# **Genomedata Documentation**

Release 1.7.4

Michael M. Hoffman

## **CONTENTS**

1	Genomedata 1.7 documentation					
	1.1	Installation	3			
	1.2	Overview	5			
	1.3	Implementation	5			
	1.4	Creation	6			
	1.5	Genomedata usage	8			
	1.6	Tips and tricks	18			
	1.7	Technical matters	19			
	1.8	Bugs	19			
	1.9	Support	19			
2 Indices and tables			21			
Ру	Python Module Index					
In						

Contents:

CONTENTS 1

2 CONTENTS

## **GENOMEDATA 1.7 DOCUMENTATION**

#### Website

http://pmgenomics.ca/hoffmanlab/proj/genomedata/

#### Author

Michael M. Hoffman <michael dot hoffman at utoronto dot ca>

#### **Organization**

Princess Margaret Cancer Centre

#### Address

Toronto Medical Discovery Tower 11-311, 101 College St, M5G 1L7, Toronto, Ontario, Canada

### Copyright

2009-2014 Michael M. Hoffman

For a broad overview, see the paper:

Hoffman MM, Buske OJ, Noble WS. (2010). The Genomedata format for storing large-scale functional genomics data. *Bioinformatics*, **26**(11):1458-1459; doi:10.1093/bioinformatics/btq164

Please cite this paper if you use Genomedata.

## 1.1 Installation

Python 3.7 (or later) and HDF5 are required before you can install Genomedata.

## 1.1.1 Installing HDF5

Ubuntu/Debian:

sudo apt-get install hdf5-tools

CentOS/RHEL/Fedora:

sudo yum -y install hdf5

OpenSUSE:

sudo zypper in hdf5 libhdf5

If HDF5 has been installed from source, set the HDF5\_DIR environment variable to the directory where it was installed.

## 1.1.2 Installing from source

Additional packages are necessary when installing from source.

Ubuntu/Debian:

sudo apt-get install libhdf5-serial-dev

CentOS/RHEL/Fedora:

sudo yum -y install hdf5-devel

OpenSUSE:

sudo zypper in hdf5-devel

## 1.1.3 Installing Genomedata

Genomedata can be installed through pip or Bioconda.

With Bioconda:

conda install genomedata

With Python 3.7 or later installed:

pip install genomedata

#### Note

The latest version of genomedata will not install with older versions of pip due to updated build dependencies. You can update your pip using the command:

pip install --upgrade pip

#### Note

Genomedata is only supported on 64 bit systems.

#### Note

The following are prerequisites:

• Linux/Unix

This software has been tested on Linux and Mac OS X systems. We would love to add support for other systems in the future and will gladly accept any contributions toward this end.

• Zlib

#### Note

For questions, comments, or troubleshooting, please refer to the *support* section.

## 1.2 Overview

Genomedata is a file format and Python API for accessing large-scale functional genomics data. The file format is both space-efficient and allows efficient random-access. The Python API works with Genomedata archives and bigWig files.

Under the surface, the Genomedata file format is *implemented* as one or more HDF5 files, but Genomedata provides a transparent interface to interact with your underlying data without having to worry about the mess of repeatedly parsing large data files or having to keep them in memory for random access. The Genomedata archives are currently write-once, although we are working to fix this.

The Genomedata hierarchy:

#### Each Genome contains many Chromosomes

#### Each Chromosome contains many Supercontigs

#### Each Supercontig contains one continuous data set

Each continuous data set is a numpy.array of floating point numbers with a column for each data track and a row for each base in the data set.

## Why have Supercontigs?

Genomic data seldom covers the entire genome but instead tends to be defined in large but scattered regions. In order to avoid storing the undefined data between the regions, chromosomes are divided into separate supercontigs when regions of defined data are far enough apart. They also serve as a convenient chunk since they can usually fit entirely in memory.

## 1.3 Implementation

Genomedata archives are implemented as one or more HDF5 files. The *API* handles both single-file and directory archives transparently, but the implementation options exist for several performance reasons.

#### Use a directory with few chromosomes/scaffolds:

- Parallel load/access
- · Smaller file sizes

#### Use a single file with many chromosomes/scaffolds:

- More efficient access with many chromosomes/scaffolds
- · Easier archive distribution

Implementing the archive as a directory makes it easier to parallelize access to the data. In particular, it makes it easy to create the archives in parallel with one chromosome on each machine. It also reduces the likelihood of running into the 2 GB file limit applicable to older applications and older versions of 32-bit UNIX. We are currently using an 81-track Genomedata archive for our research which has a total size of 18 GB, but the largest single file (chr1) is only 1.6 GB.

