

3 a) Let  $\vec{\mu}$  be mean of given sample data and  $C$  be the co-variance matrix.

unit

Let  $\vec{v}$  be a vector along the direction in which spread is maximum, then a line passing through mean & parallel to  $\vec{v}$  will give our linear approximation relationship between  $X$  &  $Y$ .

Let  $\vec{v}_1$  &  $\vec{v}_2$  be eigenvectors of  $C$  with eigenvalues  $\lambda_1$  &  $\lambda_2$ .

Let  $\lambda_1 > \lambda_2$  w.l.o.g

$\therefore$  Spread will be max. in direction in  $\vec{v}_1$ .

$$\therefore \vec{v} = \frac{\vec{v}_1}{|\vec{v}_1|}$$

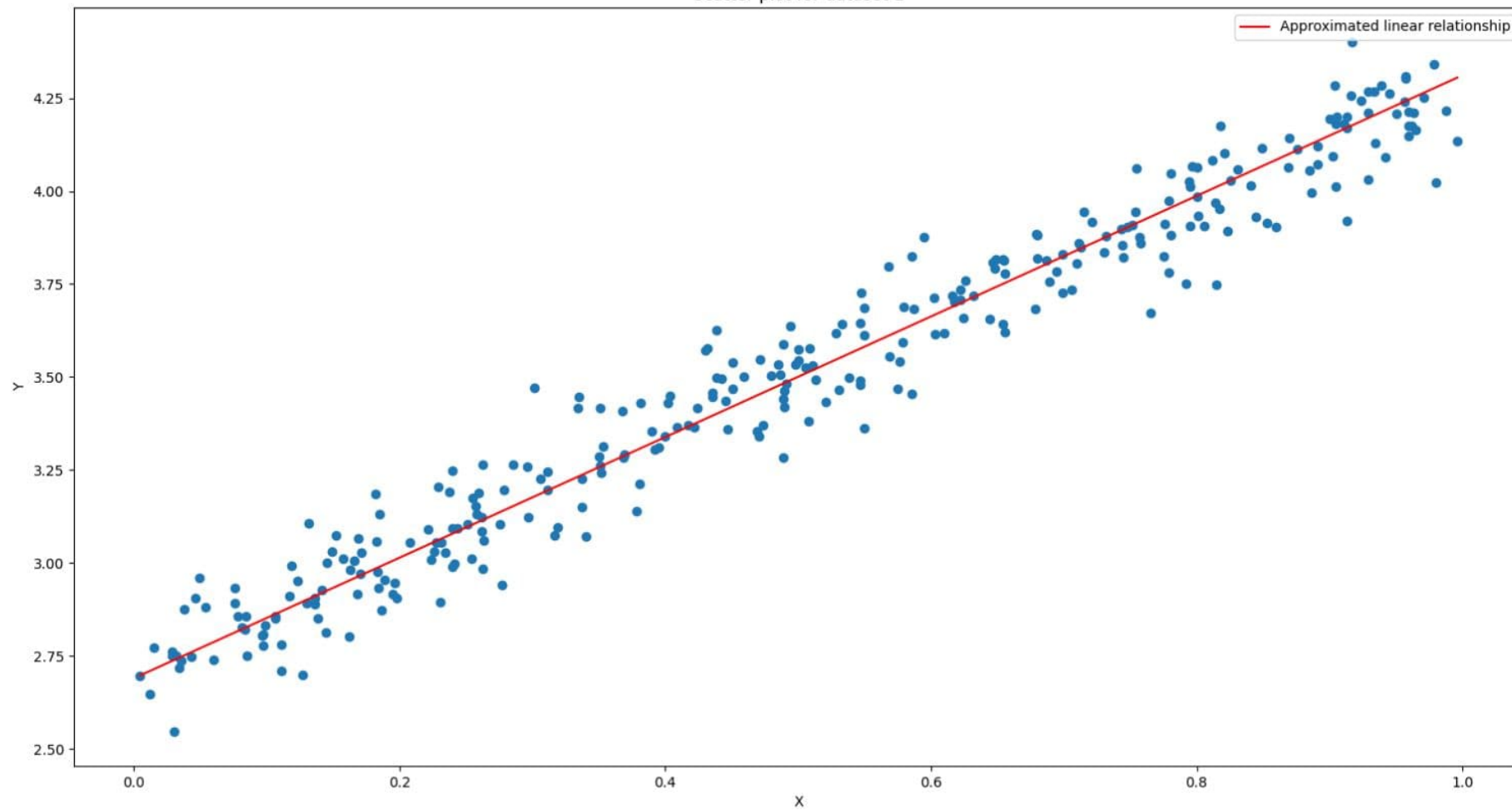
$\therefore$  line  $l$  given by

$$\vec{l} = \vec{a} + k\vec{v}$$

