

The dataChest For Dummies

Alexander Opremcak*

April 12, 2016

Contents

1	Introduction	2
2	Overview of functionality	2
2.1	Folder creation and navigation methods	2
2.1.1	ls()	2
2.1.2	mkdir()	2
2.1.3	pwd()	3
2.1.4	cd()	3
2.2	Dataset methods	3
2.2.1	createDataset()	3
2.2.2	getVariables()	5
2.2.3	addData()	5
2.2.4	getData()	5
2.2.5	getDatasetName()	6
2.2.6	openDataset()	6
2.3	Parameter methods	7
2.3.1	addParameter()	7
2.3.2	getParameter()	7
2.3.3	getParameterList()	7
3	Canonical datasets	7
3.1	1D data	8
3.1.1	Arbitrary 1D data	8
3.1.2	1D scans	9
3.2	2D data	11
3.2.1	Arbitrary 2D data	11
3.2.2	2D scans	12

*electronic address: opremcak@wisc.edu

1 Introduction

The dataChest is a python library written for storing and retrieving datasets. It enforces good user habits while remaining flexible enough to deal with data of virtually any type or shape. In this document, I will briefly describe all of the methods this class has to offer. I will then formalize our definitions of 1D and 2D datasets and provide examples of how to enter the particular kinds we discuss.

2 Overview of functionality

The dataChest has a small number of public methods which can be broken up into three categories: i) Folder Creation and Navigation methods, ii) Dataset methods, and iii) Parameter methods. The Folder Creation and Navigation methods give users basic tools for navigating and creating directories. The Dataset methods allow users to create new datasets, open existing datasets, and fetch stored data. Parameter methods allow users to attach additional information, ‘parameters’, to a dataset. In what follows, we will describe how to use the methods from each of these categories through some basic examples.

2.1 Folder creation and navigation methods

2.1.1 `ls()`

`ls()` lists the contents of the working directory. Here is how its used:

```
from dataChest.dataChest import dataChest

d = dataChest() #instantiation
contents = d.ls()
files = contents[0]
directories = contents[1]

print contents
print files
print directories
```

From here on I will assume that the **dataChest** has been imported, the object **d** above has been instantiated, and each snippet of code provided has been executed.

2.1.2 `mkdir()`

`mkdir()` allows users to make new directories. Let us see how it works:

```
directories = d.ls()[1]
print "newDir" in directories
```

```
d.mkdir("newDir")
updatedDirectories = d.ls()[1]
print "newDir" in updatedDirectories
```

Note that `mkdir()` does not automatically change your working directory when a new directory is created.

2.1.3 pwd()

`pwd()` returns the path of the working directory. This method will be used in examples that follow.

2.1.4 cd()

`cd()` is used to change the working directory. Try running this:

```
print d.pwd() #prints working directory
d.cd("newDir") #string style input
print d.pwd()
d.mkdir("newSubDir")
d.cd("..") #shortcut for bumping us back 1 directory
d.cd(["newDir", "newSubDir"]) #list style input
print d.pwd()
d.cd("") #short cut for moving to root directory
print d.pwd()
```

2.2 Dataset methods

2.2.1 createDataset()

`createDataset()` is used to create datasets. It requires users to enter a name, independent variable information, and dependent variable information. Here is a quick example of how it is used:

```
d.createDataset("voltageTimeSeries",
                [("ind1", [1], "float64", "Seconds")],
                [("dep1", [1], "float64", "Volts")]
                )
```

Lets have a look at the arguments. The first argument “voltageTimeSeries” is the name of the dataset. This string must contain valid filename characters only. The second argument `[("ind1", [1], "float64", "Seconds")]` is a list of independent variable information. If you had multiple independent variables, this list would have multiple tuple entries, one tuple per independent variable. Each tuple contains variable information, in particular a name (`"ind1"`), shape (`[1]`), type (`"float64"`), and units (`"Seconds"`). The third argument is the analog of this but for dependent variables.

