

# Assignment 1: Estimating sea-level trend using least-squares

G.H. Garrett, 4450973  
*Delft University of Technology, Delft, Zuid-Holland, 2628 HS*

This assignment gave some interesting opportunities to touch up on my passion in data analysis, providing me with the volition to write a Medium post titled “A story of Satellites, Signals and Statistics” which can be found on my Medium page (<https://medium.com/@g.h.garrett>). Furthermore the software developed within this assignment can be found on my GitHub (<https://github.com/ggarrett13>). The software will be converted into a package to provide easy installation via *python-pip* within a week of submission. I fully encourage the use of the package under the MIT License. Finally, the assignment as a whole, considering the extra documentation and handling took approximately 15 hours to complete.

## I. Estimating sea-level trend using least-squares

---

**Part a):** Estimate the sea-level trend and bias using least-squares and plot the data and your estimated trend. Write your own code for the least-squares estimation.

---

The task is broken up into steps that are required to adapt modular code for general application.

1. Create a function that solves for the unweighted vector of parameters ( $\hat{\beta}$ ) estimation depending on the given arguments, taking into consideration the three forms of the LSQ algorithm, given the information matrix ( $H$ ) and observations ( $\bar{y}$ ).

$$\hat{\beta} = \begin{cases} (H^T H)^{-1} H^T \bar{y} & \forall m > n \\ H^{-1} \bar{y} & \forall m = n \\ H^t (H H^t)^{-1} \bar{y} & \forall m < n \end{cases}$$

Where  $m$  is the quantity of observations,  $n$  is the quantity of predictors (or basis functions),  $\hat{\beta}$  is the least squares solution of the vector of parameters ( $\bar{\beta}$ ),  $H$  is the information matrix and  $\bar{y}$  is the vector of observations.

$$\bar{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{n-1} \end{pmatrix} \quad (1)$$

$$H = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{pmatrix} \quad (2)$$

$$\bar{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \quad (3)$$

2. Create a modular object (relational data structure) that represents a linear model which takes a list of basis functions as its argument and can output  $H$  given the independent variable ( $x$ ) which provides input to the basis functions.
3. Test the model analysis tool on the linear model containing trend and bias (equation 4) and plot the result.

$$y = \beta_0 + \beta_1 \cdot x \quad (4)$$

The code for the above can be found on GitHub (link found on the front page) for the most recent version. Additionally the version at the time of submission of the report can be found in the appendix. Plotting model 1, described by Equation 4 results in Figure 1,  $\hat{\beta}$  is given by Equation 5.

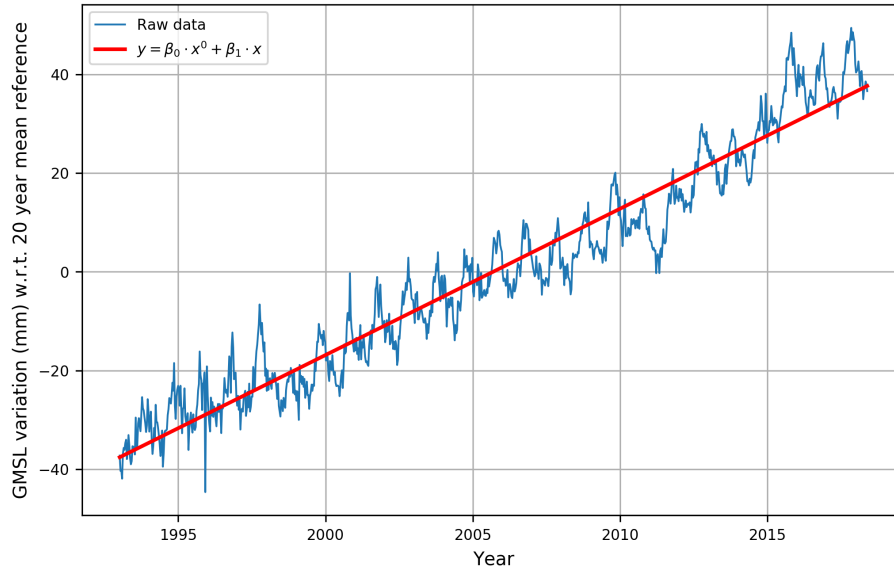


Fig. 1: Plot of model 1 ( $y = \beta_0 + \beta_1 \cdot x$ ) superimposed on the raw observations.

$$\begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} = \begin{pmatrix} -5.94 \cdot 10^3 \\ 2.96 \cdot 10^0 \end{pmatrix} \quad (5)$$

---

**Part b):** Do the same as for a) but simultaneously estimate a signal with a period of 1 year. Hint: use a combination of a sine and a cosine function. Explain in the report why it makes sense to estimate a yearly signal in sea level data, and explain if and why the estimated trend differs from a).

---

The code used for **Part a)** can be applied directly to the new equation describing the design linear model in Equation 6.

$$y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot \sin(2\pi f x) + \beta_3 \cdot \cos(2\pi f x), \quad (6)$$

where  $f = 1/T$  and the period being used to model the periodic signal is 1 year, therefore  $f = 1$ .

Plotting the predicted solution ( $\hat{y}$ ) provides Figure 2 and gives the least square solution for  $\hat{\beta}$  given by Equation 7.

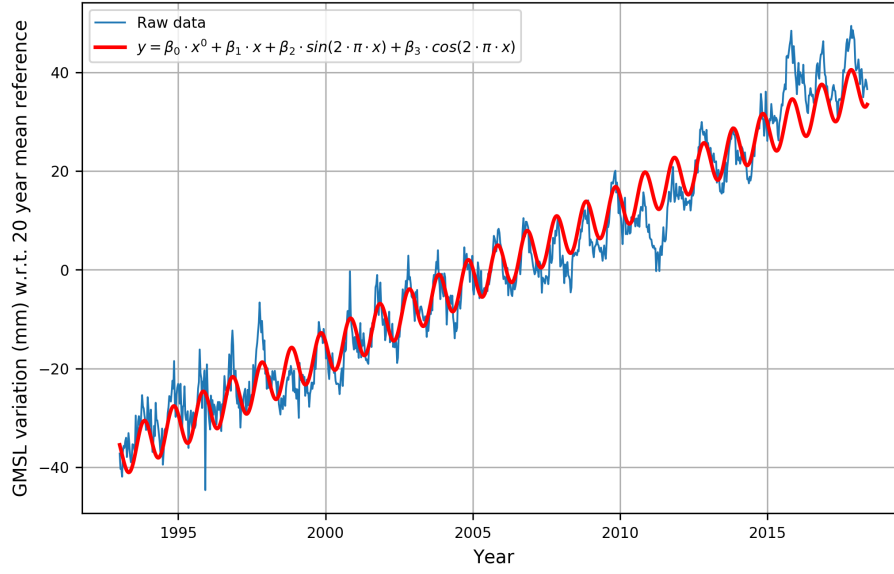


Fig. 2: Plot of model 2 ( $y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot \sin(2\pi x) + \beta_3 \cdot \cos(2\pi x)$ ) superimposed on the raw observations.

