

---

# hdfs Documentation

*Release 2.7.0*

**Author**

**Apr 23, 2023**



---

## Contents

---

<b>1</b>	<b>Installation</b>	<b>3</b>
<b>2</b>	<b>User guide</b>	<b>5</b>
2.1	Quickstart . . . . .	5
2.2	Advanced usage . . . . .	9
2.3	API reference . . . . .	12
<b>3</b>	<b>Sample script</b>	<b>25</b>
	<b>Python Module Index</b>	<b>27</b>
	<b>Index</b>	<b>29</b>



API and command line interface for HDFS.

- [Project homepage on GitHub](#)
- [PyPI entry](#)



# CHAPTER 1

---

## Installation

---

Using `pip`:

```
$ pip install hdfs
```

By default none of the package requirements for extensions are installed. To do so simply suffix the package name with the desired extensions:

```
$ pip install hdfs[avro,dataframe,kerberos]
```





## 2.1 Quickstart

This page first goes through the steps required to configure HdfsCLI's command line interface then gives an overview of the python API. If you are only interested in using HdfsCLI as a library, then feel free to jump ahead to the *Python bindings* section.

### 2.1.1 Configuration

HdfsCLI uses *aliases* to figure out how to connect to different HDFS clusters. These are defined in HdfsCLI's configuration file, located by default at `~/ .hdfscli.cfg` (or elsewhere by setting the `HDFSCLI_CONFIG` environment variable correspondingly). See below for a sample configuration defining two aliases, `dev` and `prod`:

```
[global]
default.alias = dev

[dev.alias]
url = http://dev.namenode:port
user = ann

[prod.alias]
url = http://prod.namenode:port
root = /jobs/
```

Each alias is defined as its own `ALIAS.alias` section which must at least contain a `url` option with the URL to the namenode (including protocol and port). All other options can be omitted. If specified, `client` determines which `hdfs.client.Client` class to use and the remaining options are passed as keyword arguments to the appropriate constructor. The currently available client classes are:

- `InsecureClient` (the default)
- `TokenClient`

See the *Kerberos extension* to enable the *KerberosClient* and *Custom client support* to learn how to use other client classes.

The `url` option can be configured to support High Availability namenodes when using WebHDFS, simply add more URLs by delimiting with a semicolon (;).

Finally, note the `default.alias` entry in the global configuration section which will be used as default alias if none is specified.

### 2.1.2 Command line interface

HdfsCLI comes by default with a single entry point `hdfscli` which provides a convenient interface to perform common actions. All its commands accept an `--alias` argument (described above), which defines against which cluster to operate.

#### Downloading and uploading files

HdfsCLI supports downloading and uploading files and folders transparently from HDFS (we can also specify the degree of parallelism by using the `--threads` option).

```
$ # Write a single file to HDFS.
$ hdfscli upload --alias=dev weights.json models/
$ # Read all files inside a folder from HDFS and store them locally.
$ hdfscli download export/results/ "results-$(date +%F) "
```

If reading (resp. writing) a single file, its contents can also be streamed to standard out (resp. from standard in) by using `-` as path argument:

```
$ # Read a file from HDFS and append its contents to a local log file.
$ hdfscli download logs/1987-03-23.txt - >>logs
```

By default HdfsCLI will throw an error if trying to write to an existing path (either locally or on HDFS). We can force the path to be overwritten with the `--force` option.

#### Interactive shell

The interactive command (used also when no command is specified) will create an HDFS client and expose it inside a python shell (using *IPython* if available). This makes is convenient to perform file system operations on HDFS and interact with its data. See *Python bindings* below for an overview of the methods available.

```
$ hdfscli --alias=dev

Welcome to the interactive HDFS python shell.
The HDFS client is available as `CLIENT`.

In [1]: CLIENT.list('data/')
Out[1]: ['1.json', '2.json']

In [2]: CLIENT.status('data/2.json')
Out[2]: {
  'accessTime': 1439743128690,
  'blockSize': 134217728,
  'childrenNum': 0,
  'fileId': 16389,
```

(continues on next page)

(continued from previous page)

```
'group': 'supergroup',
'length': 2,
'modificationTime': 1439743129392,
'owner': 'drwho',
'pathSuffix': '',
'permission': '755',
'replication': 1,
'storagePolicy': 0,
'type': 'FILE'
}

In [3]: CLIENT.delete('data/2.json')
Out[3]: True
```

Using the full power of python lets us easily perform more complex operations such as renaming folder which match some pattern, deleting files which haven't been accessed for some duration, finding all paths owned by a certain user, etc.

## More

Cf. `hdfscli --help` for the full list of commands and options.

