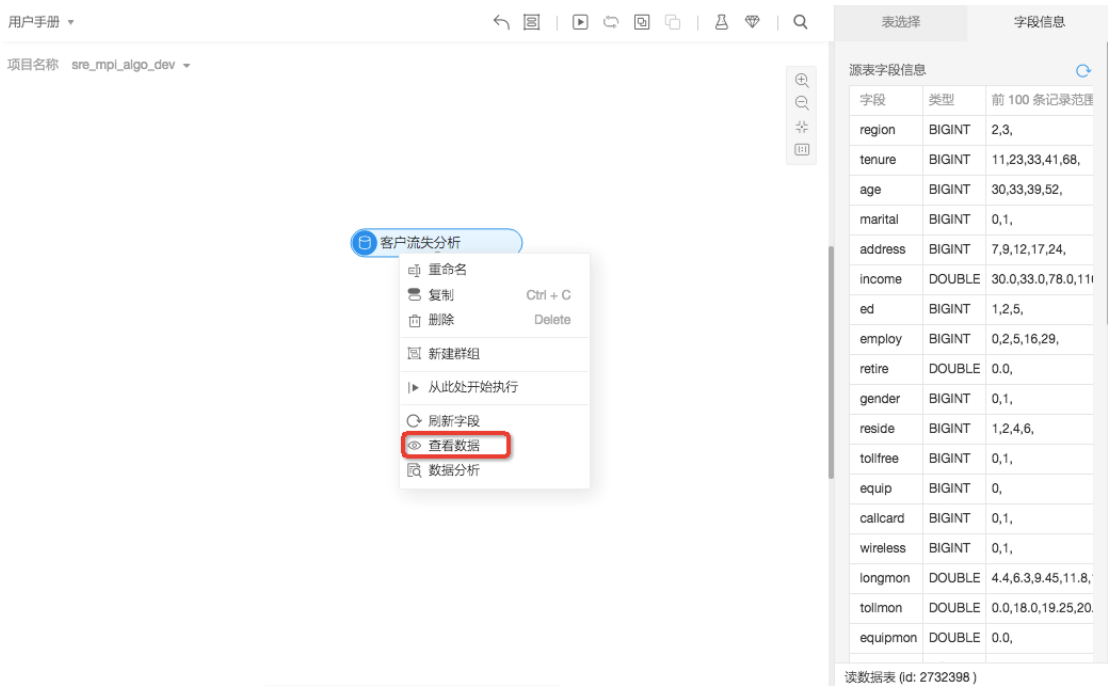
- 点击确定，组件名称显示为：“客户流失分析”，如下图：



- 切换到字段信息栏，可以查看输入表的字段名、数据类型和前 100 行数据的数值分布。如下图：

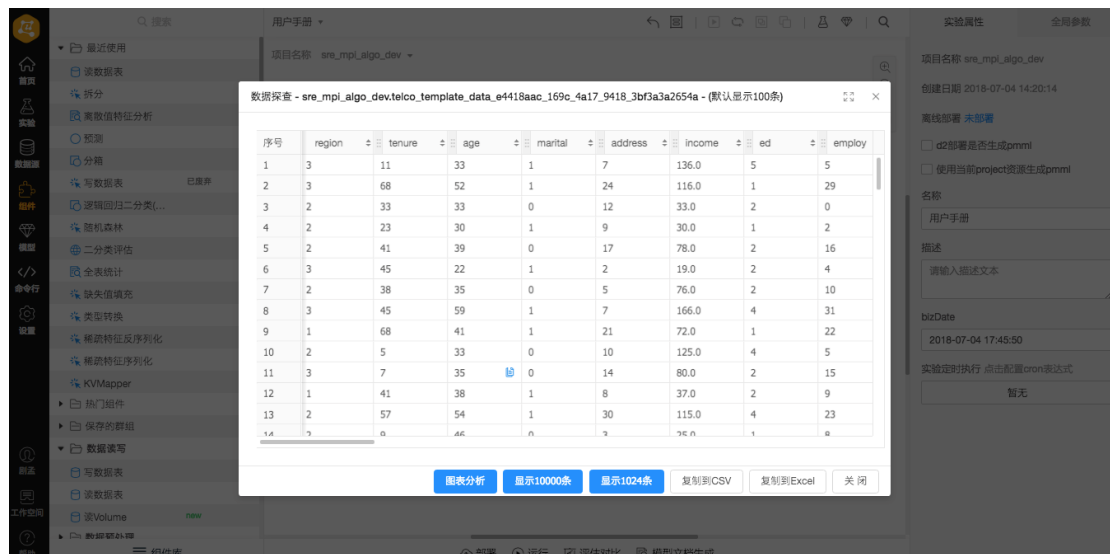


- 该案例中有 42 个字段，其中包含一些客户个人信息，例如年龄、婚姻状况、地址、收入、教育程度、行业、退休、性别、居住地和客户类别，还包含一些客户使用电信服务信息，例如使用电信服务时间，是否开通无线服务，是否开通语音信箱服务，是否开通亲情号服务，以及上月基本话费，上月长话费，上月上网费，累计基本话费，累计长话费，累计上网费等等。
- 在画布中右键点击组件可查看前 100 条详细数据。



- 点击“查看数据”，默认显示数据集的前 100 条数据，用电子表格的方式呈现，方便查看。

- 此页面提供“图表分析”、“显示 10000 条”、“显示 1024 条”、“复制到 CSV”、“复制到 Excel”等选项。如下图：

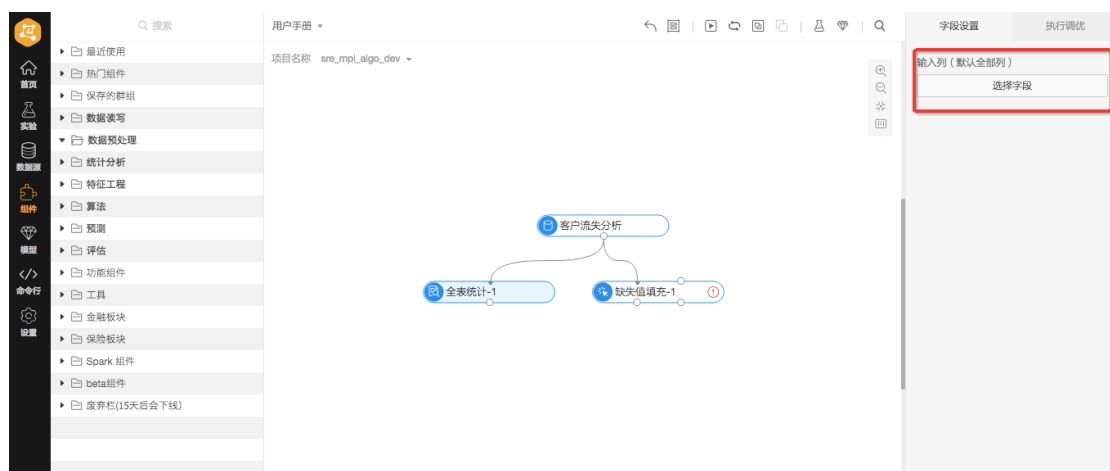


2.3. 数据探索

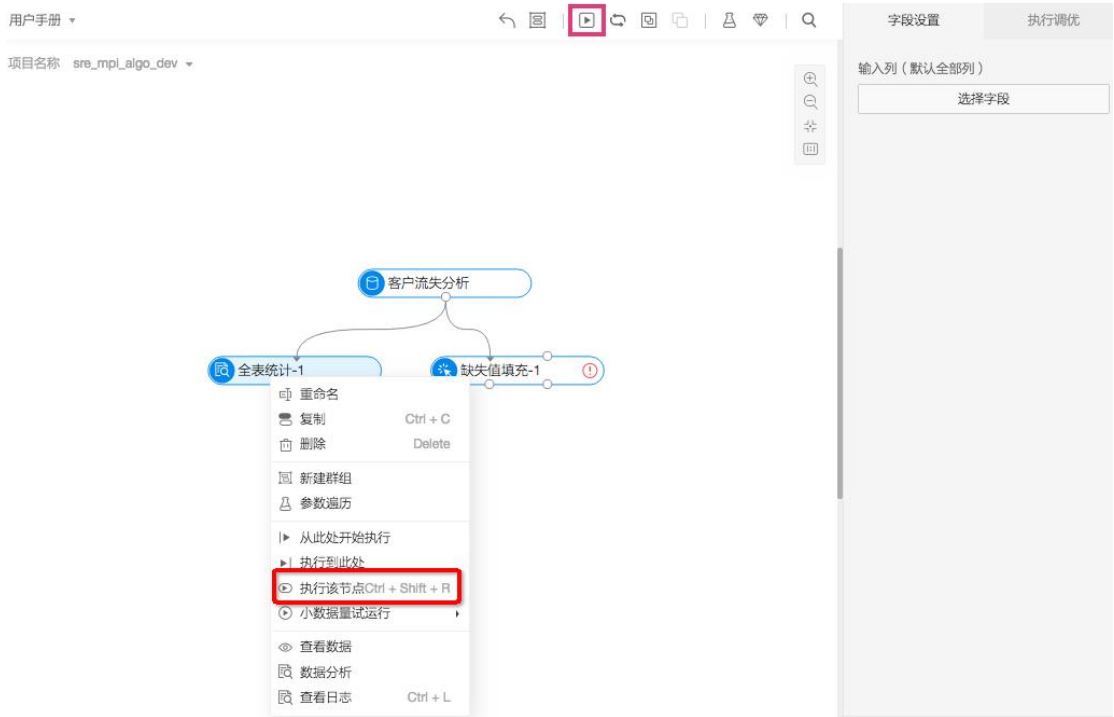
- 我们想对数据做进一步的探索，想每个字段的统计信息，包括缺省值、最大最小值、方差、偏值等等，我们使用“全表统计”组件来做分析。
- 使用搜索或从“组件—数据探索—统计分析—基本分析—全表统计”拖入画布中。
- 将“客户流失分析”和“全表统计”两个组件用线条连接。