$$\begin{pmatrix} \beta_0 & \beta_1 & \beta_2 & \beta_3 \end{pmatrix}^T = \begin{pmatrix} -5.94 \cdot 10^3 \\ 2.96 \cdot 10^0 \\ -3.85 \cdot 10^0 \\ 2.32 \cdot 10^0 \end{pmatrix} \quad (7)$$

### Observations

1. The periodic signal (of  $T = 1$  year) included in the linear model provides correct time interval predictions for the periodic signal visible within the observation data.
2. The periodic signal (of  $T = 1$  year) does not provide any predictive power for the outlying peaks (local maximums) and troughs (local minimums) seen in the observation data.

### Analysis & Evaluation

1. The yearly period can both be observed by the data and reasoned by environmental sciences as the yearly freezing and thawing of the polar ice caps due to the inclination of the Earth's axis and the effect this has due to solar radiation incident.
2. The linear model can benefit by the inclusion of additional basis functions following further analysis.

3. Within 3 significant figures, the values for  $\beta_0$  and  $\beta_1$  have remained constant during the recalculation of  $\hat{\beta}$  on Equation 6 using the algorithm for unweighted least squares.

---

**Part c):** Find out if the residuals are normally distributed. Explain your test and the conclusions in the report.

---

### Test Procedure

1. Firstly, calculate the vector of residuals for the linear model ( $\bar{\epsilon}_l$ ) given by the difference between the actual observations ( $\bar{y}$ ) and predicted values ( $\hat{y}$ ). Use an algorithm to plot a histogram of residuals with 30 bins.

$$\bar{\epsilon}_l = \bar{y} - \hat{y} = \bar{y} - H\hat{\beta} \quad (8)$$

2. Secondly compare with the result of Step 1 to a normal distribution with the mean ( $\mu$ ) and standard deviation ( $\sigma$ ) of  $\bar{\epsilon}_l$ .

$$p(\epsilon_l) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{\epsilon_l - \mu}{\sigma} \right)^2 \right] \quad (9)$$

3. Finally calculate the skewness of the discrete random variable distribution ( $\gamma_1$ ) and the kurtosis by old definition ( $\gamma_2$ ) and finally kurtosis by the new definition ( $\gamma_2'$ ).  $\gamma_1$  is calculated by:

$$\gamma_1 = \frac{E[(X - \lambda_1)^3]}{\sigma^3}, \quad (10)$$

where:

$$E[X^k] = \lambda_k = \int_{-\infty}^{\infty} x^k p(x) dx = \sum_{i=1}^m p(x_i) x_i^k, \quad (11)$$

and:

$$\mu_k = E[(X - \lambda_1)^k] = \int_{-\infty}^{\infty} (x - \lambda_1)^k p(x) dx = \sum_{i=1}^m p(x_i) (x_i - \lambda_1)^k. \quad (12)$$

Finally the old definition of kurtosis is given by Equation 13 and the new by Equation 14.

$$\gamma_2 = \frac{\mu_4}{\sigma^4} \quad (13)$$

$$\gamma'_2 = \gamma_2 - 3 \quad (14)$$

## Results

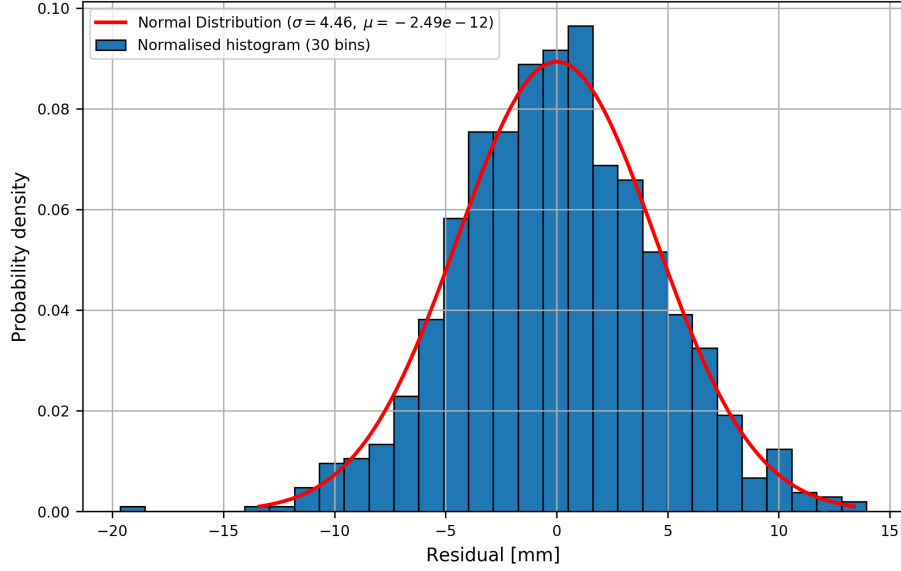


Fig. 3: Plot of model 2 ( $y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot \sin(2\pi x) + \beta_3 \cdot \cos(2\pi x)$ ) normalised residuals with ideal normal distribution superimposed.

Table 1: Tabulated values for skewness ( $\gamma_1$ ) and new kurtosis ( $\gamma'_2$ ) for model 1 and 2.

Parameter	Model 1	Model 2
$\gamma_1$	0.208	0.00183
$\gamma'_2$	-0.384	-0.0803

## Conclusions

1. The normalised histogram plot of  $\bar{\epsilon}_l$  follows the shape of the ideal normal distribution calculated closely.

2. The calculated  $\mu$  of the residuals tends towards zero which indicates that error is largely from random error and systematic error is close to negligible in the model.
3. There exists some outliers in the negative region of the residuals. This introduces some skewness in the distribution.
4. The inclusion of the yearly periodic signal reduced the skewness ( $\gamma_1$ ) by two orders of magnitude and the new kurtosis ( $\gamma_2'$ ) by one. Therefore the residuals better represent a gaussian distribution as a result of the inclusion of the periodic signal.

