

Least Squares Data Fitting

Mohd Harith Akmal Zulfaizal Fadillah, Bahari Idrus,
 Mohammad Khatim Hasan

*Centre for Artificial Intelligence (CAIT),
 Faculty of Information Science & Technology,
 Universiti Kebangsaan Malaysia, 43600 Bangi, Selangor*

Abstract

The least squares method is an approach for determining the best fit for a series of data points by minimizing the sum of the squared differences between the actual values and the approximated values from the plotted curve. The least squares method is able to approximate functions such as linear, quadratic and cubic functions as well as functions of higher degree and complexity. This report will model, fit and plot the curves for linear least squares of linear, quadratic and cubic models, multiple linear regression (MLR), least squares function approximation and nonlinear least squares using Python.

1 Introduction

The least squares method is a statistical model commonly used in regression analysis in an attempt to estimate the relationships between the dependent variable and one or more independent variables. It is used to approximate the solutions to a system of equations that minimize the sum of squared residuals which is essentially the summation of the squared difference between the y -values from the data set and the approximated model [1].

Its' applications concern mainly of regression, predictions and estimations in a wide range of fields such as finance and economy, image processing, medicine, agriculture among others. The least squares approximation is applicable whenever the dependent variable is assumed to have a causal relationship with the independent variables. Furthermore, the least squares solutions can be easily obtained from matrix linear algebra which will be discussed in the next section. However, since it is inherent in the method to minimize the sum of squared residuals, the existence of outliers can affect the solutions and line of best fit.

In the following sections, we are going to study the method of least squares for some of the general models such as linear, quadratic and cubic. We are also going to look at multiple linear regression (MLR), least squares function approximation and nonlinear least squares. After studying the mathematical foundations of those methods, we are going to model, fit and plot the approximations to relevant applicable data sets using Python and observe the difference in the different ways and modules used in fitting the data.

2 The Method of Least Squares

2.1 Linear Least Squares & Regression

A set of data points consisting of independent and dependent variables, (x_i, y_i) is plotted. A model function of some sort is chosen to best fit the data set by adjusting the parameters available. For example, we can choose the simplest linear model, $y = a_0 + a_1x$ to be approximated

to the plotted data set. The fit of the model chosen can be measured by minimizing the sum of its squared residuals which is the difference between the actual data points and the prediction from the fit approximated by the model [1],

$$F(a_i) = \sum_{i=1}^n (y_i - (a_0 + a_1 x_i))^2$$

To minimize the squared residuals, the partial derivatives of $F(a, b)$ is set to zero,

$$\frac{\partial F(a_i)}{\partial a_2} = -2 \sum_{i=1}^n (y_i - (a_0 + a_1 x_i)) = 0$$

$$\frac{\partial F(a_i)}{\partial a_1} = -2 \sum_{i=1}^n (y_i - (a_0 + a_1 x_i)) x_i = 0$$

which can be simplified to a similar system of linear equations form, $\mathbf{Ax} = \mathbf{b}$,

$$a_1 \sum_{i=1}^n x_i + n a_0 = \sum_{i=1}^n y_i$$

$$a_1 \sum_{i=1}^n x_i^2 + \sum_{i=1}^n x_i a_0 = \sum_{i=1}^n x_i y_i$$

or in its matrix form,

$$\begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum y_i x_i \end{bmatrix}$$

Solving for $\mathbf{x} = A^{-1}\mathbf{b}$, the values for a_0 and a_1 can be obtained.

This method to obtain the best fit can also work for any polynomial models by simply changing the model used in the squared residuals. Consider a quadratic model, $y = a_0 + a_1 x + a_2 x^2$,

$$F(a) = \sum_{i=1}^n (y_i - (a_0 + a_1 x + a_2 x^2))^2$$

By minimizing the squared residuals, a system of linear equations can be obtained to be solved for the parameters a_0, a_1, a_2 .

$$\begin{bmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum y_i x_i \\ \sum y_i x_i^2 \end{bmatrix}$$

From the obtained system of equations from both linear and quadratic models, it is observed that there is a pattern to the matrices A and \mathbf{b} . The system of equations to measure the parameters can be generalized to,

$$\begin{bmatrix} n & \sum x_i & \sum x_i^2 & \dots & \sum x_i^m \\ \sum x_i & \sum x_i^2 & \sum x_i^3 & \dots & \sum x_i^{m+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum x_i^m & \sum x_i^{m+1} & \sum x_i^{m+2} & \dots & \sum x_i^{m+m} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_m \end{bmatrix} = \begin{bmatrix} \sum y_i x_i^0 \\ \sum y_i x_i^1 \\ \vdots \\ \sum y_i x_i^m \end{bmatrix} \quad (1)$$

It is also possible to measure the parameter from its system of equations, $\mathbf{Ax} = \mathbf{b}$ by plugging in the values of x and y into the chosen model such that,

$$\text{Eg: Cubic model} \quad y = a_0 + a_1 x + a_2 x^2 + a_3 x^3$$

$$\begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_m & x_m^2 & x_m^3 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad (2)$$

where m is the number of rows or the number of equations available from the dataset. Since the number of rows is greater than the number of columns, it often happens that there is no exact solution for $A\mathbf{x} = \mathbf{b}$ as \mathbf{b} is not in the column space of A . Another vector is chosen such that $\mathbf{e} = \mathbf{b} - \mathbf{p}$, which is the distance between \mathbf{b} and the chosen vector is minimized while \mathbf{e} is perpendicular to the column space. The best possible solution happens to be the projection of \mathbf{b} on $C(A)$ which is $\mathbf{p} = A\hat{\mathbf{x}}$ where $\hat{\mathbf{x}}$ is the least squares solution [2]. By minimizing the squared error, \mathbf{e} ,

$$\begin{aligned} F(\mathbf{x}) &= \|\mathbf{e}\|^2 = \mathbf{e}^T \mathbf{e} \\ &= (\mathbf{b} - A\mathbf{x})^T (\mathbf{b} - \mathbf{x}) \\ &= \mathbf{b}^T \mathbf{b} - 2\mathbf{b}^T A\mathbf{x} + A^T A\mathbf{x} \end{aligned}$$

and setting the gradient to zero,

$$\begin{aligned} \frac{\partial F(\mathbf{x})}{\partial \mathbf{x}} &= -2A^T \mathbf{b} + 2A^T A\mathbf{x} = 0 \\ A^T A\mathbf{x} &= A^T \mathbf{b} \end{aligned} \quad (3)$$

the least squares solution can be solved for $\hat{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{b}$ where $A^T A$ is a much more dense invertible square matrix with nonzero entities where the dot product between the Vandermonde matrix A and its transpose A^T will result in the generalized matrix A from Equation 2,

$$A^T A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ x_1 & x_2 & x_3 & x_4 \\ x_1^2 & x_2^2 & x_3^2 & x_4^2 \\ x_1^3 & x_2^3 & x_3^3 & x_4^3 \end{bmatrix} \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ 1 & x_3 & x_3^2 & x_3^3 \\ 1 & x_4 & x_4^2 & x_4^3 \end{bmatrix} = \begin{bmatrix} n & \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 & \sum x_i^5 \\ \sum x_i^3 & \sum x_i^4 & \sum x_i^5 & \sum x_i^6 \end{bmatrix}$$

2.2 Multiple Linear Regression (MLR)

Previously, we have looked into linear least squares regression consisting of only singular independent variable. MLR concerns itself in obtaining a best fit for a dataset with multiple independent variables, (x_{i1}, x_{i2}, \dots) with respect to the dependent variable, y_i given that the independent variables are not too highly correlated with each other.