## 2.1.3 Python bindings

### Instantiating a client

The simplest way of getting a `hdfs.client.Client` instance is by using the *Interactive shell* described above, where the client will be automatically available. To instantiate a client programmatically, there are two options:

The first is to import the client class and call its constructor directly. This is the most straightforward and flexible, but doesn't let us reuse our configured aliases:

```
from hdfs import InsecureClient
client = InsecureClient('http://host:port', user='ann')
```

The second leverages the `hdfs.config.Config` class to load an existing configuration file (defaulting to the same one as the CLI) and create clients from existing aliases:

```
from hdfs import Config
client = Config().get_client('dev')
```

### Reading and writing files

The `read()` method provides a file-like interface for reading files from HDFS. It must be used in a `with` block (making sure that connections are always properly closed):

```
# Loading a file in memory.
with client.read('features') as reader:
    features = reader.read()

# Directly deserializing a JSON object.
with client.read('model.json', encoding='utf-8') as reader:
```

(continues on next page)

(continued from previous page)

```
from json import load
model = load(reader)
```

If a `chunk_size` argument is passed, the method will return a generator instead, making it sometimes simpler to stream the file's contents.

```
# Stream a file.
with client.read('features', chunk_size=8096) as reader:
    for chunk in reader:
        pass
```

Similarly, if a `delimiter` argument is passed, the method will return a generator of the delimited chunks.

```
with client.read('samples.csv', encoding='utf-8', delimiter='\n') as reader:
    for line in reader:
        pass
```

Writing files to HDFS is done using the `write()` method which returns a file-like writable object:

```
# Writing part of a file.
with open('samples') as reader, client.write('samples') as writer:
    for line in reader:
        if line.startswith('-'):
            writer.write(line)

# Writing a serialized JSON object.
with client.write('model.json', encoding='utf-8') as writer:
    from json import dump
    dump(model, writer)
```

For convenience, it is also possible to pass an iterable `data` argument directly to the method.

```
# This is equivalent to the JSON example above.
from json import dumps
client.write('model.json', dumps(model))
```

## Exploring the file system

All `Client` subclasses expose a variety of methods to interact with HDFS. Most are modeled directly after the WebHDFS operations, a few of these are shown in the snippet below:

```
# Retrieving a file or folder content summary.
content = client.content('dat')

# Listing all files inside a directory.
fnames = client.list('dat')

# Retrieving a file or folder status.
status = client.status('dat/features')

# Renaming ("moving") a file.
client.rename('dat/features', 'features')

# Deleting a file or folder.
client.delete('dat', recursive=True)
```

Other methods build on these to provide more advanced features:

```
# Download a file or folder locally.
client.download('dat', 'dat', n_threads=5)

# Get all files under a given folder (arbitrary depth).
import posixpath as psp
fpaths = [
    psp.join(dpath, fname)
    for dpath, _, fnames in client.walk('predictions')
    for fname in fnames
]
```

See the *API reference* for the comprehensive list of methods available.

## Checking path existence

Most of the methods described above will raise an *HdfsError* if called on a missing path. The recommended way of checking whether a path exists is using the *content()* or *status()* methods with a *strict=False* argument (in which case they will return *None* on a missing path).

## More

See the *Advanced usage* section to learn more.

## 2.2 Advanced usage

### 2.2.1 Path expansion

All *Client* methods provide a path expansion functionality via the *resolve()* method. It enables the use of special markers to identify paths. For example, it currently supports the *#LATEST* marker which expands to the last modified file inside a given folder.

```
# Load the most recent data in the `tracking` folder.
with client.read('tracking/#LATEST') as reader:
    data = reader.read()
```

See the method's documentation for more information.

### 2.2.2 Custom client support

In order for the CLI to be able to instantiate arbitrary client classes, it has to be able to discover these first. This is done by specifying where they are defined in the *global* section of *HdfsCLI*'s configuration file. For example, here is how we can make the *KerberosClient* class available:

```
[global]
autoload.modules = hdfs.ext.kerberos
```

More precisely, there are two options for telling the CLI where to load the clients from:

- *autoload.modules*, a comma-separated list of modules (which must be on python's path).
- *autoload.paths*, a comma-separated list of paths to python files.

Implementing custom clients can be particularly useful for passing default options (e.g. a custom `session` argument to each client). We describe below a working example implementing a secure client with optional custom certificate support.

We first implement our new client and save it somewhere, for example `/etc/hdfscli.py`.

```
from hdfs import Client
from requests import Session

class SecureClient(Client):

    """A new client subclass for handling HTTPS connections.

    :param url: URL to namenode.
    :param cert: Local certificate. See `requests` documentation for details
        on how to use this.
    :param verify: Whether to check the host's certificate.
    :param **kwargs: Keyword arguments passed to the default `Client`
        constructor.

    """

    def __init__(self, url, cert=None, verify=True, **kwargs):
        session = Session()
        if ',' in cert:
            session.cert = [path.strip() for path in cert.split(',')]
        else:
            session.cert = cert
        if isinstance(verify, basestring): # Python 2.
            verify = verify.lower() in ('true', 'yes', 'ok')
        session.verify = verify
        super(SecureClient, self).__init__(url, session=session, **kwargs)
```

We then edit our configuration to tell the CLI how to load this module and define a `prod` alias using our new client:

```
[global]
autoload.paths = /etc/hdfscli.py

[prod.alias]
client = SecureClient
url = https://host:port
cert = /etc/server.crt, /etc/key
```

Note that options used to instantiate clients from the CLI (using `hdfs.client.Client.from_options()` under the hood) are always passed in as strings. This is why we had to implement some parsing logic in the `SecureClient` constructor above.

### 2.2.3 Tracking transfer progress

The `read()`, `upload()`, `download()` client methods accept a `progress` callback argument which can be used to track transfers. The passed function will be called every `chunk_size` bytes with two arguments:

- The source path of the file currently being transferred.
- The number of bytes currently transferred for this file or `-1` to signal that this file's transfer has just finished.

