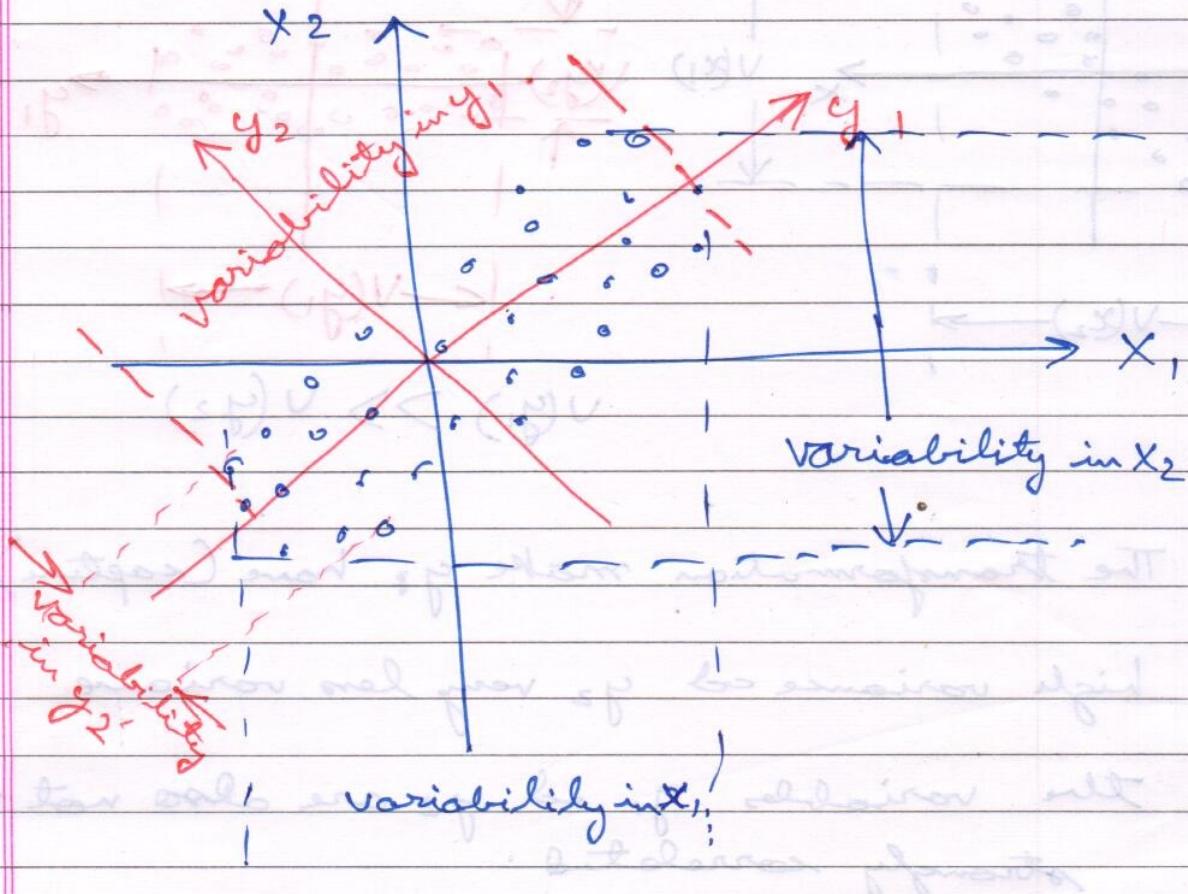


Principal Component Analysis (PCA)

- Linear projection of data so that maximum variance is captured in the lower dimensional representation of data.
- This can help in dimensionality reduction by only considering dimensions which have the highest variance.
- This is also called "Singular value decomposition"



