

$$\text{If } N=2$$

$$f_{x_1, x_2}(x_1, x_2) = f_{S_1, S_2}(\omega_1^T \underline{x} \quad \omega_2^T \underline{x}) |\omega|$$

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\omega = \begin{bmatrix} \omega_1^T \\ \omega_2^T \\ \omega_N^T \end{bmatrix}$$

$$\omega_1^T = [\omega_{11}, \omega_{12}]$$

∴

$$f_{x_1, x_2}(x_1, x_2) = f_{S_1}(\omega_1^T \underline{x}) \cdot f_{S_2}(\omega_2^T \underline{x}) |\omega|$$

$$= f_{S_1}(S_1) \cdot f_{S_2}(S_2) |\omega|$$

$$S_1 = \omega_{11} x_1 + \omega_{12} x_2$$

$$S_2 = \omega_{21} x_1 + \omega_{22} x_2$$

{ capital $S \in x$
i.e. RV's }

$$S_1 = \omega_{11} x_1 + \omega_{12} x_2$$

{ value of RV }

$$\begin{bmatrix} S_1^{(1)} \\ S_2^{(1)} \end{bmatrix} = \begin{bmatrix} \omega_{11} & \omega_{12} \\ \omega_{21} & \omega_{22} \end{bmatrix} \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \end{bmatrix}$$

In order to solve for ω , I will try to maximize $f_{x_1, x_2}(x_1, x_2)$ by choosing appropriate $\begin{bmatrix} \omega_1^T \\ \omega_2^T \end{bmatrix} = \omega$

$$f_{x_1, x_2}(x_1, x_2) = f_{S_1}(\omega_1^T \underline{x}) f_{S_2}(\omega_2^T \underline{x}) |\omega|$$

known x_1, x_2 & assumed some ω

So we have to assume some density f^n of S

The log likelihood function is

$$l(W) = \sum_{i=1}^m \left(\sum_{j=1}^N \log g'(w_j^T x^{(i)}) \right) + \log |W|$$

m = no of samples.

$$x_1 = x_1^{(1)} \dots x_1^{(m)}$$

$$x_2 = x_2^{(1)} \dots x_2^{(m)}$$

In order to maximize $l(w)$ for all w 's $w_1 \dots w_N$
we can use stochastic gradient descent.

No. of values of w to be determined are $N \times N$

Or we can minimize $-l(w)$

In order to use stochastic

we use
$$l(w) = \left(\sum_{j=1}^N \log g'(w_j^T x^{(i)}) \right) + \log |W|$$

if $N=2$

$$W_{\text{new}} = W_{\text{old}} + \alpha \nabla_W [l(w)]$$

If we consider stochastic gradient descent one
example at a time can be considered

let us consider differentiation wrt $w_{11}, w_{12}, w_{21}, w_{22}$
where $N=2$

$$\sum_{j=1}^N \log g'(w_j^T x^{(i)}) + \log |W|$$

$$= \log g(w_1^T x^{(i)}) + \log g'(w_2^T x^{(i)}) + \log |w|$$

\downarrow $s_1^{(i)}$
 \downarrow $s_2^{(i)}$

$$g(s) = \frac{1}{1 + e^{-s}}$$

$$g'(s) = g(s)(1 - g(s))$$

$$= \log [g(w_1^T x^{(i)}) (1 - g(w_1^T x^{(i)}))] + \log [g(w_2^T x^{(i)}) (1 - g(w_2^T x^{(i)}))] + \log |w|$$

Now differentiation wrt w_{11}

$$(1 - 2g(s_1^{(i)})) x_1^{(i)}$$

** The Gradient are parallel at optimum point

Unsupervised Learning

Group the cluster having similar property.

eg: Amazon & groups people according to their likes.
Segmentation in image processing is unsupervised.
label of cluster is not given

04/11/19

MIDSEM PAPER.