A directory-based Genomedata archive is not ideal for all circumstances, however, such as when working with genomes with many chromosomes, contigs, or scaffolds. In these situations, a single file implementation would be much more efficient. Additionally, having the archive as a single file allows the archive to be distributed much more easily (without tar/zip/etc).

#### Note

The default behavior is to implement the Genomedata archive as a directory if there are fewer than 100 sequences being loaded and as a single file otherwise.

Added in version 1.1: Single-file-based Genomedata archives

1.2. Overview 5

## 1.4 Creation

A Genomedata archive contains sequence and may also contain numerical data associated with that sequence. You can easily load sequence and numerical data into a Genomedata archive with the *genomedata-load* command (see command details additional details):

```
genomedata-load [-t trackname=signalfile]... [-s sequencefile]... GENOMEDATAFILE
```

This command is a user-friendly shortcut to the typical workflow. The underlying commands are still installed and may be used if more fine-grained control is required (for instance, parallel data loading or adding additional tracks later). The commands and required ordering are:

- 1. genomedata-load-seq
- 2. genomedata-open-data
- 3. genomedata-load-data
- 4. genomedata-close-data

Entire data tracks can later be replaced with the following pipeline:

- 1. genomedata-erase-data
- 2. genomedata-load-data
- 3. genomedata-close-data

Added in version 1.1: The ability to replace data tracks.

Additional data tracks can be added to an existing archive with the following pipeline:

- 1. genomedata-open-data
- 2. genomedata-load-data
- 3. genomedata-close-data

Added in version 1.2: The ability to add data tracks.

As of the current version, Genomedata archives must include the underlying genomic sequence and can only be created with *genomedata-load-seq*. A Genomedata archive can be created without any tracks, however, using the following pipeline:

- 1. genomedata-load-seq
- 2. genomedata-close-data

Added in version 1.2: The ability to create an archive without any data tracks.

Additionally, you may remove portions of data from tracks by hardmasking the specified data tracks. This can be done anytime after loading in data and unless specified otherwise will automatically close the archive as well. A track can be loaded and filtered with the following pipeline:

- 1. genomedata-open-data
- 2. genomedata-load-data
- 3. genomedata-hardmask

Added in version 1.4: The ability to hardmask tracks.

#### Note

A call to **h5repack** after *genomedata-close-data* may be used to transparently compress the data.

## 1.4.1 Example

The following is a brief example for creating a Genomedata archive from sequence and signal files.

Given the following two sequence files:

1. A text file, chr1.fa:

```
>chr1
taaccctaaccctaaccctaaccctaaccctaacccta
accctaaccctaaccctaacccta
```

2. A compressed text file, *chrY.fa.gz*:

```
>chrY
ctaaccctaaccctaaccctaaccctCTGaaagtggac
```

and the following two signal files:

1. signal\_low.wigFix:

```
fixedStep chrom=chr1 start=5 step=1
0.372
-2.540
0.371
-2.611
0.372
-2.320
```

2. signal\_high.bed.gz:

```
chrY
        0
                 12
                          4.67
chrY
        20
                 23
                          9.24
chr1
                 3
                          2.71
        1
chr1
        3
                 6
                          1.61
chr1
        6
                 24
                          3.14
```

A Genomedata archive (genomedata.test) could then be created with the following command:

```
genomedata-load -s chr1.fa -s chrY.fa.gz -t low=signal_low.wigFix \
-t high=signal_high.bed.gz genomedata.test
```

or the following pipeline:

```
genomedata-load-seq genomedata.test chr1.fa chrY.fa.gz
genomedata-open-data genomedata.test low high
genomedata-load-data genomedata.test low < signal_low.wigFix
zcat signal_high.bed.gz | genomedata-load-data genomedata.test high
genomedata-close-data genomedata.test
```

1.4. Creation 7

#### Note

chr1.fa and chrY.fa.gz could also be combined into a single sequence file with two sequences.

#### Note

If using a glob syntax for your sequence files, remember to put the glob filename in quotes to avoid having your shell expand the glob before it genomedata-load uses it (e.g. -s "chr\*.agp.gz")

#### Warning

It is important that the sequence names (*chrY*, *chr1*) in the signal files match the sequence identifiers in the sequence files exactly.