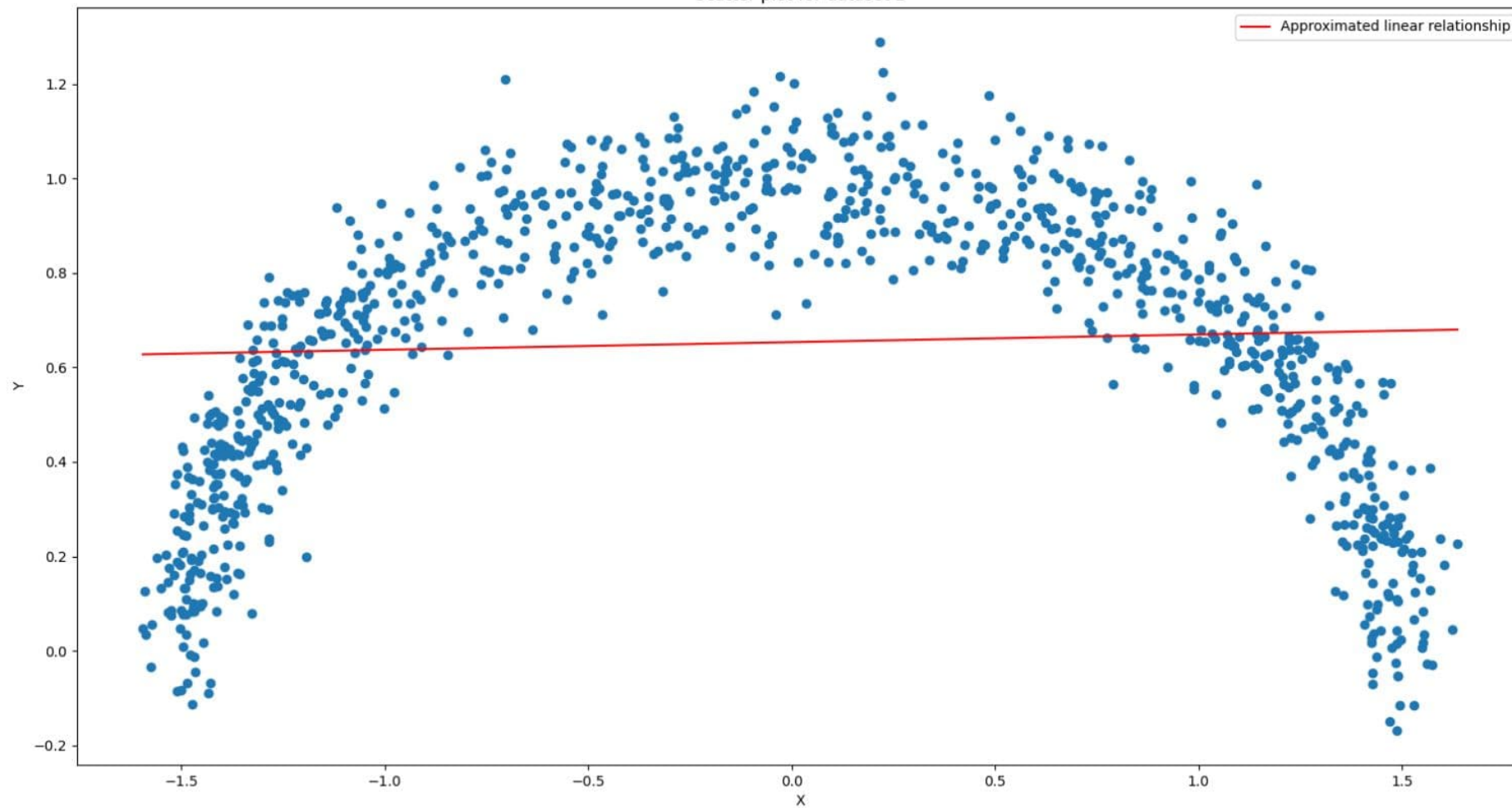
is our best approximation.



Scatter plot for dataset 1

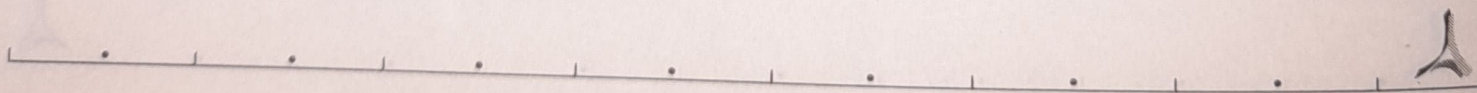


Scatter plot for dataset 2



3 For first set of points ,

we can ~~if~~ infer ~~from~~ the scatter plot that those set of points can be nicely ~~placed~~ placed around a line as ~~relation seems~~ diversions are ~~very~~ not that large from our calculated linear relationship, it is a good approximation of first set.





While for second sets of data, we can clearly infer that data points are following a curvy or non-linear pattern which is ~~very~~ not easy to put in a linear relationship hence diversions are larger and approximation is not that good compared to set 1.