Q2. $E(\hat{f}) = f$

$X(R.V)$, m (mean), σ^2 true parameters

$$M_n = \text{estimated mean} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

n = no. of samples

all x_i are iid R.V. (Assumed)

M_n taken over n examples & many times value of M_n varies but in comparable range (M_n is a R.V.)

$$E[M_n] = E\left[\frac{x_1 + x_2 + \dots + x_n}{n}\right]$$

$$= \frac{1}{n} [E[x_1] + E[x_2] + \dots + E[x_n]]$$

$$= \frac{1}{n} [m + \dots + m]$$

$$= \frac{n \cdot m}{n} = m$$

$E(\text{estimate} - \text{mean}) = \text{true mean} = \text{unbiased estimate}$

$$V_m = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\bar{x} = \frac{1}{n} [x_1 + \dots + x_n]$$

$$\text{Var}(M_n) = \text{Var}\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right)$$

X and $Y \rightarrow R.V$

$\sigma_1^2 \quad \sigma_2^2$

$$\text{var}(X+Y) = \sigma_1^2 + \sigma_2^2 \quad (\text{if they are uncorrelated})$$

$$\text{var}(M_n) = \text{var}\left(\frac{x_1}{n}\right) + \text{var}\left(\frac{x_2}{n}\right) + \dots + \text{var}\left(\frac{x_n}{n}\right)$$

$$= \frac{1}{n^2} [\sigma^2 + \dots + \sigma^2]$$

$$= \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$$

large N (sample size), gives a good estimate,
variance becomes close to 0.

$$\left\{ \begin{array}{l} V_n = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \bar{x} = \frac{1}{n} [x_1 + \dots + x_n] \\ V_n = R.V \\ \left\{ \begin{array}{l} \text{Whether } E[V_n] = \sigma^2? \\ \Rightarrow \text{Ans is no.} \end{array} \right. \end{array} \right. \rightarrow \text{Biased estimate}$$

In order to make

$$\left\{ \begin{array}{l} \text{if } V_n = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \end{array} \right.$$

$$\text{then } E[V_n] = \sigma^2$$

$\rightarrow V_n$ is an unbiased estimate of true variance.

PCA

$$y = Ax$$

$$A^{-1} = A^T$$

Parseval theorem

$$y^T y = (Ax)^T (Ax)$$

energy of transformed value

$$= x^T A^T A x$$

= energy of input value

$$= x^T x$$

energy in y is over few values.

$$\|y\|^2 = \|x\|^2$$

but in x it is distributed over all values.

08/11/19

Date: / /

Transformation matrix depends on N and not on input data.
 $N \times N$.

- DCT / DFT

DCT has better energy compaction than DFT.

If $N \times N$ is size of A & y has size of N

Ex: x : size 10,000 } with transformation y : has
 y : size 10,000 } $256 \rightarrow$ values
 $\hat{x} = A^T y$ } non zero.
 & rest values are 0.

Error: $\|\hat{x} - x\|^2$

In DCT matrix A have
 real entries.

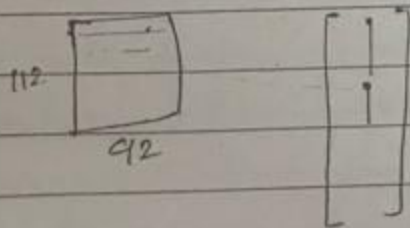
Error DCT < Error DFT

All the compressions are lossy

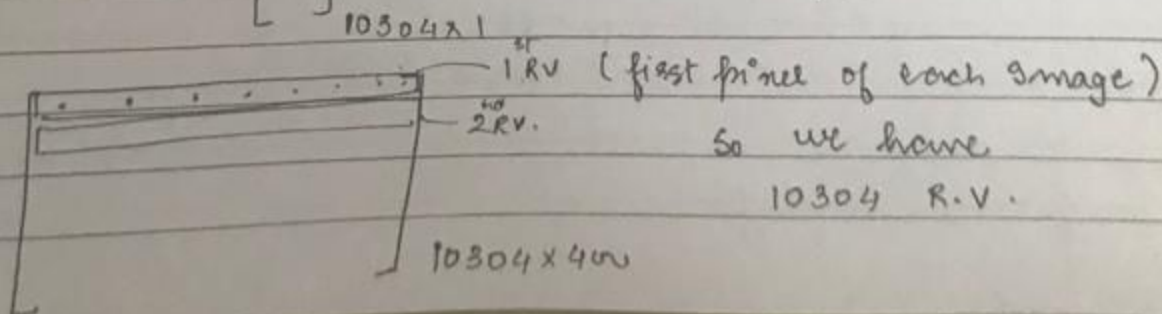
PCA Principal Component Analysis

Transformation matrix is derived from input data
 Rows of transformation matrix correspond to
 eigen vectors (orthogonal directions) obtained from
 covariance matrix of input data.

