

3.

a. $(x, y) \in \mathbb{R}^2$

$p(x, y) \Rightarrow$ Joint PDF of x, y .

\Rightarrow There can be two different methods to tackle this problem.

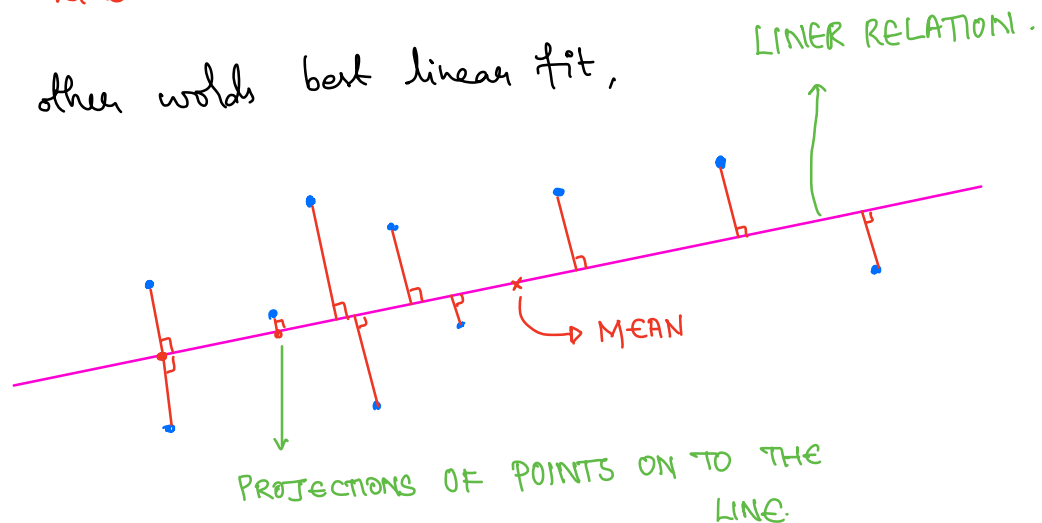
I. Using eigenvalues and eigenvectors:-

Let us say we find an eigenvector $V_{\lambda_{\max}}$ and corresponding max. eigen-value λ_{\max} .

We argue that this will give the line passing through all having $V_{\lambda_{\max}}$ as direction.

Line with LEAST SQUARED ERROR.

In other words best linear fit,



- First we show that line passes through mean
- And then we will show that PCA component is in the best direction.

Let us manipulate the sum of squared distances first :

$$e = \sum \frac{|mx + c - y|^2}{m^2 + 1} = \frac{m^2 \sum x^2 + \sum y^2 + nc^2 + 2mc \sum x - 2m \sum xy - 2c \sum y}{m^2 + 1}$$

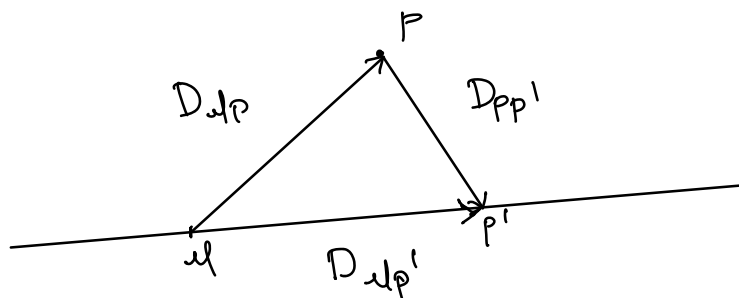
Let us first try to minimize write,

$$\frac{\partial e}{\partial c} = \frac{2cn + 2m \sum x - 2 \sum y}{m^2 + 1}$$

$$\Rightarrow \sum y = m \sum x + cn$$

$$\Rightarrow \sum \frac{y}{n} = m \cdot \frac{\sum x}{n} + c$$

\therefore Line passes through mean.



$$D_{\text{alp}}^2 = D_{\text{pp1}}^2 + D_{\text{alp}}^2$$

$$\sum D_{\text{alp}}^2 = \sum D_{\text{pp1}}^2 + \sum D_{\text{alp}}^2$$

$$\text{Var} = \sum D_{\text{pp1}}^2 + \sum D_{\text{alp}}^2$$

Hence if we fix variance if $\sum D_{\text{pp1}}^2$ is minimized, $\sum D_{\text{alp}}^2$ is obviously maximized to give that the line vector.