---

**Part d):** Do the same as b) and c) but now include an acceleration term. Explain if and why the estimated trend and the residuals differ from your answers to b) and c)

---

The equation used to model the data is as follows below with the inclusion of the acceleration term  $x^2$ . The results follow directly from the previous procedures.

$$y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot x^2 + \beta_3 \cdot \sin(2\pi x) + \beta_4 \cdot \cos(2\pi x) \quad (15)$$

## Results

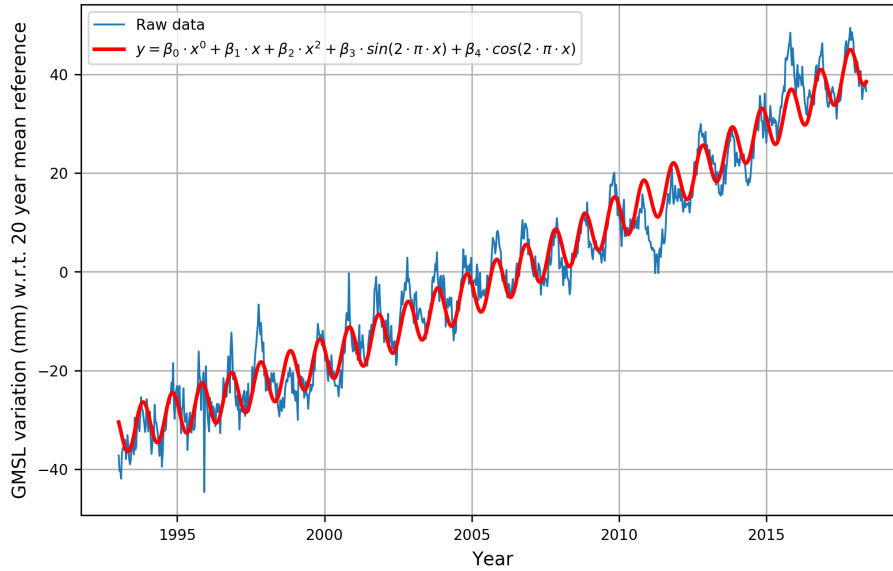


Fig. 4: Plot of model 3 ( $y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot x^2 + \beta_3 \cdot \sin(2\pi x) + \beta_4 \cdot \cos(2\pi x)$ ) superimposed on the raw observations.

$$\begin{pmatrix} \beta_0 & \beta_1 & \beta_2 & \beta_3 & \beta_4 \end{pmatrix}^T = \begin{pmatrix} 1.84 \cdot 10^5 \\ -1.87 \cdot 10^2 \\ 4.73 \cdot 10^{-2} \\ -3.96 \cdot 10^0 \\ 2.28 \cdot 10^0 \end{pmatrix} \quad (16)$$

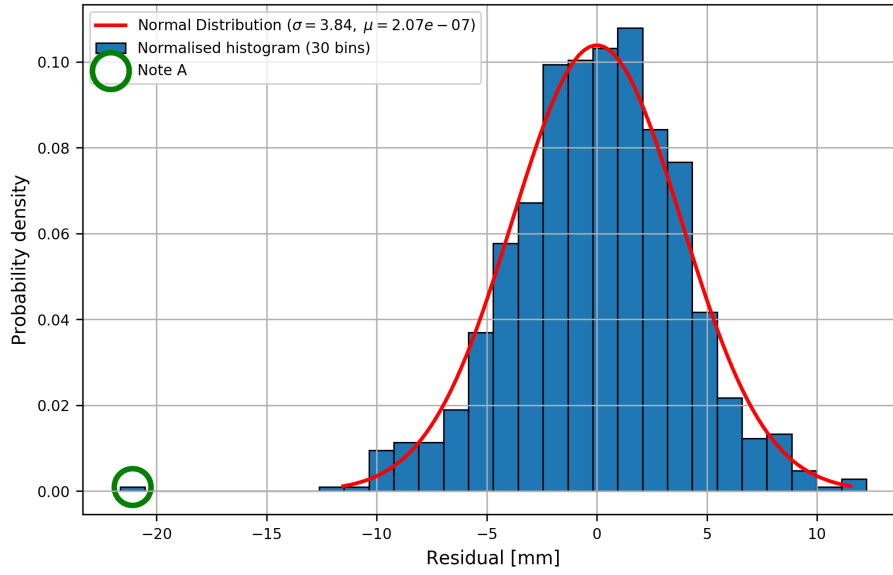


Fig. 5: Plot of model 3 ( $y = \beta_0 + \beta_1 \cdot x + \beta_2 \cdot x^2 + \beta_3 \cdot \sin(2\pi x) + \beta_4 \cdot \cos(2\pi x)$ ) normalised residuals with ideal normal distribution superimposed.

Parameter	Model 1	Model 2	Model 3
$\gamma_1$	0.208	0.00183	-0.221
$\gamma'_2$	-0.384	-0.0803	0.506

## Conclusions

1. According to the measurement of  $\gamma_1$  and  $\gamma'_2$ , the probability density function of model 3's residuals have worsened in representation of a Gaussian distribution. This could be as a result of the increased 3rd and 4th order moment about the mean caused by *Note A* in Figure 5 which has increased in negativity from model 2.
2. The standard deviation ( $\sigma$ ) has decreased significantly from model 2 to model 3. This means



that there is less variation in the residuals about the mean which is approximately equal to zero.

3. Overall it can be said that the model is a better predictor of the observable true value from the observed values, or in other words the random error has been reduced in the model.

---

**Part e):** Explain conceptually what the differences are between residuals and true errors.

---

**Residuals** Residuals are the difference between the observed values of a statistical sample and the predicted (or estimated) values according to a predictive model derived from the statistical sample set (e.g. sample mean).

**True errors** True errors are the difference between the observed values of a statistical sample and the (unobservable) true values (e.g. population mean).

---

**Part f):** Estimate the standard deviation of the trend estimate. You can now assume that standard deviation of the measurements (measurement error) can be computed from the residuals as follows:

$$\sigma_l = \frac{\epsilon \epsilon^T}{m - n} \quad (17)$$

For the symbols see section 8.5 of the lecture notes. In the report, give the equation and describe what each symbol represents and comment on the value of the standard deviation.

