

Managing Cloudera Data Science Workbench

Date published: 2020-02-28

Date modified: 2020-12-15



Legal Notice

© Cloudera Inc. 2024. All rights reserved.

The documentation is and contains Cloudera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cloudera software may be found within the documentation accompanying each component in a particular release.

Cloudera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 (“ASLv2”), the Affero General Public License version 3 (AGPLv3), or other license terms. Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cloudera software product page for more information on Cloudera software. For more information on Cloudera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cloudera reserves the right to change any products at any time, and without notice. Cloudera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cloudera.

Cloudera, Cloudera Altus, HUE, Impala, Cloudera Impala, and other Cloudera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners.

Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER’S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

Contents

Cluster Management.....	4
Cluster Monitoring with Grafana.....	4
Accessing the Grafana Dashboard.....	4
Available Grafana Dashboards.....	4
Backup and Disaster Recovery for Cloudera Data Science Workbench.....	5
Creating a Backup for Disaster Recovery for Cloudera Data Science Workbench.....	5
Monitoring/Reducing CDSW disk space.....	6
Cloudera Data Science Workbench Scaling Guidelines.....	8
Ports Used By Cloudera Data Science Workbench.....	8
Rollback Cloudera Data Science Workbench to an Older Version.....	9
Uninstalling CPM Deployments.....	9
Uninstalling RPM Deployments.....	10
RPM Deployments.....	10

Cluster Management

This section describes some common tasks and guidelines related to cluster management for Cloudera Data Science Workbench.

For a basic overview of how Cloudera Data Science Workbench fits onto a Cloudera Manager and CDH cluster, refer to the [Architecture Overview](#).

Cluster Monitoring with Grafana

Cloudera Data Science Workbench leverages Prometheus and Grafana to provide a dashboard that allows you to monitor how CPU, memory, storage, and other resources are being consumed by your CDSW deployment. Prometheus is an internal data source that is auto-populated with resource consumption data for each deployment. Grafana is the monitoring dashboard that allows you to create visualizations for resource consumption data from Prometheus.

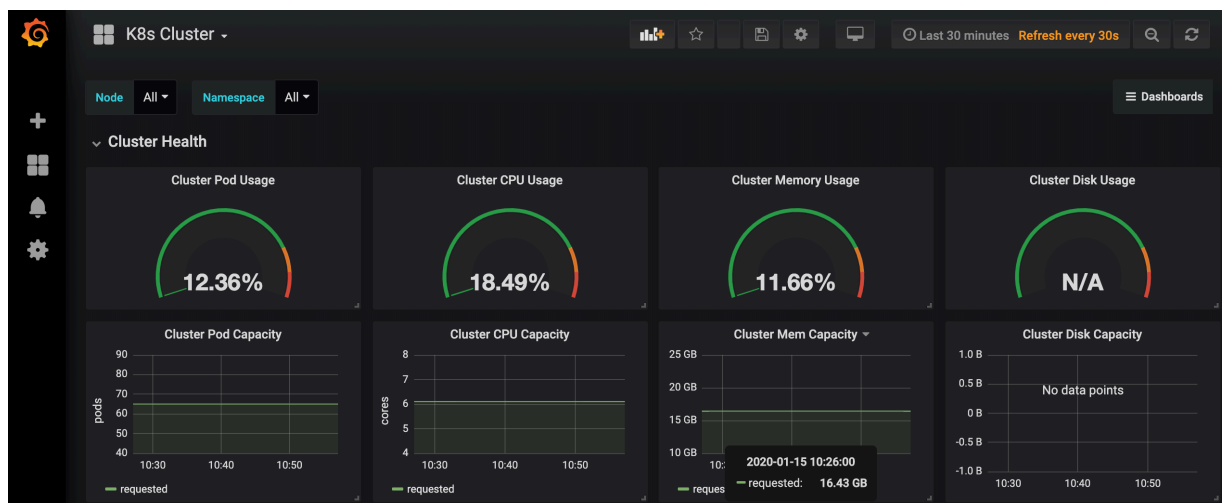
By default, CDSW provides you with three Grafana dashboards: K8s Cluster, K8s Containers, and K8s Node. You can extend these dashboards or create more panels for other metrics. For more information, see the [Grafana documentation](#).

Accessing the Grafana Dashboard

To access the Grafana dashboard for your deployment.

Procedure

1. Log into Cloudera Data Science Workbench with site administrator privileges.
2. Click Admin Overview .
3. Click the Grafana dashboard link. This will take you to the built-in Grafana server.
4. To see all the available dashboards, click Home K8 Cluster (or K8s Containers or K8s Node).



