

# Single Partition Regression Analysis

Rhys Hawkins and Malcolm Sambridge

October 4, 2012

This tutorial describes how to analyse continuous data using Monte-Carlo Markov-Chain based regression.

## Pre-requisites

This tutorial is for the Python version of the rjmcmmc library. The examples rely on the Matplotlib library for plotting. The versions used in the development of this tutorial are as follows:

- Python 2.7.1
- rjmcmmc 0.1.0
- Matplotlib 1.1.0

## The Data

For this tutorial we will use a non-trivial (in the sense that it will require a higher order polynomial to fit the function correctly) synthetic dataset with added noise.

The function that is used is an exponentially increasing sine wave over the domain 0 ... 10, ie:

$$y = e^{\frac{x}{3}} \sin \frac{2x}{3} \quad (1)$$

The actual dataset unevenly (though with fairly good coverage) samples this function and adds some Gaussian noise and these values are save to an ASCII text file. A plot of the synthetic data points with the true function is shown in Figure 1.

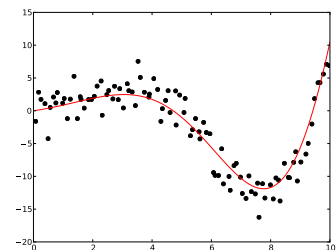


Figure 1: The Synthetic Data.

## Loading the Data

For those unfamiliar with Python, we discuss briefly the loading of data from a simple ASCII text file. The code for loading the data and doing a simple plot is shown in Listing 1.

Listing 1: Loading the Data.

```
1 #  
2 # Import the libraries we will need for analysis and plotting.  
3 #  
4 import rjmcmmc  
5 import matplotlib  
6 import matplotlib.pyplot  
7  
8 #  
9 # Open our data file which consists of one (x, y) coordinator per line
```

```

10 # separated by whitespace
11 #
12 f = open('data.txt', 'r')
13 lines = f.readlines()
14
15 x = []
16 y = []
17
18 for line in lines:
19     columns = line.split()
20     x.append(float(columns[0]))
21     y.append(float(columns[1]))
22
23 f.close()
24
25 fig = matplotlib.pyplot.figure()
26
27 matplotlib.pyplot.plot(x, y, 'ko')
28
29 fig.savefig('ch1-loading.pdf', format='PDF')
30 matplotlib.pyplot.show()
31

```

The ASCII file contains an x, y pair per line separated by a space. On line 12, we use the built in function open to open the file. On lines 15 and 16 we initialize 2 list that will contain the x and y coordinates in the file. On line 18 we loop through the lines in the file and add each x, y pair to the 2 separate lists.

From line 26 onwards, we plot the data using the matplotlib library and save the plot to a PDF file. The plot resulting from this script can be seen in Figure 2.

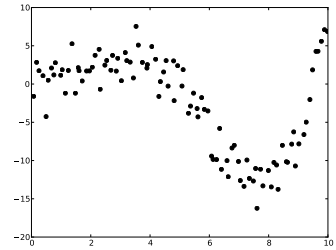


Figure 2: Noisy data and true curve used in test problem.

### Running the default analysis

For performing an regression analysis on a continuous dataset the function is called `regression_single1d`. The parameters for this function are as follows with default values shown where applicable:

*dataset* The dataset object to run the analysis on. This is an `rjmc.mc.dataset1d` object which wraps the x and y vectors you load from the file and includes individual point noise values. This is the only parameter which doesn't have a default value.

*burnin* = 10000 The number of initial samples to throw away.

*total* = 50000 The total number of samples to use for the analysis.

*max\_order* = 2 The maximum order of polynomial to use to fit the data.

*xsamples* = 100 The number of points to sample along the x direction for the curve.

*ysamples* = 100 The number of points to sample along the y directory for the statistics such as mode, median and confidence intervals. This is the number of bins for the histograms in the y direction.

*confidence\_interval* = 0.95 The confidence interval to use for minimum and maximum confidence intervals. This should be a value between 0 and 1.