2.3.1. 全表统计

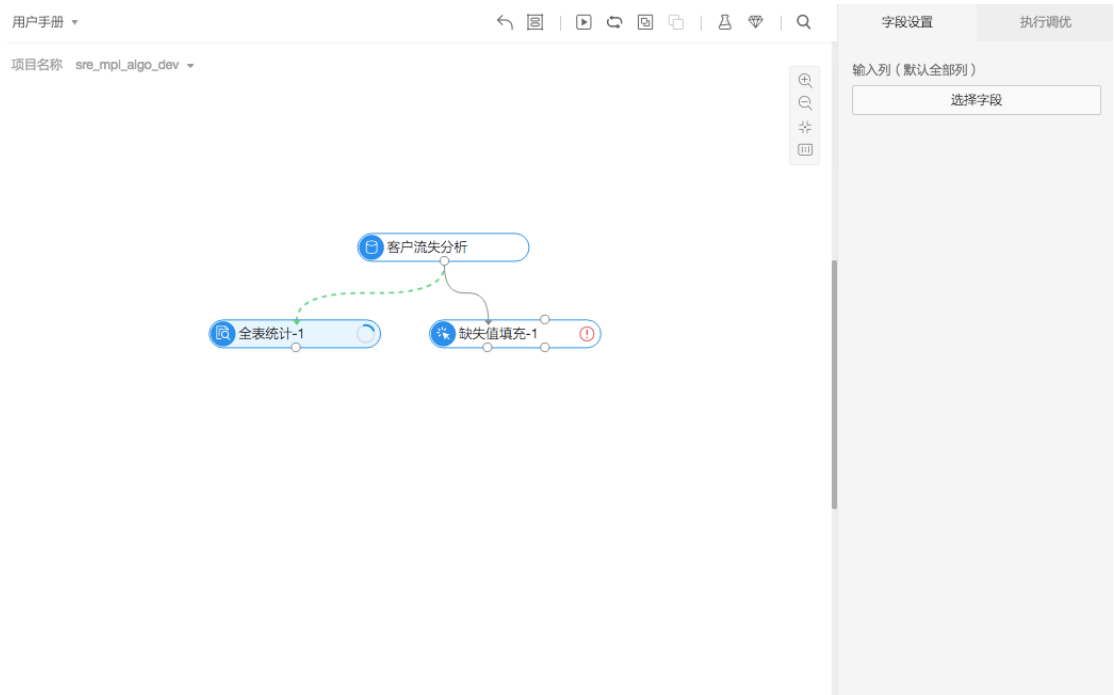
- “全表统计”组件的字段设置：默认是选中全部列，使用默认参数就可以。如下图：



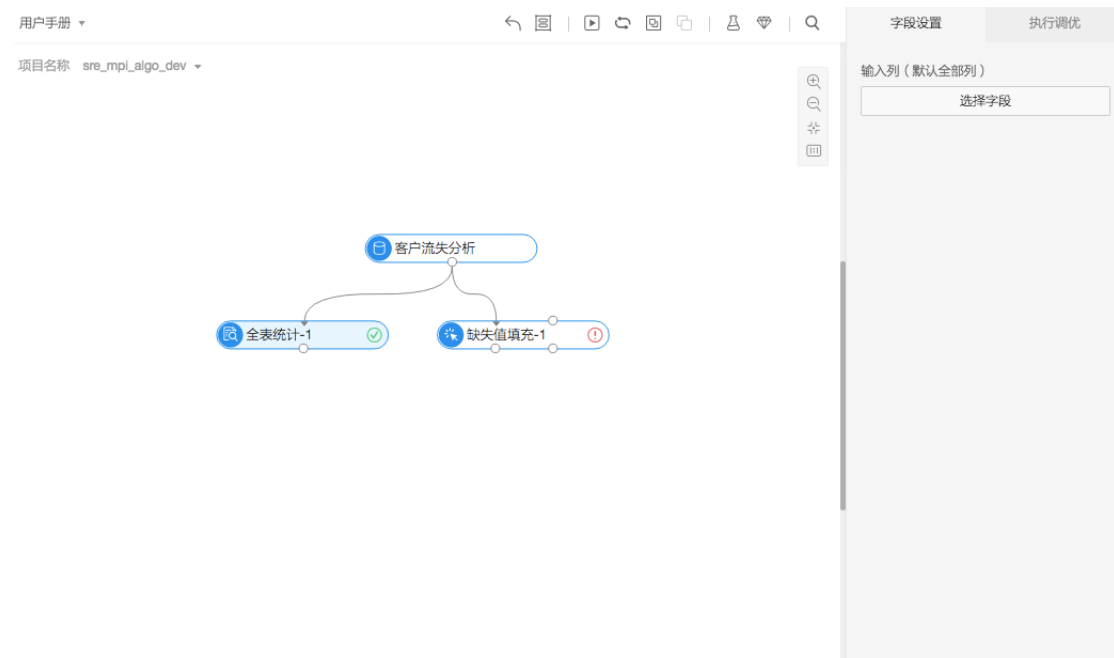
- 选中“全表统计”组件，点击右上角的“执行选择节点”；或者右键点击组件，在下拉菜单中选择“执行该节点”。如下图：



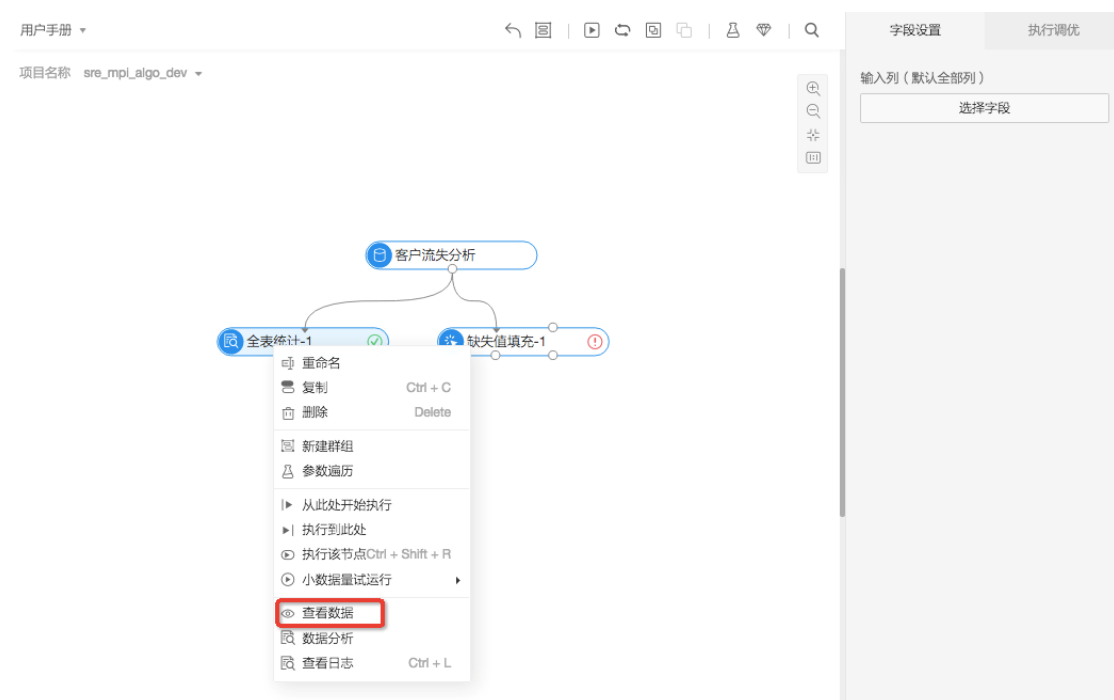
- 点击后，“全表统计”组件开始运行，运行状态如下图：



- “全表统计”组件运行成功后，状态如下图：



- 我们继续查看运行完成后的数据结果，右键点击“查看数据”，如下图：



- 运行结果中展示了数据集每个字段的统计信息：数据集的列名、字段类型、总的行数、有值的行数、缺失值行数、最小值、最大值、方差、偏度、峰度等。
- 该数据集中发现 logcard、logequi、logtoll、logwire 四个字段中有缺失值：

项目名称 sre_mpi_algo_dev

数据探查 - pai_temp_138035_2738963_1 - (默认显示100条)

序号	colname	datatype	totalcount	count	missingcount	nancount	positiveinfinitycount
18	gender	bigint	29999	29999	0	0	0
19	income	double	29999	29999	0	0	0
20	internet	bigint	29999	29999	0	0	0
21	lninc	double	29999	29999	0	0	0
22	logcard	double	29999	20339	9660	0	0
23	logequi	double	29999	11580	18419	0	0
24	loglong	double	29999	29999	0	0	0
25	logtoll	double	29999	14250	15749	0	0
26	logwire	double	29999	8880	21119	0	0
27	longmon	double	29999	29999	0	0	0
28	longten	double	29999	29999	0	0	0
29	marital	bigint	29999	29999	0	0	0
30	multiline	bigint	29999	29999	0	0	0
31	pager	bigint	29999	29999	0	0	0

图表分析

显示10000条

显示1024条

复制到CSV

复制到Excel

关闭

- 展示信息：
- 在“全表统计”的结果中，行表示数据集中的字段，列表示每个字段的统计信息。
- 具体描述如下表所示：

设置项	描述
Colname	列名：数据集中每列的名称。
Datatype	字段的存储类型。
Totalcount	总的行数。
Count	非 Null 的行数。
Missingcount	Null 的行数，缺失值数量。
Nancount	NAN（Not a Number）的数量，表示一些特殊的数值，无穷或非数值。
Positiveinfinitycount	正无穷的数量。
Negativeinfinitycount	负无穷的数量。
Min	最小值。
Max	最大值。
Mean	平均值。
Variance	方差。
Standarddeviation	标准差。
Standarderror	标准误差。

Skewness	偏度，是描述数据分布是否对称的指标。越对称的分布，skewness 越小。
Kurtosis	峰度，是用来描述数据是否长尾，如果一个分布异常点很多或者很长尾，那么其 kurtosis 也越大。
Moment2	二阶矩。
Moment3	三阶矩。
Moment4	四阶矩。
Centralmoment2	二阶中心距。
Centralmoment3	三阶中心距。
Centralmoment4	四阶中心距。
Sum	总和。
Sum2	平方和。
Sum3	立方和。
Sum4	四次方和。

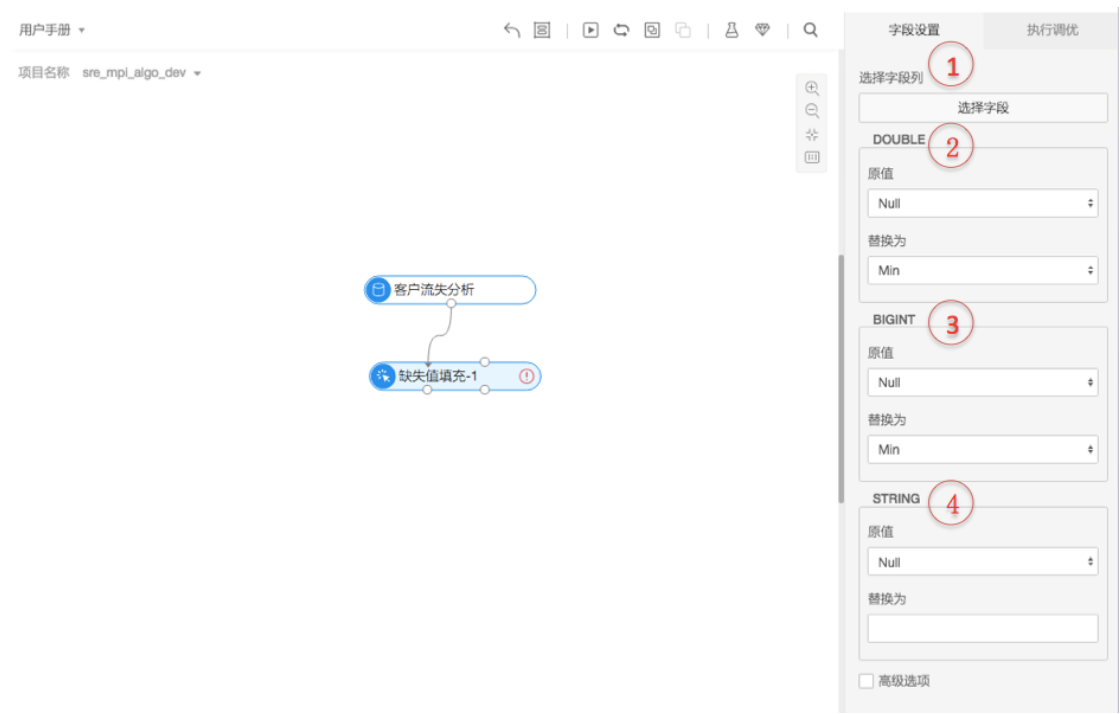
2.4. 数据预处理和统计分析

对前一步导入的数据进行预处理，并对预处理后的数据进行统计分析。

（PS: 本部分封装的是常用典型场景。但是请注意这里不能覆盖所有场景，如果需要复杂处理，请去 ETL 里写 SQL 进行）

2.4.1. 缺失值填充

- 在训练模型之前需要对上面发现的含缺失值的字段进行填充，从左侧组件栏中选中“数据预处理 - 缺失值填充”组件或搜索“缺失值填充”组件，拖入画布中，并将两个组件连线。
- 填写“缺失值填充”组件的属性信息：1）选择需要填充缺失值的列名；2）选择数值型字段的填充方法；3）选择整型字段的填充方法；4）对字符串字段的填充方法。如下图：

