

MA2252 Introduction to Computing

Lecture 13

Least Squares Regression

Sharad Kumar Keshari

School of Computing and Mathematical Sciences

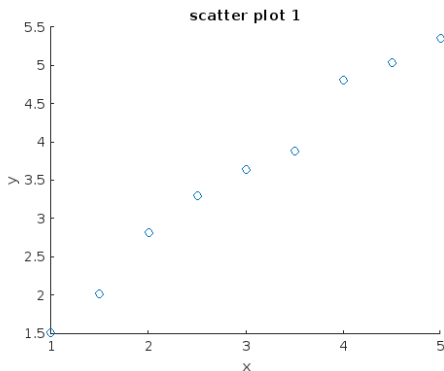
University of Leicester

At the end of lecture, students will be able to

- understand scatter plot and regression
- understand theory of Least Squares Regression
- solve basic regression problems in MATLAB

Scatter plot

A **scatter plot** plots two different sets of data using dots. Unlike line plot, the dots are not connected by a curve.



Scatter plot (contd.)

MATLAB's `scatter()` function can be used to create a scatter plot.

Example: The code below creates the scatter plot shown earlier.

```
x=1:0.5:5; %create data for vector x
s=size(x); %find size of vector x
y=x+rand(s); %create a vector y=x+'some random values'
scatter(x,y) %create scatter plot
title('scatter plot 1')
xlabel('x')
ylabel('y')
```

Demo

Regression is a statistical technique used to find the 'best-fit curve' that describes a scatter plot.

Depending on the data trend in a scatter plot, one may use

- Linear Regression curve
- Non-linear Regression curve

Regression (contd.)

Linear Regression curve example

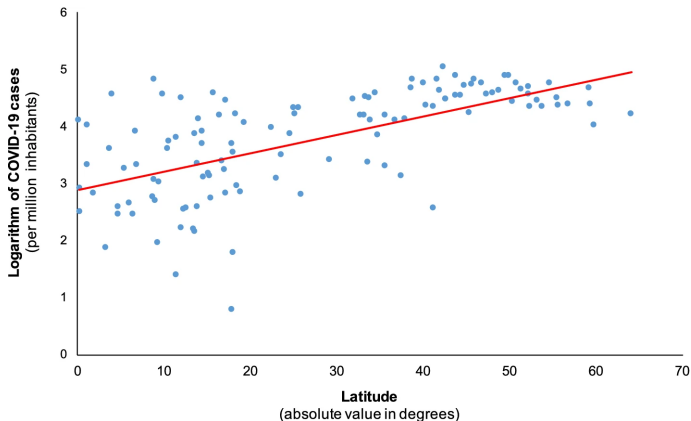


Figure: Scatter plot showing linear trend ¹

Non-linear Regression curve example

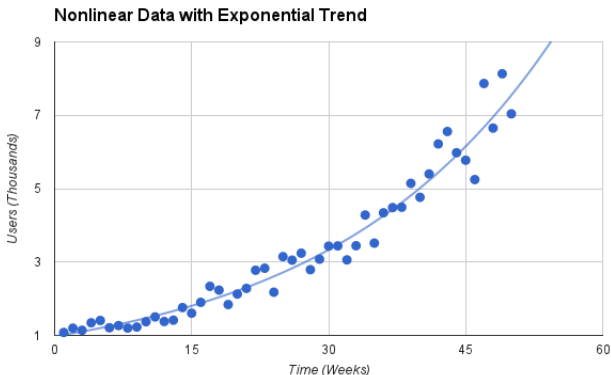


Figure: Data showing number of active users on a website with time ²

¹Chen, S., Prettner, K., Kuhn, M. et al. Climate and the spread of COVID-19. Sci Rep 11, 9042 (2021). <https://doi.org/10.1038/s41598-021-87692-z>

²<http://sam-koblenski.blogspot.com>

Regression model

A regression model provides a function to describe the relationship between one (or more) independent variables and a dependent variable.

A basic regression model is the 'Least Squares Regression model'.

Least Squares Regression

Here, the relationship between dependent data points $y_i (i = 1, 2, \dots, m)$ and independent data points x_i is modelled as

$$\hat{y}(x) = \sum_{i=1}^n \alpha_i f_i(x) \quad (1)$$

where

- $\hat{y}(x)$ is an estimation function
- α_i are parameters of estimation function
- $f_i(x)$ are linearly independent basis functions

Least Squares Regression (contd.)