## 1.5 Genomedata usage

## 1.5.1 Python interface

The data in Genomedata is accessed through the hierarchy described in *Overview*. A full *Python API* is also available.

#### Note

The Python API expects that a genomedata archive has already been created. This can be done manually via the *genomedata-load* command. Alternatively, this can be done programmatically using :\_load\_seq:*load\_seq*.

To appreciate the full benefit of Genomedata, it is most easily used as a context manager:

```
from genomedata import Genome
[...]
gdfilename = "/path/to/genomedata/archive" # or bigWig file
with Genome(gdfilename) as genome:
    [...]
```

#### Note

If Genome is used as a context manager, it will clean up any opened Chromosomes automatically. If not, the Genome object (and all opened chromosomes) should be closed manually with a call to <code>Genome.close()</code>.

## 1.5.2 Basic usage

Genomedata is designed to make it easy to get to the data you want.

Here are a few examples:

**Get arbitrary sequence** (10-bp sequence starting at chr2:1423):

```
>>> chromosome = genome["chr2"]
>>> seq = chromosome.seq[1423:1433]
```

(continues on next page)

(continued from previous page)

```
>>> seq
array([116, 99, 99, 99, 103, 103, 103, 103], dtype=uint8)
>>> seq.tostring()
'tccccggggg'
```

Get arbitrary data (data from first 3 tracks for region chr8:999-1000):

Get data for a specific track (specified data in first 5-bp of chr1):

```
>>> chromosome = genome["chr1"]
>>> data = chromosome[0:5, "sample_track"]
>>> data
array([ 47., NaN, NaN, NaN, NaN], dtype=float32)
```

Only specified data:

```
>>> from numpy import isfinite
>>> data[isfinite(data)]
array([ 47.], dtype=float32)
```

```
Note

Specify a slice for the track to keep the data in column form:

>>> col_index = chromosome.index_continuous("sample_track")

>>> data = chromosome[0:5, col_index:col_index+1]
```

## 1.5.3 BigWig differences

There are a number of minor differences between using the genomedata file format and the bigWig file format:

- There is only one track per bigWig file and it is implicitly set to the filename of the bigWig.
- Track indexing is only used to shape dimensionality of output.
- Summary statistics are taken from the bigWig file formation definition which are stored as integers. There may be some differences precision.
- Each Chromosome is represented with 1 underlying Supercontig.

#### **Command-line interface**

Genomedata archives can be created and loaded from the command line with the genomedata-load command.

## 1.5.4 genomedata-load

This is a convenience script that will do everything necessary to create a Genomedata archive. This script takes as input:

- assembly files in either FASTA (.fa or .fa.gz) format (where the sequence identifiers are the names of the chromosomes/scaffolds to create), or assembly files in AGP format (when used with --assembly). This is mandatory, despite having an option interface.
- trackname, datafile pairs (specified as trackname=datafile), where:
  - trackname is a string identifier (e.g. broad.h3k27me3)
  - datafile contains signal data for this data track in one of the following formats: WIG, BED3+1, bed-Graph, or a gzip'd form of any of the preceding
  - the chromosomes/scaffolds referred to in the datafile MUST be identical to those found in the sequence files
- the name of the Genomedata archive to create

See the *full example* for more details.

#### **Chromosome naming**

When loading sequence data that does not have a UCSC-style name (e.g. 'chr1'), you may provide a tab-delimited file that provides a mapping between naming styles to be given to the '-assembly-report' option. This file expects a commented header to provide labels to the columns below it. For example the following is a valid file for '-assembly-report':

```
# Sequence-Name GenBank-Accn RefSeq-Accn UCSC-style-name
1 CM000663.2 NC_000001.11 chr1
```

This file format style is based on the assembly reports provided by NCBI.