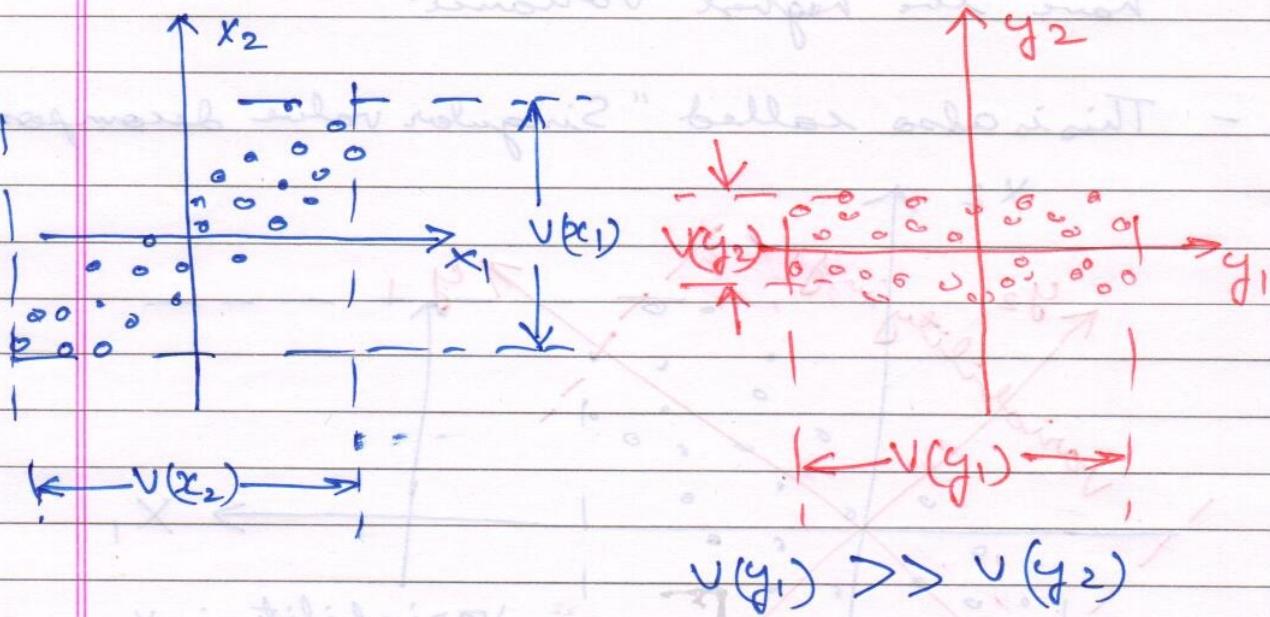
2.

X_1 and X_2 are seen to be strongly correlated, thus

$$\text{Cov}(X) = \begin{bmatrix} 1 & \rho_{12} \\ \rho_{12} & 1 \end{bmatrix}$$

where ρ_{12} is very close to 1.

Space $\xrightarrow{\text{Transformation}}$ Y space



- The transformation make y_1 have (capture) very high variance and y_2 very less variance
- the variables y_1 and y_2 are also not very strongly correlated.

In 2D space (n data points) Transformation
 X Space $\xrightarrow{\quad}$ Y Space

3.

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix} \Rightarrow Y = \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \\ \vdots & \vdots \\ y_{n1} & y_{n2} \end{bmatrix}$$

Variance : $V(y_1) > V(y_2)$

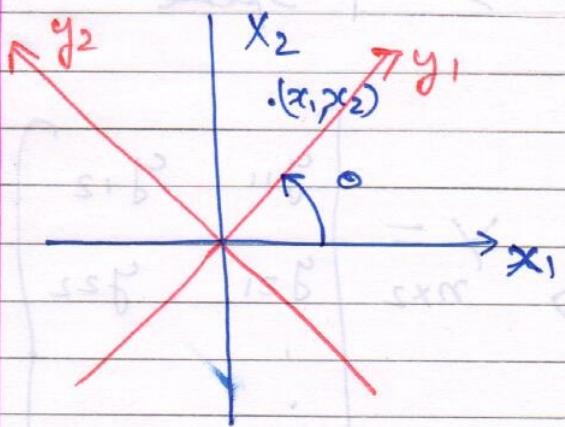
In P-dimensional Space (n data points)

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ x_{21} & x_{22} & \dots & x_{2P} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nP} \end{bmatrix} \Rightarrow Y = \begin{bmatrix} y_{11} & y_{12} & \dots & y_{1P} \\ y_{21} & y_{22} & \dots & y_{2P} \\ \vdots & \vdots & & \vdots \\ y_{n1} & y_{n2} & \dots & y_{nP} \end{bmatrix}$$

For PCA we want:

Variance in Yspace : $V(y_1) \geq V(y_2) \geq \dots \geq V(y_P)$

In 2D Space



The transformation of

X space to Y space in
2D is just a rotation.

