
PyHDFS Documentation

Release 0.1

Jing Wang

May 02, 2020

Contents

1	pyhdfs module	3
2	Indices and tables	11
	Python Module Index	13
	Index	15

For a quick introduction, see the [main README](#). For detailed documentation on available methods, see [pyhdfs.HdfsClient](#).

Contents:

WebHDFS client with support for NN HA and automatic error checking

For details on the WebHDFS endpoints, see the Hadoop documentation:

- <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-hdfs/WebHDFS.html>
- <https://hadoop.apache.org/docs/current/api/org/apache/hadoop/fs/FileSystem.html>
- <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/filesystem/filesystem.html>

```
class pyhdfs.ContentSummary (**kwargs)
    Bases: pyhdfs._BoilerplateClass
```

Parameters

- **directoryCount** (*int*) – The number of directories.
- **fileCount** (*int*) – The number of files.
- **length** (*int*) – The number of bytes used by the content.
- **quota** (*int*) – The namespace quota of this directory.
- **spaceConsumed** (*int*) – The disk space consumed by the content.
- **spaceQuota** (*int*) – The disk space quota.
- **typeQuota** (*Dict[str, TypeQuota]*) – Quota usage for ARCHIVE, DISK, SSD

```
class pyhdfs.FileChecksum (**kwargs)
    Bases: pyhdfs._BoilerplateClass
```

Parameters

- **algorithm** (*str*) – The name of the checksum algorithm.
- **bytes** (*str*) – The byte sequence of the checksum in hexadecimal.
- **length** (*int*) – The length of the bytes (not the length of the string).

```
class pyhdfs.FileStatus (**kwargs)
    Bases: pyhdfs._BoilerplateClass
```

Parameters

- **accessTime** (*int*) – The access time.
- **blockSize** (*int*) – The block size of a file.
- **group** (*str*) – The group owner.
- **length** (*int*) – The number of bytes in a file.
- **modificationTime** (*int*) – The modification time.
- **owner** (*str*) – The user who is the owner.
- **pathSuffix** (*str*) – The path suffix.
- **permission** (*str*) – The permission represented as a octal string.
- **replication** (*int*) – The number of replication of a file.
- **symlink** (*Optional[str]*) – The link target of a symlink.
- **type** (*str*) – The type of the path object.
- **childrenNum** (*int*) – How many children this directory has, or 0 for files.

exception `pyhdfs.HdfsAccessControlException` (*message: str, exception: str, status_code: int, **kwargs*)

Bases: `pyhdfs.HdfsIOException`

class `pyhdfs.HdfsClient` (*hosts: Union[str, Iterable[str]] = 'localhost', randomize_hosts: bool = True, user_name: Optional[str] = None, timeout: float = 20, max_tries: int = 2, retry_delay: float = 5, requests_session: Optional[requests.sessions.Session] = None, requests_kwargs: Optional[Dict[str, Any]] = None*)

Bases: `object`

HDFS client backed by WebHDFS.

All functions take arbitrary query parameters to pass to WebHDFS, in addition to any documented keyword arguments. In particular, any function will accept `user.name`, which for convenience may be passed as `user_name`.

If multiple HA NameNodes are given, all functions submit HTTP requests to both NameNodes until they find the active NameNode.

Parameters

- **hosts** (*list or str*) – List of NameNode HTTP host:port strings, either as `list` or a comma separated string. Port defaults to 50070 if left unspecified. Note that in Hadoop 3, the default NameNode HTTP port changed to 9870; the old default of 50070 is left as-is for backwards compatibility.
- **randomize_hosts** (*bool*) – By default randomize host selection.
- **user_name** – What Hadoop user to run as. Defaults to the `HADOOP_USER_NAME` environment variable if present, otherwise `getpass.getuser()`.
- **timeout** (*float*) – How long to wait on a single NameNode in seconds before moving on. In some cases the standby NameNode can be unresponsive (e.g. loading fsimage or checkpointing), so we don't want to block on it.
- **max_tries** (*int*) – How many times to retry a request for each NameNode. If NN1 is standby and NN2 is active, we might first contact NN1 and then observe a failover to NN1 when we contact NN2. In this situation we want to retry against NN1.

