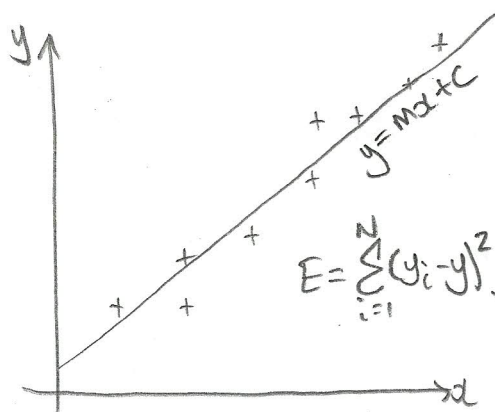


lines of best fit to (x,y) scattered data

+ Correlation



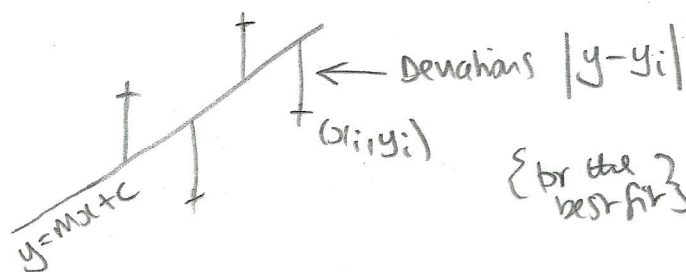
Consider a set of N pairs of x, y data
 $\{x_i, y_i\} ; i = 1 \dots N$

what is the 'best' line $y = mx + c$
 which can be associated with the data?

Define

$$E = \sum_{i=1}^N (y_i - mx_i - c)^2 \quad (1)$$

is the sum of the squares of deviation of the line of best fit from the actual y values.



To find m and c ,
 consider a Surface of $E(m, c)$

$\{ \text{for the best fit} \} \rightarrow m$ and c are the coordinates of the minimum of the Surface.

is when the overall sum of the square deviations is minimized.

Hence to find m and c we must consider solution to

$$\frac{\partial E}{\partial m} = 0 \quad \text{and} \quad \frac{\partial E}{\partial c} = 0$$

{ Partial derivatives with respect to the variables m and c correspondingly All other parameters are held constant }

$$\begin{aligned} \frac{\partial E}{\partial m} &= \sum_{i=1}^N (-x_i) (2) (y_i - mx_i - c) \\ &= 2 \sum_{i=1}^N (mx_i^2 + cx_i - x_i y_i) \end{aligned}$$

$$\therefore \text{when } \frac{\partial E}{\partial m} = 0 \Rightarrow m \sum_{i=1}^N x_i^2 + c \sum_{i=1}^N x_i - \sum_{i=1}^N x_i y_i = 0 \quad (2)$$

$$\frac{\partial E}{\partial c} = \sum_{i=1}^N (2) (-1) (y_i - mx_i - c) \quad \therefore \text{when } \frac{\partial E}{\partial c} = 0$$

$$\sum_{i=1}^N y_i - m \sum_{i=1}^N x_i - c \sum_{i=1}^N 1 = 0 \quad (3)$$

Equations (2) and (3) can be solved simultaneously to yield (m, c) which minimize E^*

(* Strictly speaking we should check that $\frac{\partial^2 E}{\partial c^2} > 0$ and $\frac{\partial^2 E}{\partial m^2} > 0$ for this to be a minima) \rightarrow we will check this later!

Firstly let us re-write (2) and (3) and define some statistical quantities

$$(2) \quad m \frac{1}{N} \sum_{i=1}^N x_i^2 + c \frac{1}{N} \sum_{i=1}^N x_i - \frac{1}{N} \sum_{i=1}^N x_i y_i = 0$$

$$\overline{x^2} = \frac{1}{N} \sum_{i=1}^N x_i^2$$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad [\text{Mean of } \{x_i\}]$$

$$\overline{xy} = \frac{1}{N} \sum_{i=1}^N x_i y_i$$

So (2): $\boxed{m \overline{x^2} + c \bar{x} - \overline{xy} = 0}$

Now (3): $\frac{1}{N} \sum_{i=1}^N y_i - m \frac{1}{N} \sum_{i=1}^N x_i - c \frac{1}{N} \sum_{i=1}^N 1 = 0$