For this analysis we are only going to use the default values and the listing is shown in Listing 2.

Listing 2: Running the Default Analysis.

```

1 #
2 # Import the libraries we will need for analysis and plotting.
3 #
4 import rjmc
5 import matplotlib
6 import matplotlib.pyplot
7
8 #
9 # Open our data file which consists of one (x, y) coordinate per line
10 # separated by whitespace
11 #
12 f = open('data.txt', 'r')
13 lines = f.readlines()
14
15 x = []
16 y = []
17
18 for line in lines:
19     columns = line.split()
20
21     x.append(float(columns[0]))
22     y.append(float(columns[1]))
23
24 f.close()
25
26 #
27 # Estimate our error standard deviation
28 #
29 sigma = 3.0
30 n = [sigma] * len(x)
31
32 #
33 # Create the rjmc dataset
34 #
35 data = rjmc.dataset1d(x, y, n)
36
37 #
38 # Run the default analysis
39 #
40 results = rjmc.regression_single1d(data)
41
42 #
43 # Retrieve the mean curve for plotting
44 #
45 xc = results.x()
46 meancurve = results.mean()
47
48 #
49 # Plot the data with black crosses and the mean with a red line
50 #
51 fig = matplotlib.pyplot.figure()
52 matplotlib.pyplot.plot(x, y, 'ko', xc, meancurve, 'r-')
53 fig.savefig('ch2-analyse.pdf', format='PDF')
54 matplotlib.pyplot.show()

```

The preamble (lines 1 . . . 24) consists of loading the file as in the previous section.

An important part of the analysis is estimating the error in the data. This is specified as a error value per data point and can be thought of a weighting as to how well the fit will attempt to fit an individual point. If the value is low, then the fit will be tight and conversely if the value is high then the fit will be loose. On lines 29 and 30 we set a value of 3.0 for all data points. Use this value for now, but try other values greater than 0.0 to see the effect.

On line 35 we construct the `dataset1d` object from the `x`, `y` and `n` lists we created. These lists must be the same length.

On line 40 we run the analysis with this `dataset1d` object. The `regression_single1d` function returns a `resultset1d` object which contains various results and diagnostics about the analysis. For this simple analysis we simply take the `x` sampling coordinates and the mean of the fits. And plot the mean with the original data points to

see how representative the mean is. This plot is shown in Figure 3.

## Order Analysis

A question we may ask is what order best represents the underlying function. We can develop some understanding of this by constraining the maximum allowable order of the fit and observing when the order histogram converges and/or the mean of the fits converges. The script to do this for this dataset is shown in Listing 3.

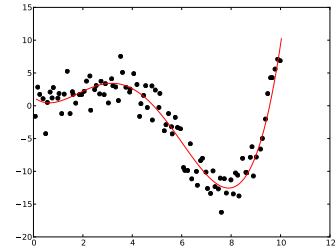


Figure 3: The Default Analysis Plot.

