
gcmstools Documentation

Release 0.1.0

Ryan Nelson

January 19, 2015

CONTENTS

1	Getting started	1
2	Installation	3
2.1	Python	3
2.2	gcmstools	4
3	Basics of working with GCMS data files	5
3.1	Set up the processing environment	5
3.2	Note on Conventions	5
3.3	AIA Files	6
3.4	Read a Data File	6
3.5	Simple plotting	7
3.6	Working with multiple data sets	8
4	Reference and Fitting	11
4.1	Non-negative Least Squares	11
5	Data Storage	15
5.1	Stored Data Tables	16
6	Calibration and Integration	17
6.1	Calibration Object	17
6.2	Calibration Information File	17
6.3	Run Calibrations and Integrations	18
7	Batch Processing	21
7.1	Complete Processing Command	22
8	Appendix A: Running Code Samples	23
8.1	Using the Command Line	23
8.2	IPython	24
8.3	Working with Text Files	27

GETTING STARTED

gcmstools is a Python package that reads some GCMS file formats and does simple fitting. The source code for this project can be found on [GitHub](#). If you are reading this in PDF format, there is [online documentation](#) as well.

This user guide is broken into a few sections.

1. *Installation*: Information about getting a Python installation up and running and installing *gcmstools*.
2. *Basic Usage*: This section covers the usage of *gcmstools* to manipulate and plot GCMS data.
3. *Referenceing and Fitting*: Incorporate reference data into your GCMS data set, and use this to manually fit the GCMS data set.
4. *Data Storage*: Create a storage container for your processed GCMS datasets. This provides a convenient storage container for your data, and is necessary for the next step in the process.
5. *Calibration*: Make a calibration file and use this to extract generate concentration information from your data.
6. *Batch Processing*: Introduces a simple function for automating this entire process. You can skip to this section if all you want to do is automate some data extractions. It is not necessary to understand the basics of data manipulation/plotting.
7. *Appendix A*: Command line basics. The examples presented in this document require a basic working knowledge of a command-line terminal interface and running Python commands from an IPython interpreter. This section covers some of the basics.
8. Appendix B: Examples. Basic data extraction and plotting examples are presented.

INSTALLATION

Gcmstools requires Python and a number of third-party packages. Below is a complete list of packages and minimum versions:

- Python ≥ 3.4 (2.x versions not supported any longer)
- Pip $\geq 6.0.6$ (might be part of new Python releases)
- Setuptools $\geq 11.3.1$ (might be part of newer Python releases)
- Numpy $\geq 1.9.1$
- Matplotlib $\geq 1.4.2$
- Pandas $\geq 0.15.2$
- IPython $\geq 2.3.1$
- netCDF4 $\geq 1.0.4$
- PyTables $\geq 3.1.1$
- Scipy $\geq 0.14.0$
- Sphinx $\geq 1.2.2$ (Optional for documentation.)
 - numfig is a Sphinx extension that is needed to autonumber figures references in the documentation.

IPython also provides a very useful advanced interactive Python interpreter, and examples in this documentation assume that you are using this environment. See the *IPython* section of [Appendix A](#) for more details.

2.1 Python

The most convenient installation method for Python other third-party packages is the all-in-one [Anaconda Python distribution](#). It combines a large number of Python packages for scientific data analysis and a program (`conda`) for managing package updates (in addition to many other advanced features). The Anaconda developers (Continuum Analytics) provide a lot of useful documentation for [installing Anaconda](#) and [using conda](#). There are other ways to install Python and its packages, but for this documentation, it will be assumed that you are using Anaconda.

Note: On Mac/Linux, Python is already part of the operating system. Do not try to install these third-party packages into the builtin Python distribution unless you really know what you are doing. You might overwrite an important file, which can cause problems for your system. Confusion between the system and Anaconda Python installation is a common source of problems for beginners, so make sure that your Anaconda Python is “activated” before running the commands in this document. (See the Anaconda documentation for more information on the activation process.)

Note: On Windows, Anaconda may not install netCDF4. In this case, you can get a prebuilt installer from [Christoph](#)

[Gohlke](#); be sure to get the Python 2.7 (“cp27”) 64-bit (“amd64”) build for the most recent version.

Learning the usage of all of these Python packages is far beyond the scope of this document. However, excellent documentation for most of the packages as well as full tutorials are [easily discovered](#).

2.2 gcmstools

To install *gcmstools* from [the main repository](#), there are two options: 1) install using `git` (recommended) or 2) download the source file and install the package.

Option 1 (recommended)

First, install the [version-control software Git](#). *gcmstools* can now be downloaded and installed with one command.

```
home>$ pip install git+https://github.com/rnelsonchem/gcmstools.git
```

The advantage here is that the same command will update your *gcmstools* installation with any any changes that have been made to the main repository.

Option 2

Download a zip file of the current state of the repository. (Look for the button shown below ([Figure 2.1](#)) at [the main repository](#).) Unzip this package wherever you’d like.

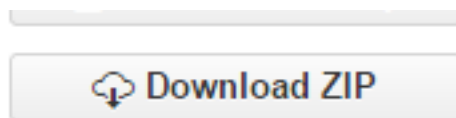


Figure 2.1: The zipfile download button.