Define $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$ $\sum_{i=1}^N 1 = N$

Hence (3): $\boxed{\bar{y} - m \bar{x} - c = 0}$

$$(2) + \bar{x}(3): \quad m \overline{x^2} - \overline{xy} + \bar{y} \bar{x} - m \bar{x}^2 = 0$$

$$\therefore m (\overline{x^2} - \bar{x}^2) = \overline{xy} - \bar{y} \bar{x}$$

$$\therefore \boxed{m = \frac{\overline{xy} - \bar{y} \bar{x}}{\overline{x^2} - \bar{x}^2}}$$

m is the ratio of two statistical quantities

Covariance of $\{x, y\}$

$$\text{Cov}[x, y] = \overline{xy} - \bar{x}\bar{y}$$

Variance of $\{x\}$

$$V[x] = \overline{x^2} - \bar{x}^2$$

Hence

$$m = \frac{\text{Cov}[x, y]}{V[x]}$$

Note:

$$\text{Cov}[x, y] = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

Now from (3): $c = \bar{y} - m\bar{x}$

in Summary, the line of best fit formed by minimizing the sum of square "y" deviations (E) is given by \rightarrow of data $\{x_i, y_i\}$

$$y = mx + c$$

$$m = \frac{\text{Cov}[x, y]}{V[x]}$$

$$c = \bar{y} - m\bar{x}$$

$$\text{Cov}[x, y] = \overline{xy} - \bar{x}\bar{y}$$

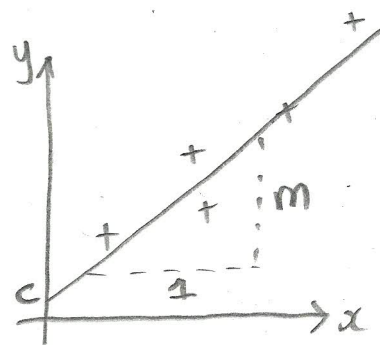
$$V[x] = \overline{x^2} - \bar{x}^2$$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\overline{x^2} = \frac{1}{N} \sum_{i=1}^N x_i^2$$

$$\overline{xy} = \frac{1}{N} \sum_{i=1}^N x_i y_i$$

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$



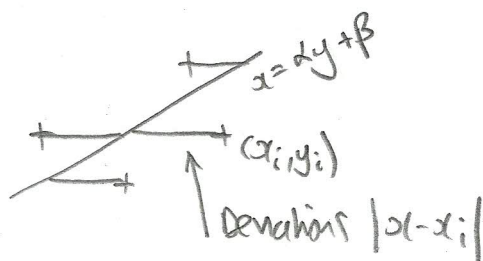
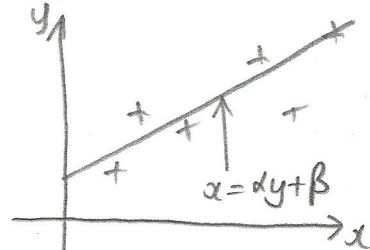
of course the entire exercise could be achieved by considering the minimum of a surface $F(\alpha, \beta)$ found from the sum of square of deviations from line of best fit

$$x = \alpha y + \beta \quad [\Rightarrow y = \frac{x}{\alpha} - \frac{\beta}{\alpha}]$$

In this case;

$$F = \sum_{i=1}^N (x_i - \alpha)^2$$

$$= \sum_{i=1}^N (x_i - \alpha y_i - \beta)^2$$



The deviation of α, β is clearly the same as above, but with $x \leftrightarrow y$ interchanged

i.e.

$$\alpha = \frac{\overline{xy} - \bar{y}\bar{x}}{\bar{y}^2 - \bar{y}^2} = \frac{\text{Cov}[x, y]}{\text{Var}[y]}$$

[dividing since $\text{Cov}[x, y] = \overline{xy} - \bar{y}\bar{x}$, $\text{Cov}[x, y] = \text{Cov}[y, x]$)

$$\beta = \bar{x} - \alpha \bar{y}$$

So for horizontal deviations

$$y = mx + c$$

$$m = \frac{\text{Var}[y]}{\text{Cov}[x, y]} \quad c = \bar{y} - m\bar{x}$$

Now, if $\{x_i, y_i\}$ are perfectly correlated, $x = \alpha y + \beta$ and $y = mx + c$ are the same line

i.e.

