

**IBM SPSS Analytic Server
V3.1.2**

管理员指南

IBM

注释

在使用本信息及其支持的产品之前，请先阅读第 21 页的『声明』中的信息。

产品信息

此版本是用于 IBM SPSS Analytic Server 的 V3.1.2 以及后续发行版和修订版，直至在新版本中另有说明为止。

目录

第 1 章 租户管理	1	版本信息	15
命名规则.	2	日志收集器	15
第 2 章 用户入门	3	常见问题	16
第 3 章 Analytic Server 作业名	5	性能调整	17
第 4 章 IBM SPSS Analytic Server 最佳 实践和建议.	7	声明	21
第 5 章 故障诊断	15	商标.	22
日志记录	15		

第 1 章 租户管理

租户提供对用户、项目和数据源的高级划分，因此无法在租户间共享对象。每个用户在其分配到的租户上下文中访问系统。

在 Analytic Server 控制台中管理租户并将用户分配到租户。能否查看“租户”页面取决于登录该控制台的用户的角色：

- 在安装期间设置的“超级用户”管理员为租户管理员。只有该用户才能创建新租户并编辑任何租户的属性。
- 具有“管理员”角色的用户可以编辑其登录的租户的属性。
- 具有“用户”角色的用户不能编辑租户属性。对这些用户隐藏“租户”页面。
- 具有“读者”角色的用户无法编辑数据源，甚至无法登录 Analytic Server 控制台。

管理员可以访问“项目”和“数据源”页面，并管理要清除和管理的任何项目或数据源。请参阅 *IBM® SPSS® Analytic Server User's Guide* 以获取更多信息。

租户列表

主“租户”页面显示表中的现有租户。只有“超级用户”管理员才能对此页面进行编辑。

- 单击租户名称以显示其详细信息和编辑其属性。
- 单击租户 URL 以在该租户的上下文中打开控制台。

注：您将从控制台注销并且需要使用租户的有效凭证来登录。

- 在搜索区域输入以对列表进行过滤，从而仅显示其名称中包含搜索字符串的租户。
- 单击新建以使用在添加新租户对话框中指定的名称创建新租户。请参阅第 2 页的『命名规则』，以了解有关可以给租户提供的名称的限制。
- 单击删除以除去所选租户。
- 单击刷新以更新列表。

个别租户详细信息

内容区域划分为多个可折叠部分。

详细信息

名称	这是一个可编辑的文本字段，用于显示租户的名称。
描述	可编辑的文本字段，您可以提供关于该租户的说明性文本。
URL	这是将提供给用户通过 Analytic Server 控制台登录租户，以及用于配置 SPSS Modeler 服务器的 URL。请参阅 <i>IBM SPSS Analytic Server Installation and Configuration Guide</i> 以获取有关配置 SPSS Modeler 的详细信息。
状态	处于活动状态的租户目前正在使用中。使租户处于不活动状态将使用户无法登录该租户，但是不会删除任何底层的信息。

主体

主体是从安装过程中设置的安全提供程序获取的用户和组。您可以将主体添加到租户，作为管理员、用户或读者。

- 在文本框中输入时，会对其名称中包含搜索字符串的用户和组进行过滤。 从下拉列表中选择**管理员**、**用户**或**读者**可以分配用户在该租户中的角色。 单击**添加参与者**以将其添加到作者列表。
- 要除去参与者，请在成员列表中选择用户或组，然后单击**除去参与者**。

度量 使您能够配置租户的资源限制。 报告租户目前使用的磁盘空间。

- 您可以为租户设置最大磁盘空间配额；当达到该限制时，不能再向该租户上的磁盘写入任何数据，直到清理出足够的磁盘空间，以使租户磁盘空间使用量低于配额。
- 您可以为租户设置磁盘空间警告级别；当超过配额时，该租户上的主体不能提交分析作业，直到清理出足够的磁盘空间，以使租户磁盘空间使用量低于配额。
- 您可以设置在该租户上一次能够运行的最大并行作业数；当超过配额时，该租户上的主体不能提交分析作业，直到目前正在运行的作业完成。
- 您可以设置数据源能够具有的最大字段数。 每当创建或更新数据源时，都会检查该限制。
- 您可以设置最大文件大小（以兆字节为单位）。 上传文件时，会检查该限制。

安装提供程序配置

使您能够指定用户认证提供程序。 缺省使用在安装和配置期间设置的缺省租户提供程序。 **LDAP** 允许您使用诸如 Active Directory 或 OpenLDAP 之类的外部 LDAP 服务器来认证用户。 为提供程序指定设置并且可选择指定过滤器设置，以控制"主体"部分中提供的用户和组。

命名规则

对于在 Analytic Server 中可给予唯一名称的一切对象（例如，数据源和项目），以下规则适用于这些名称。

- 在单个租户中，名称必须在同类型的对象中唯一。 例如，两个数据源不能同时命名为 insuranceClaims，但是一个数据源和一个项目可分别命名为 insuranceClaims。
- 名称区分大小写。 例如，insuranceClaims 和 InsuranceClaims 均被视为唯一名称。
- 名称忽略前置和后置空格。
- 以下字符在名称中无效。

~, #, %, &, *, {, }, \, :, <, >, ?, /, |, ", \t, \r, \n

第 2 章 用户入门

告知用户浏览至 `http://<host>:<port>/<context-root>/admin/<tenant>`，并输入其用户名和密码以登录到 Analytic Server 控制台。

<host>

Analytic Server 主机的地址。

<port>

Analytic Server 正在侦听的端口号。缺省情况下，此端口为 9080。

<context-root>

Analytic Server 的上下文根。缺省情况下为 `analyticserver`。

<tenant>

在多租户环境中，这是您所属的租户。在单租户环境中，缺省租户为 `ibm`。

例如，如果主机的 IP 地址为 9.86.44.232，您已创建"mycompany"租户并向其添加了用户，并且其他设置已保留为缺省值，那么用户应该导航至 `http://9.86.44.232:9080/analyticserver/admin/mycompany` 以访问 Analytic Server 控制台。

