Machine Learning

# ML Discovery & Exploration

**Date published: 2020-07-16**
**Date modified: 2022-04-11**

## CLOUDERA

# Legal Notice

# Contents

# Machine Learning Discovery and Exploration

ML Discovery and Exploration makes it simple for data scientists to get started on a data science project.

When you start a data science project, you are presented with a blank notebook, and no indication of what data sources are available on the CDP platform, or how to access them. The ML Discovery and Exploration feature automatically discovers the data sources available to you, from within the standard CML user interface.

With ML Discovery and Exploration, you can:

- Easily discover and connect to data sources (virtual warehouses or data lakes).
- Quickly set up an ML project to access a given data source.
- Start running SQL queries to explore the data.

# Prerequisites

Machine Learning Discovery and Exploration has a few prerequisites.

- You must have Admin (write) permissions to set up or edit data connections. Only read permissions are needed to list data connections or view data connection details such as the library code.
- To manually set up data connections, you must have the JDBC URL (for Impala or Hive connections) or the Data Lake external directory URL (for Spark connections). You can obtain this from the option menu for the virtual warehouse to which you want to connect.
- One or more data warehouses must be already created in the environment. For more information on how to create a data warehouse, see *Adding a New Virtual Warehouse*.
- For Spark data connections, you must use ML Runtimes, and specifically the Spark Runtime Add-on must be available.
- Hive and Impala data connections also require ML Runtimes. Legacy engines may work, but are not supported.
- For Spark data connections, you must have permissions set correctly for external access to the S3 bucket.
- For both Hive and Impala connections, the HADOOP_USERNAME and WORKLOAD_PASSWORD must be set. The HADOOP_USERNAME is set automatically by the environment. To set the WORKLOAD_PASSWORD, see *Using data connections in a project*.
- For Hive data connections, make sure that SSO is disabled.

**Related Information**

Adding a New Virtual Warehouse
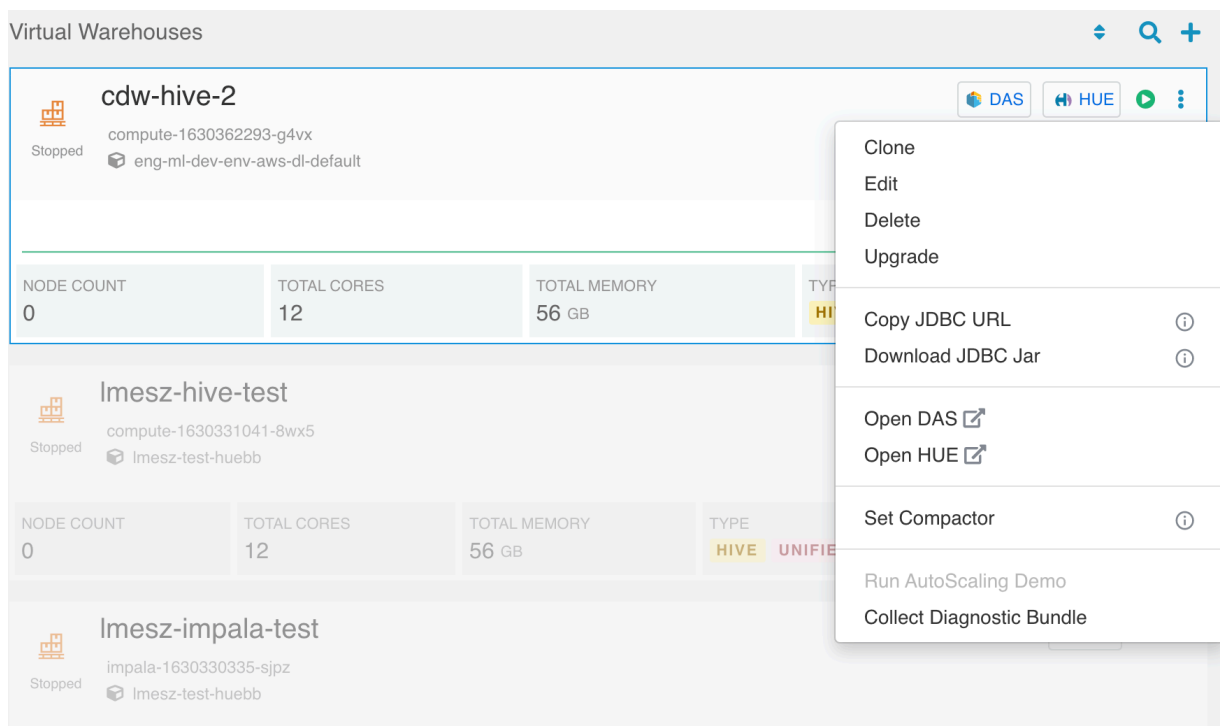
Using data connections in a project