Let us discuss shapes and types in a little more detail. While scalar data is very commonly encountered in experimental situations, it pays to allow for

other shapes of data to be stored. For example, suppose you were interested in looking at the dependence of some transmission measurement as a function of gate bias. So you fix a gate bias, take a transmission measurement, then you fix a new gate bias, take a transmission measurement and so on. At the end of the day you have a dataset that is a collection of scalars (the gate bias) and arrays (the transmission waveform). This type of data could be save by creating the following dataset:

```
d.createDataset("someName",
                [("GateBias", [1], "float64", "Seconds"),
                 ("Frequency", [1000], "float64", "Hz")],
                [("Transmission", [1000], "complex128", "dBm")])
```

where we now have a shape of [1000]. This really generalizes our notion of ‘variables’ to something that can have a non scalar shape. Also note that we have included a new data type namely "complex128". Shapes are lists which specify the dimensionality of an N-dimensional array. For example 3x3 matrix would have shape [3,3] where as a 1000 point waveform (1D array) has shape [1000]. The data type then specifies the type of data that makes up these N-dimensional arrays. Acceptable data types are given in the following table:

Table 1: Data Types

Data Types	Description#1
"bool_"	Boolean
"int8"	8-bit Signed Integer (-2^7 to $2^7 - 1$)
"int16"	16-bit Signed Integer (-2^{15} to $2^{15} - 1$)
"int32"	32-bit Signed Integer (-2^{31} to $2^{31} - 1$)
"int64"	64-bit Signed Integer (-2^{63} to $2^{63} - 1$)
"uint8"	8-bit Unsigned integer (0 to $2^8 - 1$)
"uint16"	16-bit Unsigned integer (0 to $2^{16} - 1$)
"uint32"	32-bit Unsigned integer (0 to $2^{32} - 1$)
"uint64"	64-bit Unsigned integer (0 to $2^{64} - 1$)
"int_"	Default integer type (normally either "int64" or "int32")
"float16"	Half precision float: sign bit, 5 bits exponent, 10 bits mantissa
"float32"	Single precision float: sign bit, 8 bits exponent, 23 bits mantissa
"float64"	Double precision float: sign bit, 11 bits exponent, 52 bits mantissa
"float_"	Shorthand for "float64"
"complex64"	Complex number, represented by two 32-bit floats
"complex128"	Complex number, represented by two 64-bit floats
"complex_"	Shorthand for "complex128"
"utc_datetime"	UTC DateISO String
"string"	String

The data type `"utc_datetime"` deserves some special attention. UTC, which stands for Coordinated Universal Time, is the primary time standard by which the world regulates clocks and time. Although this may seem like overkill, it was chosen to eliminate any and all ambiguity that may be introduced when it comes to storing datetimes in your datasets. Here is how you generate the current `"utc_datetime"`:

```
from datetime import datetime
utc_datetime = datetime.utcnow().isoformat()
print utc_datetime
```

Note that this has microsecond timing resolution.

2.2.2 `getVariables()`

`getVariables()` returns a list of independent and dependent variables lists as they were given when the dataset was created. Try the following:

```
varsList = d.getVariables()
indepVarsList = varsList[0]
depVarsList = varsList[1]
```

```
print indepVarsList
print depVarsList
```

2.2.3 `addData()`

`addData()` is used for adding data to a newly¹ created dataset. Assuming we ran the example above, try:

```
d.createDataset("someName",
                [("ind1", [1], "float64", "Seconds")],
                [("dep1", [1], "complex128", "Volts")]
                )
d.addData([[1.0, 2.0+1j*2.0]]) #add 1 row
d.addData([[2.0, 4.0+1j*4.0], [3.0, 6.0+1j*6.0]]) #add 2 rows
```

2.2.4 `getData()`

`getData()` is used for retrieving data from an open dataset. Whether this dataset was newly created or recently re-opened for inspection is immaterial. Assuming we carried out the example above let's run the following code

```
data = d.getData() #fetch all data
print data[0] #row 0 indices start at 0!!
print data[1] #row 1
print data[2] #row 2
```

¹Adding data to an old dataset will be discussed with the `openDataset()` method.

There may exist situations in which a user wishes to fetch data incrementally as opposed to the all or nothing style that was demonstrated above. The `getData()` method permits such behavior in a manner similar to how arrays are sliced. Here are some examples:

```
print d.getData(startIndex = 0, stopIndex = 3) #whole thing
print d.getData(startIndex = 0, stopIndex = 0) #empty set
print d.getData(startIndex = 0, stopIndex = 1) #row 0
print d.getData(startIndex = 0, stopIndex = 2) #rows 0 and 1
print d.getData(startIndex = 1, stopIndex = 3) #rows 1, and 2
print d.getData(startIndex = 2, stopIndex = 3) #row 2
```

2.2.5 `getDatasetName()`

`getDatasetName()` returns the name of the current dataset. This will be used in the next example.