Command-line usage information:

```
usage: genomedata-load [-h] [-r ASSEMBLY-REPORT] [-n NAME_STYLE] [--version] [--verbose]
→-s SEQUENCE -t NAME=FILE [-m MASKFILE] [--assembly | --sizes] [-f | -d] GENOMEDATAFILE
Create Genomedata archive named GENOMEDATAFILE by loading
specified track data and sequences. If GENOMEDATAFILE
already exists, it will be overwritten.
--track and --sequence may be repeated to specify
multiple trackname=trackfile pairings and sequence files,
respectively.
Example: genomedata-load -t high=signal.high.wig -t low=signal.low.bed.gz -s chrX.fa -s_
positional arguments:
 GENOMEDATAFILE
                       genomedata archive
optional arguments:
 -h, --help
                       show this help message and exit
 --version
                       show program's version number and exit
Chromosome naming:
 -r ASSEMBLY-REPORT, --assembly-report ASSEMBLY-REPORT
                      Tab-delimited file with columnar mappings between chromosome_
-n NAME_STYLE, --name-style NAME_STYLE
```

(continues on next page)

(continued from previous page)

```
Chromsome naming style to use based on ASSEMBLY-REPORT. Default:
→UCSC-style-name
Flags:
                       Print status updates and diagnostic messages
  --verbose
Input data:
  -s SEQUENCE, --sequence SEQUENCE
                        Add the sequence data in the specified file or files (may use
→UNIX glob wildcard syntax)
  -t NAME=FILE, --track NAME=FILE
                        Add data from FILE as the track NAME, such as: -t signal=signal.
-wig
  -m MASKFILE, --maskfile MASKFILE
                        A BED file containing regions to mask out from tracks before.
→loading
  --assembly
                        sequence files contain assembly (AGP) files instead of sequence
  --sizes
                        sequence files contain list of sizes instead of sequence
Implementation:
  -f, --file-mode
                       If specified, the Genomedata archive will be implemented as a
⇒single file, with a separate h5 group for each Chromosome. This is recommended if

→ there are a large

                        number of Chromosomes. The default behavior is to use a single_
→file if there are at least 100 Chromosomes being added.
  -d, --directory-mode If specified, the Genomedata archive will be implemented as a.
→directory, with a separate file for each Chromosome. This is recommended if there are
→a small number of
                        Chromosomes. The default behavior is to use a directory if there...
→are fewer than 100 Chromosomes being added.
```

Alternately, as described in *Overview*, the underlying Python and C load scripts are also accessible for more finely-grained control. This can be especially useful for parallelizing Genomedata loading over a cluster.

You can use wildcards when specifying sequence files, such as in <code>genomedata-load-seq -s 'chr\*.fa'</code>. You must be sure to quote the wildcards so that they are not expanded by your shell. For most shells, this means using single quotes ('chr\*.fa') instead of double quotes ("chr\*.fa").

If you aren't going to use the sequence later on, loading the assembly from an AGP file will be faster and take less memory during loading, and disk space afterward.

### 1.5.5 genomedata-load-seq

This command adds the provided sequence files to the specified Genomedata, archive creating it if it does not already exist. Sequence files should be in FASTA (.fa or .fa.gz) format. Gaps of >= 100,000 base pairs in the reference sequence, are used to divide the sequence into supercontigs. The FASTA definition line will be used as the name for the chromosomes/scaffolds created within the Genomedata archive and must be consistent between these sequence files and the data loaded later with *genomedata-load-data*. See *this example* for details.

(continued from previous page)

```
→ should be in fasta format, and a separate Chromosome will be created for each.
→definition line.
positional arguments:
  adarchive
                        genomedata archive
  seqfiles
                        sequences in FASTA format
optional arguments:
                        show this help message and exit
  -h, --help
                        show program's version number and exit
  --version
                        SEQFILE contains assembly (AGP) files instead of sequence
  -a, --assembly
  -s, --sizes
                        SEQFILE contains list of sizes instead of sequence
 -f, --file-mode
                        If specified, the Genomedata archive will be implemented as a.
⇒single file, with a separate h5 group for each Chromosome. This is recommended if __

→there are a large

                        number of Chromosomes. The default behavior is to use a single_
⇒file if there are at least 100 Chromosomes being added.
  -d, --directory-mode If specified, the Genomedata archive will be implemented as a.
→directory, with a separate file for each Chromosome. This is recommended if there are
→a small number of
                        Chromosomes. The default behavior is to use a directory if there.
→are fewer than 100 Chromosomes being added.
                       Print status updates and diagnostic messages
  --verbose
Chromosome naming:
  -r ASSEMBLY-REPORT, --assembly-report ASSEMBLY-REPORT
                        Tab-delimited file with columnar mappings between chromosome_
→naming styles.
  -n NAME_STYLE, --name-style NAME_STYLE
                        Chromsome naming style to use based on ASSEMBLY-REPORT. Default:
→UCSC-style-name
```