From the command line, navigate the newly extracted folder and use `pip` to install the package. In this case, *path-to-gcmstools-folder* is the location of the newly unzipped *gcmstools* folder. Be sure to put a dot (.) at the end of that `pip` command.

```
home>$ cd path-to-gcmstools-folder
gcmstools>$ pip install .
```

Uninstall

Uninstallation of *gcmstools* is trivial. It may be a good idea to run this command before installing updates as well to ensure that the most recent version of *gmcstools* is being installed.

```
home>$ pip uninstall gcmstools
```


BASICS OF WORKING WITH GCMS DATA FILES

3.1 Set up the processing environment

In these examples, we will run *gcmstools* from a *terminal IPython* session in a folder “gcms”, which is located in your home directory.

```
home>$ cd gcms

gcms>$ ipython
Python 3.4.1 (default, Oct 10 2014, 15:29:52)
Type "copyright", "credits" or "license" for more information.

IPython 2.3.1 -- An enhanced Interactive Python.
?                -> Introduction and overview of IPython's features.
%quickref        -> Quick reference.
help             -> Python's own help system.
object?         -> Details about 'object', use 'object??' for extra details.
```

In :

Some example files are provided with the *gcmstools* installation. These files can be moved into the current directory using the `get_sample_data` function.

```
In : from gcmstools.general import get_sample_data

In : get_sample_data()
```

This invocation copies all of the example files to the current directory. Individual file names can be passed to this function, if you only want a few data files.

```
In : get_sample_data('datasample1.CDF')
```

3.2 Note on Conventions

There are potentially many types of GCMS files; however, all file importing objects discussed in this section should have identical properties. This is important for later sections of the documentation, because fitting routines, etc., usually do not require a specific file type importer. All of the import objects are constructed with a single string input, which is the name of the file to process. This file name string can also contain path information if the file is not located in the current directory.

3.3 AIA Files

AIA, ANDI, or CDF are all related types of standard GCMS files that are all [derived from](#) the Network Common Data Format ([netCDF](#)). They may have the file extension “AIA” or “CDF”. This file type may not be the default for your instrument, so consult the documentation for your GCMS software to determine how to export your data in these formats.

To import this type of data, use the `AiaFile` object, which is located in the `gcmstools.filetypes` module.

```
In : from gcmstools.filetype import AiaFile
```

3.4 Read a Data File

First of all, you will need to import a file reader from `gcmstools.filetypes` module. In this example, we’ll use the AIA file reader, `AiaFile`; however, the results should be identical with other readers. To read a file, you can create a new instance of this object with a filename given as a string.

```
In : from gcmstools.filetype import AiaFile
```

```
In : data = AiaFile('datasample1.CDF')
Building: datasample1.CDF
```

The variable `data` now contains our processed GCMS data set. You can see its contents using [tab completion](#) in IPython.

```
In: data.<tab>
data.filename  data.intensity  data.tic      data.index   data.masses
data.filetype  data.int_extract data.index    data.times
```

Most of these attributes are data that describe our dataset. You can inspect these attributes very easily in IPython by just typing the name at the prompt.

```
In : data.times
Out:
array([0.08786667, ..., 49.8351])
```

```
In : data.tic
Out:
array([158521., ..., 0.])
```

```
In : data.filetype
Out: 'AiaFile'
```

This is a short description of these initial attributes:

- *filename*: This is the name of the file that you imported.
- *times*: A Numpy array of the times that each MS was collected.
- *tic*: A Numpy array of the total ion chromatogram intensities.
- *masses*: A Numpy array the masses that cover the data collected by the MS.
- *intensity*: This is the 2D Numpy array of raw MS intensity data. The columns correspond to the masses in the `masses` array and the rows correspond to the times in the `times` array.
- *filetype*: This is the type of file importer that was used.

The `index` and `int_extract` methods are used for finding the indices from an array and extracting integrals, respectively. Their usage is described later.

3.5 Simple plotting

Now that we've opened a GCMS data set. We can easily visualize these data using the plotting package Matplotlib. As an example, let's try plotting the total ion chromatogram. In this case, `data.times` will be our "x-axis" data, and `data.tic` will be our "y-axis" data.

```
In : import matplotlib.pyplot as plt

In : plt.plot(data.times, data.tic)
Out :
[<matplotlib.lines.Line2D at 0x7f34>]

In: plt.show()
```

This produces a pop-up window with an interactive plot, [Figure 3.1](#). (This should happen fairly quickly. However, sometimes the plot window appears behind the other windows, which makes it seem like things are stuck. Be sure to scroll through your windows to find it.) The buttons at the top of the window give you some interactive control of the plot. See the [Matplotlib documentation](#) for more information.

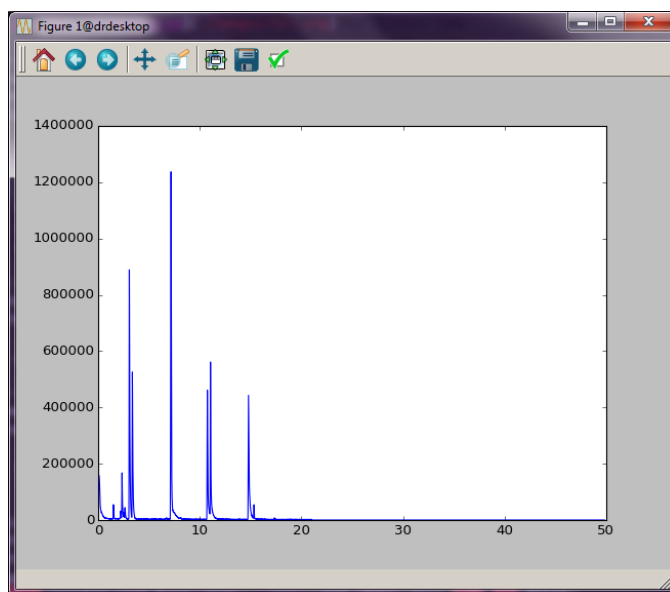


Figure 3.1: Total ion chromatogram.