A system of equations consisting of multiple independent variables can be solved by using Equation 3. Consider three independent variables of the following function,

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2^2 + \beta_3 x_1 x_2$$

The system of equations, $A\mathbf{x} = \mathbf{b}$ given that A is simply the Vandermonde matrix from Equation 2 [3],

$$\begin{bmatrix} 1 & x_{11} & x_{12}^2 & x_{11}x_{12} \\ 1 & x_{21} & x_{22}^2 & x_{21}x_{22} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{m1} & x_{m2}^2 & x_{m1}x_{m2} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

where m is the number of equations and the parameters β_i of \mathbf{x} can be solved for $\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$.

3 Least Squares Function Approximation

Given a function, $f(x)$, it can be approximated to a weighted sum of functions usually chosen to be of orthogonal functions, $\phi(x)$ where the inner product of orthogonal functions are $\langle \phi_i, \phi_j \rangle = 0$ when $i \neq j$ due to its orthogonality properties.

$$f(x) \approx \sum_{i=1}^n c_i \phi_i(x) = c_1 \phi_1(x) + c_2 \phi_2(x) + \dots + c_n \phi_n(x) = \hat{f}(x)$$

The approximation function, $\hat{f}(x)$ can best fit the original function by determining the c_i coefficients such that $\|f - \hat{f}\|^2$ is minimum.

$$\|f - \hat{f}\|^2 = \int_a^b (f(x) - \sum_{i=1}^n c_i \phi_i(x))^2 dx$$

There are generally two ways of minimizing the squared error. Firstly, by using calculus and taking the derivatives of the squared residuals with respect to c_i to get the minimum points. However, we can take advantage of the orthogonality of $\phi_i(x)$ for the minimization via linear algebra where [1],

$$\begin{aligned} \|f - \hat{f}\|^2 &= \langle f - \hat{f}, f - \hat{f} \rangle \\ &= \langle f, f \rangle - 2\langle f, \hat{f} \rangle + \langle \hat{f}, \hat{f} \rangle \\ \langle \hat{f}, \hat{f} \rangle &= \langle \sum_{i=1}^n c_i \phi_i, \sum_{j=1}^n \phi_j \rangle \\ &= c_1^2 \langle \phi_1, \phi_1 \rangle + c_2^2 \langle \phi_2, \phi_2 \rangle + \dots + c_n^2 \langle \phi_n, \phi_n \rangle \\ &= \sum_{i=1}^n c_i^2 \|\phi_i\|^2 \\ \langle f, \hat{f} \rangle &= \langle f, \sum_{i=1}^n \phi_i \rangle \\ &= \langle f, c_1 \phi_1 + c_2 \phi_2 + \dots + c_n \phi_n \rangle \\ &= c_1 \langle f, \phi_1 \rangle + c_2 \langle f, \phi_2 \rangle + \dots + c_n \langle f, \phi_n \rangle \\ &= \sum_{i=1}^n c_i \langle f, \phi_i \rangle \end{aligned}$$

The critical value of the squared error, $E(c_i)$ can then be computed for a particular coefficient, c_k ,

$$\begin{aligned} \|f - \hat{f}\|^2 &= \|f\|^2 - 2 \sum_{i=1}^n c_i \langle f, \phi_i \rangle + \sum_{i=1}^n c_i^2 \langle \phi_i, \phi_i \rangle = E(c_i) \\ \frac{\partial E}{\partial c_k} &= \frac{\partial}{\partial c_k} \|f\|^2 - 2 \sum_{i=1}^n \frac{\partial}{\partial c_k} c_i \langle f, \phi_i \rangle + \sum_{i=1}^n \frac{\partial}{\partial c_k} c_i^2 \langle \phi_i, \phi_i \rangle \\ &= -2 \langle f, \phi_k \rangle + 2c_k \langle \phi_k, \phi_k \rangle = 0 \\ c_k &= \frac{\langle f, \phi_k \rangle}{\langle \phi_k, \phi_k \rangle} \end{aligned} \quad (4)$$

Thus, the best approximation can be calculated for any given $f(x)$ with any set of orthogonal functions as [1],

$$\hat{f}(x) = \sum_{i=1}^n \frac{\langle f, \phi_i \rangle}{\langle \phi_i, \phi_i \rangle} \phi_i(x) \quad (5)$$

4 Nonlinear Least Squares

A function, $f(x)$ is considered nonlinear when it cannot be expressed as the linear combination of its coefficients such it is nonlinear in its coefficients. A function such as $y = be^{ax}$ is a nonlinear equation that can be linearized by taking the natural logarithm of both sides of the equation, $\ln y = \ln b + ax$ and is able to be approximated by ordinary least squares.

When a model chosen to be approximated to a set of data is a nonlinear function $f(x, \beta)$ that cannot be linearized, the ordinary least squares analysis does not hold true for the model. The Gauss-Newton approach is used for a nonlinear least squares analysis where the model is approximated by a linear function, commonly the first-order Taylor series [4],

$$r_i = y_i - f(x, \beta)$$

$$r_i \approx y_i - f(x_i, \beta_t) - \nabla_{\beta} f(x_i, \beta_t)[\beta - \beta_t]$$

we can redefine the Taylor's approximation to a linear function by,

$$\tilde{y}_i = y_i - f(x_i, \beta_t) \quad \tilde{x}_i = \nabla_{\beta} f(x_i, \beta_t) \quad \tilde{\beta} = \beta - \beta_t$$

$$r_i \approx \tilde{y}_i - \tilde{x}_i^T \tilde{\beta} = \tilde{y}_i - J\tilde{\beta}$$

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial \beta_1} & \cdots & \frac{\partial f_1}{\partial \beta_p} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial \beta_1} & \cdots & \frac{\partial f_n}{\partial \beta_p} \end{bmatrix}$$

where $\tilde{x}_i = J$ is the Jacobian matrix. Now that it is in the form of an ordinary least square with the solution $\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$, we can apply this solution to $\|r_i\|^2 = \|\tilde{y} - J\tilde{\beta}\|^2$ and substitute back \tilde{y} and $\tilde{\beta}$.

$$\begin{aligned} \tilde{\beta} &= (J_t^T J_t)^{-1} J_t^T \tilde{y} = (J_t^T J_t)^{-1} J_t^T r_t \\ \beta_{t+1} &= \beta_t + (J_t^T J_t)^{-1} J_t^T r_t \end{aligned} \quad (6)$$

The solutions for the regression coefficients are obtained and refined by successive approximations until the desired tolerance is achieved.