#### 1.5.6 genomedata-open-data

This command opens the specified tracks in the Genomedata archive, allowing data for those tracks to be loaded with *genomedata-load-data*.

```
usage: genomedata-open-data [-h] [-v] --trackname TRACKNAME [TRACKNAME ...]
                            [--verbose]
                            gdarchive
Open one or more tracks in the specified Genomedata archive.
positional arguments:
  gdarchive
                        genomedata archive
optional arguments:
                        show this help message and exit
  -h, --help
  -v, --version
                        show program's version number and exit
  --trackname TRACKNAME [TRACKNAME ...]
                        tracknames to open
  --verbose
                        Print status updates and diagnostic messages
```

## 1.5.7 genomedata-hardmask

This command permanently and irreversibly masks out regions from tracks in the Genomedata archive. Due to slow performance, it is not recommended for masking large genome-wide datasets. In the case of very large datasets, it is recommended you mask or filter your data first, then load the masked data with genomedata-load-data.

```
usage: genomedata-hardmask [-h] [-v] [-t TRACKNAME [TRACKNAME ...]]
                            [--hardmask OPERATOR] [--no-close] [--dry-run]
                            [--verbose]
                            maskfile qdarchive
Permanently mask TRACKNAME(s) from a genomedata archive with MASKFILE using an
optional filter operator.
positional arguments:
  maskfile
                        input mask file
  gdarchive
                        genomedata archive
optional arguments:
  -h, --help
                        show this help message and exit
                        show program's version number and exit
  -v, --version
  -t TRACKNAME [TRACKNAME ...], --trackname TRACKNAME [TRACKNAME ...]
                        Track(s) to be filtered (default: all)
                        Specify a comparison operation on a value to mask out
  --hardmask OPERATOR
                        (e.g. "lt0.5" will mask all values less than 0.5). See
                        the bash comparison operators for the two letter
                        operations (default: all values masked)
  --no-close
                        Do not close the genomedata archive after masking
                        Do not perform any masking. Useful with verbosity set
  --dry-run
                        to see what regions would be filtered
                        Print status and diagnostic messages
  --verbose
```

## 1.5.8 genomedata-load-data

This command loads data from stdin into Genomedata under the given trackname. The input data must be in one of these supported datatypes: WIG, BED3+1, bedGraph. The chromosome/scaffold references in these files must match the sequence identifiers in the sequence files loaded with *genomedata-load-seq*. See *this example* for details. A chunk-size can be specified to control the size of hdf5 chunks (the smallest data read size, like a page size). Larger values of chunk-size can increase the level of compression, but they also increase the minimum amount of data that must be read to access a single value.

BED3+1 format is interpreted the same ways as bedGraph, except that the track definition line is not required.

(continues on next page)

(continued from previous page)

Mandatory or optional arguments to long options are also mandatory or optional for any corresponding short options.

### 1.5.9 genomedata-close-data

Closes the specified Genomedata arhive.

## 1.5.10 genomedata-erase-data

Erases all data associated with the specified tracks, allowing the data to then be replaced. The pipeline for replacing a data track is:

- 1. genomedata-erase-data
- 2. genomedata-load-data
- 3. genomedata-close-data

```
usage: genomedata-erase-data [-h] [-v] --trackname TRACKNAME [TRACKNAME ...]
                             [--verbose]
                             gdarchive
Erase the specified tracks from the Genomedata archive in such a way that the
track data can be replaced (via genomedata-load-data).
positional arguments:
  gdarchive
                        genomedata archive
optional arguments:
                        show this help message and exit
  -h, --help
                        show program's version number and exit
  -v, --version
  --trackname TRACKNAME [TRACKNAME ...]
                        tracknames to erase
  --verbose
                        Print status updates and diagnostic messages
```

### 1.5.11 genomedata-info

This command displays information about a genomedata archive. Running the following command:

```
genomedata-info tracknames_continuous genomedata
```

displays the list of continous tracks. Running:

```
genomedata-info contigs genomedata
```

displays the list of contigs in BED format (0-based, half-open indexing).