Listing 3: Order Analysis.

```

1  #
2  # Import the libraries we will need for analysis and plotting.
3  #
4  import rjmc
5  import matplotlib
6  import matplotlib.pyplot
7  from mpl_toolkits.mplot3d import Axes3D, Axes3D
8  import sys
9
10 #
11 # Open our data file which consists of one (x, y) coordinator per line
12 # separated by whitespace
13 #
14 f = open('data.txt', 'r')
15 lines = f.readlines()
16
17 x = []
18 y = []
19
20 for line in lines:
21     columns = line.split()
22     x.append(float(columns[0]))
23     y.append(float(columns[1]))
24
25 f.close()
26
27 #
28 # Estimate our error standard deviation
29 #
30 sigma = 3.0
31 n = [sigma] * len(x)
32
33 #
34 # Create the rjmc dataset
35 #
36 data = rjmc.dataset1d(x, y, n)
37
38 #
39 # Run a series of analyses with varying maximum allowed order
40 #
41 results = []
42 burnin = 100
43 total = 1000
44 orderlimit = 5
45 for maxorder in range(orderlimit + 1):
46     print maxorder
47     results.append(rjmc.regression_single1d(data, burnin, total, maxorder))
48
49 colours = ['b', 'g', 'r', 'c', 'm', 'y', 'k', 'w']
50 formats = map(lambda x: x + '-', colours)
51
52 #
53 # Plot the data with black crosses the curves from each of the analyses
54 # with a different colour
55 #
56 fig = matplotlib.pyplot.figure(1)
57 ax = fig.add_subplot(111)
58
59 orders = []
60 legendtitles = []
61 for result in results:
62     order = result.order_histogram()
63     if order == None: # The max order = 0 case will return None so
64         order = [total]
65
66     ax.plot(result.x(), result.mean(), formats[len(orders)])
67
68 #
69 # Create the order histogram data (append zeros for orders not allowed
70 # in the analyses
71 #
72 legendtitles.append('Max_Order_%d' % len(orders))
73 orders.append(order + [0] * (orderlimit + 1 - len(order)))
74
75 ax.plot(x, y, 'ko')
76
77

```

```

78 legendtitles.append('Data')
79 legend = ax.legend(legendtitles, 'lower_left')
80 fig.savefig('ch3-orderanalysis.pdf', format='PDF')
81
82 #
83 # Plot a 3D bar chart showing the progression of the order histogram
84 # as the analysis maximum order is increased.
85 #
86 fig = matplotlib.pyplot.figure(2)
87 ax = Axes3D(fig)
88 xs = range(orderlimit + 1)
89 for maxorder in xs:
90     ax.bar(xs,
91           orders[maxorder],
92           zs=maxorder,
93           zdir = 'y',
94           color=colours[maxorder])
95
96 ax.set_xlabel('Order')
97 ax.set_ylabel('Maximum_Order')
98 ax.set_zlabel('Count')
99
100 fig.savefig('ch3-orderanalysishist.pdf', format='PDF')
101
102 matplotlib.pyplot.show()

```

In lines 41 ... 46 we set some parameters and run several analyses with different allowed maximum polynomial order and store the results.

In lines 51 ... 78 we plot the mean of the fits for each of the constrained analyses. This plot is shown in Figure 4. As can be seen from this plot, there is very little difference between the Max. Order 3, 4, 5 curves which implies that the data is cubic.

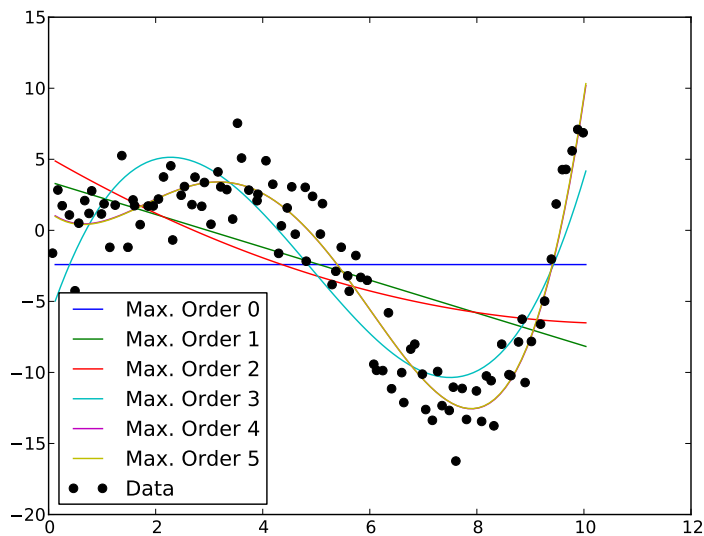


Figure 4: The Mean Curves for different maximum orders.

In lines 80 ... 99 we plot the progression of the order histogram. The order histogram reports the count of polynomials of each order where the order is randomly chosen from a probability distribution determined by the data itself. This plot is shown in Figure 5.