Below is an implementation of a toy tracker which simply outputs to standard error the total number of transferred bytes each time a file transfer completes (we must still take care to ensure correct behavior even during multi-threaded

transfers).

```

from sys import stderr
from threading import Lock

class Progress(object):

    """Basic progress tracker callback."""

    def __init__(self):
        self._data = {}
        self._lock = Lock()

    def __call__(self, hdfs_path, nbytes):
        with self._lock:
            if nbytes >= 0:
                self._data[hdfs_path] = nbytes
            else:
                stderr.write('%s\n' % (sum(self._data.values()), ))

```

Finally, note that the `write()` method doesn't expose a `progress` argument since this functionality can be replicated by passing a custom data generator (or within the context manager).

## 2.2.4 Logging configuration

It is possible to configure and disable where the CLI logs are written for each entry point. To do this, we can set the following options in its corresponding section (the entry point's name suffixed with `.command`). For example:

```

[hdfscli-avro.command]
log.level = INFO
log.path = /tmp/hdfscli/avro.log

```

The following options are available:

- `log.level`, handler log level (defaults to `DEBUG`).
- `log.path`, path to log file. The log is rotated every day (keeping a single copy). The default is a file named `COMMAND.log` in your current temporary directory. It is possible to view the currently active log file at any time by using the `--log` option at the command line.
- `log.disable`, disable logging to a file entirely (defaults to `False`).

## 2.2.5 Renaming entry points

By default the command line entry point will be named `hdfscli`. You can choose another name by specifying the `HDFSCLI_ENTRY_POINT` environment variable at installation time:

```

$ HDFSCLI_ENTRY_POINT=hdfs pip install hdfs

```

Extension prefixes will be adjusted similarly (e.g. in the previous example, `hdfscli-avro` would become `hdfs-avro`).

## 2.3 API reference

### 2.3.1 Client

WebHDFS API clients.

**class** `hdfs.client.Client` (*url*, *root=None*, *proxy=None*, *timeout=None*, *session=None*)

Bases: `object`

Base HDFS web client.

#### Parameters

- **url** – Hostname or IP address of HDFS namenode, prefixed with protocol, followed by WebHDFS port on namenode. You may also specify multiple URLs separated by semi-colons for High Availability support.
- **proxy** – User to proxy as.
- **root** – Root path, this will be prefixed to all HDFS paths passed to the client. If the root is relative, the path will be assumed relative to the user's home directory.
- **timeout** – Connection timeouts, forwarded to the request handler. How long to wait for the server to send data before giving up, as a float, or a (`connect_timeout`, `read_timeout`) tuple. If the timeout is reached, an appropriate exception will be raised. See the [requests](#) documentation for details.
- **session** – `requests.Session` instance, used to emit all requests.

In general, this client should only be used directly when its subclasses (e.g. `InsecureClient`, `TokenClient`, and others provided by extensions) do not provide enough flexibility.

**acl\_status** (*hdfs\_path*, *strict=True*)

Get `AclStatus` for a file or folder on HDFS.

#### Parameters

- **hdfs\_path** – Remote path.
- **strict** – If `False`, return `None` rather than raise an exception if the path doesn't exist.

**allow\_snapshot** (*hdfs\_path*)

Allow snapshots for a remote folder.

**Parameters** **hdfs\_path** – Remote path to a directory. If `hdfs_path` doesn't exist or does points to a normal file, an `HdfsError` will be raised. No-op if snapshotting is already allowed.

**checksum** (*hdfs\_path*)

Get a remote file's checksum.

**Parameters** **hdfs\_path** – Remote path. Must point to a file.

**content** (*hdfs\_path*, *strict=True*)

Get `ContentSummary` for a file or folder on HDFS.

#### Parameters

- **hdfs\_path** – Remote path.
- **strict** – If `False`, return `None` rather than raise an exception if the path doesn't exist.

**create\_snapshot** (*hdfs\_path*, *snapshotname=None*)

Create snapshot for a remote folder where it was allowed.



**Parameters** `hdfs_path` – Remote path to a directory. If `hdfs_path` doesn't exist, doesn't allow to create snapshot or points to a normal file, an `HdfsError` will be raised.

**:param snapshotname snapshot name; if absent, name is generated** by the server.

Returns a path to created snapshot.

**delete** (*hdfs\_path*, *recursive=False*, *skip\_trash=True*)

Remove a file or directory from HDFS.

**Parameters**

- **hdfs\_path** – HDFS path.
- **recursive** – Recursively delete files and directories. By default, this method will raise an `HdfsError` if trying to delete a non-empty directory.
- **skip\_trash** – When false, the deleted path will be moved to an appropriate trash folder rather than deleted. This requires Hadoop 2.9+ and trash to be enabled on the cluster.

This function returns `True` if the deletion was successful and `False` if no file or directory previously existed at `hdfs_path`.

**delete\_snapshot** (*hdfs\_path*, *snapshotname*)

Remove snapshot for a remote folder where it was allowed.

**Parameters** `hdfs_path` – Remote path to a directory. If `hdfs_path` doesn't exist or points to a normal file, an `HdfsError` will be raised.

**:param snapshotname snapshot name; if it does not exist, an** `HdfsError` will be raised.

**disallow\_snapshot** (*hdfs\_path*)

Disallow snapshots for a remote folder.

**Parameters** `hdfs_path` – Remote path to a directory. If `hdfs_path` doesn't exist, points to a normal file or there are some snapshots, an `HdfsError` will be raised.

