



Red Hat OpenShift Data Science self-managed 1.32

Integrating data from Amazon S3

Use data stored in an Amazon Web Services (AWS) Simple Storage Service (S3) bucket

Red Hat OpenShift Data Science self-managed 1.32 Integrating data from Amazon S3

Use data stored in an Amazon Web Services (AWS) Simple Storage Service (S3) bucket

Legal Notice

Copyright © 2023 Red Hat, Inc.

The text of and illustrations in this document are licensed by Red Hat under a Creative Commons Attribution–Share Alike 3.0 Unported license ("CC-BY-SA"). An explanation of CC-BY-SA is available at

<http://creativecommons.org/licenses/by-sa/3.0/>

. In accordance with CC-BY-SA, if you distribute this document or an adaptation of it, you must provide the URL for the original version.

Red Hat, as the licensor of this document, waives the right to enforce, and agrees not to assert, Section 4d of CC-BY-SA to the fullest extent permitted by applicable law.

Red Hat, Red Hat Enterprise Linux, the Shadowman logo, the Red Hat logo, JBoss, OpenShift, Fedora, the Infinity logo, and RHCE are trademarks of Red Hat, Inc., registered in the United States and other countries.

Linux[®] is the registered trademark of Linus Torvalds in the United States and other countries.

Java[®] is a registered trademark of Oracle and/or its affiliates.

XFS[®] is a trademark of Silicon Graphics International Corp. or its subsidiaries in the United States and/or other countries.

MySQL[®] is a registered trademark of MySQL AB in the United States, the European Union and other countries.

Node.js[®] is an official trademark of Joyent. Red Hat is not formally related to or endorsed by the official Joyent Node.js open source or commercial project.

The OpenStack[®] Word Mark and OpenStack logo are either registered trademarks/service marks or trademarks/service marks of the OpenStack Foundation, in the United States and other countries and are used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation, or the OpenStack community.

All other trademarks are the property of their respective owners.

Abstract

Learn how to use data stored in an Amazon Web Services (AWS) Simple Storage Service (S3) bucket.

Table of Contents

PREFACE	3
CHAPTER 1. PREREQUISITES	4
CHAPTER 2. CREATING AN AMAZON S3 CLIENT USING NOTEBOOK CELLS	5
CHAPTER 3. LISTING AVAILABLE AMAZON S3 BUCKETS USING NOTEBOOK CELLS	6
CHAPTER 4. LISTING FILES IN AVAILABLE AMAZON S3 BUCKETS USING NOTEBOOK CELLS	7
CHAPTER 5. DOWNLOADING FILES FROM AVAILABLE AMAZON S3 BUCKETS USING NOTEBOOK CELLS .	9
CHAPTER 6. UPLOADING FILES TO AVAILABLE AMAZON S3 BUCKETS USING NOTEBOOK CELLS	10
CHAPTER 7. ADDITIONAL RESOURCES	11

PREFACE

When working in a Jupyter Notebook, you may want to work with data stored in an Amazon Web Services (AWS) Simple Storage Service (S3) bucket. This section covers commands and procedures for working with data stored in Amazon S3.

CHAPTER 1. PREREQUISITES

- A Jupyter server running on Red Hat OpenShift Data Science.
- Access to a Amazon Web Services S3 bucket.
- Locate the **AWS Access Key ID** and **AWS Secret Access Key** for your Amazon S3 account.
- A Jupyter Notebook

CHAPTER 2. CREATING AN AMAZON S3 CLIENT USING NOTEBOOK CELLS

To interact with data in Amazon S3 buckets, you must create a local client to handle requests to that service.

Prerequisites

- Access to a Jupyter notebook server running on Red Hat OpenShift Data Science.
- Define values for the **AWS_ACCESS_KEY_ID** and **AWS_SECRET_ACCESS_KEY** environment variables when you start your notebook server, using the values from your Amazon Web Services account under **My Security Credentials**.

Procedure

1. In a new notebook cell, import the required libraries by adding the following:

```
import os
import boto3
from boto3 import session
```

2. In another new notebook cell, define the following to create your session and client.

- a. Define your credentials.

```
key_id = os.environ.get('AWS_ACCESS_KEY_ID')
secret_key = os.environ.get('AWS_SECRET_ACCESS_KEY')
```

- b. Define the client session.

```
session = boto3.session.Session(aws_access_key_id=key_id,
                                aws_secret_access_key=secret_key)
```

- c. Define the client connection.

```
s3_client = boto3.client('s3', aws_access_key_id=key_id,
                          aws_secret_access_key=secret_key)
```

Verification

- Create a new cell and run an Amazon S3 command such as the following:

```
s3_client.list_buckets()
```

A successful response includes a **HTTPStatusCode** of **200** and a list of **Buckets** similar to the following:

```
'Buckets': [{'Name': 'my-app-asdf3-image-registry-us-east-1-wbmlcvbasdfasdgvtmsmkt',
              'CreationDate': datetime.datetime(2021, 4, 21, 6, 8, 52, tzinfo=tzlocal())},
             {'Name': 'cf-templates-18rxasdfgawsb-us-east-1',
              'CreationDate': datetime.datetime(2021, 2, 15, 18, 35, 34, tzinfo=tzlocal())}]
```

CHAPTER 3. LISTING AVAILABLE AMAZON S3 BUCKETS USING NOTEBOOK CELLS

You can check which buckets you have access to by listing the buckets available to your account.