Looking at the figure, it can be seen that when the maximum permitted order is 0, all the curves sampled are 0th order as ex-

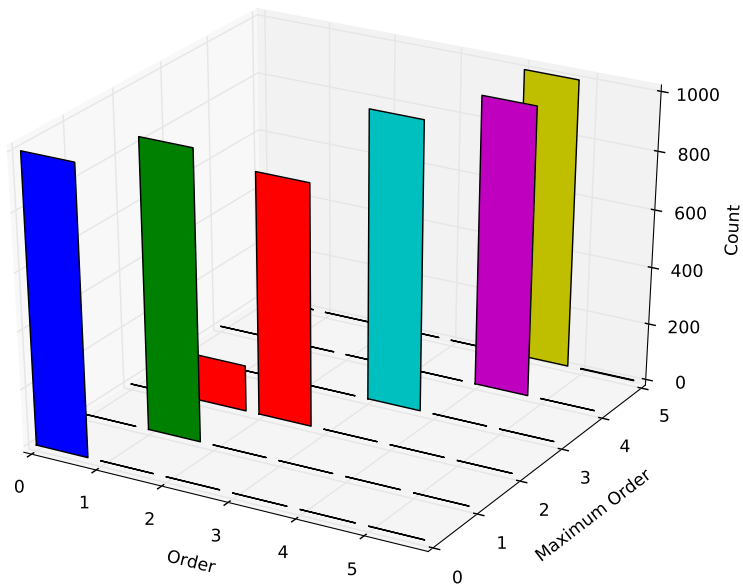


Figure 5: Posterior PDFs of the polynomial order parameter as a function of maximum allowed order.

pected. As this maximum order is increased, the distribution of orders changes until the step from 4 to 5 where there is little difference. It should be noted that if a tighter noise parameter is used, that this will change slightly as higher order polynomials will be more favoured to fit the data better.

### Confidence Intervals

So far we have only plotted the mean of the fits, however this gives us no indication of distribution of the fit (this can be thought of as the confidence of the fit). There are a number of ways in which we can look at this and one of these is to look at the curves generated during the analysis. The listing to do this is in Listing 4.

Listing 4: Confidence Intervals.

```

1  #
2  # Import the libraries we will need for analysis and plotting.
3  #
4  import rjmc
5  import matplotlib
6  import matplotlib.pyplot
7  from mpl_toolkits.mplot3d import Axes3D, Axes3D
8
9  #
10 # Open our data file which consists of one (x, y) coordinate per line
11 # separated by whitespace
12 #
13 f = open('data.txt', 'r')
14 lines = f.readlines()
15
16 x = []
17 y = []
18
19 for line in lines:
20     columns = line.split()
21     x.append(float(columns[0]))
22 
```

```

23     y.append(float(columns[1]))
24
25 f.close()
26
27 #
28 # Estimate our error standard deviation
29 #
30 sigma = 3.0
31 n = [sigma] * len(x)
32
33 #
34 # Create the rjcmc dataset
35 #
36 data = rjcmc.dataset1d(x, y, n)
37
38 #
39 # This is our callback function which samples the curves generated
40 # during the analysis
41 #
42 sample_x = None
43 sample_curves = []
44 sample_i = 0
45 sample_rate = 5
46 def sampler_cb(x, y):
47     global sample_x, sample_curves, sample_i, sample_rate
48
49     if sample_i == 0:
50         sample_x = x
51
52     if sample_i % sample_rate == 0:
53         sample_curves.append(y)
54
55     sample_i = sample_i + 1
56
57 #
58 # Run a series of analyses with varying maximum allowed order
59 #
60 results = []
61 burnin = 100
62 total = 1000
63 maxorder = 5
64 results = rjcmc.regression_single1d_sampled(data,
65                                             sampler_cb,
66                                             burnin,
67                                             total,
68                                             maxorder)
69
70 #
71 # Plot the data with black crosses, the sample curves as faint lines, and
72 # the mean as a red line
73 #
74 fig = matplotlib.pyplot.figure()
75 ax = fig.add_subplot(111)
76
77 yc = 0.5
78 yalpha = 1.0/((1.0 - yc) * float(len(sample_curves)))
79 for sy in sample_curves:
80
81     ax.plot(sample_x, sy,
82             color = str(yc),
83             alpha = yalpha,
84             linestyle = '-',
85             linewidth = 10)
86
87 ax.plot(results.x(), results.mean(), 'r-')
88 ax.plot(x, y, 'ko')
89 fig.savefig('ch4-confidence.pdf', format='PDF')
90
91 matplotlib.pyplot.show()

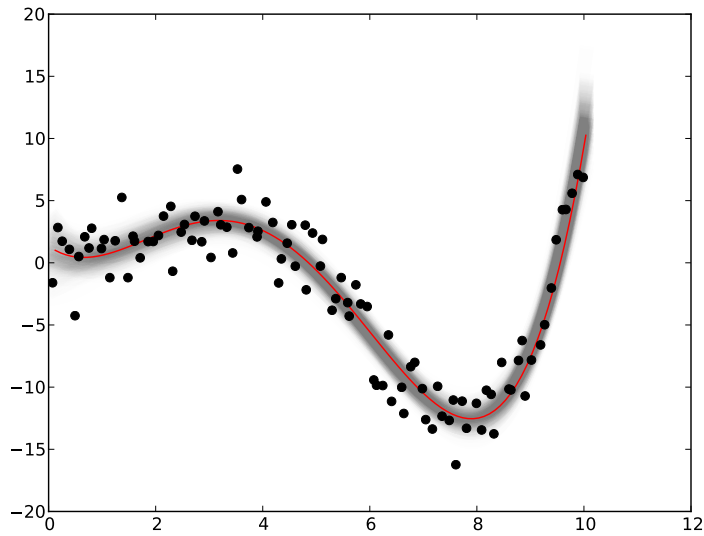
```