5 Modelling and Fitting

In this section, we are going to attempt to do least squares regression analysis on several data sets for linear, quadratic and cubic functions. In general, there are several ways to obtain the desired coefficients and best fit of the model. We are mainly going to compute and plot these parameters by using,

- i the generalized matrix from Equation 1,
- ii the Vandermonde matrix from Equation 2,
- iii built-in SciPy modules for least squares and curve fitting.

We are also going to model, fit and plot the approximated relevant functions to series of data points for multiple linear regression (MLR), least squares function approximation and nonlinear least squares.

5.1 Linear model

A simple linear function, $y = a_0 + a_1x$ is used to model a linear best fit for a set of data, x and y . The data set for this example is generated by using NumPy's random module to give the data some sort of a random noise.

A function is needed to sort the data set into the generalized matrix from Equation 1 before solving the system of equations. To obtain the generalized matrix, the arrays of x and y are passed through a function that sorts the matrix elements iteratively.

```

1 def get_gmatrices(order, x, y, n):
2     z = order + 1
3     A = np.empty(shape=(z, z))
4     for i in range(z):
5         itr = 0
6         for j in range(z):
7             if i == 0 and j == 0:
8                 A[i, j] = n
9                 itr += 1
10            else:
11                A[i, j] = sum(x ** (i + itr))
12                itr += 1
13        b = np.empty(shape=(z, 1))
14        b[0, 0] = sum(y)
15        for k in range(1, z):
16            b[k, 0] = sum(y * (x ** k))
17    return A, b

```

The function above will give output to the matrix A and \mathbf{b} which can be solved for $\mathbf{x} = A^{-1}\mathbf{b}$ to obtain the coefficients of the linear model, a_0 and a_1 , using `np.linalg.solve(A, b)`. The linear best fit is plotted with the coefficients measured in Figure 1. The sum of squared residuals as per previous section is given by $\sum(y - \hat{y})^2$ where \hat{y} is the approximated y values.

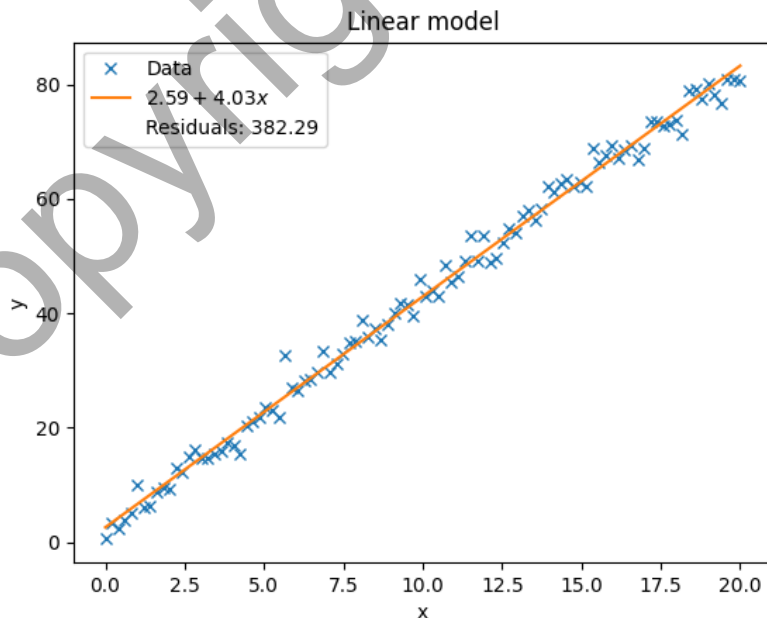


Figure 1: Linear model

5.2 Quadratic model

In addition to the generalized matrix, the coefficients for the best fit can also be obtained by using the Vandermonde matrix. For the quadratic model, the simplest quadratic function $y = a_0 + a_1x + a_2x^2$ is used. The data set for this example is generated via NumPy's random function. Sorting the system of equations in the form of $A\mathbf{x} = \mathbf{b}$, A is a 100×3 matrix where $m > n$ thus rendering solving the system impossible since there are no real solutions.

From the derivation of Equation 3, it is observed that by applying a dot product of the transpose A^T to matrix A , a much denser square matrix is obtained which will make solving the system possible. Now the system of equations can be represented as,

$$A^T A \mathbf{x} = A^T \mathbf{b}$$

A function that takes the values of x and y is constructed such that it sorts the former into a Vandermonde matrix, A and the latter into a column matrix of y values for all the data points.

```

1 def get_lamatrices(order, x, y, n):
2     power = list(range(order + 1))
3     M = x[:, np.newaxis] ** power # [1 x x^2 x^3 ...] form
4     k = y.reshape(n, 1)
5
6     return M, k

```

The function `np.newaxis` is extremely useful in constructing the Vandermonde matrix as it adds a new dimension to the array and is able to transform the elements added as seen above where every new column or dimension, the elements of said column is raised to a power. The system of equations can be solved for, in this case, $\mathbf{x} = (M^T M)^{-1} M^T \mathbf{k}$, where the coefficients are obtained and the quadratic best fit is plotted as seen in Figure 2.

```

1 c = np.linalg.inv(M.T@M)@M.T@k

```

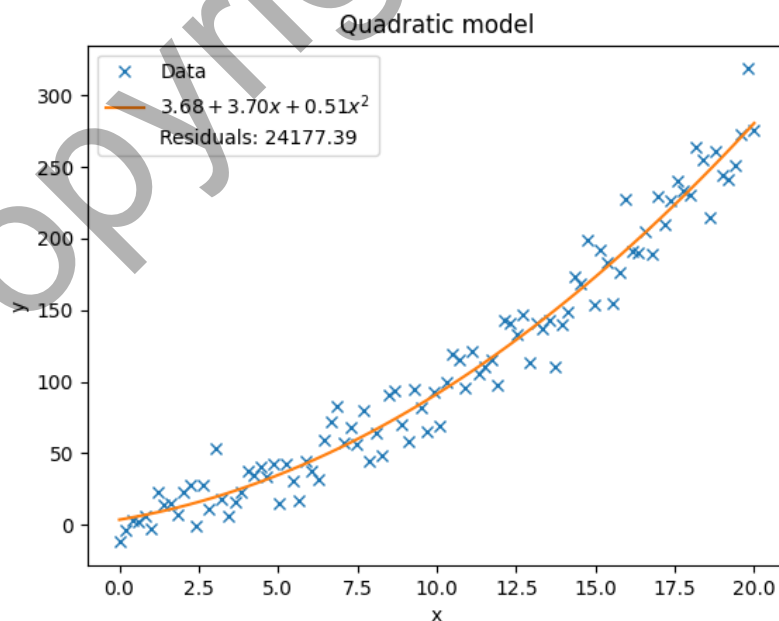


Figure 2: Quadratic model

5.3 Cubic model

As we have seen both ways of solving the system of equations to measure the coefficients of the selected model through the generalized and Vandermonde matrices, it is only apt for us to utilize the already existing algorithms built into the NumPy and SciPy module.