One drawback here is that you have to type these commands every time you want to see this plot. There is another alternative, though. You can also put all of these commands into a text file and run it with Python directly. Copy the following code into a plain text file called "tic_plot.py". (See [Working with Text Files](#) for more information on making Python program files.)

```
import matplotlib.pyplot as plt
from gcmstools.filetypes import AiaFile

data = AiaFile('datasample1.CDF')
```

```
plt.plot(data.times, data.tic)
plt.show()
```

It is common practice to do all imports at the top of a Python program. That way it is clear exactly what code is being brought into play. Run this new file using the `python` command from the terminal. Again, the plot window will appear, but you will not be able to work in the terminal until you close this window.

```
gcms>$ python tic_plot.py
```

Alternatively, you can run this program directly from IPython. This has the advantage that once the window is closed, you are dropped back into an IPython session that “remembers” all of the variables and imports that you created in your program file. See [Appendix A](#) for more information here.

```
In : %run tic_plot.py
```

3.6 Working with multiple data sets

In the example above, we opened one dataset into a variable called `data`. If you want to manipulate more than one data set, the procedure is the same, except that you will need to use different variable names for your other data sets. (Again, using `AiaFile` importer as an example, but this is not required.)

```
In : data2 = AiaFile('datasample2.CDF')
```

These two data sets can be plot together on the same figure by doing the following:

```
In : plt.plot(data.times, data.tic)
Out:
[<matplotlib.lines.Line2D at 0x7f34>]

In: plt.plot(data2.times, data2.tic)
Out:
[<matplotlib.lines.Line2D at 0x02e3>]

In: plt.show()
```

The window shown in [Figure 3.2](#) should now appear. (There is a blue and green line here that are a little hard to see in this picture. Zoom in on the plot to see the differences.)

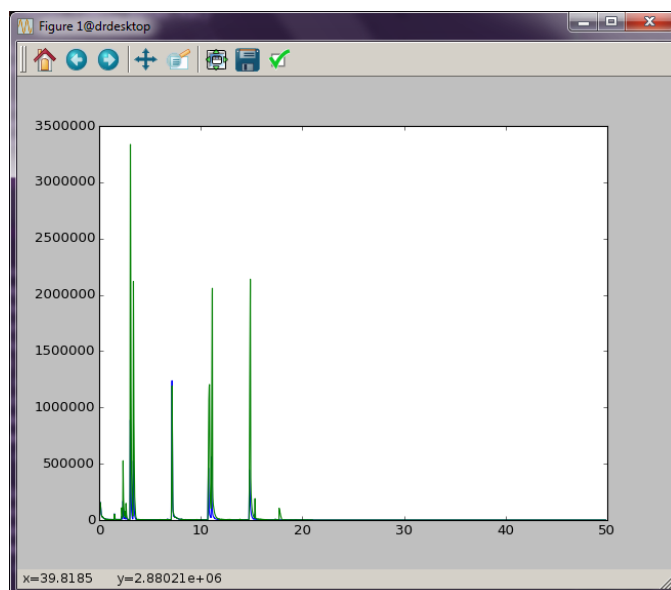


Figure 3.2: Two tic plotted together

REFERENCE AND FITTING

4.1 Non-negative Least Squares

4.1.1 Collecting References for NNLS

A series of reference spectra are required if you want to do non-negative least squares (NNLS) fitting. There are two example reference files in this repository: “ref_spec.txt” and “ref_spec2.MSL”. The MSL file is a type of MS library that can be exported by programs like the AMDIS or the NIST Mass Spectral Database. The format of both of these files is very important. MSL files are typically autogenerated by external software and may not need manual modification. The .txt file was hand generated; more information on this file format is provided in the next paragraph. A common feature of both formats, though, is that comment lines can be included by starting a line with #. This can be useful if you want to add some notes or to remove reference spectra without deleting them entirely.

Hand generated “.txt” reference files are made up of a series of reference compound information separated by blank lines. Information about the reference compounds are included using labels, and each compound must have two labels at the minimum : “NAME” and “NUM PEAK”. “NAME” will be the reference name (probably want this to be concise), and “NUM PEAKS” is followed by (at least) two space-separated columns of MS data. The first column are m/z values, and the second column are the associated intensity information. Intensities are normalized on import, so it is not necessary to do this by hand. Other labels can also be included if you would like to incorporate extra metadata about the reference compound. Each reference compound *must* be separated by a blank line. Below is a small sample of one of these files:

```
NAME:octane
FROM:www.massbank.jp
ID_NUM: JP004695
NUM PEAKS:
  42 14.07 141
  43 99.99 999
  44 2.54 25
  45 4.03 40
  53 1.58 16
  55 19.83 198
.
.
.
```

The online MS repository [massBank](http://www.massbank.jp) is a useful place to find these mass and intensity values. The data from that site is already formatted correctly for this file type.

4.1.2 Loading Reference Spectra

There are two objects located in `gcmstools.reference` for loading reference data, `TxtReference` and `MslReference`, which are used for ".txt" and ".MSL" reference files, respectively. In this example, we'll use `TxtReference`, but the other object behaves in the same manner.