Available Grafana Dashboards

This topic lists all of the available Grafana dashboards.

K8s Cluster

Provides metrics for the following:

- Overview of nodes, pods, and containers
- CPU capacity usage
- Memory capacity usage
- Pod capacity usage
- Disk capacity usage

K8s Containers

Provides metrics for the following:

- Memory usage per pod
- CPU Usage per pod
- Read/Write IOPS per pod

K8s Node

Provides metrics for the following:

- CPU usage per node
- Memory Usage per node
- Read/Write IOPS per node
- Available memory per node
- Network traffic per node

Backup and Disaster Recovery for Cloudera Data Science Workbench

All application data for Cloudera Data Science Workbench, including project files and database state, is stored on the master host at `/var/lib/cdsw`.

Given typical access patterns, it is strongly recommended that `/var/lib/cdsw` be stored on a dedicated SSD block device or SSD RAID configuration. Because application data is not replicated to HDFS or backed up by default, site administrators must enable a backup strategy to meet any disaster recovery scenarios.

Cloudera strongly recommends both regular backups and backups before upgrades and is not responsible for any data loss.

Creating a Backup for Disaster Recovery for Cloudera Data Science Workbench

Cloudera strongly recommends both regular backups and backups before upgrades and is not responsible for any data loss.

Procedure

1. Stop Cloudera Data Science Workbench.

Use one of the following steps depending on your Cloudera Data Science workbench version.


Cloudera Data Science Workbench 1.4.2 or lower

Do not stop or restart Cloudera Data Science Workbench without using the [cdsw_protect_stop_restart.sh](#) script. This is to help avoid the data loss issue detailed in TSB-346 .

Run the script on your master host and stop Cloudera Data Science Workbench (instructions below) only when instructed to do so by the script. Then proceed with step 2 of this process.

Cloudera Data Science Workbench 1.4.3 or higher

Depending on your deployment, use one of the following sets of instructions to stop the application.

- CSD - Log in to Cloudera Manager. On the Home Status tab, click  to the right of the CDSW service and select Stop from the dropdown. Wait for the action to complete.
- OR
- RPM - Run the following command on the master host:

```
cdsw stop
```

2. After stopping CDSW, and before running the following tar command, wait 2-5 minutes (depending on your disk speed) to ensure that all data from CDSW is successfully written to the disks. Otherwise the tar command may not capture all recent changes.
3. To create the backup, run the following command on the master host:

```
tar cvzf cdsw.tar.gz /var/lib/cdsw/*
```



Note: The /var/lib/cdsw directory contains all the persistent information as well as project-related information, such as database information, configurations, image details, etc. Using tar to create the backup preserves important file metadata such as file ownership, and is the recommended method to migrate the information to the new host. Other methods of copying/saving files might not preserve this information. This metadata is required for tasks such as [migrating CDSW](#) to another cluster.

4. (Optional) If needed, the following command can be used to unpack the tar bundle.

```
tar xvzf cdsw.tar.gz -C /var/lib/cdsw
```

Monitoring/Reducing CDSW disk space

If the /var/lib/cdsw disk on the CDSW master node fills up to 80% or more, then the CDSW system starts to act very poorly, with symptoms such as unable to start sessions, Kerberos tickets not working, unable to create projects, and pods just crashing. There are a number of maintenance chores that Admin teams should perform monthly on CDSW to clean up the disk storage space

About this task

Possible symptoms that CDSW storage is filling up include but are not limited to:

- unable to start sessions
- sessions randomly get killed with "Engine exited with status 1"
- unable to build models
- unable to "kinit" through Hadoop Authentication
- the UI is slow

Any of these symptoms "Go away" temporarily when CDSW is restarted. This is a major point that the problem is with the storage space, since restarting CDSW will free up a lot of temporary space.

Remove Deleted Orphan Projects

When projects are deleted from the CDSW UI, the files are left on the disk. These projects cannot be recovered in CDSW UI, since they have been removed from the database. Therefore, since nobody can view these files, they are completely safe to delete. The following command will show these files that can be safely deleted:

```
{
  for project_id in `ls /var/lib/cds/current/projects/projects/*/*`
  do
    echo "Processing $project_id"
    rows=`kubectl exec $(kubectl get pods -l role=db --no-headers=true -o custom-columns=NAME:.metadata.name) -ti -- psql -U sense -P pager=off -c "SELECT * FROM projects WHERE ID = $project_id" | grep '0 row' | wc -l`
    if [ $rows -gt 0 ]
    then
      echo -n "Project with ID $project_id has been deleted, you can archive its directory: " `ls -d /var/lib/cds/current/projects/projects/*/$project_id`;
      echo
    fi
  done
}
```

This will print out all of the projects that can safely be removed from the disk. These projects have already been deleted from the UI, so there is no risk that collaborators cannot access the files, etc.