$$y_1 = x_1 \cos \theta + x_2 \sin \theta$$

$$\text{and } y_2 = -x_1 \sin \theta + x_2 \cos \theta$$

$$\Rightarrow \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$Y = A^T X$$

$$A^T = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$$

$$\Rightarrow A = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

In 3D space

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = t \cdot A^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Rotation in 3D

$$R_3(\theta) = \begin{bmatrix} \cos\theta & \sin\theta & 0 \\ -\sin\theta & \cos\theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \text{Rotation about } x_3$$

$$R_2(\theta) = \begin{bmatrix} \cos\theta & 0 & \sin\theta \\ 0 & 1 & 0 \\ -\sin\theta & 0 & \cos\theta \end{bmatrix} \Rightarrow \text{Rotation about } x_2$$

$$R_1(\theta) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos\theta & \sin\theta \\ 0 & -\sin\theta & \cos\theta \end{bmatrix} \Rightarrow \text{Rotation about } x_1$$

General Rotation in 3D

$A^T = R = R_1(\alpha) R_2(\beta) R_3(\gamma)$ i.e. a rotation of
 γ about x_3 followed by a rotation of
 β about x_2 followed by a rotation of
 α about x_1

In general 3D rotation you can always find an axis (direction of unit length) that is not changed. If this invariant

axis is $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$ then $x = Rx$

$$\Rightarrow (R - I)x = 0$$

\uparrow
Invariant vector

In a P -dimensional Space

$$y = A^T x$$

In 2D

$$A = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} = [a_1 \ a_2]$$

$$a_1 = \begin{bmatrix} \cos\theta \\ \sin\theta \end{bmatrix} \quad a_2 = \begin{bmatrix} -\sin\theta \\ \cos\theta \end{bmatrix}$$

7.

$$\text{then } \mathbf{a}_1^T \mathbf{a}_1 = [\cos \theta \quad \sin \theta] \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}$$

$$= \begin{bmatrix} \cos \theta & \sin \theta \end{bmatrix} \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}$$

$$= 1$$

$$\text{similarly } \mathbf{a}_2^T \mathbf{a}_2 = 1$$

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

$$\Rightarrow \mathbf{A}^T \mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \mathbf{I} = \mathbf{A} \mathbf{A}^T = \mathbf{A} \mathbf{A}^{-1}$$

$$\Rightarrow \mathbf{A}^T = \mathbf{A}^{-1}$$

Thus \mathbf{A} is an orthogonal transformation

∴ in P -dimensions

$$\text{so } \left[\begin{matrix} 0 & \text{in} & 0 & \text{out} \end{matrix} \right] = \mathbf{D}^T \mathbf{D} \text{ with}$$

$$\mathbf{A} = [a_1 \ a_2 \ \dots \ a_p]$$

$$0^5 \text{ in} + 0^5 \text{ out} =$$

$$\therefore y_1 = a_1^T x = a_{11}x_1 + a_{12}x_2 + \dots + a_{1p}x_p$$

$$y_2 = a_2^T x = a_{21}x_1 + a_{22}x_2 + \dots + a_{2p}x_p$$

$$\left[\begin{matrix} 0 & \text{in} & 0 & \text{out} \end{matrix} \right] = \mathbf{A}^T \mathbf{A} :$$

$$y_j = a_j^T x = a_{j1}x_1 + a_{j2}x_2 + \dots + a_{jp}x_p$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots$$

$$y_P = a_P^T x = a_{P1}x_1 + a_{P2}x_2 + \dots + a_{PP}x_P$$

$$a_i^T a_j = \delta_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

In PCA we want this transformation

such that the variance $V(y_i)$ for y_i is as

follows:

$$V(y_1) \geq V(y_2) \geq \dots \geq V(y_P)$$

Thus, in PCA we want a transformation

matrix A^T such that we maximize the

variance of the transformed space as given

below:

$$V(y_1) \geq V(y_2) \geq \dots \geq V(y_p)$$

p -dimensional data with n number of points

This can be represented in two ways

Matrix form

$$x^p = [x_1 \ x_2 \ \dots \ x_n] \quad n \times p$$
$$x = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \quad n \times p$$