第 3 章 Analytic Server 作业名

Analytic Server 可生成 map-reduce 和 Spark 作业，可通过 Hadoop 集群的 Resource Manager 用户界面来监控这些作业。

map-reduce 作业名采用以下结构。

AS/{tenant name}/{user name}/{algorithm name}

{tenant name}

这是在其中运行作业的租户的名称。

{user name}

这是请求该作业的用户。

{algorithm name}

这是作业中的主要算法。 请注意，单个流可能生成多个 map-reduce 作业；类似地，单个流中的多个操作可包含在单个 map-reduce 作业中。

所有 map-reduce 作业都显示在 Resource Manager 用户界面中。 每个 Analytic Server 都会启动一个单独的 Spark 应用程序。 打开 Spark 应用程序的用户界面，以监视 Spark 作业（作业名称显示在描述列）。

第 4 章 IBM SPSS Analytic Server 最佳实践和建议

以下部分提供了有关数据源、集群配置和 IBM SPSS Modeler 流的 Analytic Server 最佳实践和建议。

数据源

Analytic Server 支持以下数据源类型：

- 基于文件的数据源，例如，带分隔符的固定文本和 Microsoft Excel 文件。
- 关系数据库，例如，Db2、Oracle、Microsoft SQL Server、Teradata、Postgres、Netezza、MySQL 和 Amazon Redshift。
- Hive/HCatalog 数据源，其中包含所有内置数据类型（例如，ORC 和 Parquet）以及相应的 Hive 序列化器和反序列化器实现可用于的任何定制类型。此外，还可以配置 Analytic Server 以访问 NoSQL 数据库，例如，HBase、MongoDB、Accumulo、Cassandra、Oracle NoSQL 以及相应的 Hive 存储处理程序实现可用于的其他数据库。
- 地理空间类型数据源（基于形状文件和基于地图服务）。

Hive/HCatalog 数据源上的 Analytic Server 限制

- 如果 SPSS Modeler 选择节点需要 Hive 回送，那么过滤表达式只能引用字符串类型的分区列。从 Analytic Server 3.0 开始，为以下分区列添加了数据类型支持：TINYINT、SMALLINT、INT、BIGINT。为 Hive 数据源指定的静态过滤表达式可以具有针对任何数据类型的分区列的过滤表达式。
- Analytic Server 不支持基于 Hive 视图的数据源。

集群配置 - 安全性

Kerberos 模拟

在 V3.0.1 之前，当启用 Kerberos 安全性时，Analytic Server 实例在 Analytic Server 密钥表中使用用户主体名称来认证 HDFS 操作。从 V3.0.1 开始，Analytic Server 在 Analytic Server 密钥表中利用服务主体名称，同时还使用（提出其余请求的用户的）请求用户名称来认证利用 Kerberos 模拟的 HDFS 操作。在 Kerberos 启用的集群中运行时，Analytic Server 3.0.1 或更高版本必须将模拟配置属性添加到 HDFS（或 Hive 服务配置）。对于 HDFS，必须将以下属性添加到 HDFS core-site.xml 文件：

```
hadoop.proxyuser.<analytic_server_service_principal_name> .hosts = *
hadoop.proxyuser.<analytic_server_service_principal_name> .groups = *
```

其中，<analytic_server_service_principal_name> 是在 Analytic Server 配置的 Analytic_Server_User 字段中指定的缺省 as_user 值。

如果通过 Hive/HCatalog 从 HDFS 访问数据，以下属性还必须添加到 HDFS core-site.xml 文件：

```
hadoop.proxyuser.hive.hosts = *
hadoop.proxyuser.hive.groups = *
```

Kerberos 跨域认证

Analytic Server 支持 Kerberos 跨域认证。要启用此功能部件，您必须先确保启用 KDC 跨域认证，然后将以下设置添加到 Analytic Server Ambari 配置的 **Custom analytics.cfg** 部分：

```
kerberos.user.realm.trim = true
```

集群配置 - 性能调整设置和结果

Spark 配置

Analytic Server 使用 `yarn-client` 方式与 YARN 交互并在 Hadoop 集群上运行 Spark 作业。

Analytic Server 定制配置：

- 在 Analytic Server Ambari 配置的定制 `analytics.cfg` 部分中定义 Ambari 设置。
- Cloudera 设置位于 Cloudera Manager 的 `analyticsserver-conf/config.properties` 的 **Analytic Server 高级配置片段（安全值）** 部分中。

1. 请考虑通过在 Analytic Server 定制配置中添加配置项来增大 `spark.driver.memory` 配置设置的值（如果未明确设置，那么缺省值为 1g）。例如：

```
spark.driver.memory=2g
```

2. 从包含 Spark 资源使用情况选项的以下 Analytic Server 中选择一项。

- **选项 A：静态资源分配配置**

Analytic Server 定制配置中有 3 个必须配置的参数：

```
spark.executor.instances  
spark.executor.cores  
spark.executor.memory
```

以下步骤描述如何确定参数值。

- a. 针对 CPU 和内存，建立 Analytic Server 可永久分配给 Spark 的百分比。这将产生可用于每个机器 (M) 的特定核心 (C) 数量和固定内存量。
- b. 建立每个机器都可以运行的执行者 (E) 数量。这些执行者作为单独的 Hadoop 容器（进程）在每个集群节点上运行。大于 2 的值通常比较合适，但是该值必须小于总核心数。由于为 Spark 分配的内存在这些执行者之间进行分配，因此为此参数选择一个高值将减少为每个容器分配的内存量。
- c. 建立每个执行者使用的核心数 (CE)。该值通常为 C/E（为 Spark 应用程序分配的每个机器中的核心数，除以执行者总数）。
- d. 建立用于每个执行者的内存量 (ME)。通常为 M/E。