- **retry_delay** (*float*) – How long to wait in seconds before going through NameNodes again
- **requests_session** – A `requests.Session` object for advanced usage. If absent, this class will use the default requests behavior of making a new session per HTTP request. Caller is responsible for closing session.
- **requests_kwargs** – Additional `**kwargs` to pass to requests

append (*path: str, data: Union[bytes, IO[bytes]], **kwargs*) → None
Append to the given file.

Parameters

- **data** – bytes or a file-like object
- **buffer_size** (*int*) – The size of the buffer used in transferring data.

concat (*target: str, sources: List[str], **kwargs*) → None
Concat existing files together.

For preconditions, see https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/filesystem/filesystem.html#void_concatPath_p_Path_sources

Parameters

- **target** – the path to the target destination.
- **sources** (*list*) – the paths to the sources to use for the concatenation.

copy_from_local (*localsrc: str, dest: str, **kwargs*) → None
Copy a single file from the local file system to `dest`

Takes all arguments that `create()` takes.

copy_to_local (*src: str, localdest: str, **kwargs*) → None
Copy a single file from `src` to the local file system

Takes all arguments that `open()` takes.

create (*path: str, data: Union[IO[bytes], bytes], **kwargs*) → None
Create a file at the given path.

Parameters

- **data** – bytes or a file-like object to upload
- **overwrite** (*bool*) – If a file already exists, should it be overwritten?
- **blocksize** (*long*) – The block size of a file.
- **replication** (*short*) – The number of replications of a file.
- **permission** (*octal*) – The permission of a file/directory. Any radix-8 integer (leading zeros may be omitted.)
- **buffer_size** (*int*) – The size of the buffer used in transferring data.

create_snapshot (*path: str, **kwargs*) → str
Create a snapshot

Parameters

- **path** – The directory where snapshots will be taken
- **snapshotname** – The name of the snapshot

Returns the snapshot path

create_symlink (*link: str, destination: str, **kwargs*) → None

Create a symbolic link at *link* pointing to *destination*.

Parameters

- **link** – the path to be created that points to target
- **destination** – the target of the symbolic link
- **createParent** (*bool*) – If the parent directories do not exist, should they be created?

Raises *HdfsUnsupportedOperationException* – This feature doesn't actually work, at least on CDH 5.3.0.

delete (*path: str, **kwargs*) → bool

Delete a file.

Parameters **recursive** (*bool*) – If *path* is a directory and set to true, the directory is deleted else throws an exception. In case of a file the recursive can be set to either true or false.

Returns true if delete is successful else false.

Return type bool

delete_snapshot (*path: str, snapshotname: str, **kwargs*) → None

Delete a snapshot of a directory

exists (*path: str, **kwargs*) → bool

Return true if the given path exists

get_active_namenode (*max_staleness: Optional[float] = None*) → str

Return the address of the currently active NameNode.

Parameters **max_staleness** (*float*) – This function caches the active NameNode. If this age of this cached result is less than *max_staleness* seconds, return it. Otherwise, or if this parameter is None, do a lookup.

Raises *HdfsNoServerException* – can't find an active NameNode

get_content_summary (*path: str, **kwargs*) → pyhdfs.ContentSummary

Return the *ContentSummary* of a given Path.

get_file_checksum (*path: str, **kwargs*) → pyhdfs.FileChecksum

Get the checksum of a file.

Return type *FileChecksum*

get_file_status (*path: str, **kwargs*) → pyhdfs.FileStatus

Return a *FileStatus* object that represents the path.

get_home_directory (***kwargs*) → str

Return the current user's home directory in this filesystem.

get_xattrs (*path: str, xattr_name: Union[str, List[str], None] = None, encoding: str = 'text',*

***kwargs*) → Dict[str, Union[bytes, str, None]]

Get one or more xattr values for a file or directory.

Parameters

- **xattr_name** – str to get one attribute, list to get multiple attributes, None to get all attributes.
- **encoding** – text | hex | base64, defaults to text

Returns Dictionary mapping xattr name to value. With text encoding, the value will be a unicode string. With hex or base64 encoding, the value will be a byte array.

Return type dict

list_status (*path: str, **kwargs*) → List[pyhdfs.FileStatus]

List the statuses of the files/directories in the given path if the path is a directory.

Return type list of *FileStatus* objects

list_xattrs (*path: str, **kwargs*) → List[str]

Get all of the xattr names for a file or directory.