In this script we call a slightly different function called `regression_single1d_sampled` which accepts a callback function. We define this function on lines 46 ... 55. This function accepts an `x` and `y` list which is the discretization of the current fitting polynomial being used. In this function we sample every 5th polynomial and store them.

On lines 70 ... 90 we plot the data as black dots, and the mean fit with a red line over plots of all the fits we sampled. This plot is shown in Figure 6.

The sampled fits are plotted with transparency so that where they overlap this will show increased density implying that where these sampled polynomial ensemble appears darker, we can have higher confidence that the underlying function passes through that region.

Figure 6: Sampled Confidence Intervals.



### Estimating data noise

With the hierarchical Bayesian approach we include the standard deviation of the noise on the observations as an unknown. In the above examples the noise  $\sigma$  was set to 3 units, but the actual  $\sigma$  of the data noise in Figure 1 is 2.5. Can we use the data to detect the true standard deviation of its noise? The hierarchical Bayesian sampling scheme is implemented with the script in Listing 5. Inference on the noise is implemented by introducing a new parameter,  $\lambda = \frac{\sigma_{est}}{\sigma_{true}}$ , defined as the ratio of the estimated noise to the real noise.

Listing 5: Hierarchical Bayesian sampling to determine the data noise standard deviation

```

1 #
2 # Import the libraries we will need for analysis and plotting.
3 #
4 import rjmc
5 import matplotlib
6 import matplotlib.pyplot
7 #
8 #
9 # Open our data file which consists of one (x, y) coordinate per line
10 # separated by whitespace
11 #
12 f = open('data.txt', 'r')
13 lines = f.readlines()
14
15 x = []
16 y = []
17
18 for line in lines:
19     columns = line.split()
20     x.append(float(columns[0]))
21     y.append(float(columns[1]))
22
23 f.close()
24
25 #
26 # Estimate our error standard deviation
27 #
28 #