注：所使用的执行者和核心的数量必须平衡，方法是每个执行者的内存量应该大于 $3G * CE$ 。每个执行者中的每个核心必须至少分配 3G of 内存用作存储或计算内存。

```
spark.executor.instances = <E>*N /<E> // value established in step b where N is the number of compute nodes  
spark.executor.cores = <CE> // value established in step c  
spark.executor.memory = <ME> // value established in step d
```

<code>spark.executor.cores</code>	<input type="text" value="2"/>
<code>spark.executor.instances</code>	<input type="text" value="12"/>
<code>spark.executor.memory</code>	<input type="text" value="12G"/>

图 1. 定制 `analytics.cfg` Spark 设置

- **选项 B：动态资源分配配置**

使用此选项时，YARN 分配的所有执行者将根据整个集群的实际可用资源动态增加/减少。

最小配置为：

```
spark.dynamicAllocation.enabled = true  
spark.shuffle.service.enabled = true
```

典型配置为：

```
spark.default.emitter.class = com.spss.ae.spark.ListEmitter  
spark.default.emitter.compressed = false  
spark.dynamicAllocation.enabled = true  
spark.executor.cores = 4  
spark.executor.memory = 16g  
spark.io.compression.codec = snappy  
spark.rdd.compress = true  
spark.shuffle.service.enabled = true
```

注：

- 不应使用 `spark.executor.instances = <E>`，否则将应用静态资源分配。
- 有关执行者核心和内存值的注意事项与选项 A 相同。

3. 您可以使用以下设置在 Analytic Server 定制配置中禁用 Spark 高速缓存：

```
spark.cache=false  
spark.storage.memoryFraction = 0.3
```



图 2. 定制 *analytics.cfg* Spark 高速缓存设置

使用大型 IBM SPSS Modeler 流时，不应该禁用 Spark 高速缓存。在此实例中禁用 Spark 高速缓存将导致流的运行速度变慢，但是可避免当每个执行者的指定内存量较小时可能会发生内存耗尽的情况。

JVM 配置

Ambari 设置：

1. 在 Analytic Server Ambari 配置中，设置服务器可用于本地处理的内存量。缺省值 (2 GB) 可以安全用于小到中型流，但是较高值堆大小（例如，10 GB）应该用于较大的流。

Analytic Server > 配置 > 高级 analytic-jvm-options

2. 将 `-Xmx2048M` 替换为 `-Xmx10G`，保存配置，并重新启动 Analytic Server。



图 3. 高级 *analytic-jvm-options* 设置

Cloudera 设置：

1. 在 Cloudera Manager 中，浏览到 Analytic Server 服务的配置选项卡并更新 `jvm-options` 控件以设置服务器可用于本地处理的内存量。缺省值 (2 GB) 可以安全用于小到中型流，但是较高值堆大小（例如，10 GB）应该用于较大的流。

Analytic Server 服务 > 配置 > jvm-options

2. 将 `-Xmx2048M` 替换为 `-Xmx10G`，保存配置，并重新启动 Analytic Server。

Yarn MapReduce2 配置:

- 如果您必须与用于 Analytic Server 执行的 Spark 作业并行运行 MapReduce 作业,那么必须配置 YARN 集群以便每个 YARN 容器中至少有 4 GB 内存。

Zookeeper 配置:

- Cloudera 要求您手动更新 Zookeeper 配置。有关更多信息,请参阅 https://www.cloudera.com/documentation/enterprise/5-4-x/topics/cdh_ig_zookeeper_server_maintain.html。
- 如果您使用复杂的 SPSS Modeler 流或宽数据(许多字段),那么因为 Analytic Server-Zookeeper 连接中断,您可能会遇到作业失败的问题。产生该问题的原因是 SPSS Modeler 服务器向 Analytic Server 发送了太大的程序。该问题不太可能出现在 Analytic Server 3.0 (或更高版本)中。请使用以下步骤来解决该问题:

1. 在 Ambari 控制台中,浏览到 Zookeeper 服务配置选项卡,在高级 **zookeeper-env** 下的 zookeeper-env 模板中添加以下行,然后重新启动 Zookeeper 服务。

```
export JVMFLAGS="-Xmx2048m -Djute.maxbuffer=2097152"
```

zookeeper-env template

```
export JAVA_HOME={{java64_home}}
export ZOOKEEPER_HOME={{zk_home}}
export ZOO_LOG_DIR={{zk_log_dir}}
export ZOO_PIDFILE={{zk_pid_file}}
# export SERVER_JVMFLAGS={{zk_server_heapsize}}
export JAVA=$JAVA_HOME/bin/java
export CLASSPATH=$CLASSPATH:/usr/share/zookeeper/*
export JVMFLAGS="-Xmx2048m -Djute.maxbuffer=2097152"
```

图 4. zookeeper-env 模板设置

2. 在 Ambari 控制台中,浏览到 Analytic Server 服务的配置选项卡,将以下内容添加到高级 **analytics-jvm-options**,然后重新启动 Analytic Server 服务。

```
-Djute.maxbuffer=2097152
```

content

```
erride=UTF-8 -XX:+UseParNewGC -Djute.maxbuffer=2097152
```

图 5. 高级 analytics-jvm-options 设置

注:如果问题仍然存在,请在这两个位置将 -Djute.maxbuffer 值从 2097152 增加到 4194304。

IBM SPSS Modeler 流建议

注:以下大部分建议也适用于小数据。

小数据上的原型

对流进行试验时,您通常会添加一些节点,检测传递到该点的流,也可能会添加一个节点来检出一些表格或图形输出,然后继续构建流。每当检验您的流时,您通常不能执行大数据的数据传递。

创建大数据的合适数据样本使您能够针对实际数据检验流，而不会产生在执行完整数据传递时所需的时间消耗。数据样本必须包含足够的的数据，才能成功运行您的流。例如，如果您计划对位于明尼苏达州的商店的交易进行分析，那么您的数据样本必须包含来自明尼苏达州商店的交易。