No-op if snapshotting is disallowed/never allowed.

**download** (*hdfs\_path*, *local\_path*, *overwrite=False*, *n\_threads=1*, *temp\_dir=None*, *\*\*kwargs*)

Download a file or folder from HDFS and save it locally.

**Parameters**

- **hdfs\_path** – Path on HDFS of the file or folder to download. If a folder, all the files under it will be downloaded.
- **local\_path** – Local path. If it already exists and is a directory, the files will be downloaded inside of it.
- **overwrite** – Overwrite any existing file or directory.
- **n\_threads** – Number of threads to use for parallelization. A value of 0 (or negative) uses as many threads as there are files.
- **temp\_dir** – Directory under which the files will first be downloaded when `overwrite=True` and the final destination path already exists. Once the download successfully completes, it will be swapped in.
- **\*\*kwargs** – Keyword arguments forwarded to `read()`. If no `chunk_size` argument is passed, a default value of 64 kB will be used. If a `progress` argument is passed and threading is used, care must be taken to ensure correct behavior.

On success, this method returns the local download path.

**classmethod** `from_options` (*options*, *class\_name*='Client')

Load client from options.

**Parameters**

- **options** – Options dictionary.
- **class\_name** – Client class name. Defaults to the base *Client* class.

This method provides a single entry point to instantiate any registered *Client* subclass. To register a subclass, simply load its containing module. If using the CLI, you can use the `autoload.modules` and `autoload.paths` options.

**list** (*hdfs\_path*, *status*=False)

Return names of files contained in a remote folder.

**Parameters**

- **hdfs\_path** – Remote path to a directory. If `hdfs_path` doesn't exist or points to a normal file, an *HdfsError* will be raised.
- **status** – Also return each file's corresponding *FileStatus*.

**makedirs** (*hdfs\_path*, *permission*=None)

Create a remote directory, recursively if necessary.

**Parameters**

- **hdfs\_path** – Remote path. Intermediate directories will be created appropriately.
- **permission** – Octal permission to set on the newly created directory. These permissions will only be set on directories that do not already exist.

This function currently has no return value as WebHDFS doesn't return a meaningful flag.

**parts** (*hdfs\_path*, *parts*=None, *status*=False)

Returns a dictionary of part-files corresponding to a path.

**Parameters**

- **hdfs\_path** – Remote path. This directory should contain at most one part file per partition (otherwise one will be picked arbitrarily).
- **parts** – List of part-files numbers or total number of part-files to select. If a number, that many partitions will be chosen at random. By default all part-files are returned. If `parts` is a list and one of the parts is not found or too many samples are demanded, an *HdfsError* is raised.
- **status** – Also return each file's corresponding *FileStatus*.

**read** (\*\**kws*)

Read a file from HDFS.

**Parameters**

- **hdfs\_path** – HDFS path.
- **offset** – Starting byte position.
- **length** – Number of bytes to be processed. None will read the entire file.
- **buffer\_size** – Size of the buffer in bytes used for transferring the data. Defaults the the value set in the HDFS configuration.

- **encoding** – Encoding used to decode the request. By default the raw data is returned. This is mostly helpful in python 3, for example to deserialize JSON data (as the decoder expects unicode).
- **chunk\_size** – If set to a positive number, the context manager will return a generator yielding every `chunk_size` bytes instead of a file-like object (unless `delimiter` is also set, see below).
- **delimiter** – If set, the context manager will return a generator yielding each time the delimiter is encountered. This parameter requires the `encoding` to be specified.
- **progress** – Callback function to track progress, called every `chunk_size` bytes (not available if the chunk size isn't specified). It will be passed two arguments, the path to the file being uploaded and the number of bytes transferred so far. On completion, it will be called once with `-1` as second argument.

This method must be called using a `with` block:

```
with client.read('foo') as reader:
    content = reader.read()
```

This ensures that connections are always properly closed.

---

**Note:** The raw file-like object returned by this method (when called without an encoding, chunk size, or delimiter) can have a very different performance profile than local files. In particular, line-oriented methods are often slower. The recommended workaround is to specify an encoding when possible or read the entire file before splitting it.

---

**remove\_acl** (*hdfs\_path*)

**RemoveACL\_** for a file or folder on HDFS.

**Parameters** **hdfs\_path** – Path to an existing remote file or directory. An `HdfsError` will be raised if the path doesn't exist.

**remove\_acl\_entries** (*hdfs\_path*, *acl\_spec*)

**RemoveACL\_** for a file or folder on HDFS.

**Parameters**

- **hdfs\_path** – Path to an existing remote file or directory. An `HdfsError` will be raised if the path doesn't exist.
- **acl\_spec** – String representation of an ACL spec. Must be a valid string with entries for user, group and other. For example: "user::rwx,user:foo:rw-,group::r--,other:---".

**remove\_default\_acl** (*hdfs\_path*)

**RemoveDefaultACL\_** for a file or folder on HDFS.

**Parameters** **hdfs\_path** – Path to an existing remote file or directory. An `HdfsError` will be raised if the path doesn't exist.

**.. \_RemoveDefaultAcl: REMOVEDEFAULTACL\_**

**rename** (*hdfs\_src\_path*, *hdfs\_dst\_path*)

Move a file or folder.

**Parameters**

- **hdfs\_src\_path** – Source path.

