# WeData Data Development Platform

# Practice Tutorial

# Contents

# Practice Tutorial
# Integrating WeData Platform to Execute Periodic Scheduling Tasks

Last updated：2025-04-09 14:29:15

## Background

This document introduces how to collocate Tencent Cloud Data Development and Governance Platform (Wedata) and TI-kit CLI tool to perform periodic scheduling of tasks. This feature is applicable to business scenarios that have both data governance needs and machine learning needs. You can perform unified scheduling of data development tasks and machine learning tasks on the Wedata page.

## Operation Process

You can follow the following practice process to integrate TI-ONE and Wedata platform (WeData, a one-stop data development and governance platform located in the cloud. For details, see Overview ). Currently, TI-ONE's machine learning task scheduling capability only supports the enterprise edition of Wedata in Guangzhou region.

### Step 1 Preparations

1. Create a user and a project
   Within the Wedata product, you need to first create a user and a project. For details, see Project Creation and Member and Role Management .

2. Configure a custom scheduling resource group
   To enable the TI-ONE integration feature, you need to configure the enterprise edition custom scheduling resource group first. For details, see List of Custom Scheduling Resource Groups .

### Step Two Initialize Environment

Add an execution resource group in the Wedata project space. After adding a server, you need to install Wedata Agent and TI-ONE CLI machine learning environment simultaneously. After logging in to the machine and completing the installation, you can see that the status of the resource group node is normal.



### Step Three Add Data Source

Add a data source in the Wedata project space. Add an HDFS or HIVE data source. Test connectivity. Pay attention that after creation, you need to authorize it to the projects that need to use it. Please refer to Data Source Management for data source creation.
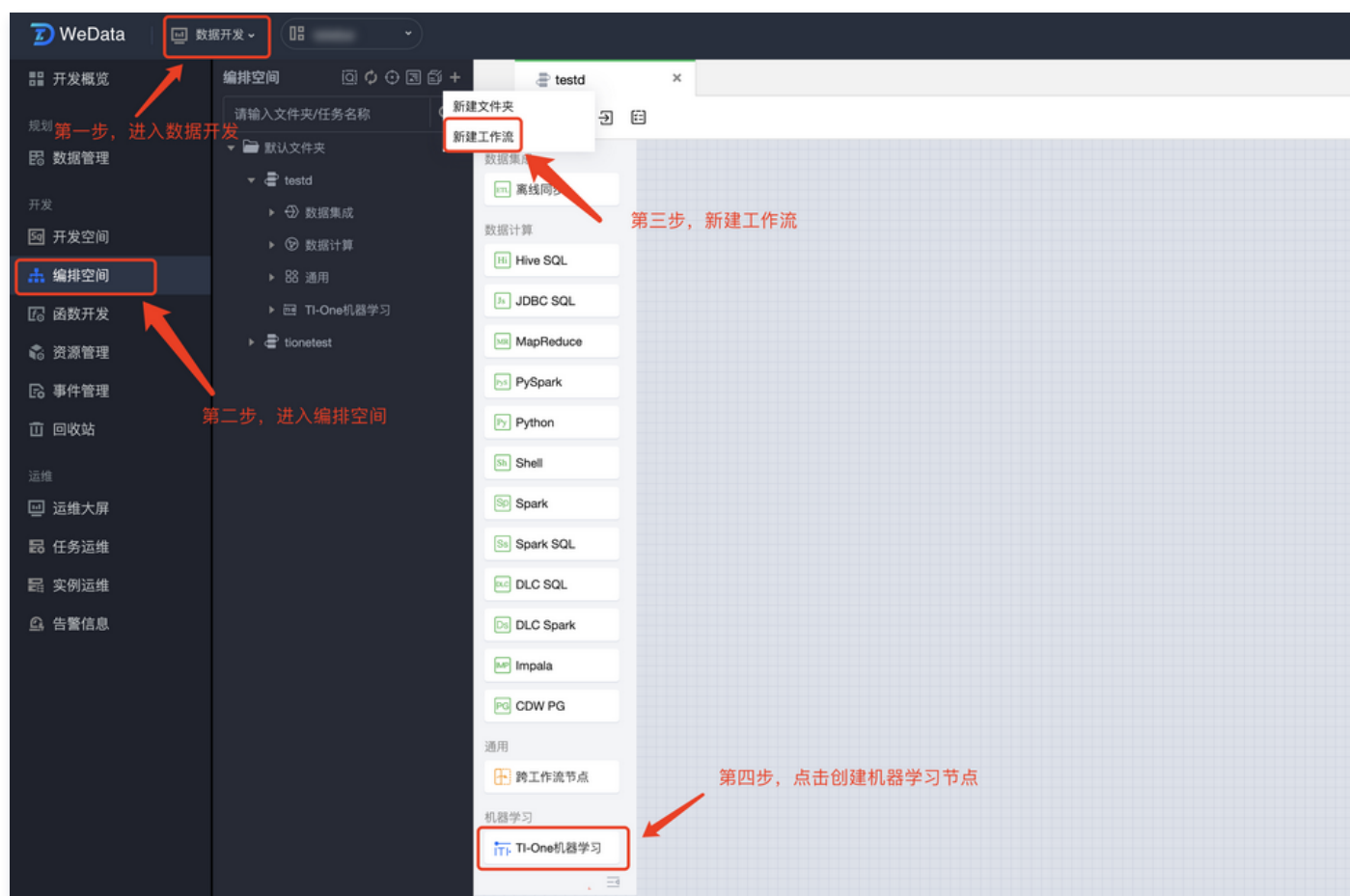
## Step Four Machine Learning Node Configuration