在采样后，您可以执行以下操作：

- 在大数据所在的集群上创建数据样本的高速缓存，或者

优点 - 简单并且不需要切换源节点

缺点 - 会话结束后，高速缓存消失

- 创建包含数据样本的新 Analytic Server 数据源，或者

优点 - 永久数据源

缺点 - 需要编辑/切换源节点

- 将数据样本下载到本地系统并创建本地数据源

优点 - 在原型设计时不消耗集群资源；使用小数据时，SPSS Modeler 客户机比 Analytic Server 更高效。

缺点 - 需要切换源节点

创建独立于源节点的输入节点和过滤节点

每个 SPSS Modeler 源节点还具有过滤节点和输入节点的组合功能。这对保持画布简化非常有用，但是在切换到不同源节点类型时存在困难。此外，这样还会掩盖正在发生输入操作和过滤操作的事实。

将过滤节点和选择节点放在尽可能靠近源节点的位置

这将减少下游操作中的记录数。

尽可能避开排序节点

Analytic Server 不支持节点中的优化，具体取决于正在排序的数据（如合并节点）。因此，中型流排序节点很少执行任何有用的操作。当排序节点后面紧跟着样本节点时，排序节点才具有值，以便获取前 N 条（或后 N 条）记录。

仅计算将使用的字段

不要计算字段，然后立即对其进行过滤。

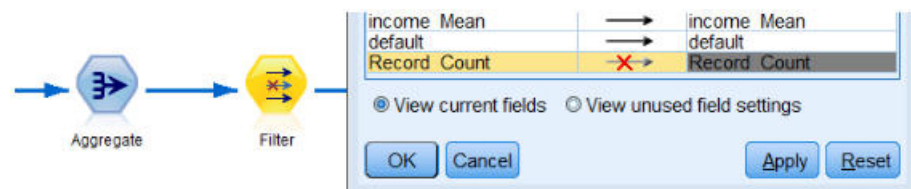


图 6. 建模器字段选项

在任何可能的时候，不执行难以理解的表达式，避免创建大量临时字段。例如，不要定义以下示例：

```
now = datetime_now()
birthdate = datetime_date(bYear, bMonth, bDay)
age = date_years_difference(birthdate, now)
```

而应该定义以下示例：

```
age = date_years_difference(datetime_date(bYear, bMonth, bDay), datetime_now())
```

在转换许多字段时，以这种方式将临时字段调入内联表达式可提高性能。

在数据源中设置存储器

更改字段的存储类型（例如，从字符串改为整数）中型流的操作可降低整体性能。在 Analytic Server 控制台中定义数据源时，您可以为字段设置存储器以避免重复执行这些转换。

在使用小数据时使用 SPSS Modeler

使用 Analytic Server 处理大数据，然后使用 SPSS Modeler 来完成对小数据的计算。

选择适当的 Analytic Server 相关流属性

配置相关流属性（工具 > 选项 > 流属性 > **Analytic Server**），确定是否允许在 Analytic Server 中退出数据处理并在 SPSS Modeler 中继续（在 Analytic Server 中不能运行节点时）。

缺省情况下，配置 SPSS Modeler 以报告错误并在此情况下停止运行。通过将设置从错误更改为警告并调整在 SPSS Modeler 中可以处理多少数据的限制，您可以绕开该错误。例如，您可以从缺省值 10000 条记录更新数据传输率（如果需要）。请注意，当查看使用 SPSS Modeler 表节点的结果时，此限制也适用。如果超出了限制，那么 SPSS Modeler 将报告数据访问超出了在流属性中设置的限制。

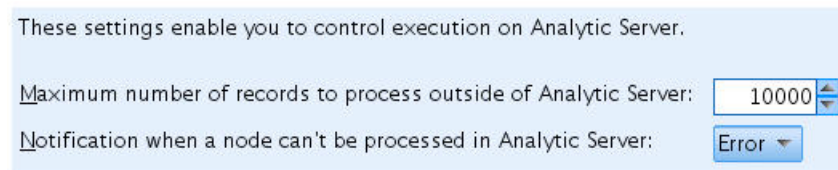


图 7. Analytic Server 设置

使用 Analytic Server 源节点

Analytic Server 可以连接到不同数据库数据源，但是 SPSS Modeler 要求所有源节点都是 Analytic Server 源节点（以使整个流作为 Analytic Server 作业运行）。对于在 Analytic Server 中运行的整个流，必须将数据库源节点切换到 Analytic Server 源节点，必须在 Analytic Server 控制台中创建 Analytic Server 数据库数据源。

考虑如何使用不受支持的节点

Analytic Server 不支持所有节点（转置节点是一个不错的示例）。要将转置操作的结果与流的其余部分合并，并使其在 Analytic Server 中运行，包含转置节点的子流应该写到使用 Analytic Server 导出节点的 Analytic Server 数据源中。然后，您可以将 Analytic Server 源节点附加到中断的流以写入到 Analytic Server 中。

注：转置操作适用于一次性或很少运行的操作，但是不应该用于例行流操作。

在运行流之前，确定流是否将在 Analytic Server 中工作

为在 Analytic Server 中运行准备流之后，选择终端节点并使用 SPSS Modeler 预览功能（工具栏上的预览运行控件）以验证运行终端节点所涉及的任何节点是否将在 Analytic Server 中工作（而不运行流）。在消息窗口中报告了问题。

将背对背的合并操作组合在一起

当一系列合并节点具有相同的密钥和连接类型时，将这些节点合并成一个节点。

合并相同的子流

在可能的情况下尝试合并相同的子流，尤其是当这些子流包含成本很高的操作（例如，合并和排序）时。SPSS Modeler 执行一次这些操作，并使用高速缓存来提高性能。在以下示例中，流与 **newField** 节点相同。