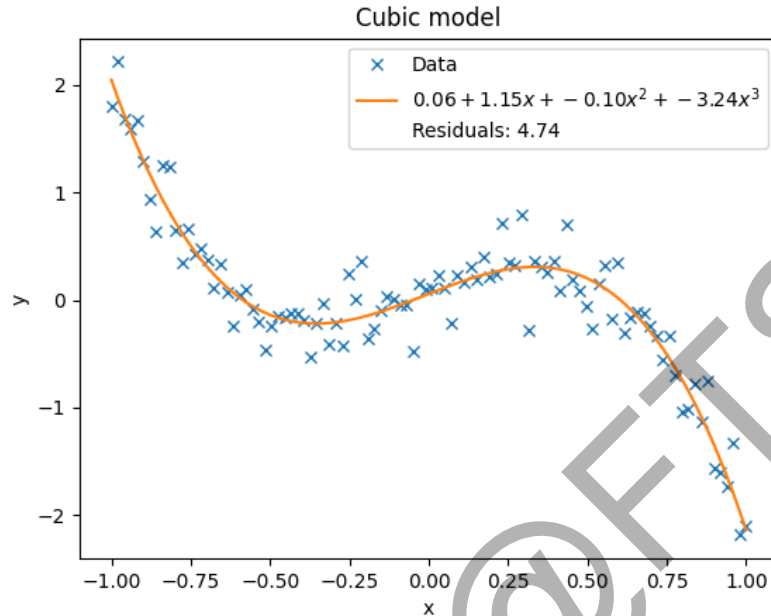


Figure 3: Cubic model using `scipy.linalg.lstsq`

The first method is `np.linalg.lstsq` and `scipy.linalg.lstsq` where it takes the Vandermonde matrix A and \mathbf{b} as its argument and yields the least squares solutions, its sum of squared residuals, rank of matrix A and singular values of A [5].

```
1 from scipy.linalg import lstsq
2
3 z, res, rank, s = lstsq(M, k)
4 print('Coefficients:', z.flatten())
5 print('Residuals', res)
```

The second method, `scipy.optimize.curve_fit` is advanced such that it uses non-linear least squares to fit a desired function to data [5, 6]. The method takes the desired function or model to be approximated and the arrays of values for the data set, x and y with some other optional arguments such as initial guesses, uncertainty for y and bounds. The method yields the optimal values of the coefficients for the model that minimized the squared residuals and the estimated covariance for the coefficients.

```
1 def func(x, a, b, c, d):
2     return a*x**3 + b*x**2 + c*x + d
3
4 x = np.linspace(0, 20, 100)
5 y = func(x, 3, 1, 4, 9) + np.random.normal(0, 1000, x.shape)
6
7 coeff, cov = scipy.optimize.curve_fit(func, x, y)
```

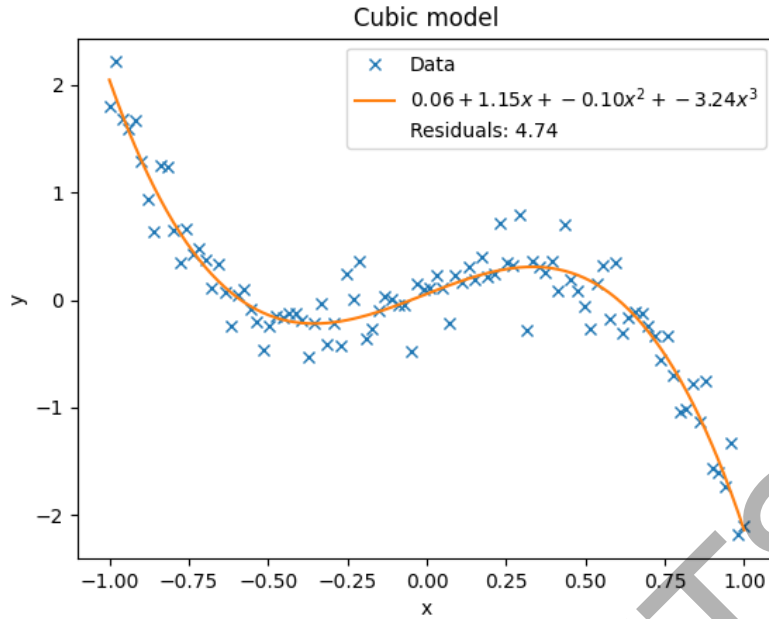



Figure 4: Cubic model using `scipy.optimize.curve_fit`

There are some other implementations that could be of use in regards to curve fitting via least squares such as NumPy's `polyfit`, SciPy's `least_squares`, `leastsq` and scikit-learn's `linear_model.LinearRegression`.

All of the models used above are polynomial models of degree 1, 2 and 3 respectively. A general polynomial regression Python code that utilizes the three aforementioned methods — the generalized matrix, the Vandermonde matrix and built-in SciPy modules — is attached as Appendix A and another Python code that uses the module scikit-learn is available in Appendix B.

5.4 Example: Global annual mean temperature

A data set of the global annual mean temperature anomalies over the course of 1880-2020 is obtained from NASA [7, 8]. By solving for the coefficients using Equation 3, the least squares polynomial regressions of degree 1, 2 and 3 are calculated and plotted against each other.

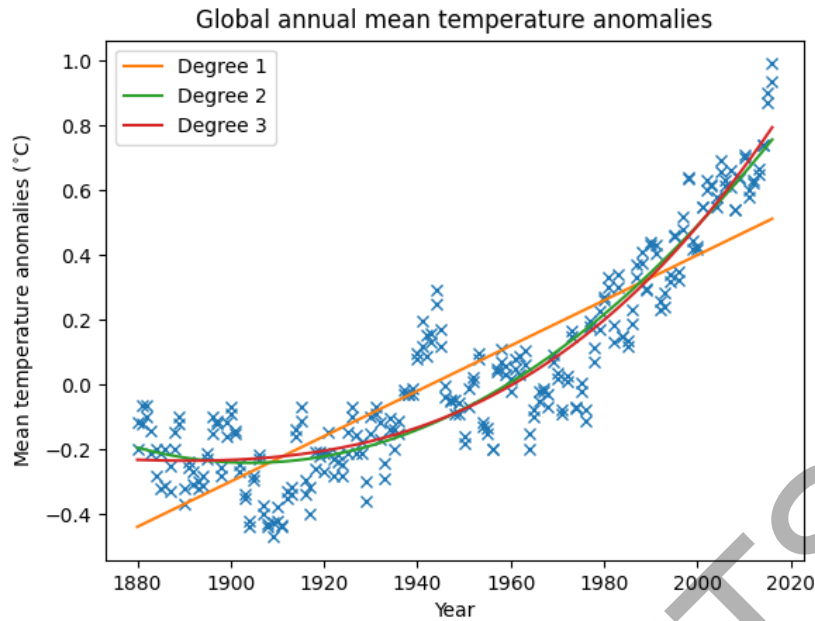


Figure 5: Polynomial models of degree 1, 2 and 3

The sum of squared residuals for degree 1, 2 and 3 are measured to be 7.0454, 3.6274 and 3.5670 respectively. The smallest sum of squared residuals is obtained when the cubic model is used for regression. By utilizing every methods and modules to obtain the least squares solutions as shown previously, it is observed that there is no significant discrepancy between the solutions for the cubic model.