Set form $X = \{x_1, x_2, \dots, x_n\}$ where

x_k is the element of the set and is a

p -dimensional feature vector (ordered p lot)

$$x_k = (x_{k1}, x_{k2}, x_{k3}, \dots, x_{kp})$$

∴ For the "kth" set

10.

$$\begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{kp} \end{bmatrix}^T = [x_{k1}, x_{k2}, \dots, x_{kp}] = x_k$$

↑
kth element of
the set representation

X vector Row representation

↓
X vector Column representation

↑
kth component of the vector

(1xP) × ... × (1xP) × (1xP)

Similarly

$$\begin{bmatrix} y_{k1} \\ y_{k2} \\ \vdots \\ y_{kp} \end{bmatrix}^T = [y_{k1}, y_{k2}, \dots, y_{kp}] = y_k$$

1xP

map vector Y

$$\begin{bmatrix} y_{k1} \\ y_{k2} \\ \vdots \\ y_{kp} \end{bmatrix} = A^T \begin{bmatrix} x_{k1} \\ x_{k2} \\ \vdots \\ x_{kp} \end{bmatrix}$$

By taking Transpose on both the sides:

11.

$$\Rightarrow \begin{bmatrix} y_{k_1} \\ y_{k_2} \\ \vdots \\ \vdots \\ y_{kp} \end{bmatrix}^T = A^T \begin{bmatrix} x_{k_1} \\ x_{k_2} \\ \vdots \\ \vdots \\ x_{kp} \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} y_{k_1} \\ y_{k_2} \\ \vdots \\ \vdots \\ y_{kp} \end{bmatrix}^T = \begin{bmatrix} x_{k_1} \\ x_{k_2} \\ \vdots \\ \vdots \\ x_{kp} \end{bmatrix}^T A$$

$$\Rightarrow y_k = x_k A$$

↑
k-th element of set representation

or row vector representation

$$\therefore Y = A^T X$$

↑ as column vector representation

$$is \text{ same as } y_k = x_k A$$

↑ as row vector representation

For PCA the X data set is not only rotated but also translated w.r.t. mean. Thus,

$$y_k = \underset{1 \times p}{(x_k - \bar{x})} \cdot \underset{1 \times p}{A} \underset{p \times p}{}$$

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k = \frac{1}{n} [\Sigma x_{k1} \ \Sigma x_{k2} \dots \ \Sigma x_{kp}]$$

$$\Rightarrow x_k = y_k A^T + \bar{x}$$

Now variance vector for Y and combination of variance and covariance of Y is

$$\text{Variance vector of } Y = [\text{var}(y_1) \ \text{var}(y_2) \ \dots \ \text{var}(y_p)]$$

$$\text{Covariance of } Y = \Sigma_y = \frac{1}{n-1} \sum_{k=1}^n y_k^T y_k$$

$A \in \mathbb{R}^{p \times p}$ as we know
 $b \times p$

$b \times 1 \quad 1 \times p$

$$\Rightarrow \Sigma_y = \frac{1}{n-1} \sum_{k=1}^n y_k^T y_k$$

$$= \frac{1}{n-1} \sum_{k=1}^n [(x_k - \bar{x}) \cdot A]^T [(x_k - \bar{x}) \cdot A]$$

$$= \frac{1}{n-1} \sum_{k=1}^n A^T (x_k - \bar{x})^T \cdot (x_k - \bar{x}) \cdot A$$

$$= A^T \left[\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^T \cdot (x_k - \bar{x}) \right] \cdot A$$

$$\Rightarrow \Sigma_y = A^T C A$$

↑ Covariance matrix of X

$$C_{ij} = \frac{1}{n-1} \sum_{k=1}^n [x_{ki} - \frac{1}{n} \sum_{l=1}^n x_{li}] \cdot [x_{kj} - \frac{1}{n} \sum_{l=1}^n x_{lj}]$$

$i, j = 1, 2, \dots, p$

Now we want to maximize Σ_y by

choosing the optimal A .