1. Enter **Data Development** > **Orchestration Space**, create a workflow. In the workflow orchestration panel, click to create a **TI-ONE Machine** learning node.

2. The Machine Learning Node in Wedata is essentially a Shell node with the Tikit execution environment for machine learning tasks installed. Users need to write Tikit commands in this node to schedule TIONE computing power for submitting training tasks.

3. After entering the node configuration page, click **Machine Learning Attribute** to configure the data source and algorithm development. Among them, for data source configuration, you can drop down to select the data source associated with the current training task (if the Machine Learning Node is upstream connected to other nodes, the upstream parent task data source can be displayed below). After dropping down, the data source ID will be displayed, which can be used for script development and training task submission.

4. Before submitting a training task, we need to prepare the training code. TIONE provides a lightweight and convenient interactive development environment, Notebook. You can click on the right to enter TIONE Notebook for code writing. (After navigating to the TIONE Notebook instance creation page, the network information of the selected data source will be carried by default. If the data source is HDFS, it will also be selected by default in the data catalog.) If the current machine learning task is associated with a certain Notebook instance, it can be directly dropped down to select. The page will display a quick jump link and the instance running status.

## Step 5 Write a Training Task Submission Command Using TICLI

1. After entering the Machine Learning Node, execute `tikit init --secretid=xxx --secretkey=xxx` for initialization before use. secretId and secretKey are the access keys of Tencent Cloud. Method for obtaining: Enter the console, click on the avatar in the upper right corner, and enter **CAM** > **API Key Management** to obtain.

2. Before using, input tikit –h to get the running modes of each command of the tikit CLI tool.

3. Submit tasks according to the currently required task type. Command testing can be run on the current shell node. After task submission, the corresponding TI-ONE task URL can be printed in the running log. You can go to the TI-ONE Console to view the training task details.

## Step 6 Submit Workflow for Periodic Scheduling

After the workflow development is completed, you can configure the workflow periodic scheduling parameters and submit the overall workflow. After submission, you can view the workflow and tasks in the **Task Ops** module. Once the periodic instance is generated, you can view the instance details on the **Instance Ops** page. For detailed operation guide related to scheduling, please refer to Task Ops.

| testtina | × | TH 模型训练 | TH test | testd | test1 | Ss 特征 | ▼ |

💾 ✈ 🔄 ⤵ ▤

数据集成

ETL 离线同步

数据计算

Hi Hive SQL

Js JDBC SQL

MR MapReduce

Ps PySpark

Py Python

Sh Shell

Sp Spark

Ss Spark SQL

DLC DLC SQL

Ds DLC Spark

MP Impala

PG CDW PG

通用

跨工作流节点

**统一调度** ×

通用设置

统一调度

历史记录

ℹ 对工作流下所有任务设置统一的调度配置，支持常规和crontab方式，常规方式可对工作流下的任务调度配置进行单独修改，crontab方式不支持对任务调度配置进行单独修改，设置后原来的配置将会被覆盖，请谨慎操作！

**调度策略**

| 配置方式 ℹ | ● 常规　　○ crontab |
| 调度周期 | ● 周期　　○ 一次性 |

| 天 ▼ | **周期说明** |

生效时间　　2022-11-01　～ 2099-12-31　📅

执行时间　　00:00　　　　　　　　　🕐

调度计划　　每天00:00执行一次

自依赖　　　○ 并行　● 无序串行　○ 有序串行

工作流自依赖 ℹ　○ 是　● 否

[重置]　[保存]