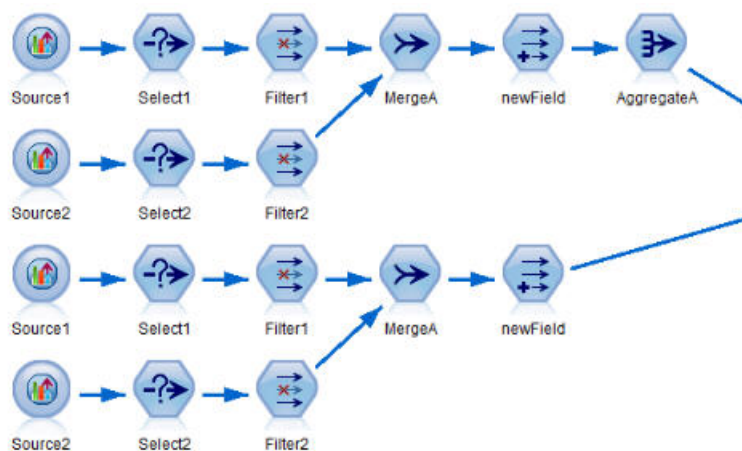


图 8. 示例流

如果按如下所示来构建子流，将更加高效（并且更方便维护）：

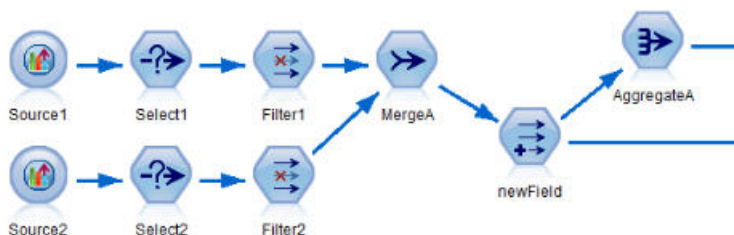


图 9. 示例流

除去额外的输入节点

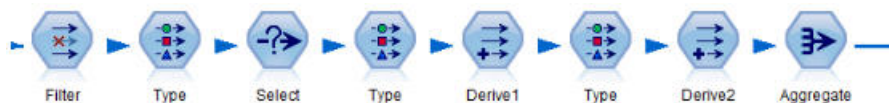


图 10. 示例流

在运行 Analytic Server 时，避开不需要的输入节点。输入节点的读取值操作启动 MapReduce 作业。如果不清除输入节点值，这通常是一次性的节省。

完全记录每个流

以下示例显示了包含许多子流的复杂流。

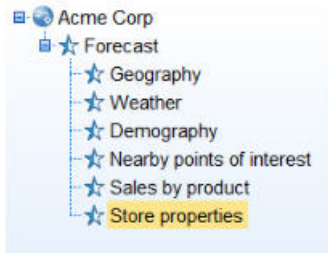


图 11. 子流示例

在此类情况下，正确命名超节点并且（像记录代码一样）记录流非常重要。清晰的注释可以向读取或维护流的其他分析人员提供宝贵的信息。例如：

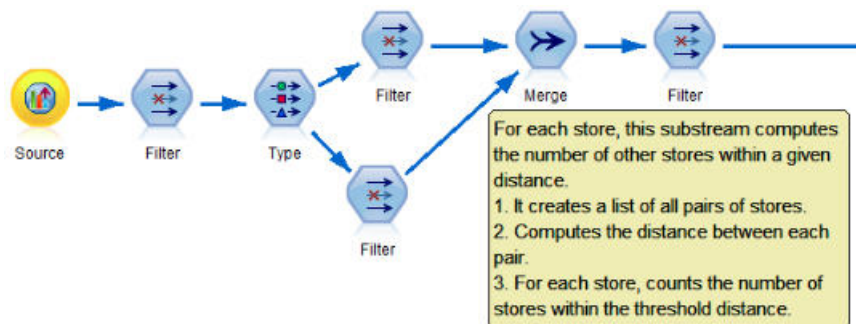


图 12. 包含注释的流示例

当开发流时，使用 SPSS Modeler 高速缓存可以快速存储中间结果

在运行 Analytic Server 的流中，通过将流的特定部分中的数据存储在 HDFS 上的临时文件（与存储在 SPSS Modeler 服务器上相反），节点高速缓存将工作。对于大数据，高速缓存运行正常并且可安全用于在 Analytic Server 上运行的流。

第 5 章 故障诊断

Analytic Server 提供了多种实用工具用于问题确定。

日志记录

Analytic Server 会创建有助于诊断问题的客户日志文件和跟踪文件。通过缺省 Liberty 安装，可以在 {AS_ROOT}/ae_wlpserver/usr/servers/aeserver/logs 目录中找到日志文件。

缺省日志记录配置将生成两个每天轮换的日志文件。

as.log

此文件包含参考警告和错误消息的高级别摘要。发生服务器错误时如果无法通过使用用户界面中显示的错误消息来解决，请首先检查此文件。

as_trace.log

此文件包含来自 ae.log 的所有条目，但添加了额外信息，这些信息主要面向 IBM 支持和开发人员，用于调试目的。

Analytic Server 使用 Apache LOG4J 作为其底层的日志记录工具。通过使用 LOG4J，可通过编辑 {AS_SERVER_ROOT}/configuration/log4j.xml 配置文件来动态调整日志记录。支持人员可能要求您执行此项操作以帮助诊断问题，或者您可能希望修改此文件以限制保留的日志文件数量。对文件进行的更改会在数秒钟内自动完成检测，因此无需重新启动 Analytic Server。

有关 log4j 和配置文件的更多信息，请参阅位于以下地址的 Apache 官方 Web 站点上的文档：<http://logging.apache.org/log4j/>。

版本信息

您可以通过检查 {AS_ROOT}/properties/version 文件夹确定所安装的 Analytic Server 版本。以下文件包含版本信息。

IBM_SPSS_Analytic_Server-*.swtag

包含详细的产品信息。

version.txt

已安装的产品版本和构建号。

日志收集器

无法通过直接复查日志文件来解决问题时，可以将所有日志捆绑在一起并发送给 IBM 支持人员。其中提供了一个实用程序来简化所有必要数据的收集。

通过使用命令 shell，运行以下命令：

```
cd {AS_ROOT}/bin
run >sh ./logcollector.sh
```

这些命令会在 {AS_ROOT}/bin 下创建一个压缩文件。该压缩文件包含所有日志文件和产品版本信息。

常见问题

本节描述了一些常见管理问题以及它们的解决方法。

运行中的流

R 作业将非英语单词转换为 Unicode

在 Cloudera 集群上，如果 Hadoop 服务器的系统编码不是 UTF-8，那么 R 将非英语单词转换为 Unicode。

1. 在 Cloudera Manager 控制台中浏览至 YARN 配置选项卡。
2. 在"NodeManager 环境高级配置片段（安全阀）"字段中添加以下设置。

```
LC_ALL=""
LANG=en_US.utf8
```