---

Using Equation 8 to calculate  $\bar{\epsilon}_l$ , and Equation 17, the standard deviation of the measurements can be obtained ( $\sigma_l$ ), where  $m$  is the quantity of observed values and  $n$  is the quantity of predictors (or basis functions). Equation 18 defines the observation covariance matrix. The leading diagonal contains the variance of each observation, and the off-diagonals contain the covariance between

observations.

$$P_{yy} = \begin{pmatrix} \sigma_{l,1}^2 & \sigma_{l,1}\sigma_{l,2} & \cdots & \sigma_{l,1}\sigma_{l,m} \\ \sigma_{l,2}\sigma_{l,1} & \sigma_{l,2}^2 & \cdots & \sigma_{l,2}\sigma_{l,m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{l,m}\sigma_{l,1} & \sigma_{l,m}\sigma_{l,2} & \cdots & \sigma_{l,m}^2 \end{pmatrix} \quad (18)$$

Typically for **independent** observations which do not influence other observations, the following equation defines the observation covariance matrix ( $P_{yy}$ ):

$$P_{yy} = \sigma_l^2 I \quad (19)$$

Following the calculation of  $P_{yy}$ , the parameter covariance matrix ( $P_{xx}$ ) can be obtained through Equation I.

$$P_{xx} = (H^t P_{yy}^{-1} H)^{-1} \quad (20)$$

The previous calculation procedure is carried out for model 3 and is seen following with all calculations made to 3 significant figures.

$$P_{yy} = \sigma_l^2 I = (14.8)^2 \cdot I_m$$

$$P_{xx} = (H^t P_{yy}^{-1} H)^{-1} = \begin{pmatrix} 1.11 \cdot 10^8 & -1.11 \cdot 10^5 & 2.76 \cdot 10^1 & -6.66 \cdot 10^1 & -2.26 \cdot 10^1 \\ -1.11 \cdot 10^5 & 1.10 \cdot 10^2 & -2.75 \cdot 10^{-2} & 6.64 \cdot 10^{-2} & 2.25 \cdot 10^{-2} \\ 2.76 \cdot 10^1 & -2.75 \cdot 10^{-2} & 6.85 \cdot 10^{-6} & -1.65 \cdot 10^{-5} & -5.63 \cdot 10^{-6} \\ -6.66 \cdot 10^1 & 6.64 \cdot 10^{-2} & -1.65 \cdot 10^{-5} & 3.17 \cdot 10^{-2} & -7.63 \cdot 10^{-5} \\ -2.26 \cdot 10^1 & 2.25 \cdot 10^{-2} & -5.63 \cdot 10^{-6} & -7.63 \cdot 10^{-5} & 3.18 \cdot 10^{-2} \end{pmatrix}$$

---

**Part g):** Explain what the off-diagonal terms of the parameter co-variance matrix represent.

---

The off-diagonal terms represent the covariance between differing parameters, which is a measure of the joint variability of two parameters. From this value we can obtain the correlation between different predictors using the following equation:

$$\rho = \frac{\sigma_{\beta_i \beta_j}}{\sigma_{\beta_i} \sigma_{\beta_j}} \quad (21)$$

---

**Part h):** Explain if you expect that the sea level data you downloaded contain other physical signals except a bias, trend, acceleration and a yearly signal. Describe how you can test this with the data you have.

---

In order to obtain better representations of the signals that are consisted within the raw data, a Fourier analysis can be carried out. The signals found can be then included in the model analysis.

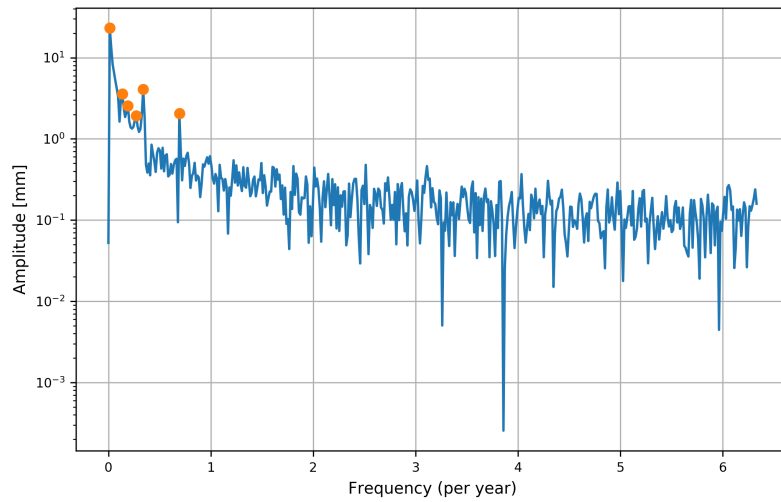


Fig. 6: Fourier analysis carried out on the data set with the 6 signals greater than 1 mm in magnitude marked.

The collection of 6 points above 1 mm in amplitude are shown in the vector seen in Equation 23.

These  $x$  values are used as new frequencies to be placed into the signal:

$$y = \beta_0 \sin(2\pi f_1 x) + \beta_1 \cos(2\pi f_1 x) \dots \beta_{m-1} \sin(2\pi f_{m/2} x) + \beta_m \cos(2\pi f_{m/2} x) \quad (22)$$

This is then added onto model 3 in order to create a new model 4. Consisting of various frequency signals.

$$(x, y) = \begin{pmatrix} (0.01359, 23.10) \\ (0.1359, 3.570) \\ (0.1902, 2.565) \\ (0.2718, 1.920) \\ (0.3397, 4.075) \\ (0.6930, 2.056) \end{pmatrix} \quad (23)$$

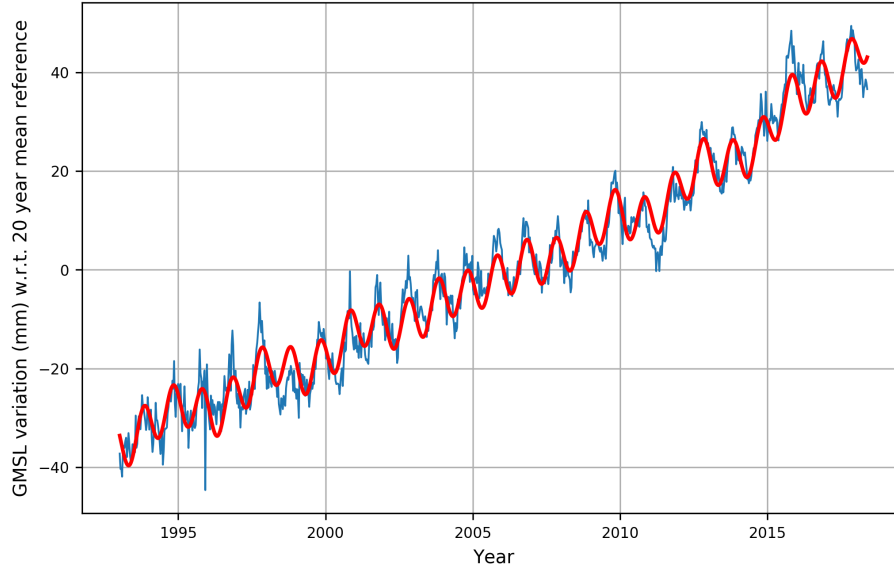


Fig. 7: Plot of model 4 superimposed on the raw observations.