$$y = mx + c$$

$$y = \frac{x - \beta}{\alpha} \quad \Rightarrow \quad m = \frac{1}{\alpha} \quad \text{and} \quad c = -\frac{\beta}{\alpha}$$

Consider

$$m = \frac{1}{\alpha} \quad \Rightarrow \quad \boxed{m\alpha = 1}$$

i.e. lines have the same gradient for vertical and horizontal deviations.

Hence if we define $\rho^2 = m^2$, the closeness of ρ to unity will tell us something useful about the degree of correlation of the data.

$$\rho^2 = \frac{\text{Cov}[x, y]}{\sqrt{V[x]}} \times \frac{\text{Cov}[x, y]}{\sqrt{V[y]}}$$

$$\rho = \frac{\text{Cov}[x, y]}{\sigma_x \sigma_y}$$

This is called the product moment correlation coefficient.

where Standard deviations

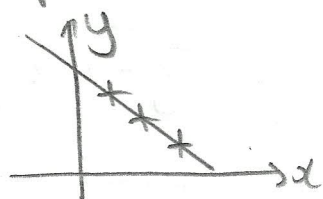
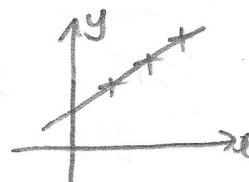
$$\sigma_x = \sqrt{V[x]}$$

$$\sigma_y = \sqrt{V[y]}$$

Note ρ could also be -ve since $\text{Cov}[x, y]$ need not be +ve.

If $\rho = +1$ the data will be perfectly correlated

If $\rho = -1$ " " " " -ve "



Ex 11.1 $\frac{\partial^2 E}{\partial m^2} = 2 \sum_{i=1}^N (x_i^2) = 2N \overline{x^2}$

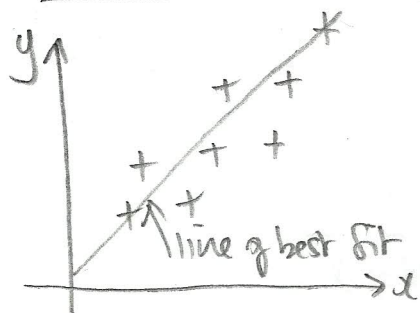
This must be > 0 so $\frac{\partial E}{\partial m} = 0$ will be a minimum

$$\frac{\partial^2 E}{\partial c^2} = \frac{\partial}{\partial c} \left\{ -2 \sum_{i=1}^N (y_i - mx_i - c) \right\}$$

$$= 2 \sum_{i=1}^N 1 = 2N$$

This is > 0 so $\frac{\partial E}{\partial c}$ will also be a minimum.

Overall linear regression and correlation Summary



$\{x_i, y_i\}$ data samples. There are N pairs.

Define

$$\begin{aligned}\bar{x} &= \frac{1}{N} \sum_{i=1}^N x_i & \bar{y} &= \frac{1}{N} \sum_{i=1}^N y_i \\ \overline{x^2} &= \frac{1}{N} \sum_{i=1}^N x_i^2 & \overline{y^2} &= \frac{1}{N} \sum_{i=1}^N y_i^2 \\ \overline{xy} &= \frac{1}{N} \sum_{i=1}^N x_i y_i\end{aligned}$$

Do this in a table \rightarrow

x	y	x^2	y^2	xy
\vdots	\vdots	\vdots	\vdots	\vdots

Then compute:

$$V[x] = \overline{x^2} - \bar{x}^2$$

$$V[y] = \overline{y^2} - \bar{y}^2$$

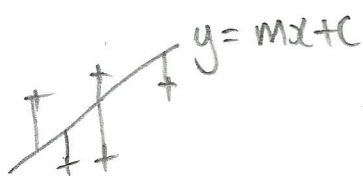
product-moment
correlation
coefficient

$$\text{Cov}[x, y] = \overline{xy} - \bar{x}\bar{y}$$

(Covariance)