PySpark 作业运行失败

确保在所有 Analytic Server 节点和所有节点管理器中部署 Spark 服务。

在启用 Kerberos 的环境下 PySpark 作业运行失败

必须先运行 kinit 命令并随后重新启动 Analytic Server，然后 PySpark 测试才能成功运行。例如：

HDP Kerberos

```
cd /etc/security/keytabs/
sudo -u as_user kinit -k -t as_user.headless.keytab as_user/lyrh1.fyre.ibm.com@IBM.COM
```

CDH Kerberos

```
cd /run/cloudera-scm-agent/process/387-analyticserver-ANALYTIC_SERVER
sudo -u as_user kinit -k -t analyticserver.keytab as_user/cdh12-1.fyre.ibm.com@IBM.COM
```

内存错误

执行者遇到内存错误后配置 YARN

所需的执行者内存超出最大阈值时将发生以下错误：

```
Caused by: com.spss.mapreduce.exceptions.JobException:
  java.lang.IllegalArgumentException: Required executor memory (1024+384 MB) is above the max
  threshold (1024 MB) of this cluster! Please increase the value of
  'yarn.scheduler.maximum-allocation-mb'.
```

以下步骤提供解决该问题所需的 YARN 配置设置。

对于 Ambari

1. 在 Ambari 用户界面中，转至 **YARN > 配置 > 设置**。
2. 将内存节点（为所有 YARN 容器分配的内存）增加到 8192MB。
3. 增加容器值：
 - 容器大小（内存）最小值增加到 682MB
 - 容器大小（内存）最大值增加的 8192MB
4. 将容器大小 (**VCores**) 最大值增加到 3。
5. 重新启动 YARN、Spark 和 Analytic Server 服务。

对于 Cloudera

1. 将 yarn.nodemanager.resource.memory-mb 增加到 8GB
 - 在 Cloudera Manager 用户界面中，转至 **YARN 服务 > 配置 > 搜索容器内存**，然后将该值增加到 8GB。
2. 在 Cloudera Manager 用户界面中，转至 **YARN 服务 > 快速链接**，然后选择动态资源池。
3. 在配置中单击每个可用池的编辑按钮，并在 **YARN** 中将最大运行应用程序数值设置为 4。

4. 重新启动 YARN、Spark 和 Analytic Server 服务。

使用 Apache Spark 2.x 的 Hadoop

- 当 Hadoop 和 Apache Spark 2.x 位于相同环境中时，大多数 forcespark 和 forcehadoop 作业将失败。Yarn 应用程序日志中显示错误：`java.lang.NoClassDefFoundError: org/apache/hadoop/fs/FSDataInputStream`。

可通过如下所示手动编辑 `/etc/spark2/conf/spark-defaults.conf` 文件解决此问题：

```
#spark.hadoop.mapreduce.application.classpath=  
#spark.hadoop.yarn.application.classpath=
```

- 在相同系统上安装两个 JDK 版本时，Cloudera 使用 JDK 1.7，而 Spark 2.x 使用 JDK 1.8。使用 Apache Spark 2.x 运行 forcespark 或 forcehadoop 作业可能导致所有作业失败，并显示以下错误消息：

```
Execution failed. Reason: org/apache/spark/api/java/function/PairFunction : Unsupported major.minor version 52.0
```

对于 Cloudera，在 Cloudera Manager 的 `server.env` 的 **Analytic Server** 高级配置片段（安全值）部分中添加以下行：

```
JAVA_HOME=/usr/java/jdk1.8.0_152
```

性能调整

本部分描述了优化系统性能的方法。

Analytic Server 是 Ambari 框架中使用其他组件（如 HDFS、Yarn 和 Spark）的组件。Hadoop、HDFS 和 Spark 的常见性能调整方法适用于 Analytic Server 工作负载。每个 Analytic Server 工作负载都不相同，因此需要根据特定部署工作负载执行调整试验。以下属性和调整提示是影响 Analytic Server 基准测试和缩放测试结果的主要更改。

当第一个作业在 Analytic Server 上运行时，该服务器启动在关闭 Analytic Server 前将处于活动状态的持久 Spark 应用程序。即使 Analytic Server 作业未主动运行，持久 Spark 应用程序都将分配并保留到在 Analytic Server 运行期间为其分配的所有集群资源。应对分配到 Analytic Server Spark 应用程序的资源量进行仔细思考。如果将所有集群资源分配给 Analytic Server Spark 应用程序，那么可以延迟或不运行其他作业。这些作业可排队等待足够的可用资源，并且这些资源将被 Analytic Server Spark 应用程序使用。

如果已配置和部署多个 Analytic Server 服务，那么每个服务实例都可以分配它自己的持久 Spark 应用程序。例如，如果部署两个 Analytic Server 服务以支持高可用性故障转移，那么您可以看到两个持久 Spark 应用程序处于活动状态，各自分配集群资源。

在某些情况下，Analytic Server 可能会启动将需要集群资源的 map reduce 作业，这样将增加复杂性。这些 map reduce 作业将需要未分配给 Spark 应用程序的资源。需要 map reduce 作业的特定组件为 PSM 模型构建。