However, we cannot choose any A ,
 since if it has to have the following
 constrain:

$$[A \cdot (\mathbb{E} - \lambda I)]^T [A \cdot (\mathbb{A}^T A - I)] = 1 - \lambda$$

To solve this optimization problem
 with a constrain, we construct a
 Lagrange function:

$$L = A^T C A - \lambda (A^T A - I)$$

and maximize the value of L .

$$\Rightarrow \frac{\partial L}{\partial A} = 0$$

$$\Rightarrow C A - \lambda A = 0$$

$$\Rightarrow (C - \lambda I) \cdot A = 0$$

This is an Eigenvalue Eigenvector Problem

The Eigenvalue problem can be traced

back to solving the following DE:

$$\frac{dy}{dt} = Ay \quad \left\{ \begin{array}{l} n \text{ linear DEs} \\ \uparrow \\ n \times n \text{ matrix} \end{array} \right.$$

This will have a solution of type

$$y(t) = e^{\lambda t} X$$

$$\Rightarrow \lambda e^{\lambda t} X = A e^{\lambda t} X$$

$$\Rightarrow \lambda X = A X$$

$$\Rightarrow (A - \lambda I) X = 0 \quad \left\{ \begin{array}{l} \text{This gives } n \\ \text{eigenvalues of } \lambda \\ \text{and } n \text{ eigenvectors} \end{array} \right.$$

$$\begin{aligned} \text{Eg. } y_1' &= 5y_1 + y_2 & \left\{ \begin{array}{l} A = \begin{bmatrix} 5 & 1 \\ 3 & 3 \end{bmatrix} \end{array} \right. \\ y_2' &= 3y_1 + 3y_2 & \end{aligned}$$

this gives two eigenvalues and two eigenvectors

$$\lambda_1 = 6 \quad x_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}; \quad \lambda_2 = 2 \quad x_2 = \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

16.

Thus general solution of the DE is

$$y(t) = \begin{bmatrix} y_1(t) \\ y_2(t) \end{bmatrix} = C_1 e^{\lambda_1 t} x_1 + C_2 e^{\lambda_2 t} x_2$$

$$\Rightarrow y(t) = C_1 e^{6t} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + C_2 e^{2t} \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

If we put the B.C.s we get

$$\text{at } t=0; y(0) = C_1 \begin{bmatrix} 1 \\ 1 \end{bmatrix} + C_2 \begin{bmatrix} 1 \\ -3 \end{bmatrix}$$

Now note if we know eigenvector for

$Ax = \lambda_1 x$ then we also know

the eigenvector for $A^2 x = \lambda_2 x$

\Rightarrow This is because $A^2 x = \lambda_2 x$

$$\Rightarrow A(\lambda_1 x) = \lambda_2 x$$

$$\Rightarrow \lambda_1^2 x = \lambda_2 x$$

\Rightarrow they have the same eigenvector

$$\text{by } \lambda_2 = (\lambda_1)^2$$

In the same manner for the equation

$$(A + cI)x = \lambda_3 x$$

$$\lambda_3 = (\lambda_1 + c)$$

and similarly

$$A^n x = \lambda_n x \text{ has same}$$

for eigenvector but $\lambda_n = (\lambda_1)^n$.

Note that in eigenvalue problem x is

an $n \times 1$ matrix (or a vector)

ie. $(A - \lambda I)x = 0$

$\begin{matrix} n \times n & n \times n & n \times 1 \end{matrix}$

whereas in PCA it has a square matrix form

$$(C - \lambda I)A = 0$$

$\begin{matrix} p \times p & p \times p & p \times p \end{matrix}$

Thus, A is a concatenation of the eigenvectors

$$\text{of } C \Rightarrow A = (a_1, a_2, \dots, a_p)$$

The variance in \hat{Y} corresponds to

the eigenvalue of $\lambda_1, \lambda_2 \dots \lambda_p$.

This is because

$$CA = \lambda A$$

$$\Rightarrow \lambda = A^T C A = \sigma_y^2$$

Thus, if we want to capture 95%

variance, then we can reduce the

number of dimensions to q such that

$$\frac{\sum_{i=1}^q \lambda_i}{\sum_{i=1}^p \lambda_i} \geq \frac{95}{100}$$

For example take $p=2$ and $n=4$ such

$$\text{that } X = \{(1,1), (2,1), (2,2), (3,2)\}$$

then $\bar{x} = \left(2, \frac{3}{2}\right)$ is the mean matrix

and the covariance matrix is

$$\text{Cov} = \frac{1}{3} \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}$$