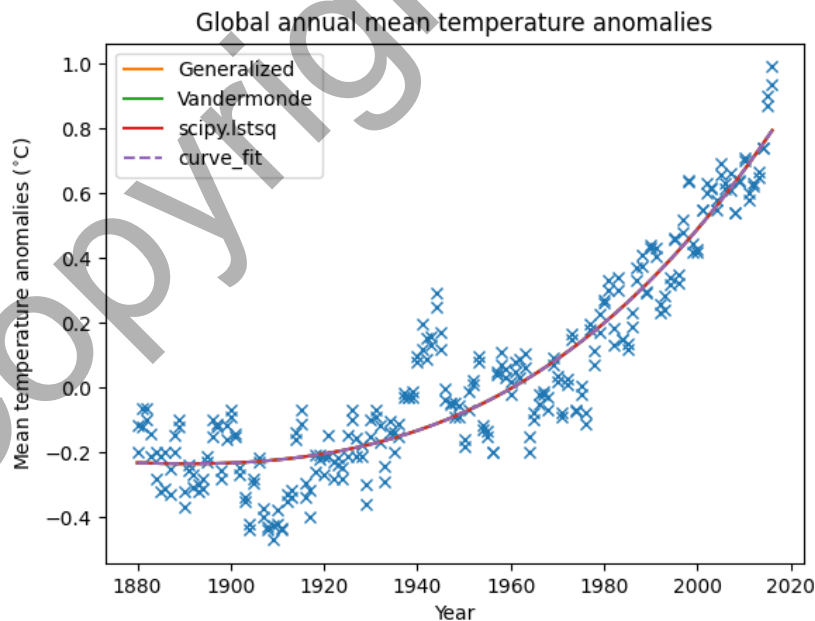


Figure 6: Cubic model

```

1 >> [-1.96857536e+03  3.17484235e+00 -1.70649924e-03  3.05672104e-07]
      Generalized
2 >> Residuals: 3.5670440372642203

```

```

3 >> [-1.96853127e+03  3.17477441e+00 -1.70646436e-03  3.05666135e-07]
      Vandermonde
4 >> Residuals: 3.5670440371850045
5 >> [-1.96847245e+03  3.17468377e+00 -1.70641782e-03  3.05658170e-07] lstsq
6 >> Residuals: 3.5670440470114393 [3.56704404]
7 >> [-1.97014946e+03  3.17726763e+00 -1.70774452e-03  3.05885183e-07]
      curve_fit
8 >> Residuals: 3.567044070463524

```

5.5 Multiple linear regression

Given a data set where there are multiple independent variables, x , the best fit can be obtained by solving a system of equation in the form of Equation 3 where A is the Vandermonde matrix and the elements is sorted based on the chosen model.

A data set of the US economy from 1976 to 1987 has several interesting variables and indicators such as the gross national product (GNP), the purchasing power of US dollar, the price of crude oil per barrel, the amount of foreign investments among others. The first two indicators, GNP (x_1) and purchasing power (x_2) are chosen against the consumer debt (in billions) as the dependent variable y and a linear model is chosen,

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

For this least squares multiple linear regression computation, scikit-learn's `linear_model.LinearRegression` method is used for the computation of the regression coefficients. The matrix A is constructed by stacking a column matrix of 1s and the respective column values of x_1 and x_2 from the dataframe [9, 10].

```

1 X = df[['PURCHASE', 'GNP']].values.reshape(-1,2)
2 y = df['CONSUMER'].values
3
4 x1 = X[:, 0]
5 x2 = X[:, 1]
6
7 x1_pred = np.linspace(min(x1), max(x1), 100)
8 x2_pred = np.linspace(min(x2), max(x2), 100)
9 x1m_pred, x2m_pred = np.meshgrid(x1_pred, x2_pred)
10 modelxs = np.array([x1m_pred.flatten(), x2m_pred.flatten()]).T
11
12 ols = linear_model.LinearRegression()
13 model = ols.fit(X, y)
14 print('c', model.coef_, model.intercept_)
15 predicted = model.predict(modelxs) # for the best fit plotting

```

The best fit is plotted where the dependent variable, y is on the z-axis and the two independent variables on x and y-axis. Based on the result of the fit, the linear regression model obtained is,

$$y = -7.855 \times 10^2 + (3.361 \times 10^2) x_1 + (2.618 \times 10^{-1}) x_2$$

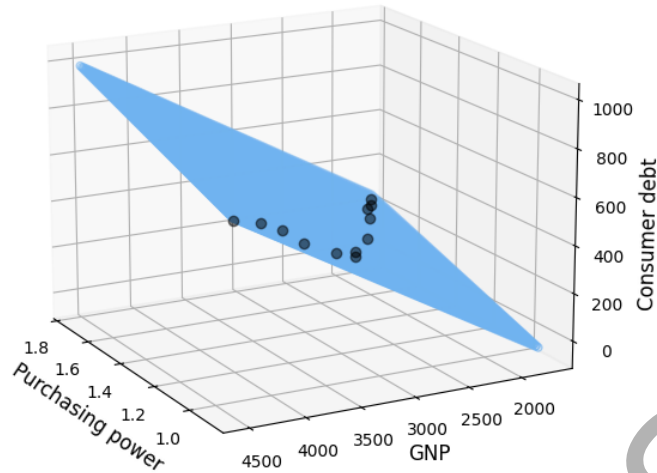


Figure 7: Gross national product (GNP) and purchasing power of US dollar against consumer debt (in billions) between 1976-1987 in USA.

As the model contains more independent variables, it will be difficult to visualize the multidimensional plot for the best fit regression but it is entirely possible to measure its regression coefficients. However, it is good practice to compute the correlation between the independent variables as to avoid multicollinearity which can introduce instability in the regression coefficients. The values for some of the collinear coefficients may be unreliable, but it is still useful for prediction of that particular model. It is a concern when we are trying to deduce a conclusion of any individual independent variable on the dependent variable [9].

5.6 Function approximation of $\sin \pi x$

Given a function $f(x) = \sin \pi x$ over $[-1, 1]$, it is possible to approximate a sum of functions, $\hat{f}(x)$ that would yield a close match to the original function $f(x)$ where,

$$\hat{f}(x) = \sum_{i=1}^n c_i \phi_i(x) = \hat{f}(x) = \sum_{i=1}^n \frac{\langle f, \phi_i \rangle}{\langle \phi_i, \phi_i \rangle} \phi_i(x)$$

The function, $\phi_i(x)$ is usually chosen to be of orthogonal functions in which case, the Legendre polynomials are chosen here [1]. Since a sinusoidal over $[-1, 1]$ resembles a cubic function, it is highly possible that a combination of the Legendre polynomials of the third order would result in a cubic function would minimize the sum of squared residuals. Hence, $\hat{f}(x)$ is chosen to be iterated until the third order of the Legendre polynomials $\phi_i(x) = P_i(x)$ for $i = 0, 1, 2, 3$ which are readily available as functions in a SciPy method, `scipy.special.legendre`.

$$P_0(x) = 1 \quad P_1(x) = x \quad P_2(x) = \frac{1}{2}(3x^2 - 1) \quad P_3(x) = \frac{1}{2}(5x^3 - 3x)$$