First, we'll need some data, and we'll use an `AiaFile` object for this example.

```
In: from gcmstools.filetypes import AiaFile

In: data = gcms.AIAFile('datasample1.CDF')
Building: datasample1.CDF
```

Next, import the reference object and create an active instance, which requires that the name of the reference file is passed into the constructor. In this example, we have a reference file called "ref_specs.txt".

```
In : from gcmstools.reference import TxtReference

In : ref = TxtReference('ref_specs.txt')

In : ref.<tab>
ref.bkg          ref.ref_build      ref.ref_file      ref.ref_meta
ref.bkg_time     ref.ref_cpds      ref.ref_mass_inten ref.ref_type
```

As you can see, several attributes have been created for this new instance. In most cases, you will not need to work with any of these yourself. To add this reference information to a GCMS data set, call the reference instance with a GCMS file object or list of objects to process several files simultaneously.

```
In : ref(data)
Referencing: datasample1.CDF

In : ref([data, otherdata1, otherdata2]) # If these other data sets exist.
Referencing: datasample1.CDF
Referencing: otherdata1.CDF
Referencing: otherdata2.CDF

In : data.<tab>
data.filename    data.index      data.masses      data.ref_meta    data.times
data.filetype    data.int_extract data.ref_array    data.ref_type    data.tic
data.intensity   data.ref_cpds
```

Several new attributes have been added to our GCMS data object. Here is a short description of each.

- *ref_cpds*: A list of reference compound names.
- *ref_array*: A 2D Numpy array of the reference mass spectra. Shape(# of ref compounds, # of masses)
- *ref_meta*: Associated meta data for each reference compound.
- *ref_type*: The name of the reference object type that was used to generate this information. (In this example, this would be "TxtReference".)

4.1.3 Fitting the data

A `Nnls` fitting object is provided in `gcmstools.fitting` for performing the non-negative least squares fit. To apply this fitting to a data set, simply call the fitting instance with a data object or list of objects.

```
In : from gcmstools.fitting import Nnls
```



```

In : fit = Nnls()

In : fit(data)
Fitting: datasample1.CDF

In : fit([data, otherdata1, otherdata2]) # If these other data sets exist.
Fitting: datasample1.CDF
Fitting: otherdata1.CDF
Fitting: otherdata2.CDF

In : data.<tab>
data.filename      data.tic           data.int_sim       data.ref_cpds
data.filetype      data.index         data.intensity     data.ref_meta
data.fits          data.int_cum       data.masses        data.ref_type
data.fittype       data.int_extract   data.ref_array     data.times

```

Again, several new attributes describing the fit have been added to our data set.

- *fittype*: A string that names the fitting object used to generate this data. (In this case, it would be “Nnls”.)
- *fits*: These are the raw fitting numbers from the NNLS routine. They do not correspond to proper integrations, so they should be used with caution.
- *int_sim*: This is a 2D numpy array of simulated GCMS curves that were generated from the fit. Shape(# of time points, # of reference compounds)
- *int_cum*: This is a cumulative summation of *int_sim*, so it has the same shape as that array. The difference between any two points in this array can be used to determine the integral over that region.

4.1.4 Plotting the Fit

You can do a quick check of how the data looks using Matplotlib. More advanced examples are presented in Appendix B. The output of the commands below is shown in [Figure 4.1](#).

```

In : import matplotlib.pyplot as plt

In : plt.plot(data.times, data.tic, 'k-', lw=1.5)
Out: [<matplotlib.lines.Line2D at 0x7f9b2905df60>]

In : plt.plot(data.times, data.int_sim)
Out:
[<matplotlib.lines.Line2D at 0x7f9b2f0df160>,
 <matplotlib.lines.Line2D at 0x7f9b29063ac8>,
 <matplotlib.lines.Line2D at 0x7f9b29063d30>,
 <matplotlib.lines.Line2D at 0x7f9b29063f98>,
 <matplotlib.lines.Line2D at 0x7f9b28fef240>,
 <matplotlib.lines.Line2D at 0x7f9b28fef4a8>,
 <matplotlib.lines.Line2D at 0x7f9b28fef710>,
 <matplotlib.lines.Line2D at 0x7f9b28faf720>]

In : plt.legend(["TIC",] + data.ref_cpds) # This isn't necessary
Out: <matplotlib.legend.Legend at 0x7f9b25a35438>

In : plt.show()

```

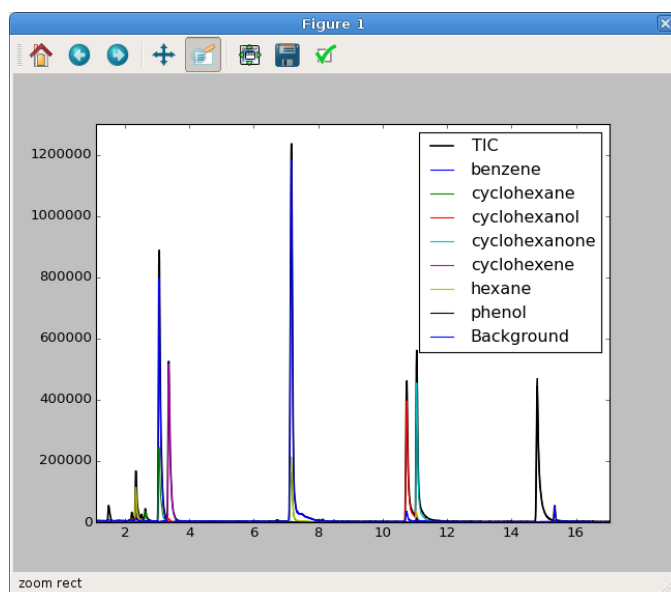


Figure 4.1: An interactive check of our fit. This has been zoomed in a little to highlight the fit and data.