2.2.6 `openDataset()`

`openDataset()` is used for opening existing datasets. Assuming we carried out the example above let's run the following:

```
oldDatasetName = d.getDatasetName() #get current dataset name
print oldDatasetName
d.createDataset("someName",
               [("ind1", [1], "float64", "Seconds")],
               [("dep1", [1], "float64", "Volts")]
               ) #create new dataset
print d.getDatasetName() #differs from oldDatasetName
d.openDataset(oldDatasetName)
print d.getDatasetName()
print d.getData() #prints data from old dataset
```

The `openDataset()` does not allow users to use `addData()` by default. If a user would like to add data² to a dataset opened with `openDataset()`, they must input an optional parameter specifying that they would like modification privileges as follows:

```
oldDatasetName = d.getDatasetName() #get current dataset name
d.createDataset("someName",
               [("ind1", [1], "float64", "Seconds")],
               [("dep1", [1], "float64", "Volts")]
               ) #create new dataset
d.openDataset(oldDatasetName, modify = True) #write/read access
d.addData([[4.0, 8.0+1j*8.0]]) #add to an old dataset with modify = True
```

²Parameters can also be added to datasets opened with `openDataset()` by the same prescription.

2.3 Parameter methods

2.3.1 addParameter()

`addParameter()` is used for appending parameters to the current dataset. Here is how it works:

```
d.createDataset("someOtherName",
                [("ind1", [1], "float64", "Seconds")],
                [("dep1", [1], "float64", "Volts")]
                ) #create new dataset
d.addParameter("Who", "Mike Jones")
d.addParameter("Time", "Goon Time")
d.addParameter("Base Temp", -13.789) #et cetera
```

If a user attempts to add a parameter that already exists, this method will throw an exception. If a user would like to overwrite an existing parameter with a new value, even new types are allowed, then they must opt for overwrite privileges using the following syntax:

```
d.addParameter("DummyParam", 11)
d.addParameter("DummyParam", 14) #raises an exception
d.addParameter("DummyParam", 13, overwrite = True)
d.addParameter("DummyParam", "blah", overwrite = True)
```

2.3.2 getParameter()

`getParameter()` is used for retrieving the value of an existing parameter. Here is how it works:

```
print d.getParameter("Who")
print d.getParameter("Time")
print d.getParameter("Base Temp")
```

2.3.3 getParameterList()

`getParameterList()` is used for retrieving a list of all available parameters. Here is how it works:

```
print d.getParameterList()
```

3 Canonical datasets

While the overview of user-end functionality may have been helpful, it will prove useful to go over some examples of how to create and add to some of the most commonly encountered types of datasets. For each type of data, I will provide motivation for the definition and then an example of how to add it to the dataChest.

3.1 1D data

By 1D data, we mean a dataset with one independent variable and m dependent variables where $m \geq 1$. Loosely speaking, this captures what it means to be a 1D dataset entirely. In what follows, I will motivate the two most commonly encountered types of 1D data.

3.1.1 Arbitrary 1D data

The most general type of 1D data has an arbitrary spacing between consecutive data points along the independent variable axis, the x -axis. By fixing the independent variable at some value, call it v , and measuring all of the m dependent quantities $Q_1(v), \dots, Q_m(v)$ of interest, we obtain a single row of data

$$\text{row} = \left[v, Q_1(v), \dots, Q_m(v) \right] \quad (1)$$

Note the size brackets that we use here to define a row. For each value of v , we repeat this process and eventually our data looks like

$$\text{data} = \left[\begin{aligned} &\left[v_0, Q_1(v_0), \dots, Q_m(v_0) \right], \\ &\left[v_1, Q_1(v_1), \dots, Q_m(v_1) \right], \dots \\ &\left[v_N, Q_1(v_N), \dots, Q_m(v_N) \right] \end{aligned} \right] \quad (2)$$

where we have $N + 1$ rows of data with $m + 1$ columns of scalar³ data in each row. Let us call data of this format **Arbitrary Type 1 Data**. Here is an example of how this data entered:

```
import numpy as np
d.createDataset("MyFavoriteTimeSeries",
               [("indepName1", [1], "float64", "Seconds")],
               [("depName1", [1], "float64", "Volts")]
               )
d.addParameter("X Label", "Time")
d.addParameter("Y Label", "Digitizer Noise")
d.addParameter("Plot Title", "Random Number Generator")
net = []
for ii in range(0, 100):
    net.append([float(ii), np.random.rand()])
d.addData(net) #add 100 rows of data at once
d.getData() #single row
```

³String data is also allowed.