- **hdfs\_dst\_path** – Destination path. If the path already exists and is a directory, the source will be moved into it. If the path exists and is a file, or if a parent destination directory is missing, this method will raise an `HdfsError`.

**rename\_snapshot** (*hdfs\_path*, *oldsnapshotname*, *snapshotname*)

Rename snapshot for a remote folder.

**Parameters** **hdfs\_path** – Remote path to a directory. If `hdfs_path` doesn't exist or points to a normal file, an `HdfsError` will be raised.

**:param oldsnapshotname snapshot name; if it does not exist**, an `HdfsError` will be raised.

**:param snapshotname new snapshot name; if it does already exist**, an `HdfsError` will be raised.

**resolve** (*hdfs\_path*)

Return absolute, normalized path, with special markers expanded.

**Parameters** **hdfs\_path** – Remote path.

Currently supported markers:

- `'#LATEST'`: this marker gets expanded to the most recently updated file or folder. They can be combined using the `'{N}'` suffix. For example, `'foo/#LATEST{2}'` is equivalent to `'foo/#LATEST/#LATEST'`.

**set\_acl** (*hdfs\_path*, *acl\_spec*, *clear=True*)

**SetAcl\_** or **ModifyAcl\_** for a file or folder on HDFS.

**Parameters**

- **hdfs\_path** – Path to an existing remote file or directory. An `HdfsError` will be raised if the path doesn't exist.
- **acl\_spec** – String representation of an ACL spec. Must be a valid string with entries for user, group and other. For example: `"user::rwx,user:foo:rw-,group::r--,other:---"`.
- **clear** – Clear existing ACL entries. If set to false, all existing ACL entries that are not specified in this call are retained without changes, behaving like **ModifyAcl\_**. For example: `"user:foo:rwx"`.

**.. \_SetAcl: SETACL\_**

**set\_owner** (*hdfs\_path*, *owner=None*, *group=None*)

Change the owner of file.

**Parameters**

- **hdfs\_path** – HDFS path.
- **owner** – Optional, new owner for file.
- **group** – Optional, new group for file.

At least one of `owner` and `group` must be specified.

**set\_permission** (*hdfs\_path*, *permission*)

Change the permissions of file.

**Parameters**

- **hdfs\_path** – HDFS path.
- **permission** – New octal permissions string of file.

**set\_replication** (*hdfs\_path*, *replication*)

Set file replication.

**Parameters**

- **hdfs\_path** – Path to an existing remote file. An `HdfsError` will be raised if the path doesn't exist or points to a directory.
- **replication** – Replication factor.

**set\_times** (*hdfs\_path*, *access\_time=None*, *modification\_time=None*)

Change remote timestamps.

**Parameters**

- **hdfs\_path** – HDFS path.
- **access\_time** – Timestamp of last file access.
- **modification\_time** – Timestamps of last file access.

**status** (*hdfs\_path*, *strict=True*)

Get `FileStatus` for a file or folder on HDFS.

**Parameters**

- **hdfs\_path** – Remote path.
- **strict** – If `False`, return `None` rather than raise an exception if the path doesn't exist.

**upload** (*hdfs\_path*, *local\_path*, *n\_threads=1*, *temp\_dir=None*, *chunk\_size=65536*, *progress=None*, *cleanup=True*, *\*\*kwargs*)

Upload a file or directory to HDFS.

**Parameters**

- **hdfs\_path** – Target HDFS path. If it already exists and is a directory, files will be uploaded inside.
- **local\_path** – Local path to file or folder. If a folder, all the files inside of it will be uploaded (note that this implies that folders empty of files will not be created remotely).
- **n\_threads** – Number of threads to use for parallelization. A value of 0 (or negative) uses as many threads as there are files.
- **temp\_dir** – Directory under which the files will first be uploaded when `overwrite=True` and the final remote path already exists. Once the upload successfully completes, it will be swapped in.
- **chunk\_size** – Interval in bytes by which the files will be uploaded.
- **progress** – Callback function to track progress, called every `chunk_size` bytes. It will be passed two arguments, the path to the file being uploaded and the number of bytes transferred so far. On completion, it will be called once with `-1` as second argument.
- **cleanup** – Delete any uploaded files if an error occurs during the upload.
- **\*\*kwargs** – Keyword arguments forwarded to `write()`. In particular, set `overwrite` to overwrite any existing file or directory.

On success, this method returns the remote upload path.

**walk** (*hdfs\_path*, *depth=0*, *status=False*, *ignore\_missing=False*, *allow\_dir\_changes=False*)

Depth-first walk of remote filesystem.

**Parameters**

- **hdfs\_path** – Starting path. If the path doesn't exist, an `HdfsError` will be raised. If it points to a file, the returned generator will be empty.
- **depth** – Maximum depth to explore. 0 for no limit.
- **status** – Also return each file or folder's corresponding `FileStatus`.
- **ignore\_missing** – Ignore missing nested folders rather than raise an exception. This can be useful when the tree is modified during a walk.
- **allow\_dir\_changes** – Allow changes to the directories' list to affect the walk. For example clearing it by setting `dirs[:] = []` would prevent the walk from entering any nested directories. This option can only be set when `status` is false.

This method returns a generator yielding tuples (`path`, `dirs`, `files`) where `path` is the absolute path to the current directory, `dirs` is the list of directory names it contains, and `files` is the list of file names it contains.