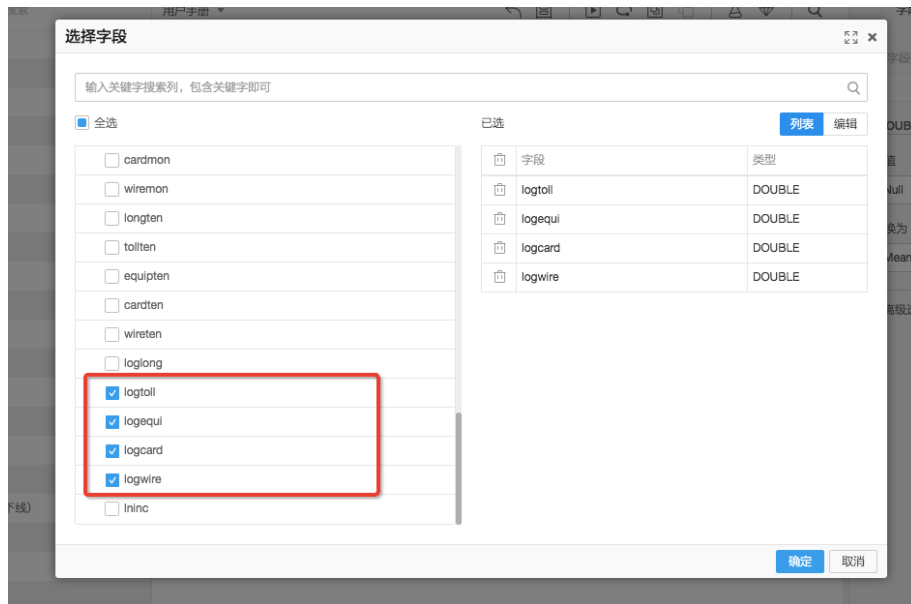


➤ 字段设置：

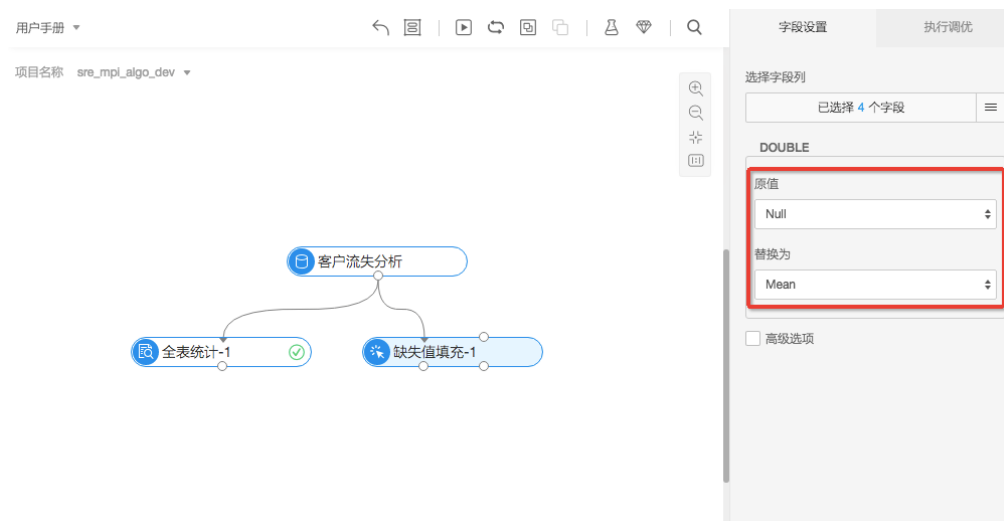
在字段设置中，可以选择字段列，设置 Double 型、Bigint 型、String 型字段的填充方法，还可以选择高级方式进行填充。
具体描述如下表所示：

设置项	描述
选择字段列	选择要填充缺失值的列。
DOUBLE	设置 Double 型字段的缺失值填充方法，原值为 Null，可替换方式有四种：Min、Max、Mean、自定义。
BIGINT	设置 Bigint 型字段的缺失值填充方法，原值为 Null，可替换方式有四种：Min、Max、Mean、自定义。
STRING	设置 String 型字段的缺失值填充方法，原值有四种类型：Null、空字符、Null 和空字符、自定义，可替换方式自定义。
高级选项	可以同时为多列进行不同方式的替换。原值，替换为三个部分组成了 config 参数，分别对应 config 参数的三个部分：列名，原值，替换值。

- 根据前面“全表统计”的分析，我们发现“logtoll”、“logequi”、“logcard”、“logwire”这四个字段有缺失值，需要进行填充，如下图：

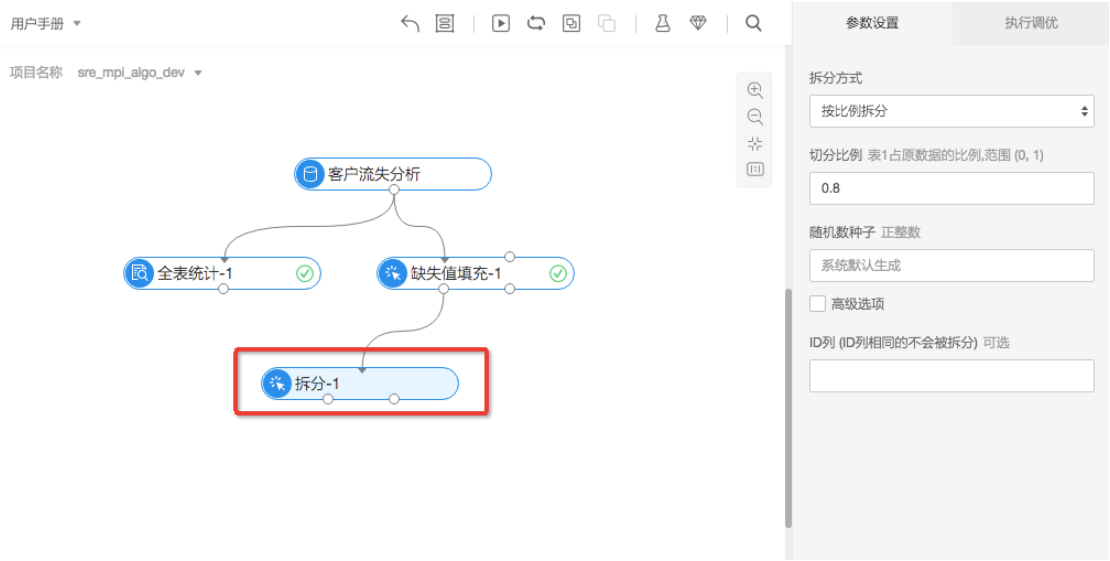


- 设置缺失值填充方法，针对数值型的变量系统提供了四种缺失值方法（Min，Max，Mean，自定义），此处用的是均值填充：

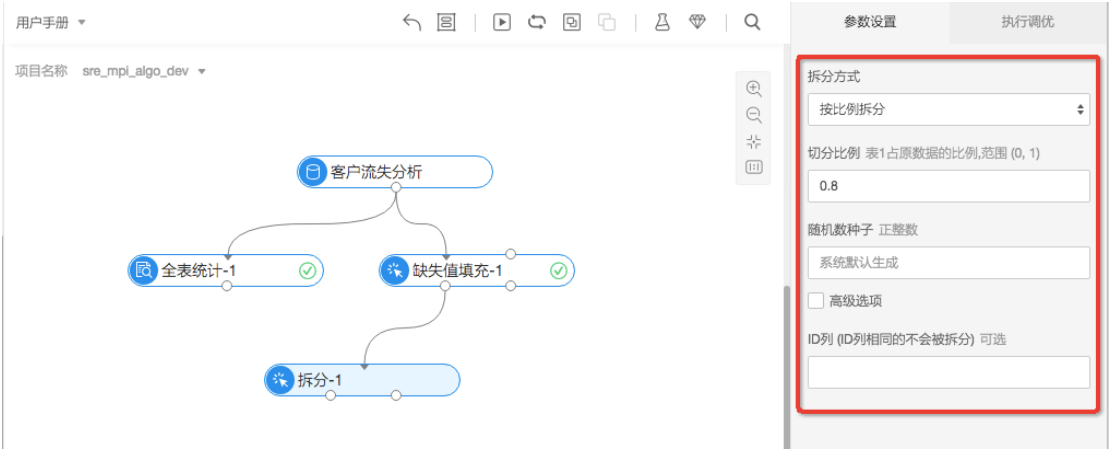


2.4.2. 数据拆分

- 训练模型之前我们对填充完缺失值的数据进行拆分，此步的目的是将数据集分成两份：一份是训练数据集，另一份是验证数据集。通过训练数据集来训练模型，通过验证数据集来验证模型的效果。
- 从左侧组件栏中选中“数据预处理 - 拆分”组件或搜索“拆分”组件，拖入画布中，用线条连接“缺失值填充”组件和“拆分”组件。如下图：



选中“拆分”组件，画布右侧显示该组件的参数设置项，如下图：

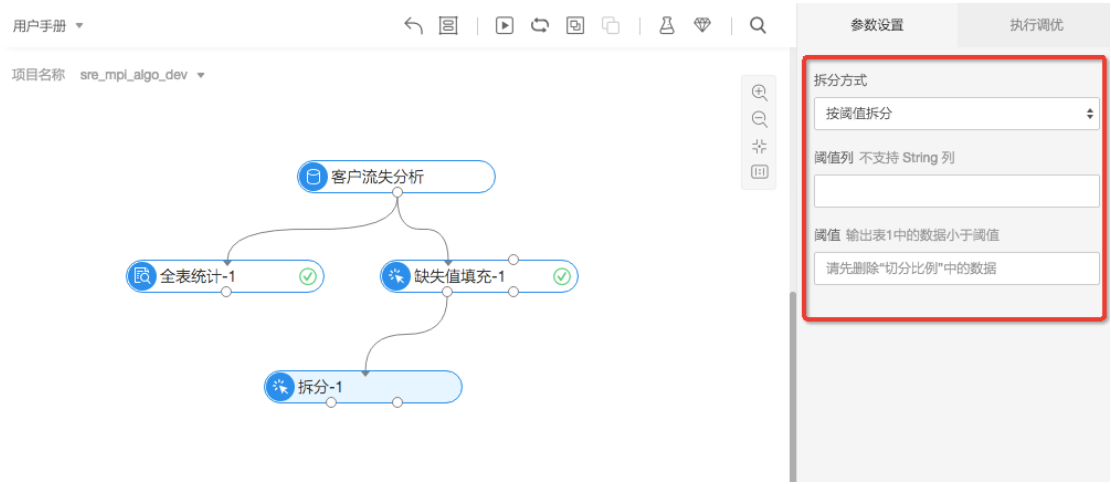


- 参数设置：
- 在参数设置中，可以设置：拆分方式、切分比例、随机数种子、ID 列。
具体描述如下表所示：

设置项	描述
拆分方式	提供了两个选择项：按比例拆分、按阈值拆分。默认是按比例拆分。
切分比例	默认是 0.8，表示 80%的数据作为训练集，20%的数据作为验证集。如果修改为 0.5，表示 50%的数据作为训练集，50%的数据作为验证集。 此处的取值范围是（0-1）的分数。
随机数种子	随机数就是就随机数种子中取出的数，系统默认生成随机数种子。如果想保证每次拆分出的数据集是相同的，那么需要手动输入一个固定的随机数种子。

ID 列	ID 列是可选项，如果设置了 ID 列，那么相同的 ID 不会被拆分。
------	-------------------------------------