节点：特征处理2

节点：模型训练
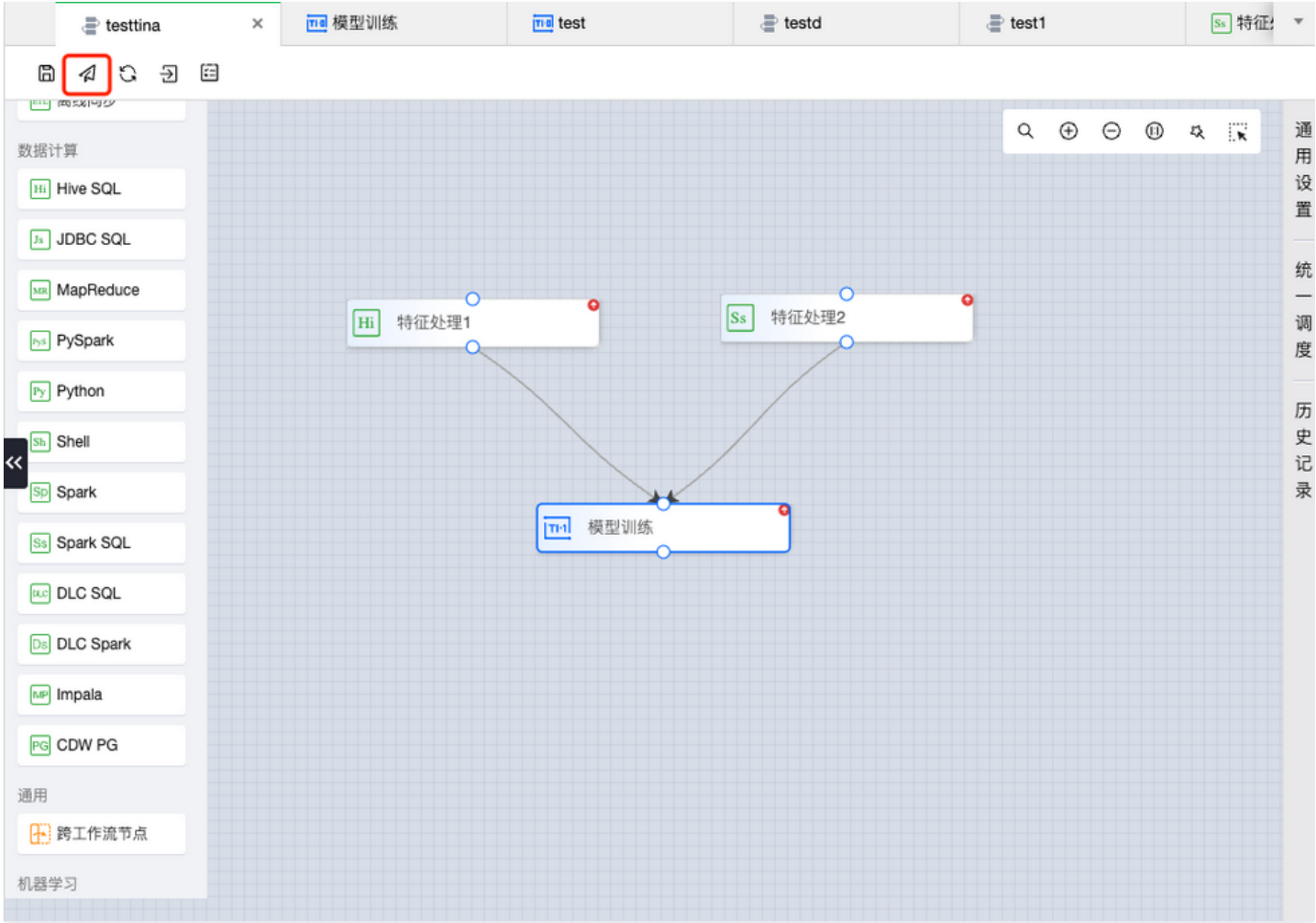
# Using Custom Templates for Data Quality

Last updated：2025-04-09 14:30:09

## Background

Tencent Cloud Data Development and Governance Platform Wedata's data quality supports creating custom templates and batch management, helping you customize table quality inspection logic based on business scenarios. This document introduces how to create rule templates through the custom template page and create detection rules for tables on the data monitoring page according to the custom rule templates.

## Operation Process



## Step 1 Preparations

1. Create a user and a project

   Within the Wedata product, you need to first create a user and a project. For detailed operation guide, check **Preparations**.