Note: The project folders are organized like /projects/projects/0/1, etc. The 0 here actually stands for "the first thousand." This means if a system has more than a thousand projects, this command might not be complete. If you have more than one thousand projects in the file system, you may have to update the above script.

After you run this command and clean up these orphaned projects, check the disk usage again to see if you are back below 70%.

Identify Huge Projects on Disk

If you still need to clean up some disk space, the next best solution is to look through the file system and find huge projects on disk. Projects can grow huge for a number of reasons. It is important to understand why so that the Admin team can educate the end users to prevent this from happening in the future. The biggest reasons for the project growing to crazy sizes are:

- Users are doing ML training on huge training data sets, and instead of using the data in HDFS, they pull massive training sets into their projects.

Resolution: Pull data from Hadoop and run your training, then delete this data. Or simply train the algorithms on the CDH cluster via spark or some other method.

- Users have created ML jobs that run on a scheduled basis and these generate output artifacts. After enough time, this can grow to substantial sizes.

Resolution: Add another job to delete this old data.

You can use the following script to list the projects in the system, along with the username, project name, email, and size. This will spit things out in a CSV format, and can be thrown into Excel and used, for example, to email the top ten users and ask them if it is possible to clean up their project sizes.

```
echo "Project Path, Size, Project Name, User name, email"
for project_id in `du -hx --max-depth=2 /var/lib/cds/current/projects/projects/ | sort -hr | awk 'NR>2 {print $2}'`
do
  echo -n $project_id ","
```

```

fileNameWithFullPath="${project_id%}/";
pid="${fileNameWithFullPath##*/}"
rows=`kubectl exec $(kubectl get pods -l role=db --no-headers=true -o custom-columns=NAME:.metadata.name) -ti -- psql -P pager=off -U sense -c "SELECT * FROM projects WHERE ID = $pid" | grep '0 row' | wc -l`
if [ $rows -gt 0 ]
then
echo -n "Project with ID $pid has been deleted, you can archive its directory: " `ls -d /var/lib/cds/current/projects/projects/*/ $pid` ; echo
else
size=`du -sh /var/lib/cds/current/projects/projects/0/$pid | awk '{print $1}'`
echo -n "$size,"
user_data=`kubectl exec $(kubectl get pods -l role=db --no-headers=true -o custom-columns=NAME:.metadata.name) -ti -- psql -U sense -t -A -c "select p.name,u.name, u.email from projects p, users u where p.user_id=u.id and p.id=$pid;" | tr '|' ','`
echo "$user_data"
fi
done

```

You should put this inside a file, chmod to 777, and run it like `get-large-projects.sh > large-projects.csv`

This can make a significant difference. In practice, the top ten projects can take up 30-40% of the free space.

Cloudera Data Science Workbench Scaling Guidelines

New hosts can be added and removed from a Cloudera Data Science Workbench deployment without interrupting any jobs already scheduled on existing hosts. Therefore, it is rather straightforward to increase capacity based on observed usage.

At a minimum, Cloudera recommends you allocate at least 1 CPU core and 2 GB of RAM per concurrent session or job. CPU can burst above a 1 CPU core share when spare resources are available. Therefore, a 1 CPU core allocation is often adequate for light workloads. Allocating less than 2 GB of RAM can lead to out-of-memory errors for many applications.

As a general guideline, Cloudera recommends hosts with RAM between 60GB and 256GB, and between 16 and 48 cores. This provides a useful range of options for end users. SSDs are strongly recommended for application data storage.



Note: At present, Kubernetes does not automatically pick the CPU/RAM that is dynamically added to the virtual machines (VM) while the VM is on. Therefore, you need to restart the CDSW service/role on the node on which you have updated the capacity.

For some data science and machine learning applications, users can collect a significant amount of data in memory within a single R or Python process, or use a significant amount of CPU resources that cannot be easily distributed into the CDH cluster. If individual users frequently run larger workloads or run workloads in parallel over long durations, increase the total resources accordingly. Understanding your users' concurrent workload requirements or observing actual usage is the best approach to scaling Cloudera Data Science Workbench.

Ports Used By Cloudera Data Science Workbench

Cloudera Data Science Workbench runs on gateway hosts in a CDH/HDP cluster. As such, Cloudera Data Science Workbench acts as a gateway and requires full connectivity to cluster services such as Impala, Spark 2, etc. Additionally, in the case of Spark 2, cluster hosts will require access to the Spark driver running on a set of random ports (20050-32767) on Cloudera Data Science Workbench hosts.