DATA STORAGE

Processed data files can be stored on-disk for later reuse using the `HDFStore` object located in the `gcmstools.datastore` module. Not only does this create a convenient storage solution for processed data sets, it is also necessary when running calibrations on a group of related data sets. This file is stored on-disk in as a [HDF file](#), which is an open-source high performance data storage container optimized for numerical data. Creation and access of this file is controlled using a combination of two Python libraries: [PyTables](#) and [Pandas](#). PyTables provides a high-level interface to the HDF file types, and Pandas is a very powerful package for working with tabular data. Both of these project have extensive documentation on their use.

A `gcmstools.HDFStore` object can be created without any arguments, and in this case, it automatically creates a file called “data.h5”. If this file already exists, it will open that file for appending or modification. If you want to have more than one data file, you can also pass in a custom data file name as the first argument to the object creation.

```
In : from gcmstools.datastore import HDFStore

In : h5 = HDFStore()

In : # Or: h5 = HDFStore('data.h5') or whatever file name you'd like
```

Added files to this storage container can be done with the `append_files` method, which can take either a single data object or a list of objects, if you have many objects to add at one time.

```
In : h5.append_files(data)
HDF Appending: datasample1.CDF

In : h5.append_files([otherdata1, otherdata2])
HDF Appending: otherdata1.CDF
HDF Appending: otherdata2.CDF
```

You can see a list of the files that are stored in this file by viewing the `files` attribute, which is a Pandas DataFrame.

```
In : h5.files
      name      filename
0  datasample1  datasample1.CDF
1  otherdata1   otherdata1.CDF
2  otherdata2   otherdata2.CDF
```

There are two name columns in this table: “name” and “filename”. The latter is the full file name as given when the `GcmsFile` object was created. This may also contain path information to that file. Keep in mind that the path information may not be correct if you’ve moved the location of this storage file. In order to efficiently store the data on disk, the “filename” is internally simplified the “name”. This simplification removes the path and suffix from “filename”, and it replaces all “.”, “-”, and spaces with “_”. In addition, if “filename” starts with a number, the prefix “num” is added to “name”. If you have two or more “filenames” that simplify to the same “name”, you will run into problems. However, if your file naming system does not produce unique filenames for different data sets, you will most certainly have more problems than just using these programs.

You can extract data from the storage file using the `extract_gcms_data` method. This function takes one argument which is the name of the dataset that you want to extract. This name can be either the simplified name or the filename (with or without the path). The extracted data is the same file object type as you stored originally.

```
In : extracted = h5.extract_gcms_data('datasample1')
```

```
In : extracted.filetype
```

```
Out: "AiaFile"
```

Several other tables are added to this file when you run data calibration, which is outlined in the next section.

5.1 Stored Data Tables

The data file may contain a number of data tables with information about the files, calibration, etc. These data tables are stored as Pandas DataFrames, and a list of currently available tables can be obtained from the `pdh5` attribute of the `HDFStore` object.

```
In : h5.pdh5
```

```
Out:
```

```
<class 'pandas.io.pytables.HDFStore'>
```

```
File path: data.h5
```

```
/calibration          frame          (shape->[6,8])
```

```
/calinput             frame          (shape->[30,9])
```

```
/datacal              frame          (shape->[49,6])
```

```
/files                frame          (shape->[1,2])
```

To view these tables, just append the table name after `pdh5`.

```
In : h5.pdh5.calibration
```

```
Out:
```

	Start	Stop	Standard	slope	intercept	r	\
Compound							
benzene	2.9	3.5	NaN	38629.931565	-367129.586850	0.998767	
phenol	14.6	15.1	NaN	30248.192619	65329.897933	0.999136	
...							

	p	stderr
Compound		
benzene	0.000052	1108.344872
phenol	0.000030	726.257380
...		

CALIBRATION AND INTEGRATION

6.1 Calibration Object

The calibration class `Calibrate` is defined in the `gcmstools.calibration` module. This object must have access to an HDF storage file that contains *all* of your data to be processed. The default object creation tries to open a file “data.h5”; however, an alternate name can also be passed in on construction to open a different HDF file.

```
In : from gcmstools.calibration import Calibrate

In : cal = Calibrate() # Opens 'data.h5' by default

In : # Or: cal = Calibrate('other.h5') if you have a different name
```

6.2 Calibration Information File

In order to integrate your GCMS data, you must first create a csv file containing all of the relevant calibration information. Again, the structure of this file is very important, so an example, “calibration.csv”, is contained with the sample data.

```
In : from gcmstools.general import get_sample_data

In : get_sample_data("calibration.csv")
```

The first row in this csv file is critical, and it must look like this:

```
Compound,File,Concentration,Start,Stop,Standard,Standard Conc
```

Each row after this describes a set of calibration information that you’d like to use for The columns of this file are as follows:

- *Compound*: The name of the compound that you are calibrating. This *must* correspond to one of the compound names (case-sensitive) used when referencing and fitting the GCMS file. See the [relevant page](#) for more information.
- *File*: This is the name of the data set that you’ll use for this compound at a particular, known concentration. Again, this filename can be the full filename (with or without the path) or the simplified name. See the [files attribute](#) section of the HDFStore docs for more info.
- *Concentration*: This is the concentration of *Compound* in *File*. This should only be a number. Do not include units. All of the concentrations should be in the same units, and keep in mind that all calibration and integration data will then be in that same unit of measurement.

- *Start* and *Stop*: These are the integration range start and stop times in minutes. Again use only numbers, no units.
- *Standard* and *Standard Conc*: If there is an internal standard used in this file, you should provide the name and concentration in these columns. Again, the standard name should have been defined when referencing your GCMS data set, and do not include units with the concentration. Make sure your concentration units are the same as the reference compound to avoid confusion.