The parameters are then found by minimising the total squared error E .

$$E = \sum_{i=1}^m (\hat{y} - y_i)^2 \quad (2)$$

Substituting (1) in (2) gives

$$E = \sum_{i=1}^m \left(\sum_{j=1}^n \alpha_j f_j(x_i) - y_i \right)^2 \quad (3)$$

E is a function of n variables namely $\alpha_j (j = 1, 2, \dots, n)$.

Least Squares Regression (contd.)

The solution for n parameters α_j which minimise the total squared error E is given as

$$\beta = \text{pinv}(A) * Y \quad (4)$$

Here,

- β is a column vector with n entries α_j
- A is a $m \times n$ matrix with entries $A(i,j) = f_j(x_i)$
- $\text{pinv}(A)$ is the pseudo-inverse of A
- Y is a column vector with m entries y_i

Least Squares Regression (contd.)

Derivation of (4): Please refer book and lecture recording

Least Squares Regression (contd.)

Example: Perform a least squares regression for the scatter plot created before using estimation function $\hat{y}(x) = \alpha_1 x + \alpha_2$.

```
openfig('scatter_plot_1.fig') %opens figure scatter_plot_1.fig
a = get(gca,'Children');
x = get(a, 'XData'); %extract x-data points
y = get(a, 'YData'); %extract y-data points
A=[x',ones(size(x'))]; %create the matrix A of basis functions
beta=pinv(A)*y'; %evaluate vector beta containing values of parameters

%plot the regression line
hold on
plot(x,beta(1)*x+beta(2))
```

Least Squares Regression (contd.)

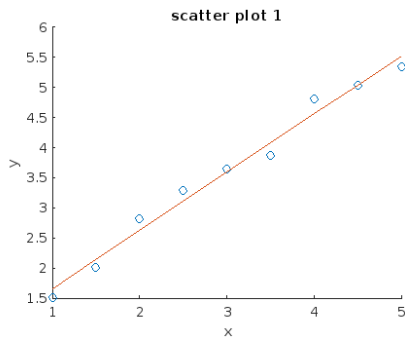


Figure: Regression line

polyfit and polyval functions

When the estimation function $\hat{y}(x)$ is a polynomial, MATLAB's polyfit and polyval functions can be used.

For vectors x and y containing x -data and y -data points respectively,

- $p = \text{polyfit}(x, y, n)$ creates a vector p of the coefficients of regression polynomial $p(x)$ of degree n
- $\text{polyval}(p, x)$ calculates the values of $p(x)$

polyfit and polyval functions (contd.)

Example: This code creates the same regression line as was shown before.

```
openfig('scatter_plot_1.fig') %opens figure scatter_plot_1.fig
a = get(gca,'Children');
x = get(a, 'XData'); %extract x-data points
y = get(a, 'YData'); %extract y-data points

p= polyfit(x,y,1); %creates coefficients of regression polynomial of degree 1
Y=polyval(p,x); %evaluates the value of polynomial at x-data points
hold on
plot(x,polyval(p,x)) %plot the regression polynomial
```

Demo

Nonlinear Estimation Functions

Sometimes, a nonlinear estimation function provides the best fit for a scatter plot. This means we require

$$\hat{y}(x) = g(\alpha_1, \alpha_2, \dots, \alpha_n, x) \quad (5)$$

where g is some nonlinear function.

In some special cases, a transformation such as

$$\tilde{y}(x) = h(\hat{y}(x)) \quad (6)$$

can linearise the equation (5) into (1).

Nonlinear Estimation Functions (contd.)

Example: Consider the estimation function

$$\hat{y}(x) = \alpha_1 e^{\alpha_2 x} \quad (7)$$

Applying the transformation

$$\tilde{y}(x) = \log(\hat{y}(x)) \quad (8)$$

converts (7) into

$$\tilde{y}(x) = \tilde{\alpha}_1 + \alpha_2 x \quad (9)$$

where we define $\tilde{\alpha}_1 = \log(\alpha_1)$. Now, least squares regression can be applied to equation (9). The parameter α_1 can be found using $\alpha_1 = e^{\tilde{\alpha}_1}$.

End of Lecture 13

Please provide your feedback [▶ here](#)