Prerequisites

- Configure an Amazon S3 client in a previous cell in the notebook. See [Creating an Amazon S3 client using notebook cells](#) for more information.

Procedure

1. Create a new notebook cell and use the **s3_client** to list available buckets.

```
s3_client.list_buckets()
```

2. You can make this list of buckets easier to read by only printing the name, rather than the full response, for example:

```
for bucket in s3_client.list_buckets()['Buckets']:
    print(bucket['Name'])
```

This returns output similar to the following:

```
my-app-asdf3-image-registry-us-east-1-wbmlcvbasdgasdgtkpt
cf-templates-18rxuasgasgvb-us-east-1
```

Additional resources

- [Creating an Amazon S3 client using notebook cells](#)
- [Amazon Web Services list buckets command reference](#)

CHAPTER 4. LISTING FILES IN AVAILABLE AMAZON S3 BUCKETS USING NOTEBOOK CELLS

You can check the files available in buckets you have access to by listing the objects in the bucket. Because buckets use object storage rather than a typical file system, object naming works differently from normal file naming. Objects in a bucket are always known by a key, which consists of the full path in the bucket plus the name of the file itself.

Prerequisites

- Configure an Amazon S3 client in a previous cell in the notebook. See [Creating an Amazon S3 client using notebook cells](#) for more information.

Procedure

1. Create a new notebook cell and list the objects in the bucket. For example:

```
bucket_name = 'std-user-bucket1'
s3_client.list_objects_v2(Bucket=bucket_name)
```

This returns several objects in the following format:

```
{'Key':
'docker/registry/v2/blobs/sha256/00/0080913dd3f10aadb34asfgsgsdgasdga072049c93606b98b
ec84adb259b424f/data',
'LastModified': datetime.datetime(2021, 4, 22, 1, 26, 1, tzinfo=tzlocal()),
'ETag': '"6e02fad2deassadfs900a4bd7344ffe"',
'Size': 4052,
'StorageClass': 'STANDARD'}
```

2. You can make this list easier to read by printing only the key rather than the full response, for example:

```
bucket_name = 'std-user-bucket1'
for key in s3_client.list_objects_v2(Bucket=bucket_name)['Contents']:
    print(key['Key'])
```

This returns output similar to the following:

```
docker/registry/v2/blobs/sha256/00/0080913dd3f10aadb34asfgsgsdgasdga072049c93606b98b
ec84adb259b424f/data
```

3. You can also filter your query to list for a specific "path" or file name, for example:

```
bucket_name = 'std-user-bucket1'
for key in s3_client.list_objects_v2(Bucket=bucket_name, Prefix='<start_of_file_path>')
['Contents']:
    print(key['Key'])
```

In the preceding example, replace **<start_of_file_path>** with your own value.

Additional resources

- [Creating an Amazon S3 client using notebook cells](#)
- [Amazon Web Services list objects command reference](#)

CHAPTER 5. DOWNLOADING FILES FROM AVAILABLE AMAZON S3 BUCKETS USING NOTEBOOK CELLS

You can download a file to your notebook server using the **download_file** method.

Prerequisites

- Configure an Amazon S3 client in a previous cell in the notebook. See [Creating an Amazon S3 client using notebook cells](#) for more information.

Procedure

1. Define the following details in a notebook cell:

- a. The bucket that the file is in. Replace **<name_of_the_bucket>** with your own value.

```
bucket_name = '<name_of_the_bucket>'
```

- b. The name of the file to download. Replace **<name_of_the_file_to_download>** with your own value.

```
file_name = '<name_of_the_file_to_download>' # Full path from the bucket
```

- c. The name that you want the file to have after it is downloaded. This can be a full path, a relative path, or just a new file name. Replace **<name_of_the_file_when_downloaded>** with your own value.

```
new_file_name = '<name_of_the_file_when_downloaded>'
```

2. Download the file, specifying the previous variables as arguments.

```
s3_client.download_file(bucket_name, file_name, new_file_name)
```



NOTE

If you want to retrieve a file as an object that you can then stream as a standard file using the `read()` method, refer to the [Amazon Web Services get object command reference](#).

Additional resources

- [Creating an Amazon S3 client using notebook cells](#)
- [Amazon Web Services download file command reference](#)

CHAPTER 6. UPLOADING FILES TO AVAILABLE AMAZON S3 BUCKETS USING NOTEBOOK CELLS

You can upload files from your notebook server to an Amazon S3 bucket by using the **upload_file** method.

Prerequisites

- Configure an Amazon S3 client in a previous cell in the notebook. See [Creating an Amazon S3 client using notebook cells](#) for more information.

Procedure

1. Define the following details in a notebook cell:

- a. The name of the file to upload. This must include the full local path to the file. Replace **<name_of_the_file_to_upload>** with your own value.

```
file_name = '<name_of_the_file_to_upload>'
```

- b. The name of the bucket to upload the file to. Replace **<name_of_the_bucket>** with your own value.

```
bucket_name = '<name_of_the_bucket>'
```

- c. The full key to use to save the file to the bucket. Replace **<full_path_and_file_name>** with your own value.

```
key = '<full_path_and_file_name>'
```

2. Upload the file, specifying the previous variables as arguments.

```
s3_client.upload_file(file_name, bucket_name, key)
```

Additional resources

- [Creating an Amazon S3 client using notebook cells](#)
- [Amazon Web Services upload file command reference](#)

CHAPTER 7. ADDITIONAL RESOURCES

- [Red Hat OpenShift Data Science documentation](#)