**write** (*hdfs\_path*, *data=None*, *overwrite=False*, *permission=None*, *blocksize=None*, *replication=None*, *buffer\_size=None*, *append=False*, *encoding=None*)  
Create a file on HDFS.

#### Parameters

- **hdfs\_path** – Path where to create file. The necessary directories will be created appropriately.
- **data** – Contents of file to write. Can be a string, a generator or a file object. The last two options will allow streaming upload (i.e. without having to load the entire contents into memory). If `None`, this method will return a file-like object and should be called using a `with` block (see below for examples).
- **overwrite** – Overwrite any existing file or directory.
- **permission** – Octal permission to set on the newly created file. Leading zeros may be omitted.
- **blocksize** – Block size of the file.
- **replication** – Number of replications of the file.
- **buffer\_size** – Size of upload buffer.
- **append** – Append to a file rather than create a new one.
- **encoding** – Encoding used to serialize data written.

Sample usages:

```
from json import dump, dumps

records = [
    {'name': 'foo', 'weight': 1},
    {'name': 'bar', 'weight': 2},
]

# As a context manager:
with client.write('data/records.jsonl', encoding='utf-8') as writer:
    dump(records, writer)

# Or, passing in a generator directly:
client.write('data/records.jsonl', data=dumps(records), encoding='utf-8')
```

```
class hdfs.client.InsecureClient (url, user=None, **kwargs)
```

Bases: `hdfs.client.Client`

HDFS web client to use when security is off.

#### Parameters

- **url** – Hostname or IP address of HDFS namenode, prefixed with protocol, followed by WebHDFS port on namenode
- **user** – User default. Defaults to the current user's (as determined by `whoami`).
- **\*\*kwargs** – Keyword arguments passed to the base class' constructor.

Note that if a session argument is passed in, it will be modified in-place to support authentication.

```
class hdfs.client.TokenClient (url, token, **kwargs)
```

Bases: `hdfs.client.Client`

HDFS web client using Hadoop token delegation security.

#### Parameters

- **url** – Hostname or IP address of HDFS namenode, prefixed with protocol, followed by WebHDFS port on namenode
- **token** – Hadoop delegation token.
- **\*\*kwargs** – Keyword arguments passed to the base class' constructor.

Note that if a session argument is passed in, it will be modified in-place to support authentication.

## 2.3.2 Extensions

The following extensions are currently available:

### Kerberos

Support for clusters using [Kerberos](#) authentication.

This extension adds a new `hdfs.client.Client` subclass, `KerberosClient`, which handles authentication appropriately with Kerberized clusters:

```
from hdfs.ext.kerberos import KerberosClient
client = KerberosClient('http://host:port')
```

To expose this class to the command line interface (so that it can be used by aliases), we add the following line inside the `global` section of `~/hdfscli.cfg` (or wherever our configuration file is located):

```
autoload.modules = hdfs.ext.kerberos
```

Here is what our earlier configuration would look like if we updated it to support a Kerberized production grid:

```
[global]
default.alias = dev
autoload.modules = hdfs.ext.kerberos

[dev.alias]
url = http://dev.namenode:port
```

(continues on next page)

(continued from previous page)

```
[prod.alias]
url = http://prod.namenode:port
client = KerberosClient
```

```
class hdfs.ext.kerberos.KerberosClient (url, mutual_auth='OPTIONAL',
                                       max_concurrency=1, root=None, proxy=None,
                                       timeout=None, session=None, **kwargs)
```

Bases: `hdfs.client.Client`

HDFS web client using Kerberos authentication.

### Parameters

- **url** – Hostname or IP address of HDFS namenode, prefixed with protocol, followed by WebHDFS port on namenode.
- **mutual\_auth** – Whether to enforce mutual authentication or not (possible values: 'REQUIRED', 'OPTIONAL', 'DISABLED').
- **max\_concurrency** – Maximum number of allowed concurrent authentication requests. This is required since requests exceeding the threshold allowed by the server will be unable to authenticate.
- **proxy** – User to proxy as.
- **root** – Root path, this will be prefixed to all HDFS paths passed to the client. If the root is relative, the path will be assumed relative to the user's home directory.
- **timeout** – Connection timeouts, forwarded to the request handler. How long to wait for the server to send data before giving up, as a float, or a (`connect_timeout`, `read_timeout`) tuple. If the timeout is reached, an appropriate exception will be raised. See the [requests](#) documentation for details.
- **session** – `requests.Session` instance, used to emit all requests.
- **\*\*kwargs** – Additional arguments passed to the underlying `HTTPKerberosAuth` class.

To avoid replay errors, a timeout of 1 ms is enforced between requests. If a session argument is passed in, it will be modified in-place to support authentication.

## Avro

Read and write [Avro](#) files directly from HDFS.

This extension enables streaming decoding and encoding of files from and to HDFS. It requires the `fastavro` library.

- `AvroWriter` writes Avro files on HDFS from python objects.
- `AvroReader` reads Avro files from HDFS into an iterable of records.

Sample usage:

It also features an entry point (named `hdfscli-avro` by default) which provides access to the above functionality from the shell. For usage examples and more information:

```
$ hdfscli-avro --help
```

```
class hdfs.ext.avro.AvroReader (client, hdfs_path, parts=None, reader_schema=None)
```

Bases: `object`

HDFS Avro file reader.