2. Create a scheduling resource group

3. Running quality inspection tasks requires creating a scheduling resource group. For detailed operation guide, check **Scheduling Resource Group**.

## Step Two Create Custom Template

1. Enter **Data Quality** > **Rule Template**, click **Custom Template,** add a template and **save**.
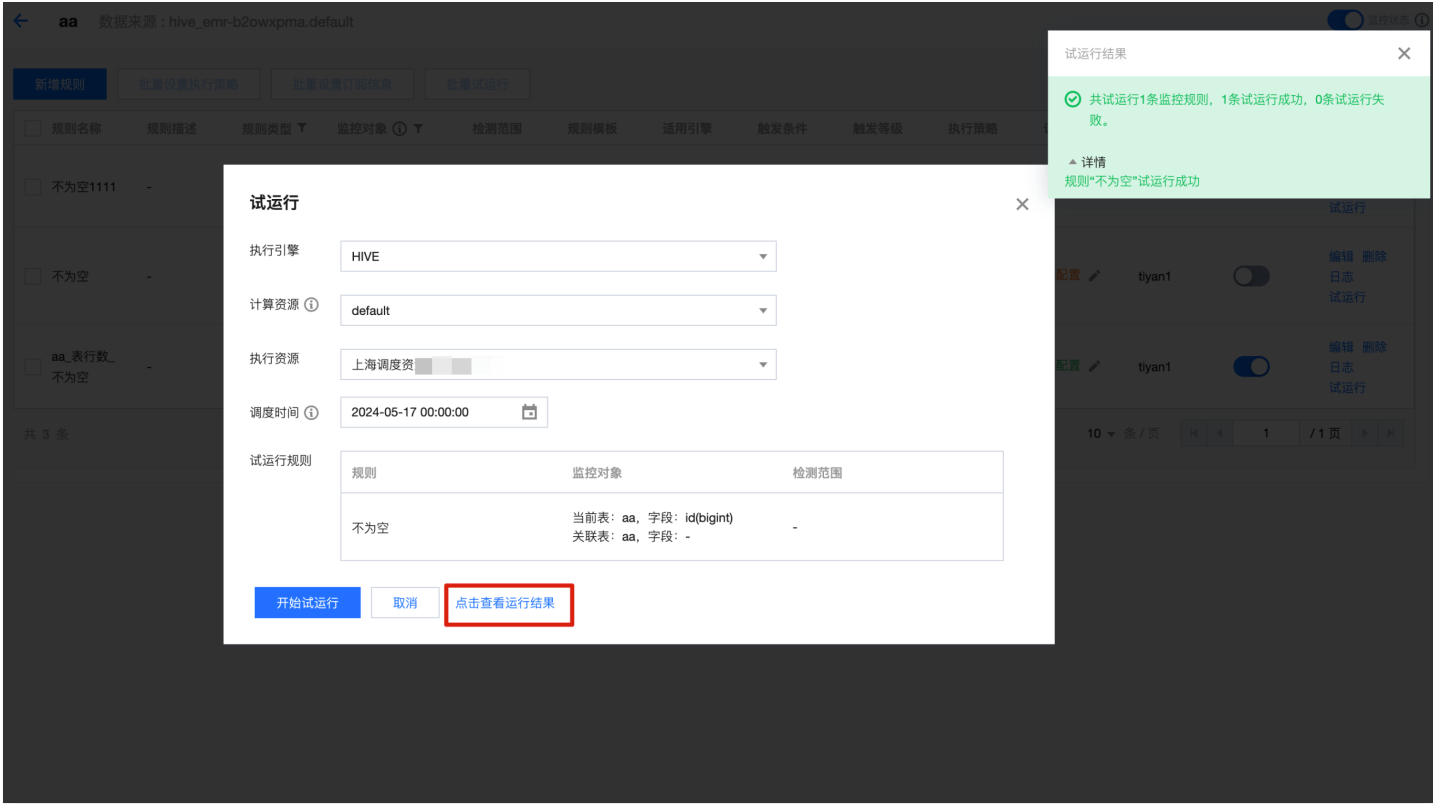
**SQL expression:**

```
select count(table1.${table_1.column_2}) AS count
from ${table_1} table1
join ${table_2} table2
on table1.${table_1.column_1} = table2.${table_2.column_1}
where table1.${table_1.column_3} >= ${param_1} and table1.${table_1.column_3} <= ${param_2}
and table2.${table_2.column_2} >= ${param_3} and table2.${table_2.column_2} <= ${param_4};
```

