

Hadoop Archives

Table of contents

1 What are Hadoop archives?	2
2 How to create an archive?	2
3 How to look up files in archives?	2

1. What are Hadoop archives?

Hadoop archives are special format archives. A Hadoop archive maps to a FileSystem directory. A Hadoop archive always has a *.har extension. A Hadoop archive directory contains metadata (in the form of `_index` and `_masterindex`) and data (`part-*`) files. The `_index` file contains the name of the files that are part of the archive and their location within the part files.

2. How to create an archive?

```
Usage: hadoop archive -archiveName name <src>* <dest>
```

`-archiveName` is the name of the archive you would like to create. An example would be `foo.har`. The name should have a *.har extension. The inputs are filesystem pathnames which work as usual with regular expressions. The destination directory would contain the archive. Note that this is a Map Reduce job that creates the archives. You would need a map reduce cluster to run this. The following is an example:

```
hadoop archive -archiveName foo.har /user/hadoop/dir1  
/user/hadoop/dir2 /user/zoo/
```

In the above example `/user/hadoop/dir1` and `/user/hadoop/dir2` will be archived in the following filesystem directory -- `/user/zoo/foo.har`. The sources are not changed or removed when an archive is created.

3. How to look up files in archives?

The archives exposes itself as a filesystem layer. So all the fs shell commands in the archives work but with a different URI. Also, note that archives are immutable. So, `rename`'s, `delete`s and `create`s return an error. URI for Hadoop Archives is

```
har://scheme-hostname:port/archivepath/fileinarchive
```

If no scheme is provided it assumes the underlying filesystem. In that case the URI would look like

```
har:///archivepath/fileinarchive
```