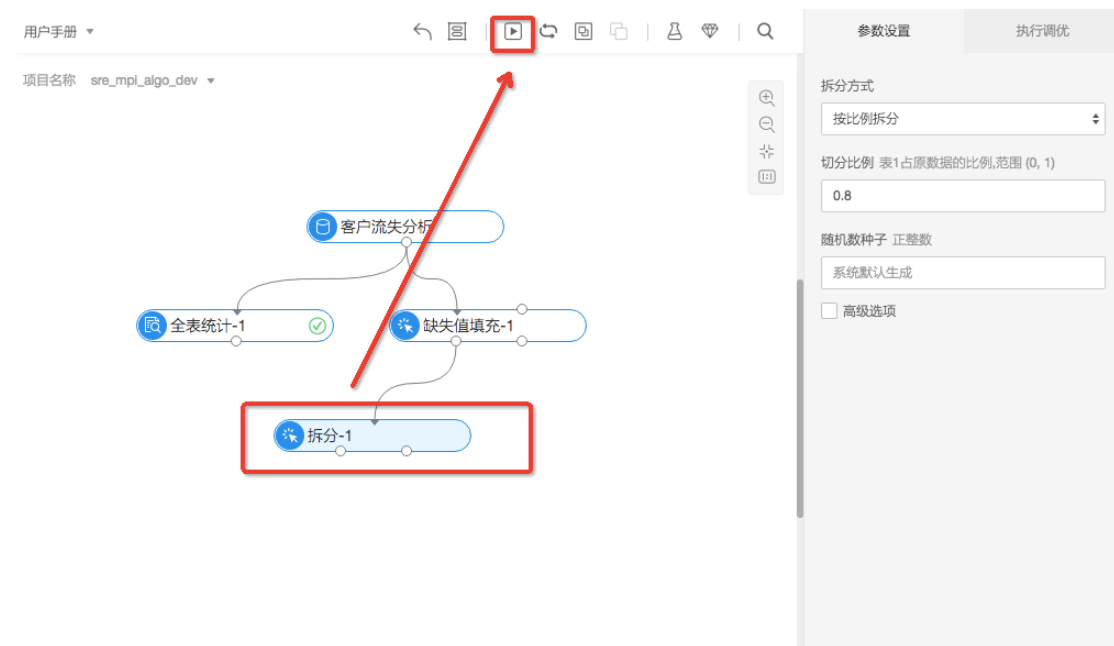
如果拆分方式选择了“按阈值拆分”，那么参数设置页面如下图：



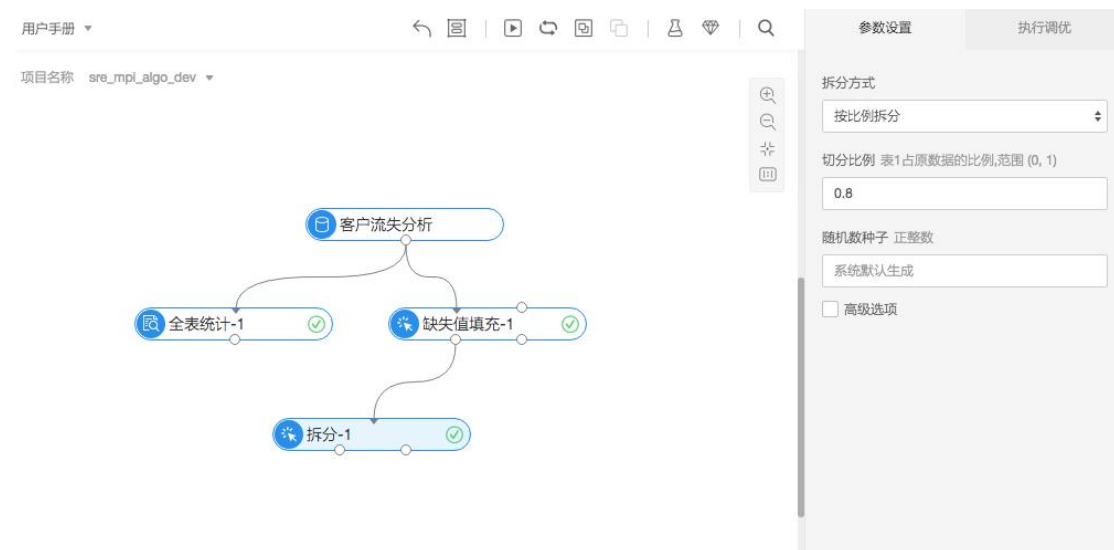
- 参数设置：
- 在参数设置中，可以设置：阈值列、阈值。
- 具体描述如下表所示：

设置项	描述
拆分方式	按阈值拆分：根据某一列的阈值进行拆分数据集。
阈值列	选择某一列作为拆分类，不能是 String
阈值	根据阈值对数据集进行拆分

设置好拆分规则之后（本案例选择的是按 0.8 比例拆分），点击右上角的“执行选择节点”，如下图：

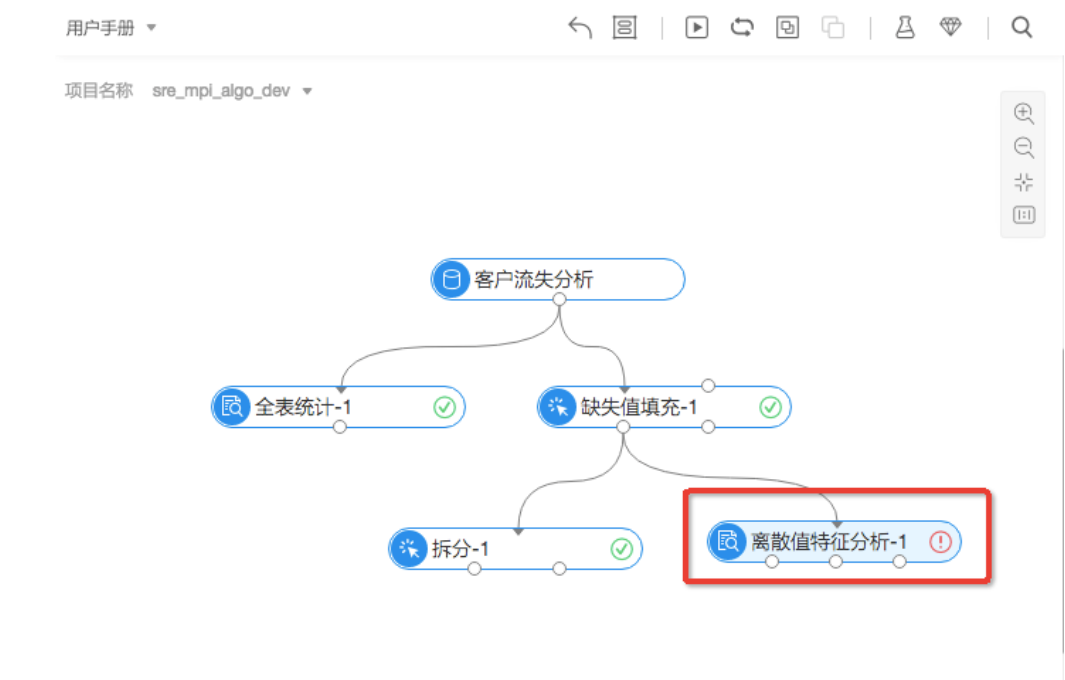


执行成功后如下图：

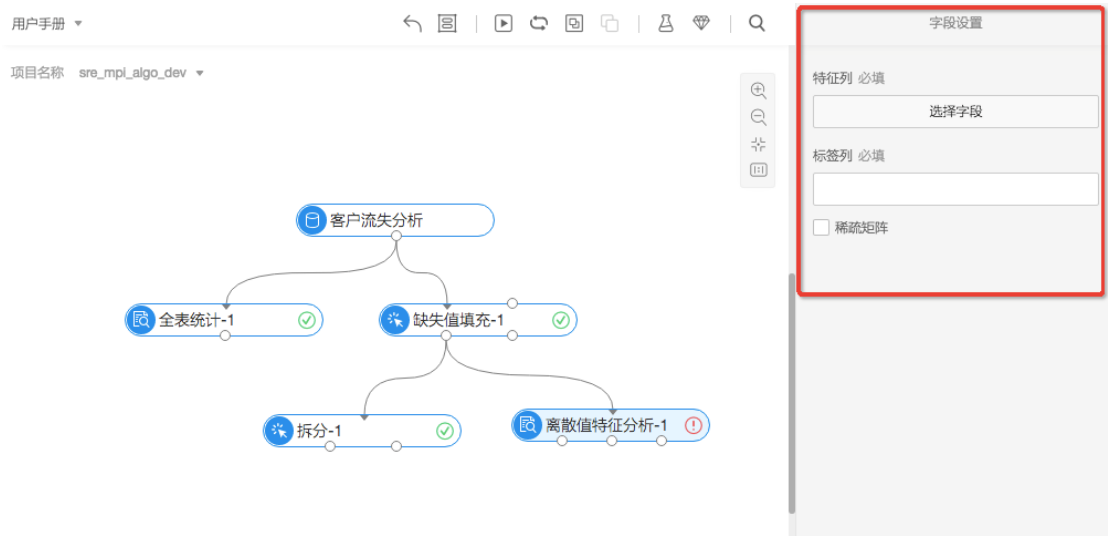


2.4.3. 离散值特征分析

- 在训练一个模型之前，我们通常会做一些简单分析来了解数据，本案例中我们用“离散值特征分析”来对数据做一个简单的分析。
- 从左侧组件栏中选择“统计分析—可视化分析—离散值特征分析”，或搜索“离散值特征分析”，拖入到画布中。
- 将“缺失值填充”和“离散值特征分析”两个组件用线条连接。如下图：



选中“离散值特征分析”组件，画布右侧显示该组件的参数设置项，如下图所示：



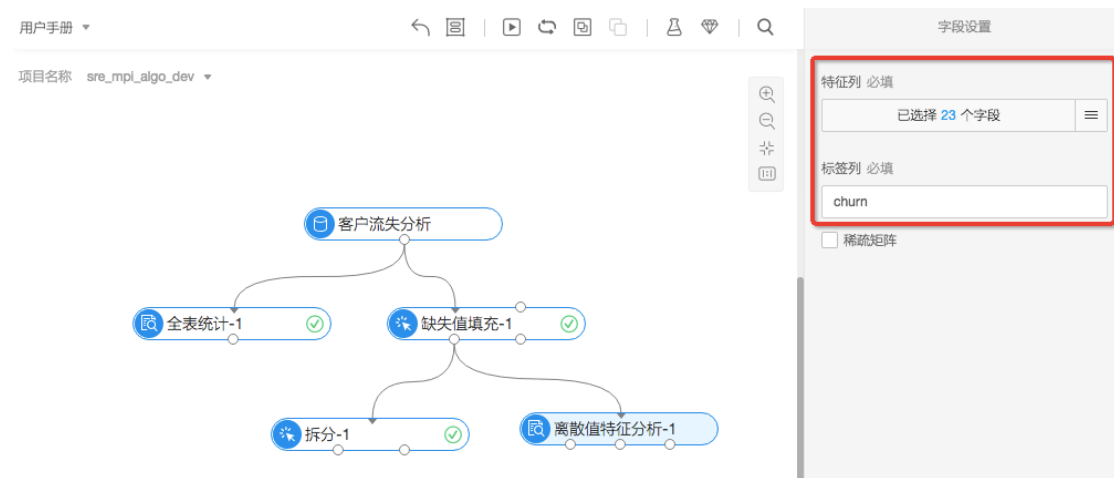
➤ 参数设置：

在字段设置中，可以设置：特征列、标签列、稀疏矩阵。

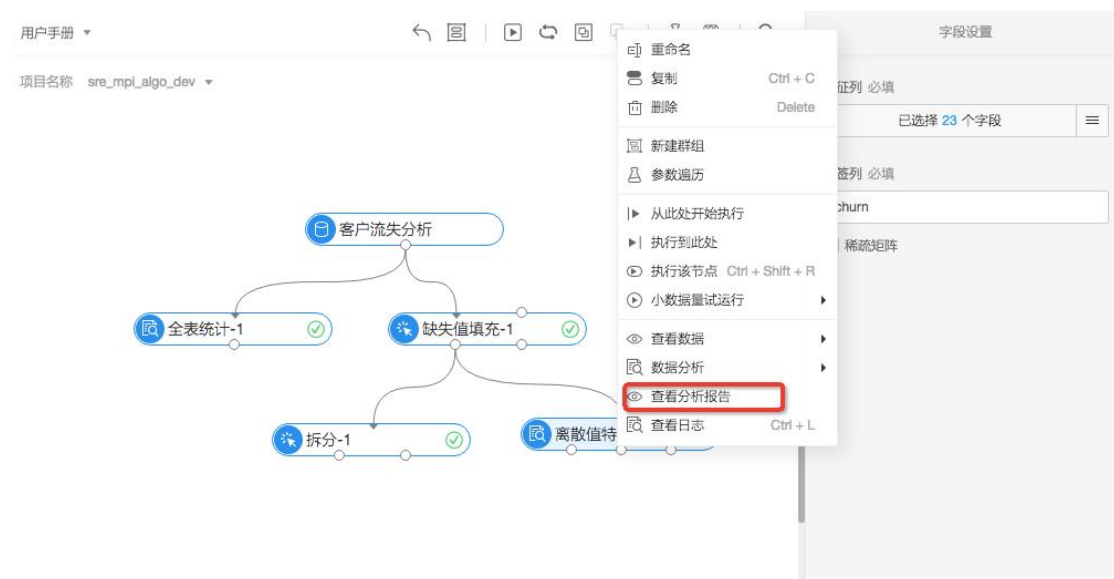
具体描述如下表所示：

设置项	描述
特征列	选择数据集中的特征列，一般是除 Label 外的其他列。
标签列	选择某一列作为 Label。
稀疏矩阵	如果输入表为 KV 格式，例如：k1:v1,k2:v2，那么特征列的输入方法为：k1,k2。