$\left(\frac{m}{\sqrt{m^2+1}}, \frac{1}{\sqrt{m^2+1}} \right)$ is the eigen-vector of maximum eigen value.

\therefore Proved.

Final algorithm to find the line :-

- ① Do the Principle component Analysis using eigen vector and values and get the vector V .
- ② Draw a line with directional vector as V and initial point as u .

③ That would be your best fit line.

$$C = \sum \begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} x & y \end{bmatrix}$$

$$D, V = \text{eig}(C)$$

$$\{V_{\text{ulmax}}\} \rightarrow \frac{V_{\text{ulmax}}(x)}{V_{\text{ulmax}}(y)} = \frac{1}{m}$$

$$y = mx + c$$

$$c = \sum \frac{y}{n} - m \sum \frac{x}{n}$$

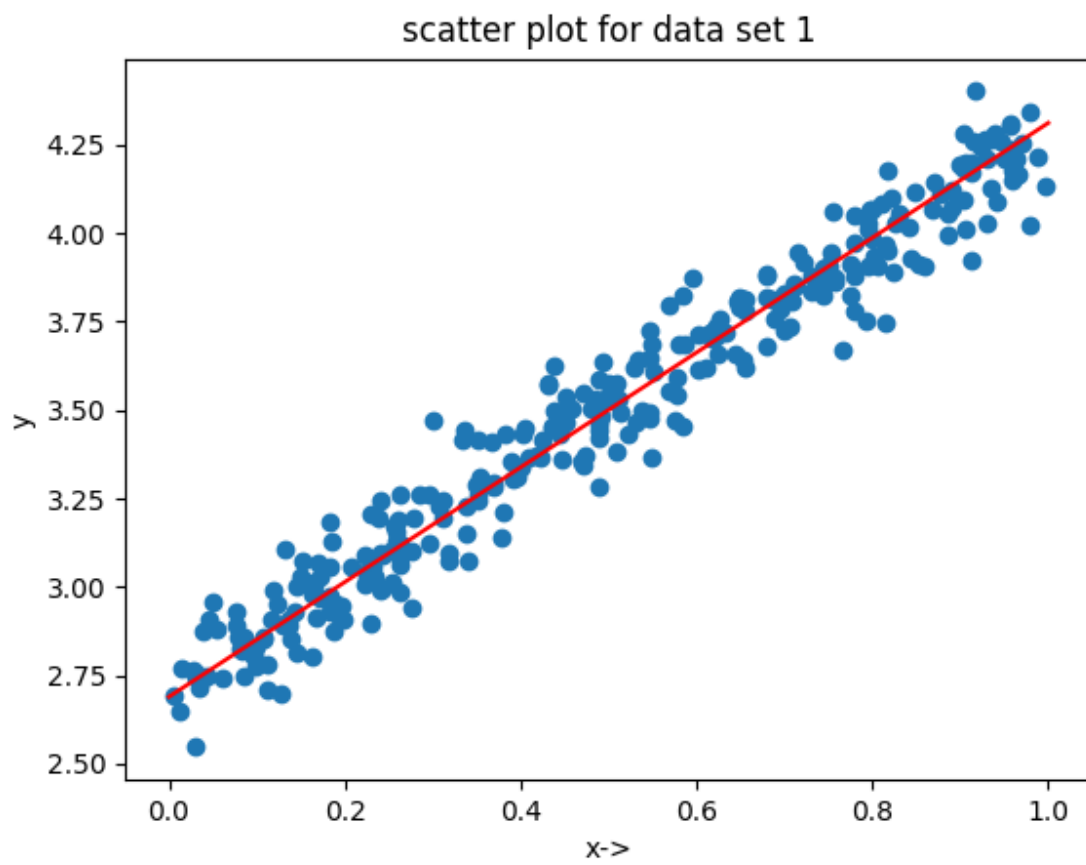
b.