```



```

29 sigma = 3.0
30 n = [sigma] * len(x)
31
32 #
33 # Create the rjcmc dataset
34 #
35 data = rjcmc.dataset1d(x, y, n)
36
37 lambda_min = 0.5
38 lambda_max = 2.0
39 lambda_std = 0.05
40
41 data.set_lambda_range(lambda_min, lambda_max)
42 data.set_lambda_std(lambda_std)
43
44 #
45 # Run the default analysis
46 #
47 results = rjcmc.regression_single1d(data)
48
49 #
50 # Retrieve the mean curve for plotting
51 #
52 xc = results.x()
53 meancurve = results.mean()
54
55 #
56 # Retrieve the results of the hierarchical
57 #
58 p = results.proposed()
59 a = results.acceptance()
60
61 print 'Lambda_Acceptance_Rate:', float(a[1])/float(p[1]) * 100.0
62
63 lh = results.lambda_history()
64
65 #
66 # Plot the data with black crosses and the mean with a red line
67 #
68 fig = matplotlib.pyplot.figure(1)
69 matplotlib.pyplot.plot(x, y, 'ko', xc, meancurve, 'r-')
70
71
72 fig = matplotlib.pyplot.figure(2)
73 matplotlib.pyplot.plot(range(len(lh)), lh)
74
75 fig = matplotlib.pyplot.figure(3)
76
77 a = matplotlib.pyplot.subplot(111)
78 lsamples = lh[10000:]
79
80 n, bins, patches = a.hist(lsamples, 100, range=(lambda_min, lambda_max))
81 a.set_title('Histogram_of_Lambda')
82 a.set_xlabel('Lambda')
83 a.set_ylabel('Count')
84
85
86 print 'Lambda_average:', sum(lsamples)/float(len(lsamples))
87 fig.savefig('ch5-hierarchical.pdf', format='PDF')
88 matplotlib.pyplot.show()

```

In this script we set up a uniform prior on  $\lambda$  over a pre-determined range and use a Gaussian distribution to perturb the  $\lambda$  values during the Markov chain. The range of the values of  $\lambda$  as well as the standard deviation of the perturbation are parameter that must be chosen. These are set in lines 37-42.

We call the function `regression_single1d` to do the work. We plot the posterior PDF of the noise  $\sigma$  as a histogram. This plot is shown in Figure 7.

The histogram shows the support of the data for a range of  $\lambda$  values. Clearly there is information in the data on the likely values of noise. Where is the peak of the histogram? How does this compare to the ratio of the estimated to true  $\sigma$ ? Usually the ability of the data to constrain noise parameters will trade-off with the model complexity given, in this case, by the order of the polynomial. You can edit the script by changing the estimated noise and rerun to see what happens.

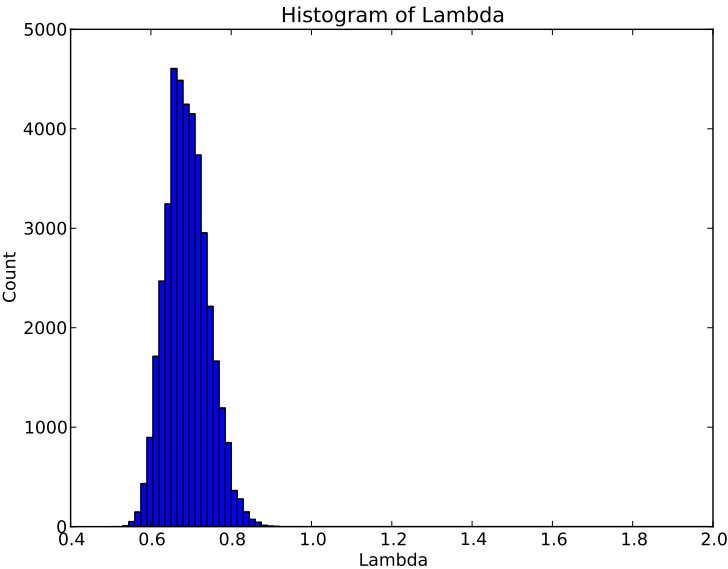


Figure 7: Posterior PDF of the noise standard deviation parameter  $\lambda$ .