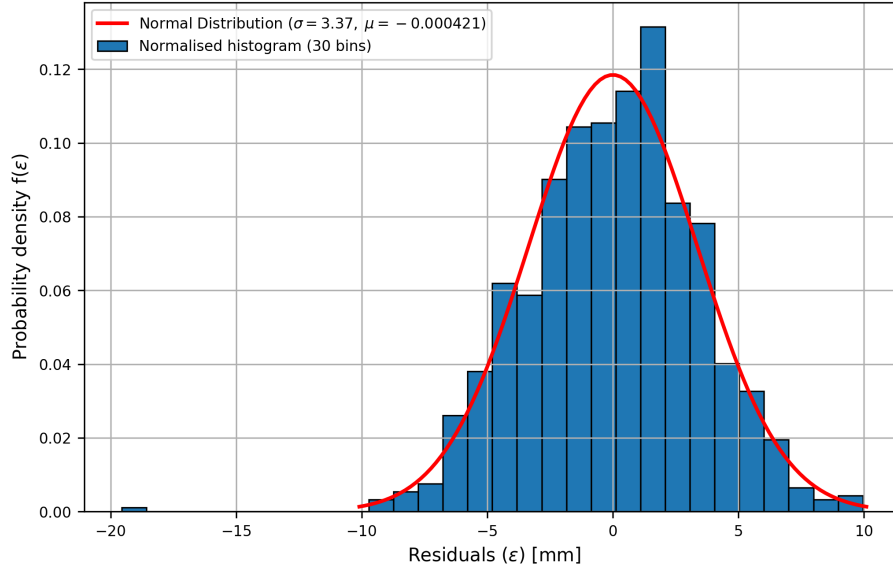


Fig. 8: Plot of model 4 normalised residuals with ideal normal distribution superimposed.

## Conclusion

1. The trend model seems to have significantly improved seen by Figure 7. Some anomalous troughs and peaks have been estimated correctly.
2. The standard deviation of the distributed residuals has decreased from 3.84 to 3.37 showing better prediction of the unobservable true values with reduction of random error.
3. Considering the quantity of extra predictive terms (or basis functions) the design model may not have significantly improved predictive capabilities when considering the amount of extra computation required has been more than doubled (due to more than double the basis functions).

---

**Part i):** Mention four different error sources that affect the global mean sea level estimate derived from a satellite altimeter, assuming that we are interested in the climate-related sea level rise.

---

#### **Error sources**

1. Potential systematic error in the satellite altimeter measurement due to imperfect calibration.
2. Further altimeter error due to ageing or error in the ground processing of data.
3. Aleatory error due to orbit solutions and gravitational models that are constantly improving.
4. The required correction for the wet troposphere which can be affected by long-term instrument drifts. These drifts may be due to internal temperature changes induced by yaw manoeuvres.

## II. Github code

### A. functions.py

```
1 import numpy as np
2
3
4 def parameter_covariance(H, Pyy=None):
5     if Pyy is not None:
6         return np.linalg.inv(np.matmul(H.T, np.matmul(np.linalg.inv(Pyy),
7             H)))
8     else:
9         return np.linalg.inv(np.matmul(H.T, H))
10
11 def unweighted_least_squares(H, y):
12     """
13     :param H: Information matrix (np.ndarray)
14     :param y: Vector of observations (np.ndarray)
15     :return: Vector of parameters (np.ndarray)
16     """
17     # m > n
18     if H.shape[0] > H.shape[1]:
19         return np.matmul(np.linalg.inv(np.matmul(H.T, H)), np.matmul(H.T,
20             y))
21     # m = n
22     elif H.shape[0] == H.shape[1]:
23         return np.matmul(np.linalg.inv(H), y)
24     # m < n
25     elif H.shape[0] < H.shape[1]:
26         return np.matmul(H.T, np.matmul(np.linalg.inv(np.matmul(H, H.T)),
27             y))
28
29 def weighted_least_squares(H, y, Py):
30     """
31     TODO: (**) Complete other forms with covariance !=(m>n).
32     :param H:
33     :param y:
34     :param Py:
35     :return:
36     """
37     # m > n
38     if H.shape[0] > H.shape[1]:
39         return np.matmul(np.linalg.inv(np.matmul(np.matmul(H.T, np.linalg.
40             inv(Py)), H)),
41             np.matmul(np.matmul(H.T, np.linalg.inv(Py)), y))
42     else:
43         raise NotImplementedError("TODO: Complete other forms with
44             covariance !=(m>n)")
45
46 def kth_order_moment_about_zero(x, k, bins=None):
47     """
48     TODO: Complete kth_order_moment_about_zero doc-strings.
49     :param x:
50     :param k:
51     :param bins:
52     :return:
53     """
54     if bins:
55         p_i, edges = np.histogram(x, density=True, bins=30)
56         x_i = [np.mean([edges[i], edges[i + 1]]) for i in range(len(edges)
57             - 1)]
58         return np.sum(np.multiply(np.power(x_i, k), p_i))
```

```

56     else:
57         return np.sum(np.multiply(np.power(x, k), x / np.sum(x)))
58
59
60 def kth_order_moment_about_mean(x, k, bins=None):
61     """
62     TODO: Complete kth_order_moment_about_mean doc-strings.
63     :param x:
64     :param k:
65     :param bins:
66     :return:
67     """
68     if bins:
69         p_i, edges = np.histogram(x, density=True, bins=30)
70         x_i = [np.mean([edges[i], edges[i + 1]]) for i in range(len(edges)
71             - 1)]
72         return np.sum(np.multiply(np.power(x_i - np.mean(x), k), p_i))
73     else:
74         return np.sum(np.multiply(np.power(x - np.mean(x), k), x / np.sum(
75             x)))