Note that the shape of each variable is [1] which is really the definition of **Arbitrary Type 1 Data** as far as the dataChest is concerned.

Alternatively we could group this data like so

$$\text{data} = \left[\begin{array}{c} [v_0, \dots, v_N], [Q_1(v_0), \dots, Q_1(v_N)], \\ \dots [Q_m(v_0), \dots, Q_m(v_N)] \end{array} \right] \quad (3)$$

Note that each column ('variable') has the same length and the implicit functional dependence of the Q_i 's as a function of index that we have assumed. Let us call data of this format **Arbitrary Type 2 Data**. Here is how it is entered:

```
res = 1e-4
timeAxis = np.arange(0.0, 1.0, res)
d.createDataset("DampedOscillations",
                [("time", [len(timeAxis)], "float64", "Seconds")],
                [("Oscillation", [len(timeAxis)], "float64", "Volts")])
d.addParameter("X Label", "Time")
d.addParameter("Y Label", "Voltage")
d.addParameter("Plot Title", "Damped Oscillations")
d.addData([ [timeAxis, np.sin(2 * np.pi * timeAxis)] ])
d.getData()
```

We use the term 'arbitrary' for these types of data because it is not necessary for the v_i 's to be related by any closed form expression as of function of array index. So there you have it. All 1D data fall under this category. So were done ... not quite!

3.1.2 1D scans

Suppose we rip a time series from an ADC with some fixed sampling rate, say 1GS/s. Then our independent axis is specified uniquely by 3 numbers, namely the initial time t_0 , the final time t_f , and the number of samples N . The time at the j^{th} point in our time series is given by

$$t[j] = t_0 + \left(\frac{t_f - t_0}{N - 1} \right) \cdot j \quad (4)$$

Let us call this style of data a **Linear 1D Scan** along the x -axis. In light of this equation, let us store our data in slightly more efficient manner

$$\text{data} = \left[\begin{array}{l} [t_0, t_f], [V_0(t_0), V_0(t_1), \dots, V_0(t_f)], \\ [V_1(t_0), V_1(t_1), \dots, V_1(t_f)], \dots \\ [V_m(t_0), V_m(t_1), \dots, V_m(t_f)] \end{array} \right] \quad (5)$$

where $[t_0, t_f]$ is simply shorthand an N dimensional array whose values are determined by equation (4). Here is how its entered:

```
length = 1e7
mu, sigma = 1, 0.1
gaussian = mu + sigma*np.random.randn(length)
t0 = 0.0
tf = 100.0
d.createDataset("LinearWaveform",
                [("indepName1", [2], "float64", "Seconds")],
                [("depName1", [int(length)], "float64", "Volts"),
                 ("depName2", [int(length)], "float64", "Volts")])
shorthandTime = [t0, tf]
d.addParameter("Scan Type", "Lin")
d.addData([[shorthandTime, gaussian, gaussian]])
d.getData()
```

As another example, suppose we are doing reflection and transmission measurements with a 2-port VNA from 20 MHz to 20 GHz. A linear frequency sweep will place a very low point density in the "low" frequency region with a majority of the points falling in the GHz region of frequency space. One way to get around this is by linearizing your frequency data on a logarithmic scale, i.e. a log sweep. Again, the independent axis, frequency, is determined by 3 numbers, that is the starting frequency f_0 , the stopping frequency f_f , and the number of points swept in frequency space N . Then the frequency at the m^{th} point is given by

$$\log(f[m]) = \log(f_0) + \left(\frac{\log(f_f) - \log(f_0)}{N - 1} \right) \cdot m \quad (6)$$

Lets call this style of data a **Logarithmic 1D Scan** along the x -axis. Following

the above paragraph our data will look like

$$\text{data} = \left[\begin{array}{l} [f_0, f_f], [S_{11}(f_0), S_{11}(f_1), \dots, S_{11}(f_f)], \\ [S_{12}(f_0), S_{12}(f_1), \dots, S_{12}(f_f)], \\ [S_{21}(f_0), S_{21}(f_1), \dots, S_{21}(f_f)], \\ [S_{22}(f_0), S_{22}(f_1), \dots, S_{22}(f_f)] \end{array} \right] \quad (7)$$

where $[f_0, f_f]$ is simply shorthand an N dimensional array whose values are determined by equation (6).