可以将以下属性配置为将资源分配给 Spark 应用程序。如果在 Spark 安装的 `spark-defaults.conf` 中进行设置，那么为在环境中运行的所有 Spark 作业对其进行分配。如果在 Analytic Server 配置中将其设置为“定制 `analytic.cfg`”部分下的定制属性，那么仅为 Analytic Server Spark 应用程序对其进行分配。

spark.executor.memory

每个执行者进程要使用的内存量。

spark.executor.instances

要启动的执行者进程的数量。

spark.executor.cores

每个执行者进程的执行者工作者线程数量。 此值应该介于 1 和 5 之间。

设置这三个主要 Spark 属性的示例。 HDFS 集群中有 10 个数据节点，每个数据节点有 24 个逻辑核心和 48 GB 内存并且只运行 HDFS 进程。 此处是为此环境配置属性的一种方法，假设您在此环境中只运行 Analytic Server 作业并期望向单个 Analytic Server Spark 应用程序的最大分配。

- 设置 `spark.executor.instances=20`。 这样将尝试在每个数据节点运行 2 个 Spark 执行者进程。
- 设置 `spark.executor.memory=22G`。 这样会将每个 Spark 执行者进程的最大堆大小设置为 22 GB，每个数据节点上分配 44 GB。 其他 JVM 和 OS 需要额外的内存。
- 设置 `spark.executor.cores=5`。 这样将为每个 Spark 执行者提供 5 个工作者线程，每个数据节点总共有 10 个工作者线程。

监视用于运行作业的 Spark UI

如果您看到溢出至磁盘，这可能会影响性能。 可能的解决方案是：

- 通过 `spark.executor.memory` 增加内存并将其分配给 Spark 执行者。
- 减少 `spark.executor.cores` 的数量。 这样将减少分配内存的并发工作线程的数量，但是还将减少并行作业的数量。
- 为 Spark 更改 Spark 执行者堆的 Spark 内存属性 `spark.shuffle.memoryFraction` 和 `spark.storage.memoryFraction` 分配百分比。

确保名称节点有足够内存

如果 HDFS 中的块数较大并继续增长，请确保名称节点堆增大以适应此增长。 这是常见的 HDFS 调整建议。

更改用于高速缓存的内存量

缺省情况下，`spark.storage.memoryFraction` 的值为 0.6。 如果数据的 HDFS 块大小为 64 MB，那么该值可以增大到 0.8。 如果输入数据的 HDFS 块大小大于 64 MB，那么仅当为每个任务分配的内存大于 2 GB 时，才可以增大此值。

模型评分的调整性能

通过执行以下步骤，您可以使用 Apache Spark 引擎提高大数据集上的模型评分作业的性能。 请注意，这些步骤不能影响集群上非 Analytic Server 服务的操作。

1. 检查是否已在集群中的每个节点上安装 `libtcmalloc_minimal.so{/version}`。

```
whereis libtcmalloc_minimal.so.*
```

2. 如果未安装 `libtcmalloc_minimal.so`，那么安装特定于操作系统的软件包（其中包含集群中每个节点上的 `libtcmalloc_minimal` 库），或手动构建和安装 `libtcmalloc_minimal`。 例如：

Ubuntu：

```
sudo apt-get install libgoogle-perftools-dev
```

Red Hat Enterprise Linux 6.x (x64)：

- a. 安装适用于 RedHat 的 EPEL 存储库（如果尚未安装）

```
wget http://dl.fedoraproject.org/pub/epel/6/x86_64/epel-release-6-8.noarch.rpm
sudo rpm -Uvh epel-release-6*.rpm
```

- b. `sudo yum install gperftools-libs.x86_64`

手动构建：

- a. 从链接 <https://github.com/gperftools/gperftools/releases> 下载 `gperftools-2.4.tar.gz`
 - b. `tar zxvf gperftools-2.4.tar.gz`
 - c. `cd gperftools-2.4`
 - d. `./configure --disable-cpu-profiler --disable-heap-profiler --disable-heap-checker --disable-debugalloc --enable-minimal`
 - e. `make`
 - f. `sudo make install`
3. 请注意已安装的库文件 `libtcmalloc_minimal.so{.version}` 的某个位置，如在一个或多个节点上运行以下命令所返回的。

```
whereis libtcmalloc_minimal.so.*
```

如果集群具有同时运行多个操作系统的节点，那么此文件可能存在多个位置。

4. 在 Ambari 控制台中，转至 Analytic Server 配置并在“定制 `analytics.cfg`”部分下，使用库位置作为值来配置主要的 `spark.executorEnv.LD_PRELOAD`。在做出此更改后，重新启动 Analytic Server 服务。例如，如果将库安装到 `/usr/lib64/libtcmalloc_minimal.so.4`，那么配置将为：

```
spark.executorEnv.LD_PRELOAD=/usr/lib64/libtcmalloc_minimal.so.4
```

如果需要多个位置，请使用空格进行分隔，如下示例所示。

```
spark.executorEnv.LD_PRELOAD=/usr/lib64/libtcmalloc_minimal.so.4 /usr/lib/libtcmalloc_minimal.so
```

如果任何节点都未将 `libtcmalloc_minimal.so` 库安装在某个已配置的位置，虽然这样不会导致错误，但是这些节点上模型评分的性能可能会降低。

Spark map 端连接

Analytic Server Spark 连接实现不支持 map 端连接功能（Spark 连接主要是 reduce 端）。当一个输入很小时，该实现不利用 map 端连接来优化连接。不利用 map 端连接会导致 Spark 作业的资源非常密集，从而最终失败。

要在运行 Analytic Server Spark map 端连接（或基于最小的 RDD 大小的本机 Spark 作业）时优化连接，您可以将 `spark.msj.maxBroadcast` 属性添加到 `analytics.cfg` 文件（SPSS Analytic Server/Configs/Custom `analytics.cfg`）或 `analytics-meta`。