**Explanation:**

- Two tables appear in the previous context: ${table_1} and ${table_2}.
  - ${table_1} indicates the primary table scanned by the monitoring rule;
  - ${table_2} refers to other tables in the same data source and database (you can also choose the primary table itself in actual use);
- Four fields of Table 1 are used, respectively:
  - ${table_1.column_1}: used for association with Table 2;
  - ${table_1.column_2}: used for result counting;
  - ${table_1.column_3}: used for filtering conditions, greater than or equal to Parameter 1, less than or equal to Parameter 2;
  - ${table_1.column_4}: represents the partition field of Table 1, which can save computing resources significantly and avoid scanning full data;
- Two fields of Table 2 are used, respectively:
  - ${table_1.column_1}: used for association with Table 1;
  - ${table_1.column_2}: used for filtering conditions, greater than or equal to Parameter 3, less than or equal to Parameter 4;
- Used 4 where parameters, which are:

○ ${param_1}: minimum value of Field 3 in Table 1 in SQL;

○ ${param_2}: maximum value of Field 3 in Table 1 in SQL;

○ ${param_3}: minimum value of Field 2 in Table 2 in SQL;

○ ${param_4}: maximum value of Field 2 in Table 2 in SQL.

- Final calculation result: the count of eligible Field 2 in Table 1, a number.

**Screenshot example:**



## Step 3 Create a Quality Rule

1. 1. Enter Data Monitoring, find the table to be monitored, and click Configure Monitoring Task.



2. Click **Add Rule**, select Custom Template for the rule type, select the newly created template, choose database and table parameters and where parameter based on the template variables, configure the trigger conditions and level, and click **Save**.
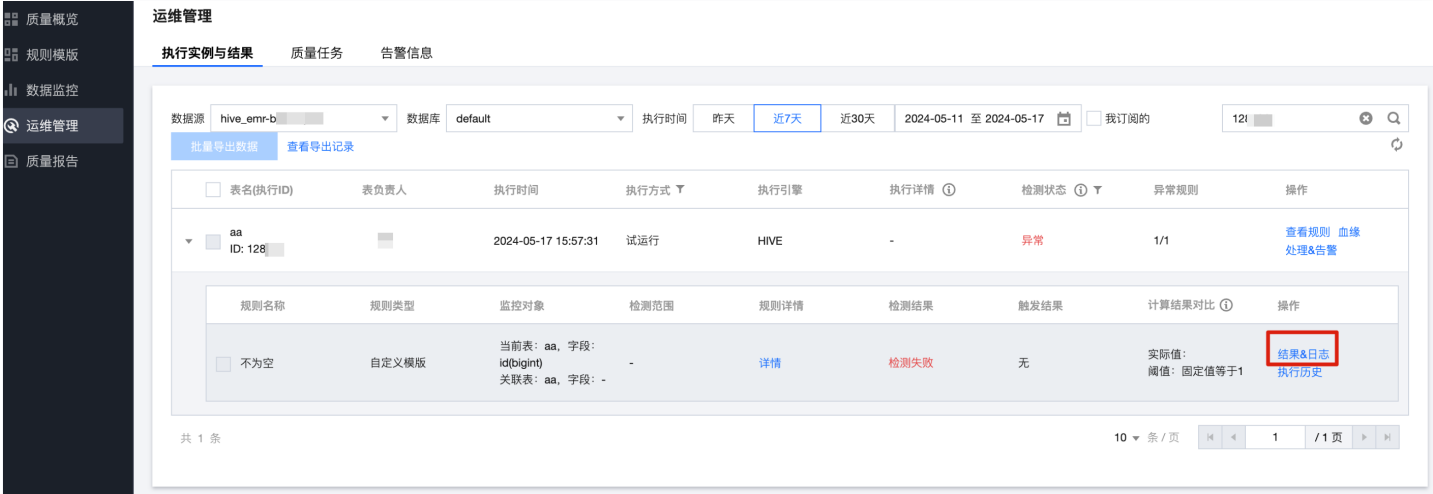
> ⚠️ **Notes:**
>
> Please first analyze what each field means before using a custom template and then map them.

## Step 4 Test-Run

1. 1. Click **trial run**, select an execution engine, computational resource, and execution resource, and select the rule just created in the validation rules.

2. 1. Click **view execution results**, navigate to **Ops management** page to view execution results.



3. 1. Click **Results & Logs** to view running logs.

Among them, EXECUTING SQL: xxxxxx prints the SQL submitted to the hive/spark/dlc engine for quality inspection.