```

## B. model\_design.py

```

1  import numpy as np
2  import sympy as sp
3  from sympy.utilities.lambdify import lambdify
4  import matplotlib.pyplot as plt
5
6
7  MOD = 0
8  SMALL_SIZE = 8 + MOD
9  MEDIUM_SIZE = 10 + MOD
10 BIGGER_SIZE = 12 + MOD
11
12 plt.rc('font', size=SMALL_SIZE) # controls default text sizes
13 plt.rc('axes', titlesize=SMALL_SIZE) # fontsize of the axes title
14 plt.rc('axes', labelsize=MEDIUM_SIZE) # fontsize of the x and y labels
15 plt.rc('xtick', labelsize=SMALL_SIZE) # fontsize of the tick labels
16 plt.rc('ytick', labelsize=SMALL_SIZE) # fontsize of the tick labels
17 plt.rc('legend', fontsize=SMALL_SIZE) # legend fontsize
18 plt.rc('figure', titlesize=BIGGER_SIZE) # fontsize of the figure title
19
20
21 class NonLinearDesignModel(object):
22     def __init__(self, basis_functions):
23         """
24         TODO: (*) Complete non-linear design model.
25         :param basis_functions: (list)
26         """
27         self._basis_functions = basis_functions
28         raise NotImplementedError("TODO: Complete non-linear design model
29             .")
30
31 class LinearDesignModel(object):
32     def __init__(self, basis_functions):
33         """
34         :param basis_functions: (list)
35         """
36         self._basis_functions = basis_functions
37
38     @property
39     def basis_functions(self):
40         return self._basis_functions
41

```



```

42     def __add__(self, other):
43         return LinearDesignModel(self.basis_functions + other.
            basis_functions)
44
45     def __radd__(self, other):
46         return LinearDesignModel(other.basis_functions + self.
            basis_functions)
47
48     def __repr__(self):
49         vec_param = ' '.join(['B' + str(i) for i in range(len(self.
            _basis_functions))])
50         vec_basis = ' '.join(self._basis_functions)
51         return 'DesignModel(y = [{]}^T [{]})'.format(vec_param, vec_basis)
52
53     def __str__(self):
54         return 'y = ' + ' + '.join(['B_{}'.format(i) + '*' + j for
55             i, j in enumerate(self.
            _basis_functions)]) + ' '
56
57     def __latex__(self):
58         return '$' + self.__str__().replace('*', '\cdot').replace('pi',
            '\pi').replace('B', '\\beta') + '$'
59
60     def information_matrix(self, x):
61         information_matrix = np.array([])
62         for basis in self.basis_functions:
63             if len(information_matrix) == 0:
64                 if (basis == '1') or (basis == 'x^0'):
65                     information_matrix = np.full((len(x), 1), 1)
66                 else:
67                     information_matrix = lambdify(sp.Symbol('x'), basis, '
            numpy')(x)
68             else:
69                 information_matrix = np.c_[information_matrix, lambdify(sp.
            Symbol('x'), basis, 'numpy')(x)]
70         return information_matrix

```

### C. data\_analysis.py

```

1  from functions import kth_order_moment_about_mean
2  from functions import kth_order_moment_about_zero
3  from functions import unweighted_least_squares
4  from functions import weighted_least_squares
5  from model_design import LinearDesignModel
6  from model_design import NonLinearDesignModel
7  import matplotlib.pyplot as plt
8  import numpy as np
9  import statistics
10 from matplotlib.pyplot import figure
11 import matplotlib.mlab as mlab
12
13
14 class DataAnalysis2D(object):
15
16     @staticmethod
17     def skewness(var, bins=None):
18         """
19         TODO: Complete skewness doc-strings.
20         :param var:
21         :param bins:
22         :return:
23         """
24         return kth_order_moment_about_mean(var, 3, bins=bins) / np.power(
            np.std(var), 3)
25

```

```

26     @staticmethod
27     def kurtosis_old(var, bins=None):
28         """
29         TODO: Complete kurtosis_old doc-strings.
30         :param var:
31         :param bins:
32         :return:
33         """
34         return kth_order_moment_about_mean(var, 4, bins=bins) / np.power(
35             np.std(var), 4)
36
37     @staticmethod
38     def kurtosis_new(var, bins=None):
39         """
40         TODO: Complete kurtosis_new doc-strings.
41         :param var:
42         :param bins:
43         :return:
44         """
45         return kth_order_moment_about_mean(var, 4, bins=bins) / np.power(
46             np.std(var), 4) - 3
47
48     def __init__(self, design_model: LinearDesignModel, x: np.ndarray, y:
49         np.ndarray, Py: np.ndarray = None):
50         """
51         :param design_model: Design model to be analysed (
52             LinearDesignModel/NonLinearDesignModel)
53         :param x: Independent variable associated with observations (np.
54             ndarray)
55         :param y: Vector of observations (np.ndarray)
56         """
57         self._x = x
58         self._y = y
59         self._Py = Py
60         self._design_model = design_model
61
62     @property
63     def Py(self):
64         try:
65             return self._Py
66         except AttributeError:
67             raise AttributeError("Matrix Py has not been provided, please
68                 set using Py setter.")
69
70     @Py.setter
71     def Py(self, arg):
72         self._Py = arg
73
74     def unweighted_least_squares(self):
75         """
76         :return: Unweighted least squares solution of the vector of
77             parameters (np.ndarray)
78         """
79         return unweighted_least_squares(self._design_model.
80             information_matrix(self._x), self._y)
81
82     def unweighted_prediction(self, x):
83         """
84         :param x: Independent variable associated with observations (np.
85             ndarray)
86         :return: Unweighted predictions using the design model (np.ndarray
87             )
88         """
89         return np.matmul(self._design_model.information_matrix(x), self.