声明

本信息是为在美国国内供应的产品和服务而编写的。可以从 IBM 获取本资料的其他语言版本。但是，您必须拥有该语言的产品或产品版本副本，才能访问对应语言的资料。

IBM 可能在其他国家或地区不提供本文档中讨论的产品、服务或功能特性。有关您所在区域当前可获得的产品和服务的信息，请向您当地的 IBM 代表咨询。任何对 IBM 产品、程序或服务的引用并非意在明示或默示只能使用 IBM 的产品、程序或服务。只要不侵犯 IBM 的知识产权，任何同等功能的产品、程序或服务，都可以代替 IBM 产品、程序或服务。但是，评估和验证任何非 IBM 产品、程序或服务的操作，由用户自行负责。

IBM 可能已拥有或正在申请与本文档内容有关的各项专利。提供本文档并不意味着授予用户使用这些专利的任何许可。您可以用书面形式将许可查询寄往：

IBM Director of Licensing
IBM Corporation
North Castle Drive, MD-NC119
Armonk, NY 10504-1785
US

有关双字节 (DBCS) 信息的许可查询，请与您所在国家或地区的 IBM 知识产权部门联系，或用书面方式将查询寄往：

Intellectual Property Licensing
Legal and Intellectual Property Law
IBM Japan Ltd.
19-21, Nihonbashi-Hakozakicho, Chuo-ku
Tokyo 103-8510, Japan

International Business Machines Corporation"按现状"提供本出版物，不附有任何种类的（无论是明示的还是默示的）保证，包括但不限于默示的有关非侵权、适销和适用于某种特定用途的保证。某些管辖区域在某些交易中不允许免除明示或默示的保证。因此本条款可能不适用于您。

本信息可能包含技术方面不够准确的地方或印刷错误。本信息将定期更改；这些更改将编入本信息的新版本中。IBM 可以随时对本出版物中描述的产品和/或程序进行改进和/或更改，而不另行通知。

本信息中对任何非 IBM Web 站点的引用都只是为了方便起见才提供的，不以任何方式充当对那些 Web 站点的保证。那些 Web 站点中的资料不是 IBM 产品资料的一部分，使用那些 Web 站点带来的风险将由您自行承担。

IBM 可以按它认为适当的任何方式使用或分发您所提供的任何信息而无须对您承担任何责任。

本程序的被许可方如果要了解有关程序的信息以达到如下目的：(i) 使其能够在独立创建的程序和其它程序（包括本程序）之间进行信息交换，以及 (ii) 使其能够对已经交换的信息进行相互使用，请与下列地址联系：

IBM Director of Licensing
IBM Corporation

North Castle Drive, MD-NC119
Armonk, NY 10504-1785
US

只要遵守适当的条件和条款，包括某些情形下的一定数量的付费，都可获得这方面的信息。

本文档中描述的许可程序及其所有可用的许可资料均由 IBM 依据 IBM 客户协议、IBM 国际程序许可协议或任何同等协议中的条款提供。

此处引用的性能数据和客户示例仅用于描述目的。实际性能可能因特定配置和操作条件而异。

涉及非 IBM 产品的信息可从这些产品的供应商、其出版说明或其他可公开获得的资料中获取。IBM 没有对这些产品进行测试，也无法确认其性能的精确性、兼容性或任何其他关于非 IBM 产品的声明。有关非 IBM 产品性能的问题应当向这些产品的供应商提出。

关于 IBM 未来方向或意向的声明都可随时更改或收回，而不另行通知，它们仅仅表示了目标和意愿而已。

所有 IBM 的价格均是 IBM 当前的建议零售价，可随时更改而不另行通知。经销商的价格可与此不同。

本信息仅用于规划的目的。在所描述的产品上市之前，此处的信息会有更改。

本信息包含在日常业务操作中使用的数据和报告的示例。为了尽可能完整地说明这些示例，示例中可能会包括个人、公司、品牌和产品的名称。所有这些名称都是虚构的，若实际人员或企业与此相似，纯属巧合。

版权许可证：

本信息包含在日常业务操作中使用的数据和报告的示例。为了尽可能完整地说明这些示例，示例中可能会包括个人、公司、品牌和产品的名称。所有这些名称都是虚构的，若实际人员或企业与此相似，纯属巧合。

凡这些实例程序的每份拷贝或其任何部分或任何衍生产品，都必须包括如下版权声明：

© (贵公司的名称) (年)。此部分代码是根据 IBM Corp. 公司的样本程序衍生出来的。

© Copyright IBM Corp. (输入年份)。All rights reserved.

商标

IBM、IBM 徽标和 ibm.com 是 International Business Machines Corp.，在全球许多管辖区域的商标或注册商标。其他产品和服务名称可能是 IBM 或其他公司的商标。最新的 IBM 商标列表可以在 Web 上的 "Copyright and trademark information" 中获取，地址为：www.ibm.com/legal/copytrade.shtml。

Adobe、Adobe 徽标、PostScript 以及 PostScript 徽标是 Adobe Systems Incorporated 在美国和/或其他国家或地区的注册商标或商标。

IT Infrastructure Library 是 Central Computer and Telecommunications Agency 的注册商标，该企业现已成为 Office of Government Commerce 的一部分。

Intel、Intel 徽标、Intel Inside、Intel Inside 徽标、Intel Centrino、Intel Centrino 徽标、Celeron、Intel Xeon、Intel SpeedStep、Itanium 和 Pentium 是 Intel Corporation 或其子公司在美国和其他国家或地区的商标或注册商标。

Linux 是 Linus Torvalds 在美国和@3B72其他国家或地区的注册商标。

Microsoft、Windows、Windows NT 以及 Windows 徽标是 Microsoft Corporation 在美国和/或其他国家或地区的商标。

ITIL 是一个注册商标，是 Minister for the Cabinet Office 的共同体注册商标，并且已在 U.S. Patent and Trademark Office 进行注册。

UNIX 是 The Open Group 在美国和/或其他国家或地区的注册商标。

Cell Broadband Engine 是 of Sony Computer Entertainment, Inc. 在美国和/或其他国家或地区的商标并且在当地许可证下使用。

Linear Tape-Open、LTO、LTO 徽标、Ultrium 和 Ultrium 徽标是 HP、IBM Corp 和 Quantum 在美国和其他国家或地区的商标。



Printed in China