This command generates a tab-delimited file containing chromosome name and sizes, suitable for use as a UCSC "chrom sizes" file:

```
genomedata-info sizes genomedata
```

## 1.5.12 genomedata-query

Prints data from a genomedata archive, for the track TRACKNAME, on CHROM, in the region BEGIN-END (0-based, half-open indexing). Intended as a convenience function only; this is much slower than the Python interface, so it should not be used for large regions.

```
usage: genomedata-query [-h] [-v] gdarchive trackname chrom begin end
print data from genomedata archive in specified trackname and coordinates
positional arguments:
  gdarchive
                 genomedata archive
  trackname
                 track name
  chrom
                 chromosome name
  begin
                 chromosome start
  end
                 chromosome end
optional arguments:
                 show this help message and exit
  -h, --help
  -v, --version show program's version number and exit
```

### **Python API**

The Genomedata package is designed to be used from a variety of scripting languages, but currently only exports the following Python API.

```
class genomedata.Genome(filename, *args, **kwargs)
```

The root level of the genomedata object hierarchy.

If you use this as a context manager, it will keep track of any open Chromosomes and close them (and the Genome object) for you later when the context is left:

```
with Genome("/path/to/genomedata") as genome:
   chromosome = genome["chr1"]
[...]
```

If not used as a context manager, you are responsible for closing the Genomedata archive once you are done:

```
>>> genome = Genome("/path/to/genomedata")
>>> chromosome = genome["chr1"]
[...]
>>> genome.close()
```

```
__init__(filename, *args, **kwargs)
```

Create a Genome object from a genomedata archive.

#### **Parameters**

- **filename** (*string*) the root of the Genomedata object hierarchy. This can either be a .genomedata file that contains the entire genome or a directory containing multiple chromosome files.
- \*args args passed on to open\_file if single file or to Chromosome if directory
- \*\*kwargs keyword args passed on to open\_file if single file or to Chromosome if directory

#### Example:

```
>>> genome = Genome("./genomedata.ctcf.pol2b/")
>>> genome
Genome("./genomedata.ctcf.pol2b/")
    [...]
>>> genome.close()
>>> genome = Genome("./cat_chipseq.genomedata", mode="r")
    [...]
>>> genome.close()
```

## \_\_iter\_\_()

Return next chromosome, in sorted order, with memoization.

#### Example:

```
for chromosome in genome:
    print chromosome.name
    for supercontig, continuous in chromosome.itercontinuous():
        [...]
```

```
__getitem__(name)
```

Return a reference to a chromosome of the given name.

#### **Parameters**

**name** (*string*) – name of the chromosome (e.g. "chr1" if chr1.genomedata is a file in the Genomedata archive or chr1 is a top-level group in the single-file Genomedata archive)

#### Returns

Chromosome

#### Example:

```
>>> genome["chrX"]
<Chromosome 'chrX', file='/path/to/genomedata/chrX.genomedata'>
>>> genome["chrZ"]
KeyError: 'Could not find chromosome: chrZ'
```

#### add\_track\_continuous(trackname)

Add a new track

The Genome object must have been created with :param mode:="r+". Behavior is undefined if this is not the case.

Currently sets the dirty bit, which can only be erased with genomedata-close-data

#### close()

Close this Genomedata archive and any open chromosomes

If the Genomedata archive is a directory, this closes all open chromosomes. If it is a single file, this closes that file. This should only be used if Genome is not a context manager (see *Genome*). The behavior is undefined if this is called while Genome is being used as a context manager.

#### erase\_data(trackname)

Erase all data for the given track across all chromosomes

The Genome object must have been created with :param mode:="r+". Behavior is undefined if this is not the case.

Currently sets the dirty bit, which can only be erased with genomedata-close-data

#### property format\_version

Genomedata format version

None means there are no chromosomes in it already or there is no information available.

#### index continuous(trackname)

Return the column index of the trackname in the continuous data.

#### Parameters

```
trackname (string) – name of data track
```

#### Returns

integer

This is used for efficient indexing into continuous data:

```
>>> col_index = genome.index_continuous("sample_track")
>>> data = genome["chr3"][100:150, col_index]
```

although for typical use, the track can be indexed directly:

```
>>> data = genome["chr3"][100:150, "sample_track"]
```

#### property isopen

Return a boolean indicating if the Genome is still open

#### property maxs

Return a vector of the maximum value for each track.

#### Returns

numpy.array

#### property means

Return a vector of the mean value of each track.

#### **Returns**

numpy.array

#### property mins

Return the minimum value for each track.