Return type list

listdir (*path: str, **kwargs*) → List[str]

Return a list containing names of files in the given path

makedirs (*path: str, **kwargs*) → bool

Create a directory with the provided permission.

The permission of the directory is set to be the provided permission as in `setPermission`, not `permission&~umask`.

Parameters **permission** (*octal*) – The permission of a file/directory. Any radix-8 integer (leading zeros may be omitted.)

Returns true if the directory creation succeeds; false otherwise

Return type bool

open (*path: str, **kwargs*) → IO[bytes]

Return a file-like object for reading the given HDFS path.

Parameters

- **offset** (*long*) – The starting byte position.
- **length** (*long*) – The number of bytes to be processed.
- **buffer_size** (*int*) – The size of the buffer used in transferring data.

Return type file-like object

remove_xattr (*path: str, xattr_name: str, **kwargs*) → None

Remove an xattr of a file or directory.

rename (*path: str, destination: str, **kwargs*) → bool

Renames Path src to Path dst.

Returns true if rename is successful

Return type bool

rename_snapshot (*path: str, oldsnapshotname: str, snapshotname: str, **kwargs*) → None

Rename a snapshot

set_owner (*path: str, **kwargs*) → None

Set owner of a path (i.e. a file or a directory).

The parameters owner and group cannot both be null.

Parameters

- **owner** – user
- **group** – group

set_permission (*path: str, **kwargs*) → None

Set permission of a path.

Parameters `permission` (*octal*) – The permission of a file/directory. Any radix-8 integer (leading zeros may be omitted.)

set_replication (*path: str, **kwargs*) → bool

Set replication for an existing file.

Parameters `replication` (*short*) – new replication

Returns true if successful; false if file does not exist or is a directory

Return type bool

set_times (*path: str, **kwargs*) → None

Set access time of a file.

Parameters

- **modificationtime** (*long*) – Set the modification time of this file. The number of milliseconds since Jan 1, 1970.
- **accesstime** (*long*) – Set the access time of this file. The number of milliseconds since Jan 1 1970.

set_xattr (*path: str, xattr_name: str, xattr_value: Optional[str], flag: str, **kwargs*) → None

Set an xattr of a file or directory.

Parameters

- **xattr_name** – The name must be prefixed with the namespace followed by .. For example, `user.attr`.
- **flag** – CREATE or REPLACE

walk (*top: str, topdown: bool = True, onerror: Optional[Callable[[pyhdfs.HdfsException], None]] = None, **kwargs*) → Iterator[Tuple[str, List[str], List[str]]]

See `os.walk` for documentation

exception `pyhdfs.HdfsDSQuotaExceededException` (*message: str, exception: str, status_code: int, **kwargs*)

Bases: `pyhdfs.HdfsQuotaExceededException`

exception `pyhdfs.HdfsException`

Bases: `Exception`

Base class for all errors while communicating with WebHDFS server

exception `pyhdfs.HdfsFileAlreadyExistsException` (*message: str, exception: str, status_code: int, **kwargs*)

Bases: `pyhdfs.HdfsIOException`

exception `pyhdfs.HdfsFileNotFoundException` (*message: str, exception: str, status_code: int, **kwargs*)

Bases: `pyhdfs.HdfsIOException`

exception `pyhdfs.HdfsHadoopIllegalArgumentException` (*message: str, exception: str, status_code: int, **kwargs*)

Bases: `pyhdfs.HdfsIllegalArgumentException`

exception `pyhdfs.HdfsHttpException` (*message: str, exception: str, status_code: int, **kwargs*)

Bases: `pyhdfs.HdfsException`

The client was able to talk to the server but got a HTTP error code.

Parameters

- **message** – Exception message

- **exception** – Name of the exception
- **javaClassName** – Java class name of the exception
- **status_code** (*int*) – HTTP status code
- **kwargs** – any extra attributes in case Hadoop adds more stuff

exception `pyhdfs.HdfsIOException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsHttpException`

exception `pyhdfs.HdfsIllegalArgumentException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsHttpException`

exception `pyhdfs.HdfsInvalidPathException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsHadoopIllegalArgumentException`

exception `pyhdfs.HdfsNSQuotaExceededException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsQuotaExceededException`

exception `pyhdfs.HdfsNoServerException`
 Bases: `pyhdfs.HdfsException`

The client was not able to reach any of the given servers

exception `pyhdfs.HdfsPathIsNotEmptyDirectoryException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsIOException`