CODE \Rightarrow q3-b.py

GRAPH \Rightarrow Scatter_3b.png

The procedure given in the previous question has been applied in the given code.

Plotted the scatter plot and also the line which shows the linear relation between x and y variables that we got along the eigenvector with maximum eigen value.



SCATTER PLOT FOR DATASET 1

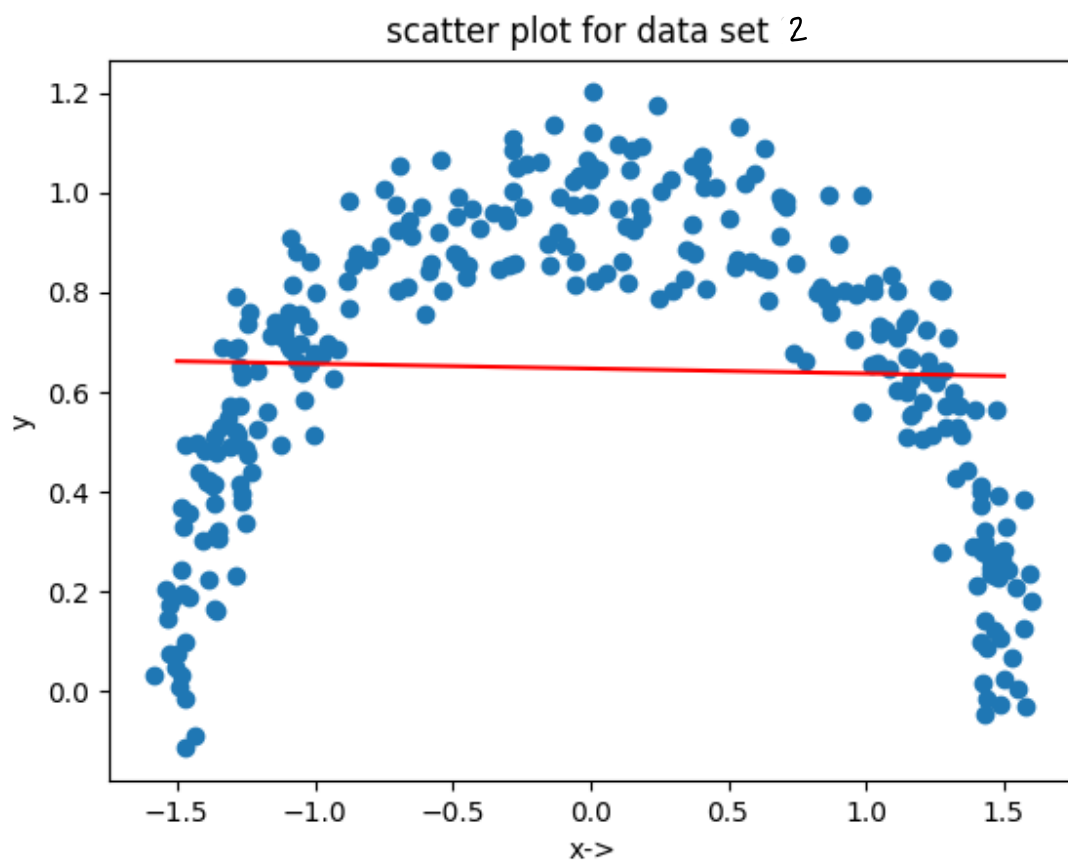
C.

CODE \Rightarrow q3-C.py
q3-C.ipynb

IMAGE \Rightarrow scatter-3C.png

OBSERVATIONS:

- Data in set 2 does not belong to variables which are related by a linear relation
- This approximation works best when data follows nearly linear trend with slight deviation and errors
- But whenever the relation exists the PCA analysis gives the relation which is very accurate.
- Like the set 1 case the scatter plot is linear (mostly) but in set 2 case the scatter plot is non-linear, so the approximate in set 2 case is not valid.
- So the linear relation we get from the PCA analysis doesn't imply that there is a relation. But if the relation exists then we can get it by this PCA analysis.



SCATTER PLOT FOR DATASET 2