To calculate for the coefficient, c_i , the inner product of $\langle f, \phi_i \rangle$ and $\langle \phi_i, \phi_i \rangle$ can be measured by using the integration method available through SciPy, `scipy.integrate.quad` [5].

```
1 func = lambda x: np.sin(np.pi*x)
2 func2 = lambda x: 0.5(3*x**2-1)
3 func_product = lambda x: func(x)*func2(x) # Multiplying two functions
```

```

4 a, b = -1, 1
5
6 integral_f = quad(func, a, b)
7 integral_comb = quad(func_product, a, b)

```

After being able to measure all the coefficients, the approximated function $\hat{f}(x)$ can be obtained by summing all the product of coefficients and Legendre polynomials and plotted against the actual original function.

```

1 def y(order, x, coeff):
2     polyn = 0
3     for i in range(order + 1):
4         polyn += coeff[i] * legendre(i)(x)
5     return polyn

```

$$\hat{f}(x) = (0.955)x + (4.213 \times 10^{-18})x^2 + (-1.158)x^3$$

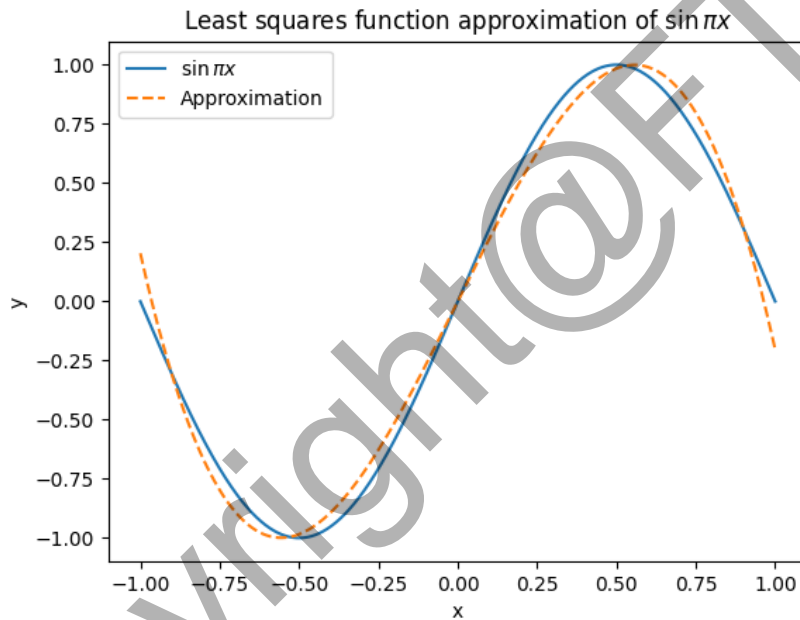


Figure 8: $\sin \pi x$ approximation using Legendre polynomials

5.7 Nonlinear regression

Given a nonlinear function to be approximated with,

$$f(x) = \frac{ax}{b + e^x}$$

there are three components of importance, β_t , J_t and r_t , that are needed for the optimal least squares solutions,

$$\beta_{t+1} = \beta_t + (J_t^T J_t)^{-1} J_t^T r_t$$

Since solving a nonlinear least squares regression is done iteratively, an initial guess condition for β needs to be defined along with a function that computes the Jacobian matrix and the column matrix of the residuals r_t .

The Jacobian matrix is essentially a matrix of partial derivatives of the function with respect to each of the coefficients. These partial derivatives can be calculated numerically by using the finite difference method where h is a small change in x [1],

$$\frac{\partial f}{\partial a} = \frac{f(x+h) - f(x-h)}{2h}$$

```

1 def Jacobian(f, x, a, b):
2     h = 1e-6
3     pdv_a = (f(x, a+h, b)-f(x, a-h, b))/(2*h)
4     pdv_b = (f(x, a, b+h)-f(x, a, b-h))/(2*h)
5     return np.column_stack([pdv_a, pdv_b])

```

The Gauss-Newton equation can be made into a function that takes the nonlinear function, the values of the data set, the initial guess values for a and b and the number of iterations.

```

1 def GaussNewton(f, x, y, a0, b0, iterations): # Gauss-Newton
2     tol = 1e-12
3     g_i = g = np.array([a0, b0])
4     for itr in range(iterations):
5         g = g_i
6         J = Jacobian(f, x, g[0], g[1])
7         r = y - f(x, g[0], g[1])
8         g_i = g + np.linalg.inv(J.T@J)@J.T@r
9         if np.linalg.norm(g-g_i) < tol:
10            print(np.linalg.norm(g-g_i), 'test')
11            break
12    return g

```

By calling the Gauss-Newton function for a certain data set and guess values, it will return the values of the coefficients that can be plotted.

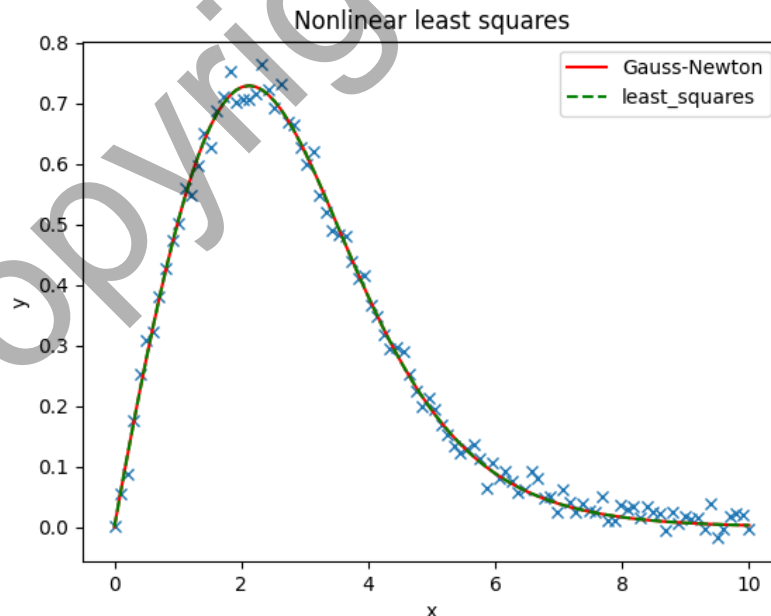


Figure 9: Nonlinear least squares of function $f(x) = \frac{ax}{b+e^x}$.

6 Conclusion

The least squares method is one of the most widely used and reliable approach in regression analysis with the ability to approximate data set with a wide range of functions. In this report, we have shown several models that can be approximated to such as the linear, quadratic and cubic models. The least square solutions to these models have been computed using the various methods we have discussed and shown no significant variation of the coefficients and the residuals between the methods – generalized matrix, Vandermonde matrix and built-in modules. We have also studied, modelled and fitted other types of regression such as multiple linear regression (MLR), least squares function approximation and nonlinear regression. This work is part of the underlying study for the project ” *Implementing quantum algorithm for least squares data fitting through IBM-Q*”.

Acknowledgements

This work is related to the ongoing research by Ts. Dr. Bahari Idrus and Assoc. Professor. Dr. Mohammad Khatim Hasan funded by GUP-2020-061.