**Parameters**

- **client** – *hdfs.client.Client* instance.
- **hdfs\_path** – Remote path.
- **parts** – Part-files to read, when reading a distributed file. The default is to read all part-files in order. See *hdfs.client.Client.parts()* for details.
- **reader\_schema** – Schema to read the data as. If specified, it must be compatible with the writer's schema (the default).

The contents of the file will be decoded in a streaming manner, as the data is transferred. This makes it possible to use on files of arbitrary size. As a convenience, the content summary object of the remote file is available on the reader's `content` attribute.

Usage:

```
with AvroReader(client, 'foo.avro') as reader:
    schema = reader.writer_schema # The remote file's Avro schema.
    content = reader.content # Content metadata (e.g. size).
    for record in reader:
        pass # and its records
```

**schema**

Get the underlying file's schema.

The schema will only be available after entering the reader's corresponding `with` block.

**writer\_schema**

Get the underlying file's schema.

The schema will only be available after entering the reader's corresponding `with` block.

```
class hdfs.ext.avro.AvroWriter(client, hdfs_path, schema=None, codec=None,
                               sync_interval=None, sync_marker=None, metadata=None,
                               **kwargs)
```

Bases: object

Write an Avro file on HDFS from python dictionaries.

**Parameters**

- **client** – *hdfs.client.Client* instance.
- **hdfs\_path** – Remote path.
- **schema** – Avro schema. If not specified, the writer will try to infer it from the first record sent. There are however limitations regarding what can be inferred.
- **codec** – Compression codec. The default is 'null' (no compression).
- **sync\_interval** – Number of bytes after which a block will be written.
- **sync\_marker** – 16 byte tag used for synchronization. If not specified, one will be generated at random.
- **metadata** – Additional metadata to include in the container file's header. Keys starting with 'avro.' are reserved.
- **\*\*kwargs** – Keyword arguments forwarded to *hdfs.client.Client.write()*.

Usage:

```
with AvroWriter(client, 'data.avro') as writer:
    for record in records:
        writer.write(record)
```

**schema**

Avro schema.

**write** (*record*)

Store a record.

**Parameters** **record** – Record object to store.

Only available inside the `with` block.

## Dataframe

Read and write [Pandas](#) dataframes directly from HDFS.

This extension requires both the `avro` extension and `pandas` to be installed. Currently only Avro serialization is supported.

`hdfs.ext.dataframe.read_dataframe` (*client*, *hdfs\_path*)

Read dataframe from HDFS Avro file.

**Parameters**

- **client** – `hdfs.client.Client` instance.
- **hdfs\_path** – Remote path to an Avro file (potentially distributed).

`hdfs.ext.dataframe.write_dataframe` (*client*, *hdfs\_path*, *df*, **\*\*kwargs**)

Save dataframe to HDFS as Avro.

**Parameters**

- **client** – `hdfs.client.Client` instance.
- **hdfs\_path** – Remote path where the dataframe will be stored.
- **df** – Dataframe to store.
- **\*\*kwargs** – Keyword arguments passed through to `hdfs.ext.avro.AvroWriter`.

## 2.3.3 Configuration

Command line interface configuration module.

This module provides programmatic access to `HdfsCLI`'s configuration settings. In particular it exposes the ability to instantiate clients from aliases (see `Config.get_client()`).

**class** `hdfs.config.Config` (*path=None*, *stream\_log\_level=None*)

Bases: `ConfigParser.RawConfigParser`

Configuration class.

**Parameters**

- **path** – path to configuration file. If no file exists at that location, the configuration parser will be empty. If not specified, the value of the `HDFSCLI_CONFIG` environment variable is used if it exists, otherwise it defaults to `~/hdfscli.cfg`.

- **stream\_log\_level** – Stream handler log level, attached to the root logger. A false-ish value will disable this handler. This is particularly useful with the `catch()` function which reports exceptions as log messages.

On instantiation, the configuration object will attempt to load modules defined in the `autoload` global options (see *Custom client support* for more information).

**get\_client** (*alias=None*)  
Load HDFS client.

**Parameters** **alias** – The client to look up. If not specified, the default alias be used (`default.alias` option in the `global` section) if available and an error will be raised otherwise.

Further calls to this method for the same alias will return the same client instance (in particular, any option changes to this alias will not be taken into account).

**get\_log\_handler** (*command*)  
Configure and return log handler.

**Parameters** **command** – The command to load the configuration for. All options will be looked up in the `[COMMAND.command]` section. This is currently only used for configuring the file handler for logging. If logging is disabled for the command, a *NullHandler* will be returned, else a `TimedRotatingFileHandler`.

**class** `hdfs.config.NullHandler` (*level=0*)  
Bases: `logging.Handler`

Pass-through logging handler.

This is required for python <2.7.

Initializes the instance - basically setting the formatter to `None` and the filter list to empty.

**emit** (*record*)  
Do nothing.

`hdfs.config.catch` (*\*error\_classes*)  
Returns a decorator that catches errors and prints messages to `stderr`.

**Parameters** **\*error\_classes** – Error classes.

Also exits with status 1 if any errors are caught.

## 2.3.4 Utilities

Common utilities.

**class** `hdfs.util.AsyncWriter` (*consumer*)  
Bases: `object`

Asynchronous publisher-consumer.