#### Returns

numpy.array

#### property num\_datapoints

Return the number of datapoints in each track.

#### Returns

a numpy.array vector with an entry for each track.

#### property num\_tracks\_continuous

Returns the number of continuous data tracks.

#### property sums

Return a vector of the sum of the values for each track.

#### Returns

numpy.array

#### property sums\_squares

Return a vector of the sum of squared values for each track's data.

#### Returns

numpy.array

## property tracknames\_continuous

Return a list of the names of all data tracks stored.

#### property vars

Return a vector of the variance in the data for each track.

#### Returns

numpy.array

## 1.6 Tips and tricks

If you find yourself creating many Genomedata archives on the same genome, it might be useful to save a copy of an archive after you load sequence, but before you load any data. Obviously, you can only do this if you use the fine-grained workflow of *genomedata-load-seq*, *genomedata-open-data*, *genomedata-load-data*, and *genomedata-close-data*.

## 1.7 Technical matters

## 1.7.1 Chunking and chunk cache overhead

Genomedata uses an HDF5 data store. The data is stored in chunks. The chunk size is 10,000 bp and one data track of 32-bit single-precision floats, which makes the chunk 40 kB. Each chunk is gzip compressed so on disk it will be smaller. To read a single position you have to read its entire chunk off of the disk and then decompress it. There is a tradeoff here between latency and throughput. Larger chunk sizes mean more latency but better throughput and better compression.

The only disk storage overhead is that compression is slightly less efficient than compressing the whole binary data file when you break it into chunks. This is far outweighed by the efficient random access capability. If you have different needs, then it should be possible to change the chunk shape (genomedata.CONTINUOUS\_CHUNK\_SHAPE) or compression method (genomedata.\_util.FILTERS\_GZIP).

The memory overhead is dominated by the chunk cache defined by PyTables. On the version of PyTables we use, this is 2 MiB. You can change this by setting tables.parameters.CHUNK\_CACHE\_SIZE.

## 1.8 Bugs

There is currently an interaction between Genomedata and PyTables that can result in the emission of Performance-Warnings when a Genomedata file is opened. These can be ignored. We would like to fix these at some point.

## 1.9 Support

To stay informed of **new releases**, subscribe to the moderated **genomedata-announce** mailing list (mail volume very low):

https://listserv.utoronto.ca/cgi-bin/wa?A0=genomedata-announce-l

For **discussion and questions** about the use of the Genomedata system, there is a **genomedata-users** mailing list:

https://listserv.utoronto.ca/cgi-bin/wa?A0=genomedata-l

For issues related to the use of Genomedata on **Mac OS X**, please use the above mailing list or contact Jay Hesselberth <jay dot hesselberth at ucdenver dot edu>.

If you want to **report a bug or request a feature**, please do so using our issue tracker:

https://github.com/hoffmangroup/genomedata/issues/

For other support with Genomedata, or to provide feedback, please write contact the authors directly. We are interested in all comments regarding the package and the ease of use of installation and documentation.

1.7. Technical matters 19

## **CHAPTER**

## TWO

## **INDICES AND TABLES**

- genindex
- modindex
- search

## **PYTHON MODULE INDEX**

g

genomedata, 16

24 Python Module Index

## **INDEX**

```
Symbols
                                                   S
__getitem__() (genomedata.Genome method), 16
                                                   sums (genomedata.Genome property), 18
__init__() (genomedata.Genome method), 16
                                                   sums_squares (genomedata.Genome property), 18
__iter__() (genomedata.Genome method), 16
Α
                                                   tracknames_continuous (genomedata.Genome prop-
add_track_continuous()
                              (genomedata.Genome
                                                            erty), 18
        method), 17
                                                   V
C
                                                   vars (genomedata.Genome property), 18
close() (genomedata.Genome method), 17
E
erase_data() (genomedata.Genome method), 17
F
format_version (genomedata.Genome property), 17
G
Genome (class in genomedata), 16
genomedata
    module, 16
index_continuous() (genomedata.Genome method),
isopen (genomedata.Genome property), 18
M
maxs (genomedata.Genome property), 18
means (genomedata.Genome property), 18
mins (genomedata.Genome property), 18
module
    genomedata, 16
Ν
num_datapoints (genomedata.Genome property), 18
num_tracks_continuous (genomedata.Genome prop-
        erty), 18
```