Since 1D scans are really just specific cases of arbitrary 1D datasets, you may wonder why this discussion is justified. However if the distinction being made here seems artificial, consider the fact that data returned from most ADCs, VNAs, and Spectrum Analyzers contains only the y values, the x -axis is built in software and returned for your convenience.

3.2 2D data

By 2D dataset, we mean a dataset with 2 independent variables and m dependent variables. The example that jumps to mind is 2D bias sweep over some rectangular grid but there are many others.

3.2.1 Arbitrary 2D data

Following the section on Arbitrary 1D Data, the most general type of 2D data has an arbitrary spacing between consecutive points with respect to each of the dependent variables. Meaning you could move along a spiral, raster over some arbitrary 2D shape, or sweep over a rectangle. Suppose we can vary two independent quantities, v and p . We are interested in studying how the m dependent quantities $Q_1(v, p), \dots, Q_m(v, p)$ vary with v and p . For each combination $(v, p) \in \{(v_0, p_0), (v_1, p_1), \dots, (v_N, p_N)\}$, we measure each of the Q_i 's. The data will look something like

$$\text{data} = \left[\begin{array}{l} [v_0, p_0, Q_1(v_0, p_0), \dots, Q_m(v_0, p_0)], \\ [v_1, p_1, Q_1(v_1, p_1), \dots, Q_m(v_1, p_1)], \dots \\ [v_n, p_n, Q_1(v_n, p_n), \dots, Q_m(v_n, p_n)] \end{array} \right] \quad (8)$$

Again let's call data of this format **Arbitrary Type 1 Data** where we now have two independent quantities, v and p , instead of one.

Alternatively we could format it as

$$\text{data} = \left[\begin{aligned} &[v_0, v_1, \dots, v_n], [p_0, p_1, \dots, p_n], \\ &[Q_1(v_0, p_0), Q_1(v_1, p_1), \dots, Q_1(v_n, p_n)], \\ &[Q_2(v_0, p_0), Q_2(v_1, p_1), \dots, Q_2(v_n, p_n)], \dots \\ &[Q_m(v_0, p_0), Q_m(v_1, p_1), \dots, Q_m(v_n, p_n)] \end{aligned} \right] \quad (9)$$

Again we call data of this format **Arbitrary Type 2 Data**.

3.2.2 2D scans

While sweeping around in 2D space in an arbitrary manner is fun, it sometimes makes sense to sweep in a more systematic way. Suppose we are interested in rastering over some rectangular grid, namely all pairs of points (v, p) such that $v \in \{v_0, v_1, \dots, v_m\}$ and $p \in \{p_0, p_1, \dots, p_n\}$ where each of these sets follows a formula similar to equation (4) or (6). By measuring each of the Q_i 's for all pairs, we map out a 2D surface over some rectangular domain that we decide upon in software. Lets fix a point along the v -axis and perform a 1D Scan along the p -axis, then move to the next point along the v -axis and repeat, we obtain data that looks like

$$\text{data} = \left[\begin{aligned} &[v_0, [p_0, p_n], [Q_1(v_0, p_0), Q_1(v_0, p_1), \dots, Q_1(v_0, p_n)]] \\ &[v_1, [p_0, p_n], [Q_1(v_1, p_0), Q_1(v_1, p_1), \dots, Q_1(v_1, p_n)]], \dots \\ &[v_m, [p_0, p_n], [Q_1(v_m, p_0), Q_1(v_m, p_1), \dots, Q_1(v_m, p_n)]] \end{aligned} \right] \quad (10)$$

where $[p_0, p_n]$ is compact notation for a linear or logarithmic scan. Let us call this a **2D Scan** with scan direction along the y -axis.

Alternatively, we could have fixed a point along the p -axis and performed a 1D scan along the v -axis, then stepped to the next point on the p axis and so on which would give the data

$$\text{data} = \left[\begin{aligned} &[v_0, v_m], p_0, [Q_1(v_0, p_0), Q_1(v_1, p_0), \dots, Q_1(v_m, p_0)] \\ &[v_0, v_m], p_1, [Q_1(v_0, p_1), Q_1(v_1, p_1), \dots, Q_1(v_m, p_1)], \dots \\ &[v_0, v_m], p_n, [Q_1(v_0, p_n), Q_1(v_1, p_n), \dots, Q_1(v_m, p_n)] \end{aligned} \right] \quad (11)$$

Again we call this a **2D Scan** with scan direction along the x -axis.

These two forms encapsulate how nearly all buffered 2D scans are actually carried out with instruments. Fix one independent variable, sweep the other while measuring some quantity, step and repeat.