# Set up a data connection manually

Data connections are automatically discovered. You can also set up a data connection manually, which works across CDP environments.

**Procedure**

1. Log into the CDP web interface and navigate to the Data Warehouse service.
2. In the Data Warehouse service, select Virtual Warehouses in the left navigation panel.

**3.** Select the options menu for the warehouse you want to access, and select Copy JDBC URL.



**4.** In  Site Administration Data Connections , select New Connection.

**5.** Enter the connection name. You cannot have duplicate names for data connections within a workspace or within a given project.

**6.** Select the connection type:

    **a.** Hive Virtual Warehouse

    **b.** Impala Virtual Warehouse

**7.** Enter the JDBC URL for the data connection. You can obtain this from the virtual warehouse in Cloudera Data Warehouse.

**8.** (Optional) Enter the Virtual Warehouse Name. This is the name of the warehouse in Cloudera Data Warehouse.

**Results**
The data connection is available to users by default. To change availability, click the Available switch.

# Data connection management

There are a few things to keep in mind about data connections.

- Manage data connections in a workspace

  At the workspace level, you can check the data connections that are available in a workspace. In  Project Settings Data Connections , check that your desired data source is present. You can also set the availability for any discovered connections, if necessary.

- Manage data connections in a project

  At the project level, you can see all data connections that were made available in the workspace when the project was created. The available connections can then further be marked as unavailable if so desired. You can update any changes to the connections that were made at the workspace level by clicking Sync with Workspace. Any changes made here only apply to your project.

- Data connection availability

  Keep in mind these two scenarios for setting data connection availability.

  1. If a workspace data connection is marked Unavailable, and you then create a project, the data connection will not appear in the project. If the connection is then changed to Available, and then the Sync with Workspace button is clicked, the connection will appear in the project.
  2. If a workspace data connection is marked Available, and you then create a project, the connection shows up. If the workspace data connection is then toggled to Unavailable, and you click Sync with Workspace in the project, the data connection will remain available in the project.

# Using data connections in a project

As a data scientist, you can manage the data connections within a project, and connect your code to the data.

Ensure your environment is set up as follows:

- The project must be configured to use ML Runtimes.
- To use a Spark data connection, make sure you have the correct permissions set in the s3 bucket.
- Configure the WORKLOAD_PASSWORD environmental variable.

You need to set an environmental variable for the password (WORKLOAD_PASSWORD). This environmental variable is needed to access a Hive or Impala warehouse.

To set the password:

1. Go to User Settings Environment .
2. Under the Reserved Environment Variables set your WORKLOAD_PASSWORD.
3. Click Save.



# Using data connection snippets

As a data scientist, you can connect your project to data with a data connection snippet.

**Procedure**

1. Select New Project.
2. Enter the project name.
3. Start a new session. To use a Spark connection, make sure to select an ML Runtime engine.
4. In the Connection Code Snippet pane, select the data connection to use.