You can add extra columns to this table without penalty, in case you need to add additional information to this table. You can also comment out lines by starting a line with a # character. This is useful if you want to ignore a bad data point without completely removing the line from the calibration file.

6.3 Run Calibrations and Integrations

6.3.1 Calibration Curves

The calibration curves can be generated using the `curvegen` method of the `Calibrate` object. This function must be called with the name of your calibration file. In this example, that filename is “calibration.csv”.

```
In : cal.curvegen('calibration.csv')
Calibrating: benzene
Calibrating: phenol
...
```

```
In :
```

This process creates two new tables as attributes to your calibration object, `calinput` and `calibration`. The former table is simply your input csv information with columns appended for the concentrations (“conc”) and integrals (“integral”) used for generating the calibration curve. If no internal standard is defined, then “conc” will be the same as the compound concentration you used in the input file. If an internal standard was defined, then “conc” and “integral” will be these values divided by the corresponding internal standard values. These tables also stored in the HDF file as well, if you want to check them at a later date.

The calibration table contains all of the newly created calibration curve information, such as slope, intercept, r value, etc.

```
In : cal.calibration
Out:
```

	Start	Stop	Standard	slope	intercept	r	\
Compound							
benzene	2.9	3.5	NaN	38629.931565	-367129.586850	0.998767	
phenol	14.6	15.1	NaN	30248.192619	65329.897933	0.999136	
...							

	p	stderr
Compound		
benzene	0.000052	1108.344872
phenol	0.000030	726.257380
...		

6.3.2 Plotting Calibrations

By default, no plots are generated for these calibrations. There are a couple of ways to get some plots of the calibration data.

1. `cal.curvegen('calibration.csv', calfolder='cal', picts=True)` : This method will auto generate pictures for all of the calibration compounds and place them in a folder defined by the keyword argument `calfolder`. This argument is optional, if you don't mind the default folder name of "cal". Be careful! This will delete this folder before generating new plots, so if this folder exists, make sure it is clear of important data.
2. `cal.curveplot('benzene')` : This method will generate a plot of the benzene calibration information and save it to the current folder. There are several keyword arguments to this function:
 - `folder='.'` : This sets the folder where the picture will be saved. By default it is the current directory.
 - `show=False` : Change this value to `True` if you want an interactive plot window to be displayed. Default is `False`.
 - `save=True` : Save the calibration plot to the folder.

If both `save` and `show` are set to `False`, nothing will happen.

Of course, this function must be done after a call to `curvegen`, but it does provide a method to look at calibration data from an previously processed HDF file without rerunning the calibration.

6.3.3 Integrating Data

Once the calibration curves have been generated, you can integrate all of the remaining data in the HDF file using the `datagen` method of the `Calibrate` object.

```
In : cal.datagen()
Processing: datasample1.CDF
Processing: otherdata1.CDF
Processing: otherdata2.CDF
...
```

After processing, another data table attributed (`datacal`) is created and saved to the HDF file.

```
In : cal.datacal
Out:
```

	benzene	phenol	...
name			
datasample1	4239.070627	58.336917	...
otherdata1	5475.778519	20.401981	...
otherdata2	4355.094930	19.171877	...
...			

6.3.4 Plotting Integrals

By default, no plots are generated for the integrals. If you'd like to see plots of the integrals, there are a couple of methods.

1. `cal.datagen(datafolder='data', picts=True)` : This method will auto generate pictures for all of the calibration compounds and place them in a folder defined by the keyword argument `datafolder`. This argument is optional, if you don't mind the default folder name of "data". Be careful! This will delete this folder before generating new plots, so if this folder exists, make sure it is clear of important data.
2. `cal.dataplot('benzene', 'datasample1')` : This method will generate a plot of the benzene integral for 'datasample1' and save it to the current folder. There are several keyword arguments to this function:
 - `folder='.'` : This sets the folder where the picture will be saved. By default it is the current directory.

- `show=False` : Change this value to `True` if you want an interactive plot window to be displayed. Default is `False`.
- `save=True` : Save the calibration plot to the folder.

If both `save` and `show` are set to `False`, nothing will happen.

Of course, this function must be done after a call to `datagen`, but it does provide a method to look at calibration data from an previously processed HDF file without rerunning the calibration and data integration functions.

BATCH PROCESSING

All of the previous steps can be automated using the `proc_data` function located in the `gcmstools.general` module.

```
In : from gcmstools.general import proc_data
```

This function only has two required arguments: 1) the path to the folder that contains *all* of the GCMS files and 2) the name of the data file that you'd like to generate. In this example, our data is in the folder “data” and our processed data file is called “data.h5”.

```
In : proc_data('data/', 'data.h5')
... # Lots of stuff will get printed at this point.
```

The `data_proc` function accepts numerous keyword arguments to allow some process control.