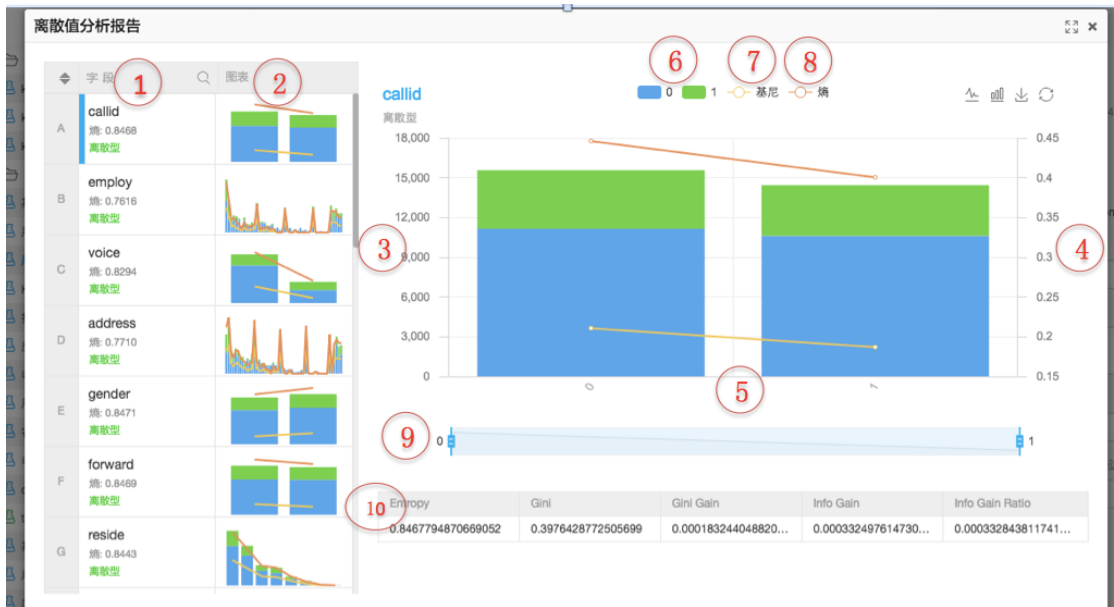
本案例中我们选择了除 Label 外的所有 Bigint 字段，设置 Churn 作为标签，输入表并非 KV 格式，所以不选稀疏矩阵。如下图：



设置完成后，运行此节点，运行成功后想查看结果，右键点击组件，选择“查看分析报告”，如下图：



点击“查看分析报告”，浮层展示结果，如下图：



- 1、字段：做离散值特征分析的字段。
- 2、图表：离散值特征分析的缩略图表。
- 3、图表纵坐标（左）：符合条件的记录数。
- 4、图表纵坐标（右）：基尼系数和熵值。
- 5、图表横坐标：Callid 的取值，0 表示无呼叫保持服务、1 表示有呼叫保持服务。
- 6、标签 Label：0 表示无流失，1 表示有流失。
- 7、基尼系数：根据洛伦茨曲线找出了判断分配平等程度的指标。
- 8、熵：代表随机变量不确定度的度量，即发生概率越高的事件，其所携带的信息熵越低。
- 9、拖动滑动条来选择离散值的范围
- 10、熵值、基尼系数、基尼增益、信息增益、信息增益比。

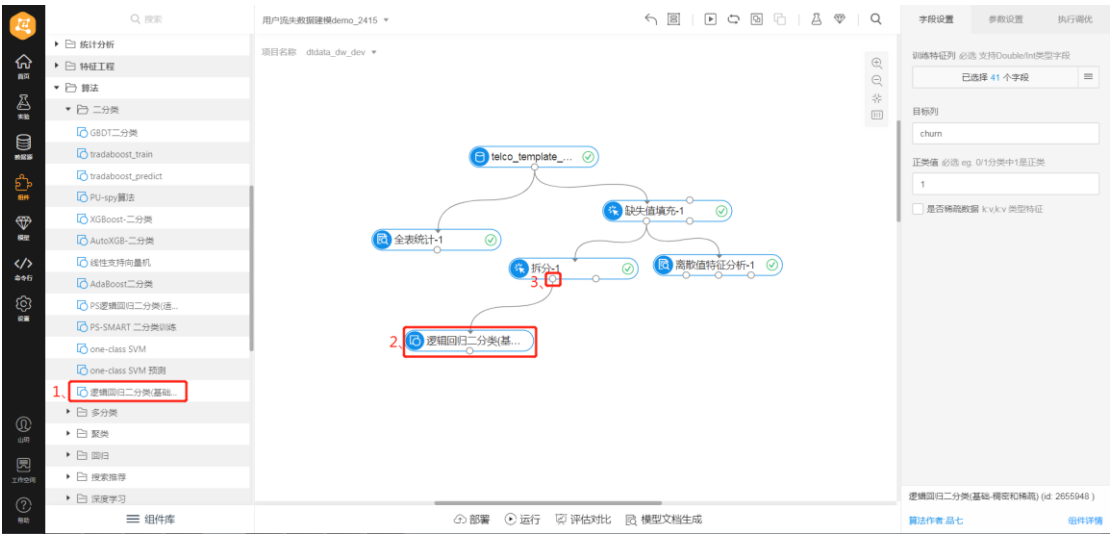
2.5. 算法建模

2.5.1. 逻辑回归二分类

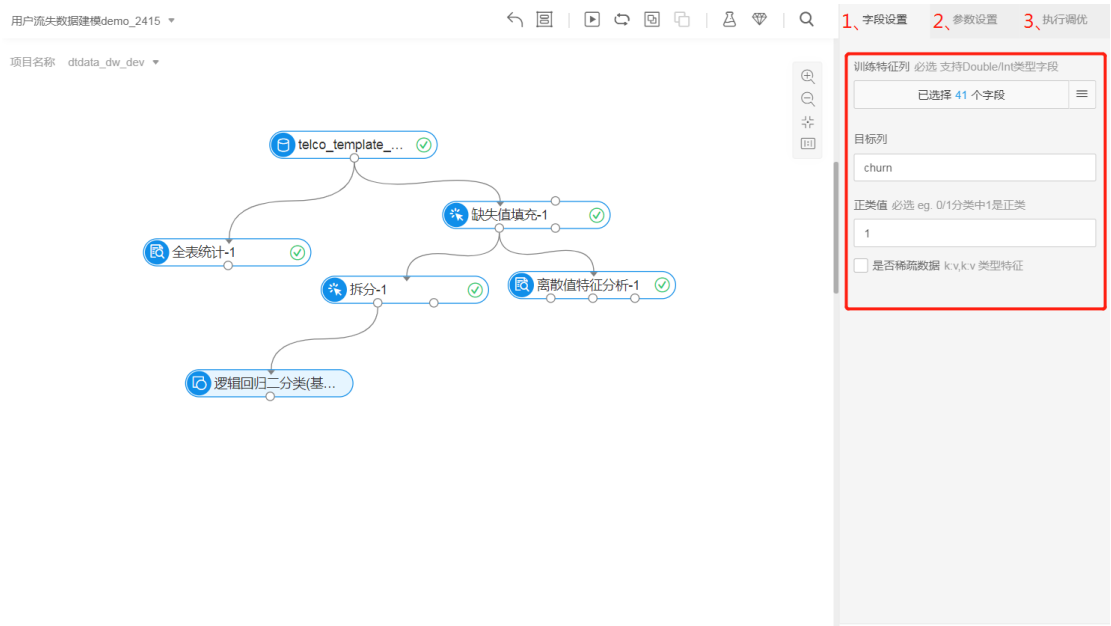
算法简介：一种线性二分类模型，基于样本特征的加权以及 logistic 函数，实现对样本分类的概率预测。训练过程中可以通过设置 L1, L2 正则项的系数控制非零参数的规模或数值幅度，从而使模型具有更好的泛化能力。

本实验使用二分类模型进行建模，从“算法-->二分类”栏目，将“逻辑回归二分类”组件，拖拽进入面板。

将拆分组件的第一个输出端点连接到逻辑回归二分类组件上。如下图所示：



选中“逻辑回归二分类”组件，画布右侧的参数栏显示它的参数配置项。分为 字段设置、参数设置、执行调优。



➤ 字段设置：

在字段设置中，选择训练特征列、目标列、正类值、是否稀疏数据。

具体描述如下表所示：

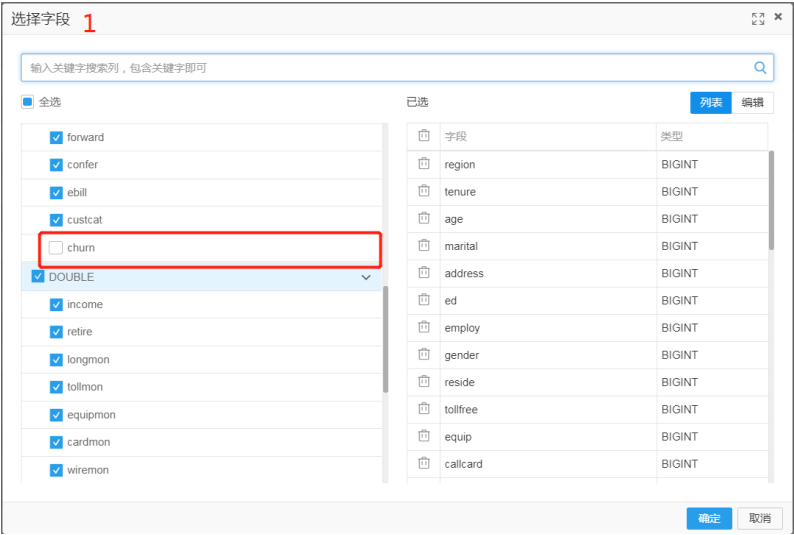
设置项	描述
训练特征列	用来训练模型的特征列，除标签列、权重列以外的其他列

目标列	数据的标签列，即 label												
正类值	目标样本的标记，如“1”表示用户的状态为流失，“0”表示用户的状态正常。所以在“正类值”中填入目标样本的标签，“1”												
是否稀疏数据	<div>数据是否具有 key-value 特征。具有 key-value 特征的数据如下：</div> <table><tr><td colspan="3">kv</td></tr><tr><td colspan="3">a:0 b:1 c:2</td></tr></table> <div>平台解析后变为，黑体的为字段名（key），值为字段值（value）：</div> <table><tr><td>1</td><td>2</td><td>3</td></tr><tr><td>a</td><td>b</td><td>c</td></tr></table> <div>数据经过 kv 转换后，能有效的减少存储空间。</div>	kv			a:0 b:1 c:2			1	2	3	a	b	c
kv													
a:0 b:1 c:2													
1	2	3											
a	b	c											