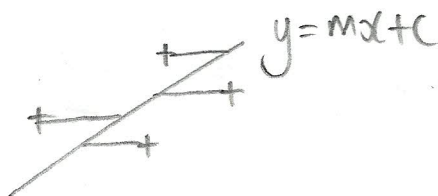
$$\rho = \frac{\text{Cov}[x, y]}{\sqrt{V[x]V[y]}}$$

lines of best fit are then:



Vertical fit

$$m = \frac{\text{Cov}[x, y]}{V[x]} \quad c = \bar{y} - m\bar{x}$$

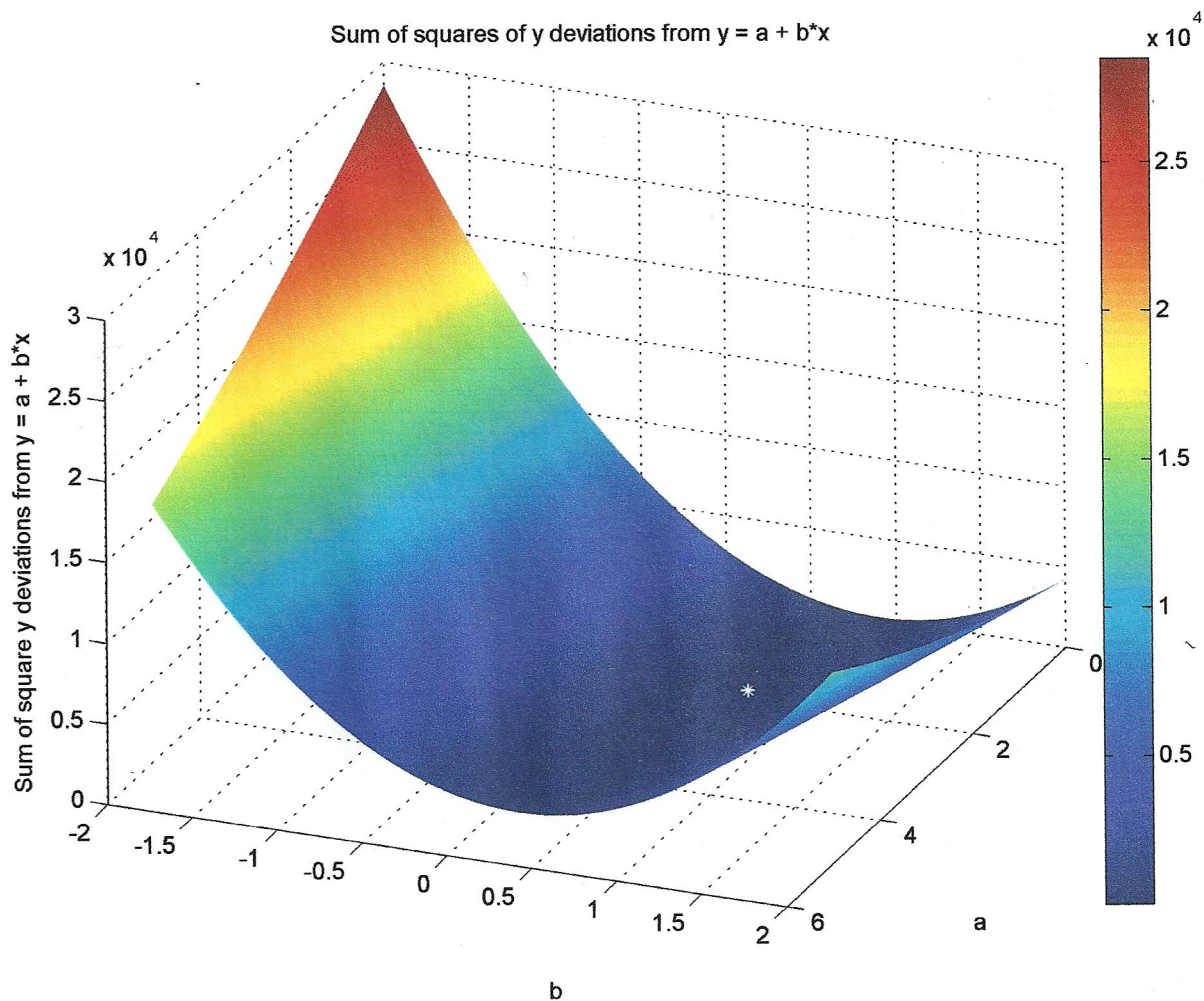
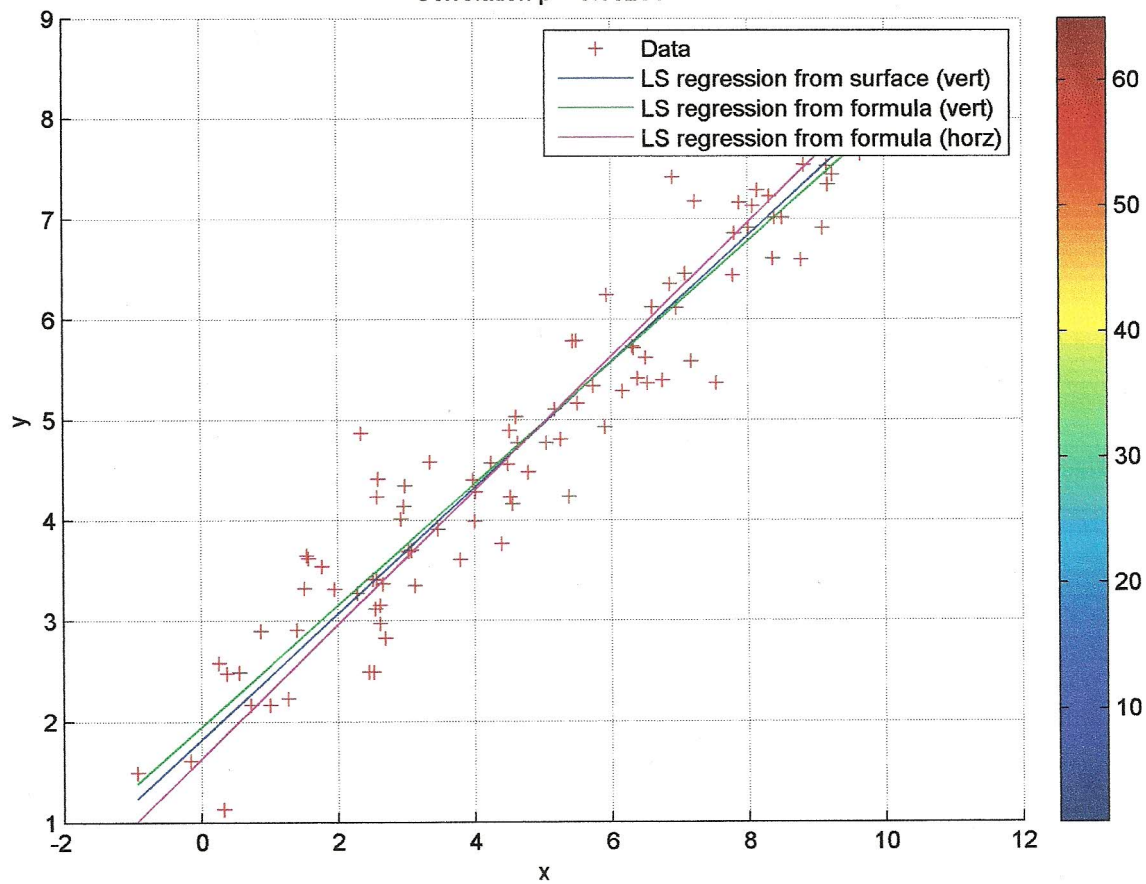


Horizontal fit

$$m = \frac{V[y]}{\text{Cov}[x, y]} \quad c = \bar{y} - m\bar{x}$$

+ve correlation

Least squares linear regression
 Vertical: $y = 1.9461 + 0.60346x$.
 Horizontal: $y = 1.6304 + 0.66584x$.
 Correlation $p = 0.95201$



-ve correlation

Least squares linear regression
 Vertical: $y = 2.0275 + -0.80835x$.
 Horizontal: $y = 2.2835 + -0.86039x$.
 Correlation $p = -0.96929$