- `filetype='aia'` : This flag can be used to control the type of GCMS objects used for processing the data in your folder. The default value `aia` uses the `AiaFile` object. See *Basics of working with GCMS data files* for detailed information.
- `reffile=None` : Pass the name of a reference file for referencing your GCMS data. The default is `None`, so no referencing will be done. Otherwise, the reference object will be determined by the file extension. For example, `reffile='ref_specs.txt'` will create a `TxtReference` object using the file “ref_specs.txt” (which must exist of course). See *Reference and Fitting* for more information.
- `fittype=None` : Set the fitting type to use for fitting the GCMS data. See *Reference and Fitting* for detailed information. The valid choices are:
 - `'nnls'` for non-negative least squares fitting.
- `calfile=None` : Pass in the name of a reference csv file to generate calibration curves and integrate the data. For example, `calfile='calibration.csv'` will calibrate your data using the information in the csv file “calibration.csv”. See *Calibration and Integration* for detailed information, especially on the expected structure of the csv file.
- `picts=False` : Set this argument to true if you want to generate pictures of your calibration curves and data fits. If you have a lot of data files, this can be very slow. It is possible to view these plots after processing. See *Calibration and Integration* for detailed information.
- `chunk_size=4` : This sets the number of files that will be processed at any given time. This keeps the total number of opened GCMS files to a minimum, which is important if you have a large number of files to process. You probably don't need to change this.
- `multiproc=False` : Setting this argument to `True` will use IPython's parallel machinery to run the processing using multiple cores. Before using this command, you must start an IPython node cluster from the command line. Open a new terminal, and type the following command:

```
home>$ ipcluster start -n 2
```

This will start a cluster of two (`-n 2`) nodes. You can have up to one node per core on your processor. Setting this number greater than the number of cores will result in a degradation of performance. To stop the node cluster, type `Ctrl-C` from the terminal where you started the cluster. Or else, you can stop it from another terminal using the following command:

```
home>$ ipcluster stop
```

This only works for processing the files, not for generating plots. So plotting your calibration data will still be very slow for a lot of data files.

See [IPython's parallel documentation](#) for more information.

- This function can also accept any keyword arguments for any file type, reference, fitting, and calibration objects. See their documentation for more information.

7.1 Complete Processing Command

A complete version of this processing command might look like the following.

```
In : proc_data('data/', 'data.h5', filetype='aia', reffile='ref_spects.txt',  
             fittype='nnls', calfile='calibration.csv', multiproc=True)  
### Lot's of stuff gets printed  
...
```

The HDF file “data.h5” contains all of your data. See the [Data Storage, Calibration and Integration](#) documentation and Appendix B for more information on how to view and plot these data.

APPENDIX A: RUNNING CODE SAMPLES

8.1 Using the Command Line

Running the code samples in this documentation requires a rudimentary knowledge of the command-line terminal (command prompt on Windows). The terminal can seem very “texty” and confusing at first; however, with a little practice it gets to be fairly intuitive and *efficient*. There are many tutorials online, for example, [The Command Line Crash Course](#). However, A few basic command line concepts are covered here, for reference.

When you start a terminal, you will be presented with a window to type commands. In this documentation, the terminal command prompt will be denoted as `home>$`, where “home” indicates the current folder where the commands will be executed and “>\$” is just a separator. These things do not need to be typed when entering commands. A similar format is common in a lot of online documentation. Some commands generate output. The output will occur after the command prompt, but will not be preceded by a command prompt symbol.

The most important thing you’ll want to be able to do is move to different directories (i.e. folders). To do this there are a couple of useful terminal commands: “change directory” (`cd`) and “present working directory” (`pwd`). When you open the terminal, you will usually start in your “home” directory. This will probably be “/Users/username” on Mac, “/home/username” on Linux, or “C:\Users\username” on Windows. To move to a different directory, use the `cd` command; to find out the location of the current folder, use `pwd`. Here’s an example.

```
home$> pwd
/home/username/

home$> cd folder1

folder1$> pwd
/home/username/folder1
```

The second command here moved active directory to the folder *folder1*. There are also a few special directory shortcuts:

- `~` This refers to the home directory.
- `..` (Double dot) This refers to the parent directory of the current directory.
- `.` (Single dot) This refers to the current directory.
- `\` or `/` Separators to combine directory names. The first works on Linux/Mac (and in IPython, see below), the second is required on Windows.

Here’s these shortcuts in action.

```
home$> cd folder1/folder2

folder2$> pwd
/home/username/folder1/folder2
```

```
folder2>$ cd .

folder2>$ pwd
/home/username/folder1/folder2

folder2>$ cd ../../..

home>$ pwd
/home/username

home>$ cd folder1/folder2

folder2>$ cd ~

home>$ pwd
/home/username

home>$ cd folder1/folder2

folder2>$ cd ~/folder3

folder3>$ pwd
/home/username/folder3
```

The other important command is `ls`, which lists the contents of the current directory. (In Windows, the equivalent command is `dir`.)

```
folder3>$ ls
file1 file2 folder4
```

8.2 IPython

8.2.1 Start IPython

One of Python's strengths as a data analysis language is its interactive interpreter. This mode is accessed from a terminal by typing `python`.

```
home>$ python
Python 3.4.1 (default, Oct 10 2014, 15:29:52)
[GCC 4.7.3] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

Any command that is typed at the `>>>` prompt is treated as Python code, and executed appropriately. That means you can interactively write and explore code in this manner.

```
>>> 2 + 2
4
>>> print('Hello World')
Hello World
```

The default Python interpreter is very limited, which is why IPython was developed. IPython is an advanced Python interpreter, which has several advanced features like autocompletion and introspection, just to name two. Over the years, this project has grown substantially, and in addition to a terminal based interpreter, there is now a GUI version and a very cool web-based Notebook as well. To learn more about the other features, consult the [IPython documentation](#).

IPython is started from the terminal using the `ipython` command:

```
home>$ ipython
Python 3.4.1 (default, Oct 10 2014, 15:29:52)
Type "copyright", "credits" or "license" for more information.