5. Select Copy Code.
6. Paste the snippet into your code.

7. Edit the username and password, if necessary. (This is not necessary for the Spark Data Lake connection, or if the WORKLOAD_PASSWORD environmental variable is set up for the user.)

8. Enter the data connection name.

9. Uncomment and edit the SQL query.

10. Select Run.

### Results
The results of the SQL query display in the command window.

### What to do next
When you have finished exploring the available data sources, you can select Don't show me this again for this project on the Connection Code Snippet pane, and it will no longer display when you start a session in the project.

# Managing default and backup data connections

As a data scientist, you may want to set default and backup data connections (virtual warehouses). This makes it easy to manage the case where the default data source (virtual warehouse or data lake) becomes unavailable, for example.

### Before you begin

You need Admin privileges to perform this task. Here, the data connection names of default and backup are used as examples.

### Procedure

1. In  Site Administration Data connections , click Edit to change the name of the default data connection to default.

2. Change another connection to backup.

3. In Project Settings, click Sync with Workspace to update the names of the data connections. If two different data connections have the same name, an error occurs during synchronization.

4. If the default data connection becomes unavailable, the Workspace Administrator can go to the Data Connections tab, and rename the connections.

   For example, after changing default to unavailable, change backup to default. Projects that use library code for their connection to the default data connection will continue to operate, because it is now using the new default connection. Note: you need to click Sync with Workspace at the project level to get the updated connections.

# API Permissions For Projects

The tasks available to collaborators depends on their level of access.

A collaborator with read-only access is able to do the following:

- List data connections
- Get data connection details

> **Note:** The list of data connections is unavailable in the UI to viewer-only users because the Setting tab in a project is hidden for viewer collaborators. However, they can use curl to submit api requests to list the data connections.

A collaborator with write access can:

- Create data connections
- Edit data connections
- Delete data connections

- List data connections
- Get data connection details

# Troubleshooting: 401 Unauthorized

Problem: Session returns a 401: Unauthorized error.

When you are in a session and try to run the code for a Hive or Impala connection, the session returns a 401: Unauthorized HTTP error.

```
Modify next 2 lines to update your credentials

> USERNAME = os.getenv('HADOOP_USER_NAME')
> PASSWORD = os.getenv('WORKLOAD_PASSWORD')
> conn = cmldata.getConnection({
        'CONNECTION_NAME': CONNECTION_NAME,
        'USERNAME': USERNAME,
        'PASSWORD': PASSWORD
    })
HttpError: HTTP code 401: Unauthorized
HttpError                          Traceback (most recent call last)
<ipython-input-1-a69c5e2af83b> in <module>
      2          'CONNECTION_NAME': CONNECTION_NAME,
      3          'USERNAME': USERNAME,
----> 4          'PASSWORD': PASSWORD
      5      })

~/.local/lib/python3.7/site-packages/cmldata.py in getConnection(properties)
    105      fmt_codesnippet = _getConnectionSnippet(properties)
    106      scope = {}
--> 107      exec(fmt_codesnippet, scope)
    108      return scope["conn"]

<string> in <module>

<string> in getCursor(self)

/usr/local/lib/python3.7/site-packages/impala/hiveserver2.py in cursor(self, user, configuration, convert_type
    127          log.debug(' cursor(): getting new session handle')
```

Solution: You need to set your workload password. Ensure you are using the correct credentials.

# Troubleshooting: Existing connection name

Problem: When the user attempts to sync data connections, an error message displays, stating the crn or name is a duplicate.

Solution: This indicates a project connection (one that is not copied from the workspace) has the same name or crn as a workspace connection. To resolve this, you need to change the name or crn of the data connection at the project level.

# Troubleshooting: Existing connection name

Problem: When the user attempts to sync data connections, an error message displays, stating the crn or name is a duplicate.

Solution: This indicates a project connection (one that is not copied from the workspace) has the same name or crn as a workspace connection. To resolve this, you need to change the name or crn of the data connection at the project level.