Firewall restrictions must be disabled across Cloudera Data Science Workbench and CDH/HDP cluster hosts. Internally, the Cloudera Data Science Workbench master and worker hosts require full connectivity with no firewalls. Externally, end users connect to Cloudera Data Science Workbench exclusively through a web server running on the master host, and therefore do not need direct access to any other internal Cloudera Data Science Workbench or CDH services.

This information has been summarized in the following table.

Components	Details
Communication with the CDH / HDP cluster	CDH / HDP -> Cloudera Data Science Workbench The CDH/HDP cluster must have access to the Spark driver that runs on Cloudera Data Science Workbench hosts, on a set of randomized ports in the range, 20050-32767.
	Cloudera Data Science Workbench -> CDH / HDP As a gateway service, Cloudera Data Science Workbench must have access to all the ports used by CDH and Cloudera Manager .
Communication with the Web Browser	The Cloudera Data Science Workbench web application is available at port 80. HTTPS access is available over port 443.

Rollback Cloudera Data Science Workbench to an Older Version

All stateful data for Cloudera Data Science Workbench is stored in the `/var/lib/cds` directory on the Master host. The contents of this directory are forward compatible, which is what allows for upgrades. However, they are not backward compatible. Therefore, to rollback Cloudera Data Science Workbench to a previous version, you must have a backup of the `/var/lib/cds` directory, taken prior to the last upgrade.

About this task

In general, the steps required to restore a previous version of Cloudera Data Science Workbench are:

Procedure

1. Depending on your deployment, either uninstall the RPM or deactivate the current CDSW parcel in Cloudera Manager.
2. On the master host, restore the backup copy you have of `/var/lib/cds`. Note that any changes made after this backup will be lost.
3. Install a version of Cloudera Data Science Workbench that is equal to or greater than the version of the `/var/lib/cds` backup.


Uninstalling CPM Deployments

This topic describes how to uninstall Cloudera Data Science Workbench from a cluster.

About this task

Procedure

1. Log in to the Cloudera Manager Admin Console.

2. On the Home > Status tab, click  to the right of the CDSW service and select Stop from the dropdown and confirm that you want to stop the service.





Note: On Cloudera Data Science Workbench 1.4.0 (and lower), do not stop or restart Cloudera Data Science Workbench without using the [cdsw_protect_stop_restart.sh](#) script. This is to help avoid the data loss issue detailed in TSB-346.

3. (Strongly Recommended) On the master host, backup the contents of the /var/lib/cdsw directory. This is the directory that stores all your application data.

To create the backup, run the following command on the master host:

```
tar cvzf cdsw.tar.gz /var/lib/cdsw/*
```

4. Go back to Cloudera Manager and click Hosts Parcels .
5. Go to the CDSW parcel and click Deactivate.
6. Select Deactivate Only from the list.
7. Click the  to the right of the Activate button and select Remove From Hosts.
8. Click OK to confirm.
9. Go back to the Home Status page and click  to the right of the CDSW service. Select Delete from the dropdown and confirm that you want to delete the service.
10. Remove all your user data that is stored in /var/lib/cdsw on the master host from the deployment. This step will permanently remove all user data.

```
sudo rm -Rf /var/lib/cdsw
```

For more details on how to uninstall Cloudera Manager and its components, see [Uninstalling Cloudera Manager and Managed Software](#).

Uninstalling RPM Deployments

This topic describes how to uninstall Cloudera Data Science Workbench from a cluster.

RPM Deployments

Procedure

1. Run the following command on the master host to stop Cloudera Data Science Workbench.

```
cdsw stop
```



Note: On Cloudera Data Science Workbench 1.4.0 (and lower), do not stop or restart Cloudera Data Science Workbench without using the [cdsw_protect_stop_restart.sh](#) script. This is to help avoid the data loss issue detailed in TSB-346.

2. (Strongly Recommended) On the master host, backup the contents of the /var/lib/cdsw directory. This is the directory that stores all your application data.

To create the backup, run the following command on the master host:

```
tar cvzf cdsw.tar.gz /var/lib/cdsw/*
```

3. To uninstall Cloudera Data Science Workbench, run the following commands on the master host and all the worker hosts.

```
cdsw stop  
yum remove cloudera-data-science-workbench
```

4. Remove all your user data that is stored in /var/lib/cdsw on the master host from the deployment. This step will permanently remove all user data.

```
sudo rm -Rf /var/lib/cdsw
```