exception `pyhdfs.HdfsQuotaExceededException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsIOException`

exception `pyhdfs.HdfsRemoteException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsIOException`

exception `pyhdfs.HdfsRetriableException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsIOException`

exception `pyhdfs.HdfsRuntimeException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsHttpException`

exception `pyhdfs.HdfsSecurityException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsHttpException`

exception `pyhdfs.HdfsSnapshotException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsIOException`

exception `pyhdfs.HdfsStandbyException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsIOException`

exception `pyhdfs.HdfsUnsupportedOperationException` (*message: str, exception: str, status_code: int, **kwargs*)
 Bases: `pyhdfs.HdfsHttpException`

class `pyhdfs.TypeQuota` (***kwargs*)
 Bases: `pyhdfs._BoilerplateClass`

Parameters

- **consumed** (*int*) – The storage type space consumed.
- **quota** (*int*) – The storage type quota.

CHAPTER 2

Indices and tables

- `genindex`
- `modindex`
- `search`

p

pyhdfs, 3

A

append() (*pyhdfs.HdfsClient method*), 5

C

concat() (*pyhdfs.HdfsClient method*), 5
 ContentSummary (*class in pyhdfs*), 3
 copy_from_local() (*pyhdfs.HdfsClient method*), 5
 copy_to_local() (*pyhdfs.HdfsClient method*), 5
 create() (*pyhdfs.HdfsClient method*), 5
 create_snapshot() (*pyhdfs.HdfsClient method*), 5
 create_symlink() (*pyhdfs.HdfsClient method*), 5

D

delete() (*pyhdfs.HdfsClient method*), 6
 delete_snapshot() (*pyhdfs.HdfsClient method*), 6

E

exists() (*pyhdfs.HdfsClient method*), 6

F

FileChecksum (*class in pyhdfs*), 3
 FileStatus (*class in pyhdfs*), 3

G

get_active_namenode() (*pyhdfs.HdfsClient method*), 6
 get_content_summary() (*pyhdfs.HdfsClient method*), 6
 get_file_checksum() (*pyhdfs.HdfsClient method*), 6
 get_file_status() (*pyhdfs.HdfsClient method*), 6
 get_home_directory() (*pyhdfs.HdfsClient method*), 6
 get_xattrs() (*pyhdfs.HdfsClient method*), 6

H

HdfsAccessControlException, 4
 HdfsClient (*class in pyhdfs*), 4
 HdfsDSQuotaExceededException, 8

HdfsException, 8
 HdfsFileAlreadyExistsException, 8
 HdfsFileNotFoundException, 8
 HdfsHadoopIllegalArgumentException, 8
 HdfsHttpException, 8
 HdfsIllegalArgumentException, 9
 HdfsInvalidPathException, 9
 HdfsIOException, 9
 HdfsNoServerException, 9
 HdfsNSQuotaExceededException, 9
 HdfsPathIsNotEmptyDirectoryException, 9
 HdfsQuotaExceededException, 9
 HdfsRemoteException, 9
 HdfsRetriableException, 9
 HdfsRuntimeException, 9
 HdfsSecurityException, 9
 HdfsSnapshotException, 9
 HdfsStandbyException, 9
 HdfsUnsupportedOperationException, 9

L

list_status() (*pyhdfs.HdfsClient method*), 7
 list_xattrs() (*pyhdfs.HdfsClient method*), 7
 listdir() (*pyhdfs.HdfsClient method*), 7

M

makedirs() (*pyhdfs.HdfsClient method*), 7

O

open() (*pyhdfs.HdfsClient method*), 7

P

pyhdfs (*module*), 3

R

remove_xattr() (*pyhdfs.HdfsClient method*), 7
 rename() (*pyhdfs.HdfsClient method*), 7
 rename_snapshot() (*pyhdfs.HdfsClient method*), 7

S

`set_owner()` (*pyhdfs.HdfsClient method*), 7
`set_permission()` (*pyhdfs.HdfsClient method*), 7
`set_replication()` (*pyhdfs.HdfsClient method*), 8
`set_times()` (*pyhdfs.HdfsClient method*), 8
`set_xattr()` (*pyhdfs.HdfsClient method*), 8

T

`TypeQuota` (*class in pyhdfs*), 9

W

`walk()` (*pyhdfs.HdfsClient method*), 8