在本案例中，churn 字段表示用户是否会流失。Churn 为 1 表示用户会流失，为 0 表示用户不会流失。

点击“选择字段”按钮，选择除标签（churn）以外的其余字段。点击“目标列”的输入框，输入标签列的字段名，系统将进行模糊匹配，用户点击选择目标列。





➤ 参数设置：
在参数设置中，设置逻辑回归二分类的算法参数。具体描述如下表所示：

设置项	是否必选	可选范围	默认值	描述
正则项	可选	'l1'、'l2'、 'None'	l1	在逻辑回归函数中加入正则项。正则项的类别有：L1 正则、L2 正则。
最大迭代次数	可选	-	100	指定 L-BFGS 的最大迭代次数
正则系数	可选，正值类型为 None 时，此值无效	-	1	正则项的系数。 <ul style="list-style-type: none">L1 正则项系数，主要用于控制非零参数的规模，设置的越小，非零的参数越多。对训练集的拟合效果会越好。L2 正则项系数，主要用于控制参数的数值幅度。设置的越大，参数的变化幅度越小，不容易拟合到现有的数据集。
最小收敛误差	必选	-	1.0e-06	L-BFGS 的终止条件，即两次迭代之间 log-likelihood 的差

本案例中选择正则项的类型为“None”，即在函数后，不追加正则项。其他系数如下图所示：

字段设置

参数设置

执行调优

正则项 可选

None

最大迭代次数 可选

100

正则系数 可选 正则类型为None时此值无效

1

最小收敛误差

0.000001

- 执行调优：
- 在执行调优中，设置核数目和每个核的内存大小（MB）。具体描述如下表所示：

设置项	是否必选	可选范围	默认值	描述
核数目	可选	视集群的具体情况而定	默认自动调整	内存中为该实验分配的内存核数目
每个核的内存大小（MB）	可选	视集群的具体情况而定	默认自动调整	每个内存核的内存大小

本案例中核数目和每个核的内存大小都设置为“默认自动调整”。

字段设置

参数设置

执行调优

核数目

默认自动调整

每个核的内存大小(MB)

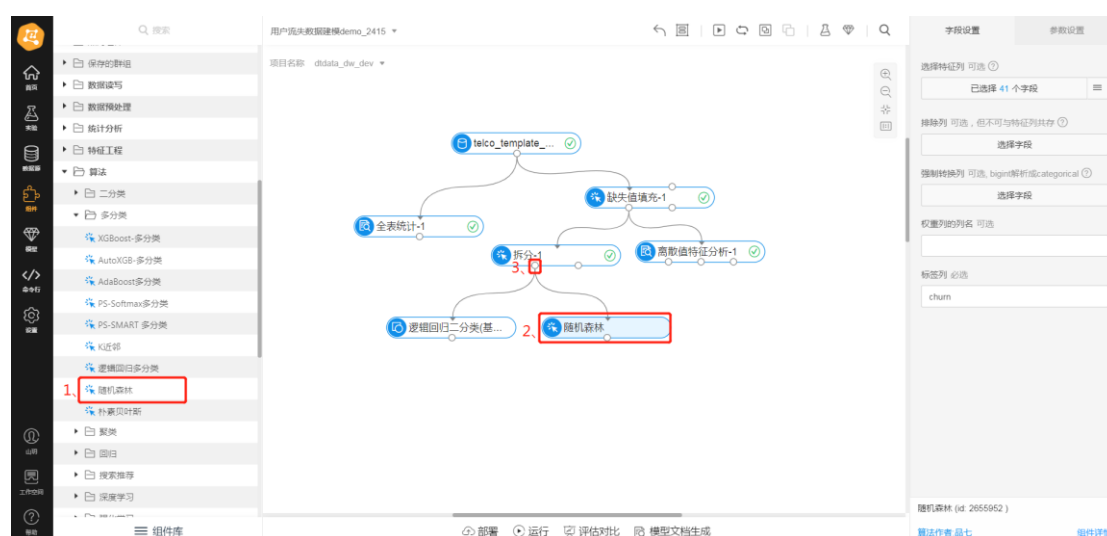
默认自动调整

2.5.2. 随机森林

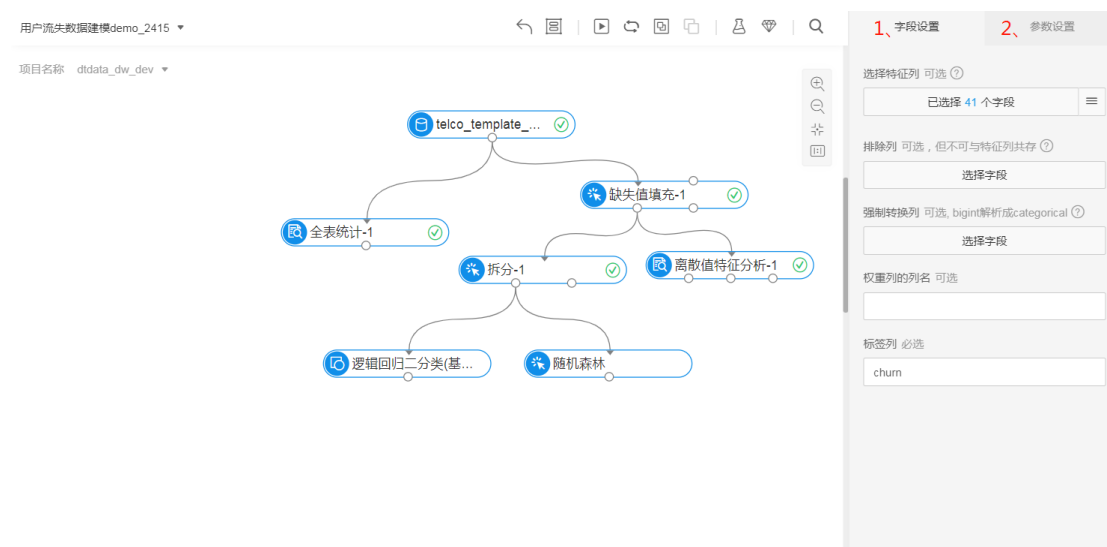
算法简介：随机森林是一个包含多个决策树的分类器。采用多个决策树的平等投票机制。在分类问题中，随机森林整体的输出结果是票数最多的分类选项；在回归问题中，随机森林的输出是所有决策树输出的平均值。

本实验使用两个算法进行建模，建模完成后将评估对比两种算法的模型效果。在这步操作中，我们选用“随机森林”算法。从“算法-->多分类”栏目，将“随机森林”组件，拖拽进入面板。

将拆分组件的第一个输出端点连接到随机森林组件上。如下图所示：



选中“随机森林”组件，画布右侧的参数栏显示它的参数配置项。分为 字段设置、参数设置。



➤ 字段设置：

在字段设置中，选择训练特征列、目标列、正类值、是否稀疏数据。

具体描述如下表所示：

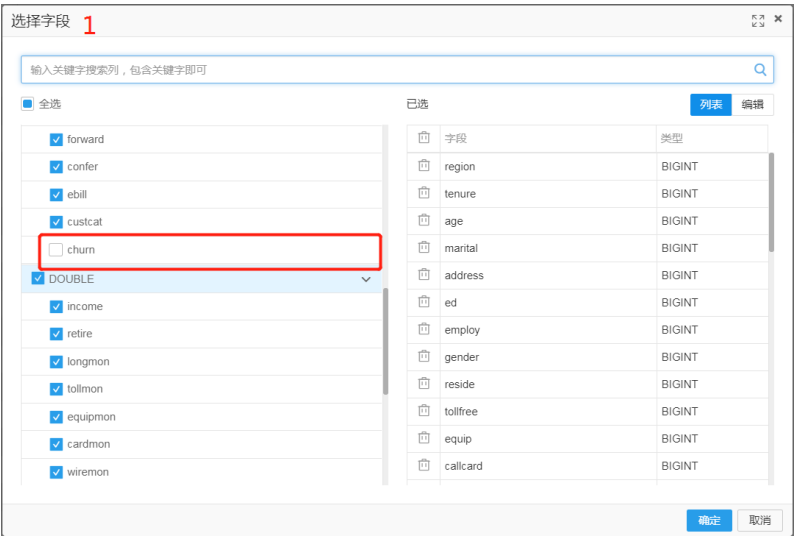
设置项	是否必选	描述
训练特征列	必选	用来训练模型的特征列，除标签列、权重列以外的其他列
排除列	可选	用来排除不参加训练的列。不能与特征列共存
强制转换列	可选	因为在随机森林算法中，需要解析特征是离散型还是连续型。所以此处设置强制转换列，将输入此列的字段，由 bigint 转换为 categorical 。 此外，对于其他字段类型的默认解析规则如下： <ul style="list-style-type: none"> • string、boolean、datetime 类型的列解析为离散类型 • double、bigint 类型的列解析为连续类型 • 若有将 bigint 解析为 categorical 的情况，通过参数 forceCategorical 指定
权重列的列名	可选	记录特征权重的列的字段名
标签列	必选	数据的标签列，即 label

在本案例中，我们只设置训练特征列和标签列。

churn 字段表示用户是否会流失。**Churn** 为 1 表示用户会流失，为 0 表示用户不会流失。

点击“选择字段”按钮，选择除标签（**churn**）以外的其余字段。点击“目标列”的输入框，输入标签列的字段名，系统将进行模糊匹配，用户点击选择目标列。

The image shows two screenshots of the 'Field Settings' (字段设置) interface. The left screenshot shows the 'Select Feature Column' (选择特征列) step, where a red '1' is next to the 'Select Field' (选择字段) button. The right screenshot shows the 'Select Target Column' (选择目标列) step, where a red '2' is next to the 'churn' field in the dropdown menu.



➤ 参数设置：

在参数设置中，设置随机森林的算法参数。具体描述如下表所示：

设置项	是否必选	可选范围	默认值	描述
森林中树的个数	必选	正整数， (0, 1000]	100	指定随机森林中要创建的树的个数
单颗树的算法在森林中的位置	可选	可选，默认算法在森林中均分	-	<p>由于决策树存在多种算法形式，该设置用于指定每个树的算法类型。</p> <ul style="list-style-type: none">如果用户不需要指定，则输入 None，算法在森林中均分。如果用户需要指定，则输入格式如下：比如有 n 棵树，用户输入 a,b，表示[0,a) 是 ID3, [a,b) 是 cart， [b,n) 是 C4.5。 <p>例如：在一个拥有 5 棵树的森林中， [2, 4]表示[0,2)为 ID3 算法， [2,3)为 cart 算法， [4,5)为 c4.5 算法。</p>
单颗树随机特征数	可选	正整数， [1,N],N 为 feature 数	log2N	指定每颗树的特征数。
叶节点数据的最小个数	必选	正整数	2	指定决策树的叶节点数据的最小个数
叶节点	可选	[0,1]	0.0	指定叶节点数据个数占父节点的最小比例。当比例小于输入值时，即停止该父

数据个数占父节点的最小比例				节点的分裂。
单颗树的最大深度	可选	$[1, \infty)$	∞	指定单颗树的最大深度
单颗树输入的随机数据个数	可选	$(1000, 1000000]$	默认 100000	指定森林中单颗树输入的随机数据的个数
随机数种子	可选	$[0, \infty)$	0	单颗树的样本和特征需要随机选取。该设置项指定随机计算中需要的随机数种子。

本案例中，参数设置如下图所示：

字段设置

参数设置

森林中树的个数 正整数，范围 (0, 1000]

100

单颗树的算法在森林中的位置 可选 ?

