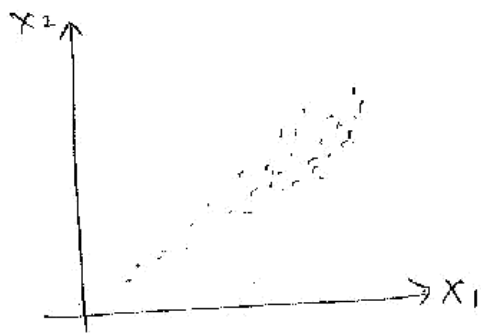# Principal Component Analysis ( PCA )

PCA helps reduce high-dimensional data into something that can be explained in fewer dimensions and gain an understanding of data.
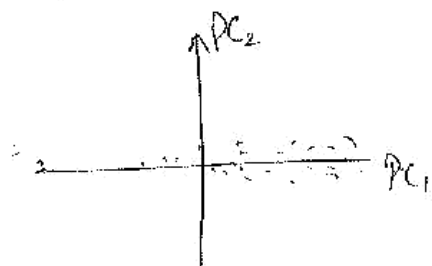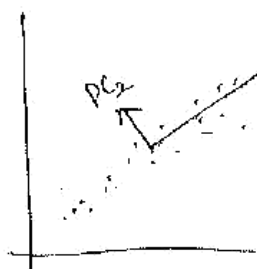Since our suspicion is that in our intresting data set not all the measures/features are independent but there exist correlations/structure/patterns.

Assume 2D data as shown in graph below



The data is not random as $x_1$ increases $x_2$ also increase and they are positively co-rrelated.

Now assume that we have 1-D to explain this data. Then the most of the variation in the data is along the direction of $PC_1$. This is the intution behind PCA.



This is the Principal Component direction of data. Hence if we transform the data on to my new basis of $PC_1$, we can actually get the same data. What PCA is doing is this transformation. In this case it is rotation. This is in 2-D what about D-Dimensions.
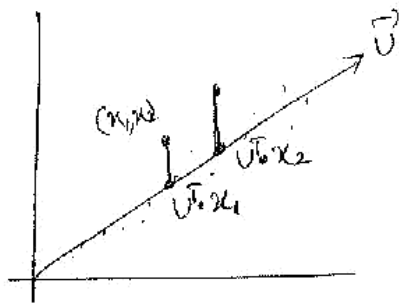
Assume that following X data Matrix & a vector U.

$$X = \begin{bmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nD} \end{bmatrix} \rightarrow \text{first feature} \qquad U = \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_D \end{bmatrix}$$

Where N samples are collected in D Dimensions



we need to project all the original data points on to a generic vector $\vec{U}$, and find the variance after projecting all points onto $\vec{U}$.

The variance of all the projected points is

$$\frac{1}{N} \sum_{i=1}^{n} \left( U^T . x_n - U^T . \bar{x} \right)^2$$

$$\frac{1}{N} \sum_{i=1}^{n} \left( U^T (x_n - \bar{x}) \right)^2 \qquad \begin{aligned} (A.B)^2 &= (A.B)(A.B)^T \\ &= (A.B) B^T A^T \end{aligned}$$

$$\frac{1}{N} \sum_{i=1}^{n} U^T (x_n - \bar{x}) \left( U^T (x_n - \bar{x}) \right)^T$$

$$\frac{1}{N} \sum_{i=1}^{n} U^T (x_n - \bar{x}) \left( (x_n - \bar{x}) u \right) \qquad \begin{aligned} &\text{since } U^T \text{ has} \\ &\text{nothing to do} \\ &\text{with data points} \\ &\text{we can rewrite this} \\ &\text{as} \end{aligned}$$

$$U^T \left( \frac{1}{N} \sum_{i=1}^{n} (x_n - \bar{x})(x_n - \bar{x}) \right) u \qquad \text{———} \quad ①$$

This quantity can be Evalued as

If $X = \begin{bmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nD} \end{bmatrix}$ and $\bar{x}$ is $\bar{x}_1, \bar{x}_2, \bar{x}_3 \ldots \bar{x}_D$

$\downarrow$

mean of 1st feature

$$x_n - \bar{x} = \begin{bmatrix} x_{n1} - \bar{x}_1 \\ x_{n2} - \bar{x}_2 \\ \vdots \\ x_{nD} - \bar{x}_D \end{bmatrix}$$

$$\left(x_n - \bar{x}\right)^T = \begin{bmatrix} x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \cdots & x_{nD} - \bar{x}_D \end{bmatrix}$$

$$\left(x_n - \bar{x}\right)\left(x_n - \bar{x}\right)^T = \begin{bmatrix} (x_{n1} - \bar{x}_1)^2 & (x_{n1} - \bar{x}_1)(x_{n2} - \bar{x}_2) & \cdots & (x_{n1} - \bar{x}_1)(x_{nD} - \bar{x}_D) \\ (x_{n2} - \bar{x}_2)(x_{n1} - \bar{x}_1) & (x_{n2} - \bar{x}_2)^2 & & (x_{n2} - \bar{x}_2)(x_{nD} - \bar{x}_D) \\ \vdots & & & \vdots \\ \vdots & & & \vdots \end{bmatrix}$$

Substituting this in ①

$$U^T \left[ \frac{1}{N} \sum_{i=1}^{n} \begin{bmatrix} (x_{n1} - \bar{x}_1)^2 & (x_{n1} - \bar{x}_1)(x_{n2} - \bar{x}_2) & \cdots & (x_{n1} - \bar{x}_1)(x_{nD} - \bar{x}_D) \\ \vdots & & & \vdots \end{bmatrix} \right]$$

Applying this matrix over all data points gives us

$$\begin{bmatrix} \frac{1}{N} \sum_{i=1}^{n} (x_{n1} - \bar{x}_1)^2 & \frac{1}{N} \sum_{i=1}^{n} (x_{n1} - \bar{x}_1)(x_{n2} - \bar{x}_2) & \cdots & \cdots \\ \frac{1}{N} \sum_{i=1}^{n} (x_{n2} - \bar{x}_2)(x_{n1} - \bar{x}_1) & \frac{1}{N} \sum_{i=1}^{n} (x_{n2} - \bar{x}_2)^2 & \cdots & \\ \vdots & & & \end{bmatrix}$$

This resultant matrix is a Co-variance matrix of the term

$$\begin{pmatrix} \text{Var } x_{n1} & \text{Cov}(x_{n1}, x_{n2}) & \cdots & & \text{Cov}(x_{n1}, x_{nD}) \\ \text{Cov}(x_{n2}, x_{n1}) & \text{Var } x_{n2} & \cdots \end{pmatrix}$$

Once the Co-Variance matrix is derived take the Eigen Value decomposition

$$Eign\ (\text{co-variance matrix})$$

This decomposition gives a set of Eigen values $-\lambda$ and Eigen Vectors $- W$

In this decomposition $\lambda$ specifies the variance preserved after projecting original data on the vectors in W.

Order the Eigen vectors by magnitude of the $\lambda$ and pick the highest $\lambda$ and its corresponding Eigen Vector. based on how many principal components you want according to the problem.