References

- [1] J.F. Epperson. *An Introduction to Numerical Methods and Analysis*. Wiley, 2013. ISBN: 9781118367599. URL: <https://books.google.com.my/books?id=3101AgAAQBAJ>.
- [2] Gilbert Strang. *Introduction to Linear Algebra*. Fourth. Wellesley, MA: Wellesley-Cambridge Press, 2009. ISBN: 9780980232714 0980232716 9780980232721 0980232724 9788175968110 8175968117.
- [3] Nathaniel E. Helwig. *Multiple Linear Regression*. Jan. 2017.
- [4] Geof H. Givens and Jennifer A. Hoeting. *Computational Statistics*. eng. 2nd ed. Hoboken, N.J: Wiley, 2013. ISBN: 9780470533314.
- [5] Eric Jones, Travis Oliphant, Pearu Peterson, et al. *SciPy: Open source scientific tools for Python*. 2001–. URL: <http://www.scipy.org/>.
- [6] Eric Ayars. *Computational Physics With Python*. California State University, 2013.
- [7] . *GISS Surface Temperature Analysis (GISTEMP), version 4*. publisher: NASA Goddard Institute for Space Studies. 2021. URL: <https://data.giss.nasa.gov/gistemp/>.
- [8] Nathan J. L. Lenssen et al. “Improvements in the GISTEMP Uncertainty Model”. en. In: *Journal of Geophysical Research: Atmospheres* 124.12 (June 2019), pp. 6307–6326. ISSN: 2169-897X, 2169-8996. DOI: 10.1029/2018JD029522. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1029/2018JD029522> (visited on 04/14/2021).
- [9] Eric Kim. *Multiple Linear Regression and Visualization in Python*. Nov. 2019. URL: https://aegis4048.github.io/mutiple_linear_regression_and_visualization_in_python.
- [10] F. Pedregosa et al. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.

A Polynomial regression

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 from scipy.linalg import lstsq
5 from scipy.optimize import curve_fit
6
7 def func(order, x, *a):
8     polyn = 0
9     i = 0
10    for p in range(order + 1):
11        polyn += a[i] * x ** p
12        i += 1
13    return polyn
14
15 def get_lamatrices(order, x, y, n):
16     power = list(range(order + 1))
17     M = x[:, np.newaxis] ** power # [1 x x^2 x^3 ...] form
18     k = y.reshape(n, 1)
19     return M, k
20
21 def get_gmatrices(order, x, y, n):
22     z = order + 1
23     A = np.empty(shape=(z, z))
24     for i in range(z):
25         itr = 0
26         for j in range(z):
27             if i == 0 and j == 0:
28                 A[i, j] = n
29                 itr += 1
30             else:
31                 A[i, j] = sum(x ** (i + itr)) # Generalized matrix form for
32                 A
33                 itr += 1
34
35     b = np.empty(shape=(z, 1))
36     b[0, 0] = sum(y)
37     for k in range(1, z):
38         b[k, 0] = sum(y * (x ** k))
39
40     return A, b
41
42 def get_error(y, fh):
43     return np.sum((y - fh) ** 2)
44
45 df = pd.read_csv('annualglobaltemp.csv')
46 xdata = df[['Year']].to_numpy().reshape(-1).T # Down one dimension,
47         transpose
48 ydata = df[['Mean']].to_numpy().reshape(-1).T
49
50 plt.plot(df['Year'], df['Mean'], 'x')
51
52 n = xdata.size
53 x, y = xdata.astype('longdouble'), ydata
54 xlin = np.linspace(min(x), max(x), 1500)
55 orderlst = [3]
56
57 def fcv(x, *s):
58     return s[0] + s[1]*x + s[2]*x**2 + s[3]*x**3

```



```

57
58 for order in orderlst:
59     A, b = get_lamatrices(order, x, y, n)
60     M, k = get_gmatrices(order, x, y, n)
61
62     d = np.linalg.solve(M, k)
63     d_err = get_error(y, func(order, x, *d.flatten()))
64     print(d.flatten(), 'Generalized')
65     print('Residual:', d_err)
66
67     c = np.linalg.inv(A.T @ A) @ A.T @ b
68     print(c.flatten(), 'Vandermonde')
69     print('Residuals:', get_error(y, func(order, x, *c.flatten())))
70
71     z, res, rank, s = lstsq(A, b) # SVD decomposition
72     print(z.flatten(), 'lstsq')
73     z_err = get_error(y, func(order, x, *z.flatten()))
74     print('Residuals:', z_err, res)
75
76     p, _ = curve_fit(fcv, x, y, p0=[1, 1, 1, 1])
77     y_cv = fcv(xlin, *p)
78     print(p, 'curve_fit')
79     print('Residuals:', get_error(y, fcv(x, *p)))
80
81     #print(np.linalg.norm(y-func(order, x, *c.flatten())))
82
83     plt.plot(xlin, func(order, xlin, *d.flatten()), label=r'Generalized')
84     plt.plot(xlin, func(order, xlin, *c.flatten()), label=f'Vandermonde')
85     plt.plot(xlin, func(order, xlin, *z.flatten()), label=r'scipy.lstsq')
86     plt.plot(xlin, y_cv, '--', label=r'curve_fit')
87
88 plt.legend()
89 plt.xlabel('Year')
90 plt.ylabel('Mean temperature anomalies ( $^{\circ}$ C)')
91 plt.title('Global annual mean temperature anomalies')
92 plt.show()

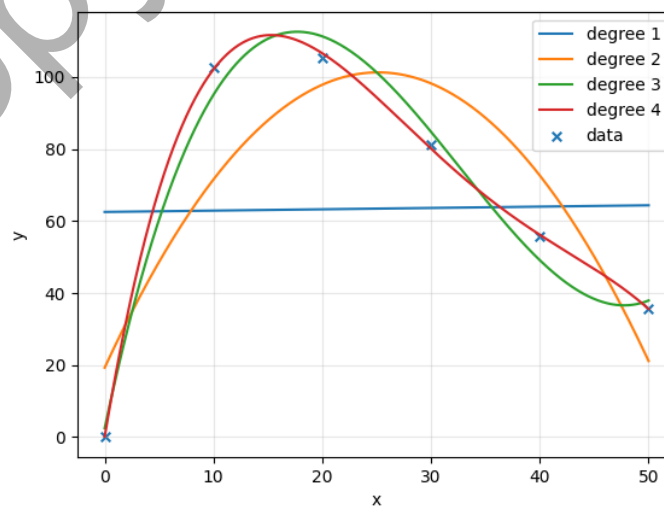
```

B Polynomial regression – sci-kit learn

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 from sklearn.linear_model import LinearRegression
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.pipeline import make_pipeline
7
8 x_plot = np.linspace(0, 50, 1000)
9
10 x = np.array([0, 10, 20, 30, 40, 50], dtype=np.longdouble)
11 y = np.array([0, 102.6903, 105.4529, 81.21744, 55.6016, 35.6859])
12
13 # Column matrix
14 X = x[:, np.newaxis]
15 X_plot = x_plot[:, np.newaxis]
16
17 plt.scatter(x, y, s=30, marker='x', label='data')
18 degree1st = [1, 2, 3, 4]
19
20 for degree in degree1st:
21     model = make_pipeline(PolynomialFeatures(degree), LinearRegression(
22         fit_intercept=False))
23     model.fit(X, y)
24     polyn = model.named_steps['linearregression']
25     coeff = polyn.coef_
26     print('Coefficients:', coeff)
27     print('R^2:', model.score(X, y))
28     y_plot = model.predict(X_plot)
29     plt.plot(x_plot, y_plot, label=f'degree {degree}')
30
31 plt.grid(alpha=0.3)
32 plt.xlabel('x')
33 plt.ylabel('y')
34 plt.legend()
35 plt.show()