IPython 2.3.1 -- An enhanced Interactive Python.
?                -> Introduction and overview of IPython's features.
%quickref        -> Quick reference.
help             -> Python's own help system.
object?         -> Details about 'object', use 'object??' for extra details.

In [1]: 2 + 2
Out[1]: 4

In [2]: print('Hello World')
Hello World
```

The `In [#] :` prompt now takes the place of `>>>` in the regular Python interpreter. In addition, certain types of output are preceded by an `Out[#] :` prompt. The numbers in brackets help you to determine the order that commands are processed. For this documentation, though, the numbers will be stripped for clarity, e.g. `In :` and `Out :`. If you see these prompts, you should know that the commands are being run in an IPython session.

8.2.2 Autocompletion and Introspection

The take home message of this section is *use the Tab key a lot!* It will make you much more productive.

Two very nice aspect of the IPython interpreter are autocompletion and object introspection. Both of these will make use of the Tab key on your keyboard; in code snippets, this key will be denoted as `<tab>`, which means you should press the Tab key rather than typing it out. To see these two operations in action, we can first create a new string object.

```
In : my_string = 'Hello World'

In : print(my_string)
Hello World
```

To determine the methods available to a string object, we can use IPython's object introspection.

```
In : my_string.<tab>
my_string.capitalize    my_string.isidentifier  my_string.rindex
my_string.casefold      my_string.islower       my_string.rjust
my_string.center        my_string.isnumeric     my_string.rpartition
my_string.count         my_string.isprintable   my_string.rsplit
my_string.encode        my_string.isspace       my_string.rstrip
my_string.endswith      my_string.istitle       my_string.split
my_string.expandtabs    my_string.isupper       my_string.splitlines
my_string.find          my_string.join          my_string.startswith
my_string.format        my_string.ljust         my_string.strip
my_string.format_map    my_string.lower         my_string.swapcase
my_string.index         my_string.lstrip        my_string.title
my_string.isalnum       my_string.maketrans     my_string.translate
my_string.isalpha       my_string.partition     my_string.upper
my_string.isdecimal     my_string.replace       my_string.zfill
my_string.isdigit       my_string.rfind
```

As you can see, there are many, many things that you can do with this string object. IPython can also use the Tab key to autocomplete long names for variables, path strings, etc. Here's an example:

```
In : my_string.is<tab>
my_string.isalnum      my_string.isidentifier  my_string.isspace
my_string.isalpha      my_string.islower      my_string.istitle
my_string.isdecimal    my_string.isnumeric    my_string.isupper
my_string.isdigit      my_string.isprintable
```

```
In : my_string.isi<tab>
```

Notice that when you type `tab` here IPython automatically expands this to `my_string.isidentifier`. This works for path strings as well.

Note: It should be pointed out that tab completion also works on the regular command line terminal interface as well.

8.2.3 Magic Commands

IPython has a number of special commands that make its interpreter behave much like a command-line terminal. These commands, called Magic Commands, are preceded by `%` or `%%`. The [magic command documentation](#) covers many of them, but a few that are useful to the examples in this document are discussed here.

The magics `%cd`, `%pwd`, and `%ls` serve the exact same purpose as in the terminal. Another very useful magic is `%run`. This command executes a Python program file from inside the IPython session, and in addition to executing the code, it also loads the data and variables into the current IPython session. This is best explained by example. Create a new folder called `folder1` in your home directory. Create the file `test.py` in `folder1` and paste the following code into that file. (See [Working with Text Files](#) for some information on text files and Python programs.)

```
var1 = 7
var2 = "Hello World"
var3 = var1*var2
```

Now let's start up IPython and run this new program.

```
home>$ ipython
Python 3.4.1 (default, Oct 10 2014, 15:29:52)
Type "copyright", "credits" or "license" for more information.

IPython 2.3.1 -- An enhanced Interactive Python.
?          -> Introduction and overview of IPython's features.
%quickref  -> Quick reference.
help       -> Python's own help system.
object?    -> Details about 'object', use 'object??' for extra details.
```

```
In : %pwd
/home/username
```

```
In : %cd folder1
/home/username/folder1
```

```
In : %ls
test.py
```

```
In : %run test.py
```

```
In :
```

At this point, it seems like nothing has happened; however, the variable that we defined in our file “`test.py`” are now contained in our IPython session. Assuming that the following IPython code is the same session as above.

```
In : var1  
Out: 7
```

```
In : var3  
Out: Hello WorldHello WorldHello WorldHello WorldHello  
WorldHello WorldHello World
```

As you can see, this is a very powerful way to save your work for later or to run code that is fairly repetitive.

8.2.4 Notebook Interface

Todo.

8.3 Working with Text Files

There are many instances where you will need to work with plain text files, including when writing Python programs. Plain text files are *not* word processing documents (e.g. MS Word), so you will want to use a dedicated text editor. Another source of problems for beginners is that leading white space in Python programs is important. For these reasons, a dedicated Python text editor can be very useful for beginners. Anaconda is bundled with [Spyder](#), which has a builtin text editor. The Anaconda FAQ has [information on running Spyder](#) on your system. Spyder is actually a full development environment, so it can be very intimidating for beginners. Don't worry! The far left panel is the text editor, and you can use that without knowing what any of the other panels are doing. Some internet searches will reveal other text editors if you'd prefer something smaller. (Do *not* use MS Notepad.)

The ".py" suffix for Python programs can be important. On Windows, however, file extensions are not shown by default, which makes them difficult to modify. In these cases, you may inadvertently create a file with the extension ".py.txt", which will not behave as you expect. Consult the internet for ways to show file extensions on a Windows machine.