```

```

unweighted_least_squares())
80
81 def unweighted_residuals(self):
82     """
83     :return: Unweighted residuals between observations and unweighted
            prediction (np.ndarray)
84     """
85     return self._y - self.unweighted_prediction(self._x)
86
87 def weighted_least_squares(self):
88     """
89     :return: Unweighted least squares solution of the vector of
            parameters (np.ndarray)
90     """
91     return weighted_least_squares(self._design_model.
            information_matrix(self._x), self._y, Py=self.Py)
92
93 def weighted_prediction(self, x):
94     """
95     :param x: Independent variable associated with observations (np.
            ndarray)
96     :return: Weighted predictions using the design model (np.ndarray)
97     """
98     return np.matmul(self.weighted_least_squares(), self._design_model
            .information_matrix(x))
99
100 def weighted_residuals(self):
101     """
102     :return: Weighted residuals between observations and weighted
            prediction (np.ndarray)
103     """
104     return self._y - self.weighted_prediction(self._x)
105
106 def plot_model(self, plot_type='matplotlib', show_save=('show'), title
= None, ylabel=None, xlabel=None, name=None,
107               type='unweighted', legend=True):
108     """
109     TODO: (*) Finish doc-strings for plot_model.
110     :param plot_type:
111     :param show_save:
112     :param title:
113     :param ylabel:
114     :param xlabel:
115     :param name:
116     :param type:
117     :return:
118     """
119     if plot_type is 'matplotlib':
120         figure(num=None, figsize=(8, 5), dpi=300, facecolor='w',
            edgecolor='k')
121         plt.plot(self._x, self._y, label='Raw data', linewidth=1)
122         if title:
123             plt.title(title)
124         plt.grid()
125         smoothed_x = np.linspace(self._x[0], self._x[-1], 1000)
126         if type is 'weighted':
127             predicted_y = np.matmul(self._design_model.
                information_matrix(smoothed_x),
128                                     self.weighted_least_squares())
129         elif type is 'unweighted':
130             predicted_y = np.matmul(self._design_model.
                information_matrix(smoothed_x),
131                                     self.unweighted_least_squares())
132         else:

```

```

133         raise SystemError("{} type not recognised, please use <
           weighted> or <unweighted>.".format(type))
134     plt.plot(smoothed_x, predicted_y,
135              label=self._design_model._latex__() + ' ' + str(type
           ),
136              linestyle='-',
137              linewidth=2,
138              color='red')
139     if legend:
140         plt.legend()
141     plt.ylabel(ylabel)
142     plt.xlabel(xlabel)
143     if 'show' in show_save:
144         plt.show()
145     if 'save' in show_save:
146         plt.savefig('plots/' + str(name) + '.png', bbox_inches='
           tight')
147
148     def plot_residual_hist(self, plot_type='matplotlib', show_save=('show
           '), title=None, ylabel=None, xlabel=None,
149                           name=None, bins=30, type='unweighted'):
150         '''
151         TODO: (*) Finish doc-strings for plot_residual_hist.
152         :param plot_type:
153         :param show_save:
154         :param title:
155         :param ylabel:
156         :param xlabel:
157         :param name:
158         :param bins:
159         :return:
160         '''
161         if type is 'unweighted':
162             residuals = self.unweighted_residuals()
163         elif type is 'weighted':
164             residuals = self.weighted_residuals()
165         else:
166             raise SystemError("{} type argument not recognised. Please use
           weighted or unweighted".format(type))
167         mu = statistics.mean(residuals)
168         sigma = statistics.stdev(residuals)
169         x = np.linspace(mu - 3 * sigma, mu + 3 * sigma, 100)
170         if plot_type is 'matplotlib':
171             figure(num=None, figsize=(8, 5), dpi=300, facecolor='w',
           edgecolor='k')
172             if title:
173                 plt.title(title)
174             plt.grid()
175             plt.hist(residuals,
176                      bins=bins,
177                      density=True,
178                      label='Normalised histogram ({0} bins)'.format(bins),
179                      edgecolor='black',
180                      linewidth=0.8)
181             plt.plot(x, mlab.normpdf(x, mu, sigma),
182                      label='Normal Distribution ({0:sigma={0:.3g}}'.format(
           sigma) + ',\;\mu={0:.3g}$)'.format(mu),
183                      color='red', linewidth='2')
184
185         # DEVELOPMENT FOR ANOMALIES
186         # p, ed = np.histogram(residuals, bins=30, density=True)
187         # mid = [np.mean([ed[i], ed[i + 1]]) for i in range(len(ed) -
           1)]
188         # plt.scatter(mid[0], p[0], s=400, facecolors='none',

```

```

189         edgecolors='g', linewidths=3, label='Note A')
190     plt.legend()
191     plt.ylabel(ylabel)
192     plt.xlabel(xlabel)
193     if 'show' in show_save:
194         plt.show()
195     if 'save' in show_save:
196         plt.savefig('plots/' + str(name) + '.png', bbox_inches='
            tight')

```

#### D. gmst\_handler.py

```

1 import numpy as np
2 import pandas as pd
3
4 GMSL_HDR = 50
5 GMSL_COLUMNS = ['altt', 'mfc', 'year', 'n_obs', 'n_wobs', 'gmsl', '
    gmsl_std', 'smoothed', 'gmsl_gia', 'gmsl_gia_std',
6                 'gmsl_gia_smthed_20', 'gmsl_gia_smthed_sig']
7 GMSL_FILENAME = 'example_dataset/GMSL_TPJAOS_4.2_199209_201805.txt'
8
9
10 def txt_to_array(relative_path, skip_header):
11     return np.genfromtxt(relative_path, skip_header=skip_header)
12
13
14 def array_to_dataframe(array, columns):
15     return pd.DataFrame(array, columns=columns)

```

#### E. example\_analysis.py

```

1 from example.gmst_handler import *
2 from model_design import LinearDesignModel
3 from data_analysis import DataAnalysis2D
4 import scipy.fftpack
5 from scipy.signal import argrelextrema
6 import matplotlib.pyplot as plt
7 from matplotlib.pyplot import figure
8
9 if __name__ == '__main__':
10     data_array = txt_to_array(GMSL_FILENAME, skip_header=GMSL_HDR)
11     data_frame = array_to_dataframe(data_array, columns=GMSL_COLUMNS)
12     refined_data_frame = data_frame[['year', 'gmsl']]
13
14     #####
15
16     # MODEL 1
17     #####
18
19     # model 2 involving a trend, bias and signal
20     model_1 = LinearDesignModel(basis_functions=['x^0', 'x'])
21
22     # instantiate a data analysis model with the given data and linear
23     model
24     model_1_analysis = DataAnalysis2D(model_1, refined_data_frame['year'].
25         values, refined_data_frame['gmsl'].values)
26
27     # plot the trend + bias model for the data set
28     model_1_analysis.plot_model(
29         ylabel='GMSL variation (mm) w.r.t. 20 year mean reference',
30         xlabel='Year',
31         show_save=('save'),
32         name='model1_plot')
33
34     # plot the residual probability density function