**Parameters** **consumer** – Function which takes a single generator as argument.

This class can be used to transform functions which expect a generator into file-like writer objects. This can make it possible to combine different APIs together more easily. For example, to send streaming requests:

```
import requests as rq

with AsyncWriter(lambda data: rq.post(URL, data=data)) as writer:
    writer.write('Hello, world!')
```

**flush()**

Pass-through implementation.

**seekable()**

Implement file-like method expected by certain libraries.

fastavro relies on it in python 3.

**tell()**

No-op implementation.

**write(chunk)**

Stream data to the underlying consumer.

**Parameters chunk** – Bytes to write. These will be buffered in memory until the consumer reads them.

**exception** `hdfs.util.HdfsError(message, *args, **kwargs)`

Bases: `exceptions.Exception`

Base error class.

**Parameters**

- **message** – Error message.
- **args** – optional Message formatting arguments.

`hdfs.util.temppath(*args, **kws)`

Create a temporary path.

**Parameters dpath** – Explicit directory name where to create the temporary path. A system dependent default will be used otherwise (cf. `tempfile.mkstemp`).

Usage:

```
with temppath() as path:
    pass # do stuff
```

Any file or directory corresponding to the path will be automatically deleted afterwards.

## CHAPTER 3

---

### Sample script

---

More examples can be found in the `examples/` folder on GitHub.



## h

`hdfs.client`, 12  
`hdfs.config`, 22  
`hdfs.ext.avro`, 20  
`hdfs.ext.dataframe`, 22  
`hdfs.ext.kerberos`, 19  
`hdfs.util`, 23





**A**

`acl_status()` (*hdfs.client.Client method*), 12  
`allow_snapshot()` (*hdfs.client.Client method*), 12  
`AsyncWriter` (*class in hdfs.util*), 23  
`AvroReader` (*class in hdfs.ext.avro*), 20  
`AvroWriter` (*class in hdfs.ext.avro*), 21

**C**

`catch()` (*in module hdfs.config*), 23  
`checksum()` (*hdfs.client.Client method*), 12  
`Client` (*class in hdfs.client*), 12  
`Config` (*class in hdfs.config*), 22  
`content()` (*hdfs.client.Client method*), 12  
`create_snapshot()` (*hdfs.client.Client method*), 12

**D**

`delete()` (*hdfs.client.Client method*), 13  
`delete_snapshot()` (*hdfs.client.Client method*), 13  
`disallow_snapshot()` (*hdfs.client.Client method*), 13  
`download()` (*hdfs.client.Client method*), 13

**E**

`emit()` (*hdfs.config.NullHandler method*), 23

**F**

`flush()` (*hdfs.util.AsyncWriter method*), 23  
`from_options()` (*hdfs.client.Client class method*), 14

**G**

`get_client()` (*hdfs.config.Config method*), 23  
`get_log_handler()` (*hdfs.config.Config method*), 23

**H**

`hdfs.client` (*module*), 12  
`hdfs.config` (*module*), 22  
`hdfs.ext.avro` (*module*), 20  
`hdfs.ext.dataframe` (*module*), 22  
`hdfs.ext.kerberos` (*module*), 19

`hdfs.util` (*module*), 23  
`HdfsError`, 24

**I**

`InsecureClient` (*class in hdfs.client*), 18

**K**

`KerberosClient` (*class in hdfs.ext.kerberos*), 20

**L**

`list()` (*hdfs.client.Client method*), 14

**M**

`makedirs()` (*hdfs.client.Client method*), 14

**N**

`NullHandler` (*class in hdfs.config*), 23

**P**

`parts()` (*hdfs.client.Client method*), 14

**R**

`read()` (*hdfs.client.Client method*), 14  
`read_dataframe()` (*in module hdfs.ext.dataframe*), 22  
`remove_acl()` (*hdfs.client.Client method*), 15  
`remove_acl_entries()` (*hdfs.client.Client method*), 15  
`remove_default_acl()` (*hdfs.client.Client method*), 15  
`rename()` (*hdfs.client.Client method*), 15  
`rename_snapshot()` (*hdfs.client.Client method*), 16  
`resolve()` (*hdfs.client.Client method*), 16

**S**

`schema` (*hdfs.ext.avro.AvroReader attribute*), 21  
`schema` (*hdfs.ext.avro.AvroWriter attribute*), 22  
`seekable()` (*hdfs.util.AsyncWriter method*), 24  
`set_acl()` (*hdfs.client.Client method*), 16

`set_owner()` (*hdfs.client.Client method*), 16  
`set_permission()` (*hdfs.client.Client method*), 16  
`set_replication()` (*hdfs.client.Client method*), 16  
`set_times()` (*hdfs.client.Client method*), 17  
`status()` (*hdfs.client.Client method*), 17

## T

`tell()` (*hdfs.util.AsyncWriter method*), 24  
`temp_path()` (*in module hdfs.util*), 24  
`TokenClient` (*class in hdfs.client*), 19

## U

`upload()` (*hdfs.client.Client method*), 17

## W

`walk()` (*hdfs.client.Client method*), 17  
`write()` (*hdfs.client.Client method*), 18  
`write()` (*hdfs.ext.avro.AvroWriter method*), 22  
`write()` (*hdfs.util.AsyncWriter method*), 24  
`write_dataframe()` (*in module hdfs.ext.dataframe*),  
22  
`writer_schema` (*hdfs.ext.avro.AvroReader attribute*),  
21