Thus, the eigenvalues are

$$\lambda_1 = 0.8727$$

$$\text{and } \lambda_2 = 0.1273$$

This corresponds to

$$a_1 = \begin{bmatrix} -0.85065 \\ -0.52571 \end{bmatrix}$$

$$a_2 = \begin{bmatrix} 0.52571 \\ -0.85065 \end{bmatrix}$$

$$\text{and } A = [a_1 \ a_2]$$

$$\Rightarrow A = \begin{bmatrix} -0.85065 & 0.52571 \\ -0.52571 & -0.85065 \end{bmatrix}$$

\therefore Transformation of X by A^T rotation

matrix would lead to Y as given by

$$Y = A^T(X - \bar{X}) \quad \left| \begin{array}{l} \text{In column vector representation} \\ \text{In row vector representation} \end{array} \right. \quad y_k = (x_k - \bar{x}) \cdot A$$

$$\Rightarrow X = \{(1, 1), (2, 1), (2, 2), (3, 2)\}$$

will transform to

$$y_{n=1} = \begin{bmatrix} -0.85065 & -0.52571 \\ 0.52571 & -0.85065 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ 1 & -1.5 \end{bmatrix} \begin{bmatrix} x - \bar{x} \end{bmatrix}$$

$$= \begin{bmatrix} 1.1135 \\ -0.10039 \end{bmatrix}$$

$$y_{n=2} = \begin{bmatrix} A^T \end{bmatrix} \begin{bmatrix} 2 & -2 \\ 1 & -1.5 \end{bmatrix} = \begin{bmatrix} 0.26286 \\ 0.42533 \end{bmatrix}$$

similarly $y_{n=3} = \begin{bmatrix} -0.26286 \\ -0.42533 \end{bmatrix}$ and $y_{n=4} = \begin{bmatrix} -1.1135 \\ 0.10039 \end{bmatrix}$

$\therefore X = \{(1, 1), (2, 1), (2, 2), (3, 2)\}$ is transformed to

$$Y = \{(1.1135, -0.10039), (0.26286, 0.42533), (-0.26286, -0.42533), (-1.1135, 0.10039)\}$$

Without dimensionality Reduction

$$X \xrightarrow{A^T} Y \quad A^T = \begin{bmatrix} -0.85065 & -0.525 \\ -0.52571 & -0.8506 \end{bmatrix}$$

$$\{(1,1), (2,1), (2,2), (3,2)\} \quad \left\{ \begin{array}{l} (1.1135, -0.10039), \\ (0.26286, 0.42533), \\ (-0.26286, -0.42533), \\ (-1.1135, -0.10039) \end{array} \right.$$

Mean X

$$= (2, 1.5)$$

Mean Y

$$= (0, 0)$$

Covariance of X

$$= \begin{bmatrix} \frac{2}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

Covariance of Y

$$= \begin{bmatrix} 0.87266 & 0 \\ 0 & 0.12732 \end{bmatrix}$$

$$= \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

Covariance is a diagonal matrix

$$\{(1,1), (2,1), (2,2), (3,2)\}$$

$$A \quad Y$$

The inverse transform results in the same X data.

multibled dimensionality reduction

With dimensionality reduction

Since the variance of y_1 is 0.87266

and y_2 is 0.12732, thus if we take only

y_1 dimension then we get 0.87266 i.e.

87% of variance.

Thus we do not take $A = [a_1 \ a_2]$

but instead $A' = [a_1 \ 0] = \begin{bmatrix} -0.85065 \\ -0.52571 \end{bmatrix}$

$$\Rightarrow X \xrightarrow{(A')^T} Y'$$

$$(A')^T = \begin{bmatrix} -0.85065 & -0.52571 \\ 0 & 0 \end{bmatrix}$$

$$X \xrightarrow{\quad\quad\quad} Y'$$

$$X = \begin{cases} (1, 1) \\ (2, 1) \\ (2, 2) \\ (3, 2) \end{cases} \quad Y' = \begin{cases} (1.11351, 0) \\ (0.26286, 0) \\ (-0.26286, 0) \\ (-1.11351, 0) \end{cases}$$

The inverse transform of Y' with A' gives

$$X' = \begin{cases} (1.0528, 0.91482) \\ (1.7764, 1.36181) \\ (2.2236, 1.63819) \\ (2.9472, 2.08538) \end{cases}$$