格式: 2,3。详见问号中描述

单棵树随机特征数 范围 [1, N]，N为feature数

可选, 默认 log2N

叶节点数据的最小个数 正整数

2

叶节点数据个数占父节点的最小比例 [0,1]

0

单颗树的最大深度 可选 范围 [1, ∞)

默认 无穷

单颗树输入的随机数据个数 (1000, 1000000]

100000

随机数种子 可选，范围[0, ∞)，默认0

0

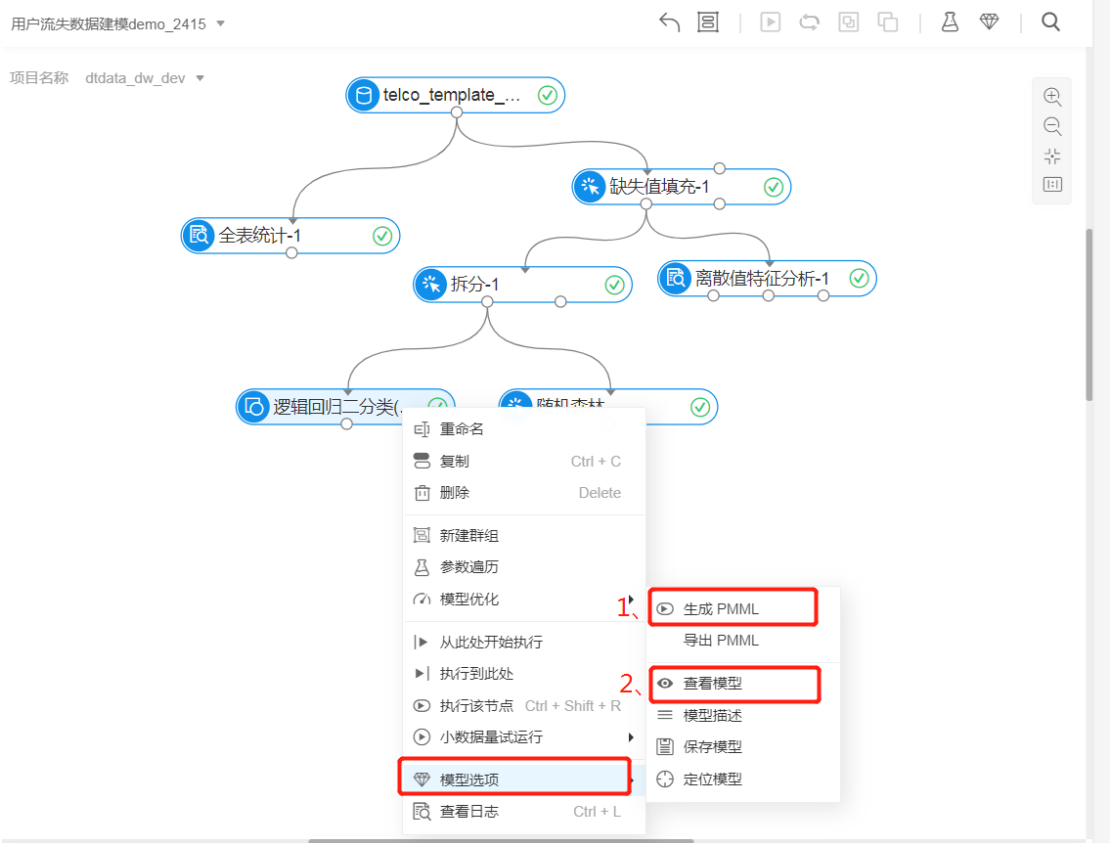
2.6. 模型预测及评估

2.6.1. 查看模型

运行下图中的 DAG，运行成功后可以得到模型文件。选择算法组件，右键点击，在展开的右键菜单中，选择“生成 PMML”，平台将生成该算法对应的 PMML 文件。

在右键菜单中，点击“查看模型”，即可查看模型的详细信息。

右键菜单描述详见 3.x.x 章节。



逻辑回归二分类

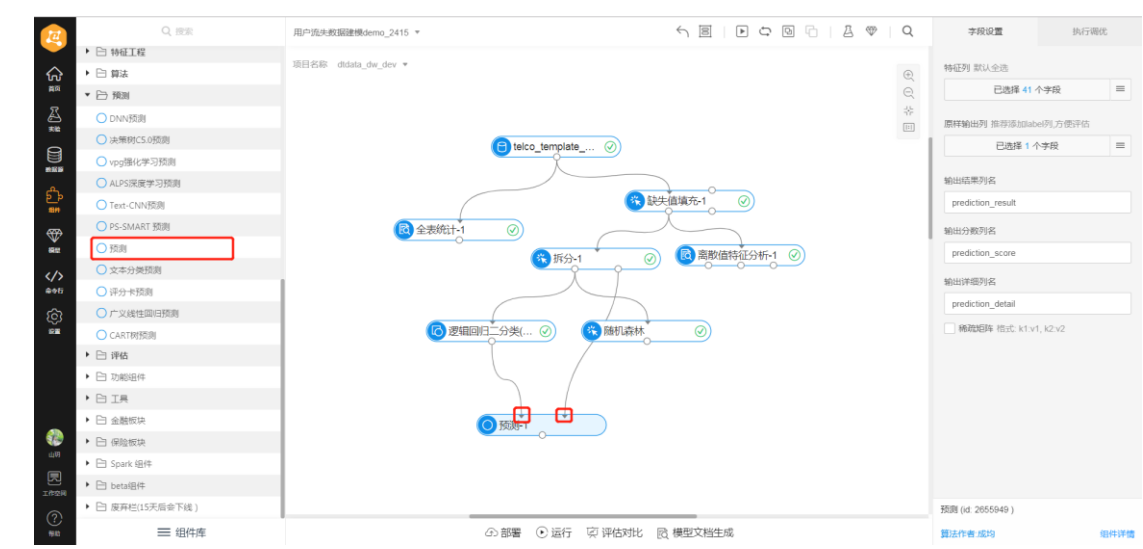
在输入数据为稀疏的时候，不显示 weight 全是 0 的特征

字段名 ▲	权重	
	1 ▲	0 ▲
region	-0.01964648914714521	-
tenure	-0.03740506169914636	-
age	-0.01341764934039284	-
marital	0.03686722810932691	-
address	-0.01766192950862205	-
ed	0.1506115101456351	-
employ	-0.014690362958871	-
gender	0.03008620939892213	-

保存到ODPS: pai_lr_coefficient_xlab_m_logisticregres_2655948_v0 保存 关闭

2.6.2. 预测

训练好的模型，可直接用于数据的预测。拖入“预测”组件，将其与数据连接起来。预测组件第一个输入节点为模型，第二个输入节点为测试集。



- 字段设置：
- 在字段设置中，选择特征列、原样输出列、输出结果列名、输出分数列名、输出详细列名、系数矩阵。

具体描述如下表所示：

设置项	默认值	描述										
特征列	-	用来测试模型的特征列										
原样输出列	-	将测试集中的某些字段，原样输出在结果表中。 一般选择数据的标签列，即 label										
输出结果列名	prediction_result	预测结果列										
输出分数列名	prediction_score	预测结果概率得分； 仅模型为二分类时有效										
输出详细列名	prediction_detail	每个类别的预测概率得分； 仅模型为二分类时有效										
稀疏矩阵	-	<p>输入的测试数据是否具有 key-value 特征。具有 key-value 特征的数据如下：</p> <table><tr><td colspan="2">kv</td></tr><tr><td colspan="2">a:0 b:1 c:2</td></tr></table> <p>平台解析后变为，黑体的为字段名（key），值为字段值（value）：</p> <table><tr><td>1</td><td>2</td><td>3</td></tr><tr><td>a</td><td>b</td><td>c</td></tr></table> <p>数据经过 kv 转换后，能有效的减少存储空间。</p>	kv		a:0 b:1 c:2		1	2	3	a	b	c
kv												
a:0 b:1 c:2												
1	2	3										
a	b	c										

本案例中，字段设置输入如下：

字段设置

执行调优

特征列 默认全选

已选择 41 个字段

≡

原样输出列 推荐添加label列,方便评估

已选择 1 个字段

≡

输出结果列名

prediction_result

输出分数列名

prediction_score

输出详细列名

prediction_detail

☐ 稀疏矩阵 格式: k1:v1, k2:v2

➤ 执行调优：

在执行调优中，设置核数目和每个核的内存大小（MB）。具体描述如下表所示：

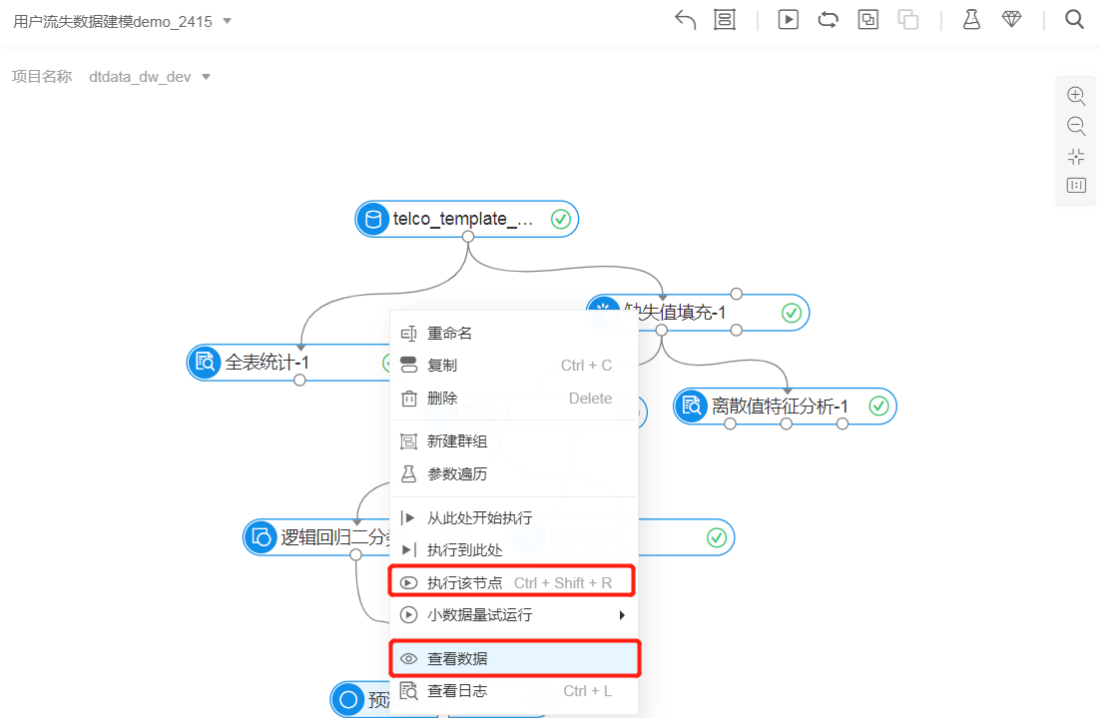
设置项	是否必选	可选范围	默认值	描述
核数目	可选	视集群的具体情况而定	默认自动调整	内存中为该实验分配的内存核数目
每个核的内存大小（MB）	可选	视集群的具体情况而定	默认自动调整	每个内存核的内存大小

本案例中核数目和每个核的内存大小都设置为“默认自动调整”。



在本案例中，用户点击特征列的“选择字段”按钮，选择除标签（churn）以外的其余字段；点击“目标列”的输入框，输入 **churn**，系统将进行模糊匹配，点击选择目标列。