```

```

31 model_1_analysis.plot_residual_hist(
32     ylabel='Probability density f($\epsilon$)',
33     xlabel='Residuals ($\epsilon$) [mm]',
34     show_save=('save'),
35     name='model1_residuals')
36
37 # print the least square solution
38 print('Vector of parameters: ', model_1_analysis.
    unweighted_least_squares())
39
40 # print the skewness of the residual distribution
41 print('Skewness: ', model_1_analysis.skewness(model_1_analysis.
    unweighted_residuals(), 30))
42
43 # print the kurtosis of the residual distribution
44 print('Kurtosis: ', model_1_analysis.kurtosis_new(model_1_analysis.
    unweighted_residuals(), 30))
45
46 #####
47 # MODEL 2
48 #####
49
50 # model 2 involving a trend, bias and signal
51 model_2 = LinearDesignModel(basis_functions=['x^0', 'x', 'cos(2*pi*x)
52     ', 'sin(2*pi*x)'])
53
54 # instantiate a data analysis model with the given data and linear
55 # model
56 model_2_analysis = DataAnalysis2D(model_2, refined_data_frame['year'].
57     values, refined_data_frame['gmsl'].values)
58
59 # plot the trend + bias model for the data set
60 model_2_analysis.plot_model(
61     ylabel='GMSL variation (mm) w.r.t. 20 year mean reference',
62     xlabel='Year',
63     show_save=('save'),
64     name='model2_plot'
65 )
66
67 # plot the residual probability density function
68 model_2_analysis.plot_residual_hist(
69     ylabel='Probability density f($\epsilon$)',
70     xlabel='Residuals ($\epsilon$) [mm]',
71     show_save=('save'),
72     name='model2_residuals'
73 )
74
75 # print the least square solution
76 print('Vector of parameters: ', model_2_analysis.
    unweighted_least_squares())
77
78 # print the skewness of the residual distribution
79 print('Skewness: ', model_2_analysis.skewness(model_2_analysis.
    unweighted_residuals(), 30))
80
81 # print the kurtosis of the residual distribution
82 print('Kurtosis: ', model_2_analysis.kurtosis_new(model_2_analysis.
    unweighted_residuals(), 30))
83
84 #####
85 # MODEL 3

```

```

83 #####
84 # model 3 involving a trend, bias and signal
85 model_3 = LinearDesignModel(basis_functions=['x^0', 'x', 'x^2', 'cos
      (2*pi*x)', 'sin(2*pi*x)'])
86
87 # instantiate a data analysis model with the given data and linear
      model
88 model_3_analysis = DataAnalysis2D(model_3, refined_data_frame['year'].
      values, refined_data_frame['gmsl'].values)
89
90 # plot the trend + bias model for the data set
91 model_3_analysis.plot_model(
92     ylabel='GMSL variation (mm) w.r.t. 20 year mean reference',
93     xlabel='Year',
94     show_save=('save'),
95     name='model3_plot'
96 )
97
98 # plot the residual probability density function
99 model_3_analysis.plot_residual_hist(
100     ylabel='Probability density f($\epsilon$)',
101     xlabel='Residuals ($\epsilon$) [mm]',
102     show_save=('save'),
103     name='model3_residuals'
104 )
105
106 # print the least square solution
107 print('Vector of parameters: ', model_3_analysis.
      unweighted_least_squares())
108
109 # print the skewness of the residual distribution
110 print('Skewness: ', model_3_analysis.skewness(model_3_analysis.
      unweighted_residuals(), 30))
111
112 # print the kurtosis of the residual distribution
113 print('Kurtosis: ', model_3_analysis.kurtosis_new(model_3_analysis.
      unweighted_residuals(), 30))
114
115 #####
116 # FOURIER ANALYSIS
117 #####
118 x = refined_data_frame['year'].values
119 y = refined_data_frame['gmsl'].values
120
121 # Number of sample points
122 N = len(x)
123
124 # sample spacing
125 T = (y[-1]-y[0]) / N
126
127 # fast fourier transform
128 yf = scipy.fftpack.fft(y)
129 xf = np.linspace(0.0, 1.0 / (2.0 * T), N / 2)
130
131 refined = 2.0 / N * np.abs(yf[:N // 2])
132 idx = argrelextrema(refined, np.greater)
133
134 figure(num=None, figsize=(8, 5), dpi=300, facecolor='w', edgecolor='k
      ')
135 plt.grid()
136 # fig, ax = plt.subplots()

```

```

137 # plt.plot(xf, np.log(2.0 / N * np.abs(yf[:N // 2])))
138 print(refined[idx][refined[idx] >= 1.0])
139 print(xf[idx][refined[idx] >= 1.0])
140
141 plt.semilogy(xf, 2.0 / N * np.abs(yf[:N // 2]))
142 plt.semilogy(xf[idx][refined[idx] >= 1.0], refined[idx][refined[idx]
143 >= 1.0], linewidth=0, marker='o')
144 plt.ylabel('Amplitude [mm]')
145 plt.xlabel('Frequency (per year)')
146 plt.savefig('plots/fourier_analysis.png')
147 # plt.show()
148 #####
149
150 # MODEL 4
151 #####
152
153 f = [0.01358823, 0.13588228, 0.19023519, 0.27176456, 0.3397057,
154 0.69299964]
155 _basis_functions = ['x^0', 'x', 'x^2', 'cos(2*pi*x)', 'sin(2*pi*x)']
156 for freq in f:
157     _basis_functions += ['cos(2*pi*x*{})'.format(str(freq)), 'sin(2*pi
158 *x*{})'.format(str(freq))]
159
160 print(_basis_functions)
161
162 # model 4 involving a trend, bias and signal
163 model_4 = LinearDesignModel(basis_functions=_basis_functions)
164
165 # instantiate a data analysis model with the given data and linear
166 model
167 model_4_analysis = DataAnalysis2D(model_4, refined_data_frame['year'].
168 values, refined_data_frame['gmsl'].values)
169
170 # plot the trend + bias model for the data set
171 model_4_analysis.plot_model(
172     ylabel='GMSL variation (mm) w.r.t. 20 year mean reference',
173     xlabel='Year',
174     show_save=('save'),
175     name='model4_plot',
176     legend=False
177 )
178
179 # plot the residual probability density function
180 model_4_analysis.plot_residual_hist(
181     ylabel='Probability density f($\epsilon$)',
182     xlabel='Residuals ($\epsilon$) [mm]',
183     show_save=('save'),
184     name='model4_residuals',
185 )
186
187 # print the least square solution
188 print('Vector of parameters: ', model_4_analysis.
189 unweighted_least_squares())
190
191 # print the skewness of the residual distribution
192 print('Skewness: ', model_4_analysis.skewness(model_4_analysis.
193 unweighted_residuals(), 30))
194
195 # print the kurtosis of the residual distribution
196 print('Kurtosis: ', model_4_analysis.kurtosis_new(model_4_analysis.
197 unweighted_residuals(), 30))

```