```



C Least squares function approximation

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.special import legendre
4 from scipy.integrate import quad
5
6 func = lambda x: np.sin(np.pi*x) # Original function to be approximated
7 orth = [lambda x: 1, lambda x: x, lambda x: x**2-(1/3), lambda x: x**3-(3/5)
8         *x] # Test
9
10 def c(f, phi2, a, b):
11     num = quad(f, a, b) # Integral of f, limits a to b
12     dnm = quad(phi2, a, b)
13     return num[0]/dnm[0] # c_i = \frac{\langle f, \phi_i \rangle}{\langle \phi_i, \phi_i \rangle}
14
15 def y(order, x, coeff):
16     polyn = 0
17     for i in range(order + 1):
18         polyn += coeff[i] * legendre(i)(x) # \hat{f}(x) = \sum_{i=1}^n c_i \phi_i(x)
19     return polyn
20
21 def sqerror(y, fhat):
22     return np.sum((y-fhat)**2)
23
24 def main():
25     order = 3
26     a, b = -1, 1
27     coeff = []
28     xlin = np.linspace(a, b, 1000)
29
30     for i in range(order+1):
31         phi = legendre(i) # Orthogonal function -- Legendre polynomials
32
33         f = lambda x: func(x)*phi(x) # Combining 2 functions
34         phi2 = lambda x: phi(x)*phi(x)
35
36         coeff.append(c(f, phi2, a, b))
37
38     print(coeff)
39     print('Squared error:', sqerror(func(xlin), y(order, xlin, coeff)))
40
41     plt.plot(xlin, func(xlin), label=r'$\sin\{\pi x\}$')
42     plt.plot(xlin, y(order, xlin, coeff), '--', label=r'Approximation')
43     plt.legend()
44     plt.title('Least squares function approximation of $\sin\{\pi x\}$')
45     plt.show()
46
47 if __name__ == '__main__':
48     main()

```

D Multiple linear regression

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn import linear_model
5 import seaborn as sn
6
7 df = pd.read_csv('usecons7687.csv')
8
9 """
10 CRUDE = dollars/barrel crude oil
11 INTEREST = % interest on ten yr. U.S. treasury notes
12 FOREIGN = foreign investments/billions of dollars
13 DJIA = Dow Jones industrial average
14 GNP = GNP/billions of dollars
15 PURCHASE = purchasing power U.S. dollar (1983 base)
16 CONSUMER = consumer debt/billions of dollars
17
18 data from https://college.cengage.com/mathematics/brase/
19   understandable_statistics/7e/students/datasets/mlr/frames/frame.html
20 """
21 X = df[['PURCHASE', 'GNP']].values.reshape(-1,2)
22 y = df['CONSUMER'].values
23
24 x1 = X[:, 0]
25 x2 = X[:, 1]
26
27 x1_pred = np.linspace(min(x1), max(x1), 100)
28 x2_pred = np.linspace(min(x2), max(x2), 100)
29 x1m_pred, x2m_pred = np.meshgrid(x1_pred, x2_pred)
30 modelxs = np.array([x1m_pred.flatten(), x2m_pred.flatten()]).T
31
32 ols = linear_model.LinearRegression()
33 model = ols.fit(X, y)
34 print(model.coef_)
35 predicted = model.predict(modelxs)
36 print(predicted)
37
38 r2 = model.score(X, y)
39
40 plt.style.use('default')
41 fig = plt.figure(figsize=(8, 8))
42 ax = fig.add_subplot(111, projection='3d')
43
44 ax.plot(x1, x2, y, color='k', zorder=15, linestyle='none', marker='o', alpha
45         =0.5)
46 ax.scatter(x1m_pred.flatten(), x2m_pred.flatten(), predicted, facecolor
47            =(0,0,0,0), s=20, edgecolor='#70b3f0')
48 ax.set_xlabel('Purchasing power', fontsize=12)
49 ax.set_ylabel('GNP', fontsize=12)
50 ax.set_zlabel('Consumer debt', fontsize=12)
51 ax.view_init(elev=28, azim=120)
52 fig.tight_layout()
53 plt.show()
54
55 corrMatrix = df.corr(method='pearson')
56 sn.heatmap(corrMatrix, annot=True)
57 plt.show()

```

E Nonlinear least squares

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.optimize import least_squares
4
5 def func(x, a, b):
6     return a*x/(b+np.exp(x))
7
8 def Jacobian(f, x, a, b):
9     h = 1e-6
10    pdv_a = (f(x, a+h, b)-f(x, a-h, b))/(2*h) # Finite difference: pdvF = (F
11    (x+h)-F(x-h))/(2h)
12    pdv_b = (f(x, a, b+h)-f(x, a, b-h))/(2*h)
13    return np.column_stack([pdv_a, pdv_b])
14
15 def GaussNewton(f, x, y, a0, b0, iterations): # Gauss-Newton
16    tol = 1e-12
17    g_i = g = np.array([a0, b0])
18    for itr in range(iterations):
19        g = g_i
20        J = Jacobian(f, x, g[0], g[1])
21        r = y - f(x, g[0], g[1])
22        g_i = g + np.linalg.inv(J.T@J)@J.T@r # \beta_{i+1} = \beta{i} + (J_i
23        ^T J_i)^{-1} J_i^T r_i
24        if np.linalg.norm(g-g_i) < tol:
25            break
26    return g
27
28 def r(a, x, y):
29    return a[0]*x/(a[1]+np.exp(x)) - y
30
31 def main():
32    x = np.linspace(0, 10, 50)
33    y = func(x, 6, 9) + np.random.normal(0, 0.02, x.shape) # + Noise
34
35    a, b = GaussNewton(func, x, y, 2.5, 0.6, 50)
36    print(a, b)
37    y_hat = func(x, a, b)
38
39    a_guess = np.array([2, 3])
40    res = least_squares(r, a_guess, args=(x, y)) # Scipy's least squares
41    module
42    res_cauchy = least_squares(r, a_guess, loss='cauchy', f_scale=0.02, args
43    =(x, y)) # Cauchy loss function
44    print(res.x)
45    y_hatscp = func(x, res.x[0], res.x[1])
46
47    plt.plot(x, y, 'x')
48    plt.plot(x, y_hat, label=r'Gauss-Newton')
49    plt.plot(x, y_hatscp, label=r'least_squares')
50    plt.legend()
51    plt.show()
52
53 if __name__ == '__main__':
54    main()

```