选中预测组件，右键点击，在展开的菜单中，点击“执行该节点”。运行成功后，在右键菜单中点击“查看数据”，即可查看模型在测试集上得出的打分。



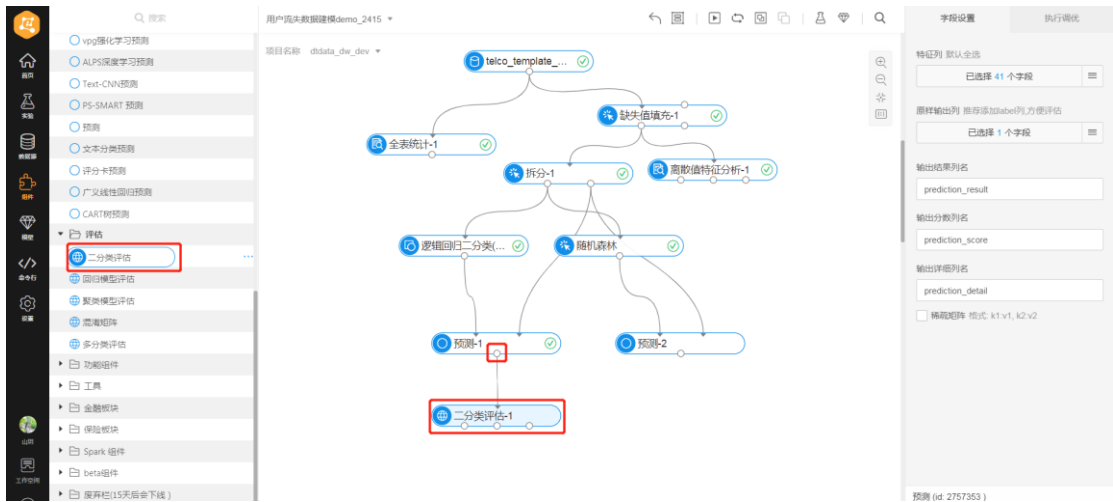
数据探查 - pai_temp_135628_2655949_1 - (默认显示100条)

序号	churn	prediction_result	prediction_score	prediction_detail
1	1	0	0.5676012090563383	{ "0": 0.5676012090563383, "1": ...
2	0	0	0.9769259448405181	{ "0": 0.9769259448405181, "1": ...
3	1	0	0.7677910521396251	{ "0": 0.7677910521396251, "1": ...
4	0	0	0.894274653354411	{ "0": 0.894274653354411, "1": 0...
5	0	0	0.9140403192469182	{ "0": 0.9140403192469182, "1": ...
6	0	0	0.7546431373694349	{ "0": 0.7546431373694349, "1": ...
7	0	0	0.8362015403310399	{ "0": 0.8362015403310399, "1": ...
8	1	0	0.809269115178162	{ "0": 0.809269115178162, "1": 0...
9	0	0	0.6272576978900392	{ "0": 0.6272576978900392, "1": ...
10	0	1	0.5332721111534128	{ "0": 0.4667278888465872, "1": ...
11	1	1	0.7109696941808096	{ "0": 0.2890303058191904, "1": ...
12	1	1	0.7521247654067211	{ "0": 0.2478752345932789, "1": ...
13	0	1	0.5118082925302373	{ "0": 0.4881917074697627, "1": ...

图表分析 显示10000条 显示1024条 复制到CSV 复制到Excel 关闭

2.6.3. 二分类评估

模型预测完毕后，对模型的测试效果进行评估。投入“二分类评估”组件，将预测结果输入到“二分类评估”组件中。



- 字段设置：
- 在字段设置中，选择特征列、原样输出列、输出结果列名、输出分数列名、输出详细列名、系数矩阵。

具体描述如下表所示：

设置项	是否必选	默认值	描述
原始标签列列名	必选	-	测试集中原始标签列的列名

分数列列名	必选	prediction_score	测试集中的预测分数列的列名
正样本的标签值	必选	1	目标样本的标记，如“1”表示用户的状态为流失，“0”表示用户的状态正常。所以在“正类值”中填入目标样本的标签，“1”
计算 KS、PR 等指标时按等频分成多少个桶	必选	1000	计算 KS、PR 等指标时按等频分成多少个组
权重列列名	可选	-	指定测试集中的特征权重列，仅支持 double 和 bigint
分组列列名	可选	-	分组 ID 列，通过不同的分组 ID 区分数据所属的分组，并对各个分组的数据分别计算相关评估指标

本案例中，字段设置输入如下：

字段设置

执行调优

原始标签列列名

churn

分数列列名

prediction_score

正样本的标签值

1

计算KS,PR等指标时按等频分成多少个桶

1000

权重列列名 可选，仅支持double和bigint

分组列列名 可选，仅支持string类型 ?

☐ 高级选项

➤ 执行调优：

在执行调优中，设置核数目和每个核的内存大小（MB）。具体描述如下表所示：

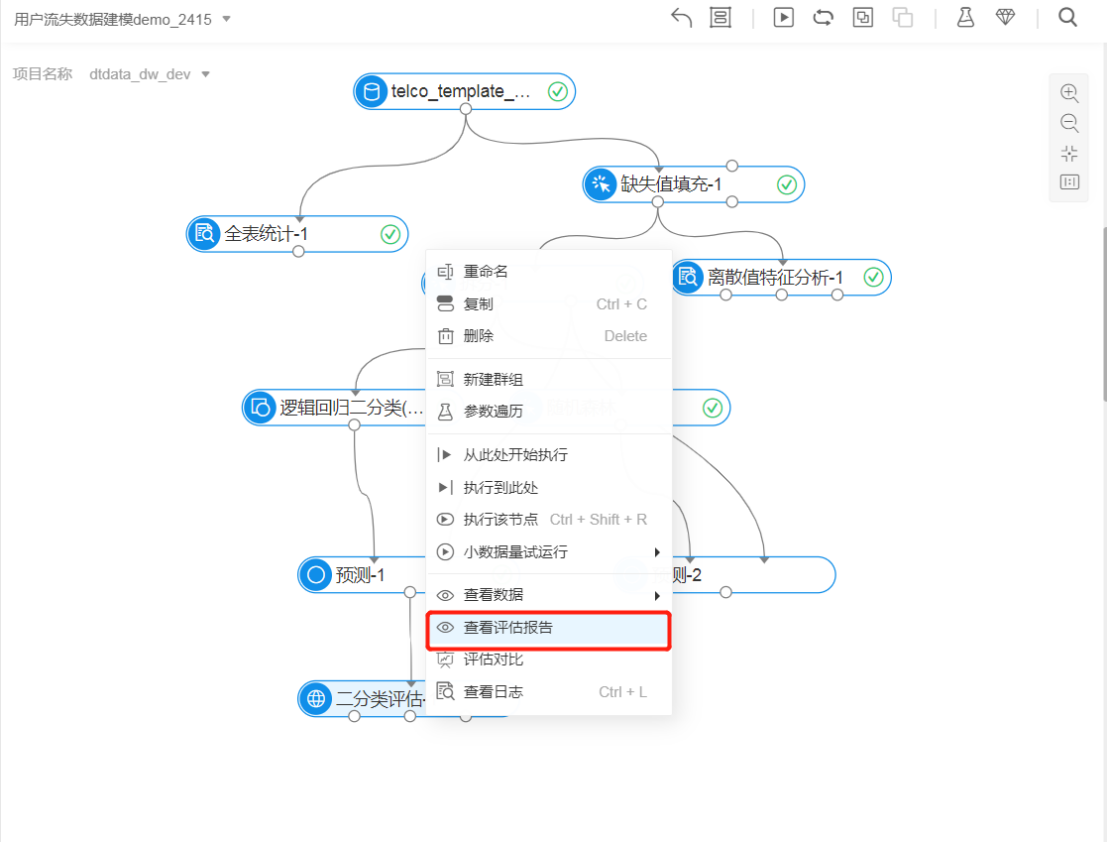
设置项	是否必选	可选范围	默认值	描述
核数目	可选	视集群的具体情况	默认自	内存中为该实验分配的内存核数目

		而定	动调整	
每个核的内存大小（MB）	可选	视集群的具体情况而定	默认自动调整	每个内存核的内存大小

本案例中核数目和每个核的内存大小都设置为“默认自动调整”。



运行该组件，运行成功后，在右键菜单中选中“查看评估报告”。



针对测试集的结果，从 ROC、K-S 曲线、Lift 曲线、PR 曲线四个维度衡量了模型的预测能力，评估报告如下图所示：



评估指标说明：

某池塘有 1400 条鲤鱼，300 只虾，300 只鳖。现在以捕鲤鱼为目的。撒一大网，逮着了 700 条鲤鱼，200 只虾，100 只鳖。那么，评估指标分别如下：

$$\text{正确率} = 700 / (700 + 200 + 100) = 70\%$$

$$\text{召回率} = 700 / 1400 = 50\%$$

$$\text{F 值} = 70\% * 50\% * 2 / (70\% + 50\%) = 58.3\%$$

由此可见，正确率是评估捕获的成果中目标成果所占得比例；召回率，顾名思义，就是从关注领域中，召回目标类别的比例；而 F 值，则是综合这二者指标的评估指标，用于综合反映整体的指标。

评估指标描述详见 x.x.x 章 术语说明。

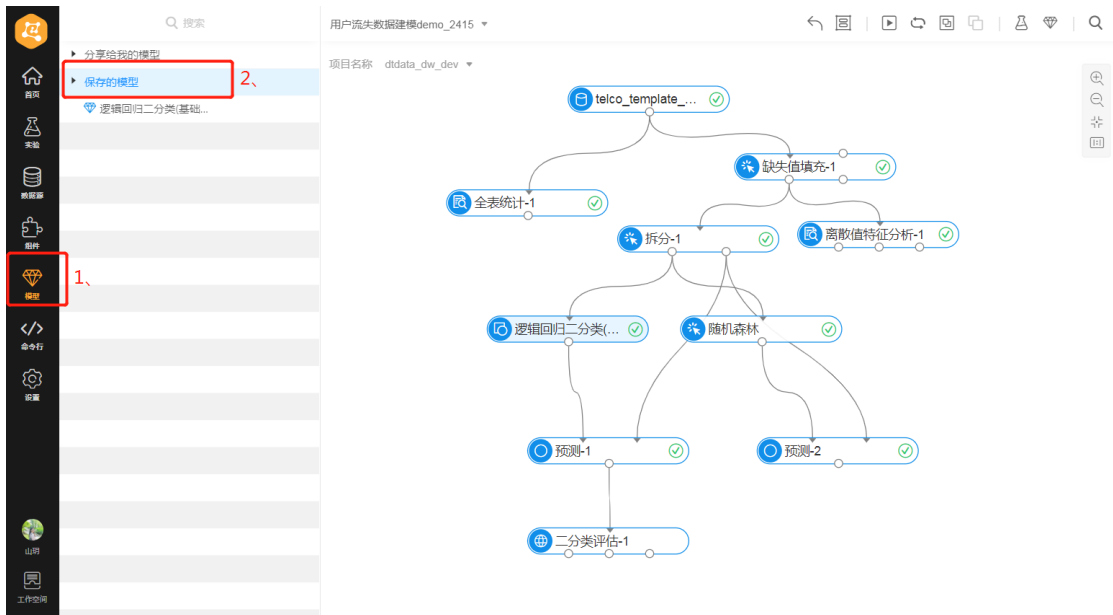
2.7. 结果存储

2.7.1. 保存模型

模型训练、测试完成后，可以保存该模型。选中训练好的模型，右键点击，选择“模型选项—保存模型”。



在左侧的“模型”一级菜单中，查找保存的模型。



3. 深度学习

3.1 TensorFlow-GPU 组件

该部分仅用于 GPU 敏感性模型（比如 tensorflow 模型），而且是独占型任务，其他不合适 GPU 来跑的操作（比如特征预处理）将导致操作时间延长，并影响其他同学使用，请谨慎使用！

相关操作方式请详见：

<https://openclub.alipay.com/read.php?tid=7869&fid=97>

- 1、蚂蚁平台用户指南-ETL数据开发平台（在该平台进行数据查看、处理、变量特征等工作）

download:蚂蚁平台用户指南-ETL数据开发平台.pdf

- 2、蚂蚁平台用户指南-机器学习平台（在该平台进行数据建模、训练、保存等工作）

download:蚂蚁平台用户指南-机器学习平台.pdf

download:蚂蚁平台用户指南-机器学习平台（NLP题）-TensorFlowGPU组件使用说明2.pdf

3.2 其他深度学习组件

点击右键，查看操作详情即可