First row correspond to highest value of eigen value
 & second row to second highest or second Principal
 Component.



Input data matrix
 for 400 images.
 Size 10304×400



So we have

10304 R.V.

$$\text{cov. matrix} = 10304 \times 10304$$

$$\text{each eigen vector size} = 10304 \times 1$$

$$\text{no. of eigen value} = 10304$$

Suppose we retain 100 eigen vectors.

$$Y = AX$$

↳ rows of A are the eigen vectors of covariance matrix of X. $= C_X$

08/11/19

$$X = A^T Y$$

$$X = \begin{bmatrix} \text{columns} \\ \text{of} \\ C_X \end{bmatrix} \begin{bmatrix} Y \end{bmatrix} \begin{matrix} 10304 \times 400 \\ 100 \end{matrix}$$

~~10304~~ 10304×100

100 values as

$$X = 10304 \times 400$$

$$C_X = 10304 \times 10304$$

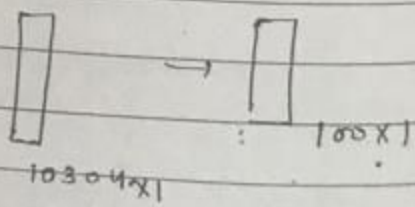
now by retaining only 100 eigen vector from 10304 eigen vectors

as you retain more & more you get better performance (lower error)

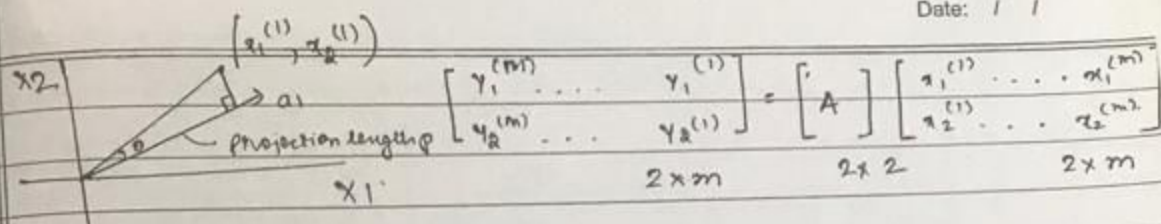
$$Y = AX \begin{matrix} 10304 \times 400 \\ 10304 \times 10304 \\ 10304 \times 400 \end{matrix}$$

↓ after reduction

$$Y_{100 \times 400} = \begin{pmatrix} 100 \times 10304 \end{pmatrix} \cdot (10304 \times 400)$$



Dimensionality Reduction



$$a_1 = \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix}$$

a_1 is the vector ^{on} which projected.

m : no of examples.

$x^{(1)} \dots x^{(m)}$

\downarrow
 $x_1^{(1)}, x_2^{(1)}$

$$\underline{x^{(1)}} \cdot \underline{a_1} = |\underline{x^{(1)}}| |\underline{a_1}| \cos \theta$$

$$\cos \theta = \frac{p}{|\underline{x^{(1)}}|}$$

$$\underline{x^{(1)}} \cdot \underline{a_1} = |\underline{a_1}| p$$

$$p = \frac{\underline{x^{(1)}} \cdot \underline{a_1}}{|\underline{a_1}|}$$

If $a_1 = \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix}$ has unit norm then

$$p = \underline{x^{(1)}} \cdot \underline{a_1}$$

$$\therefore p = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} \end{bmatrix} \cdot \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix}$$

$$p = a_{11} x_1^{(1)} + a_{12} x_2^{(1)} = y_1^{(1)}$$

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

projection of $x^{(1)} \cdot a_2 = y_2^{(1)}$

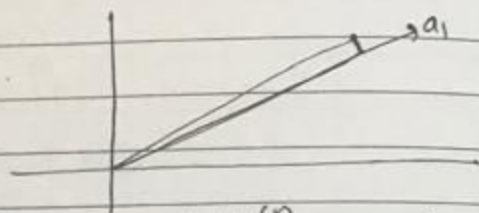
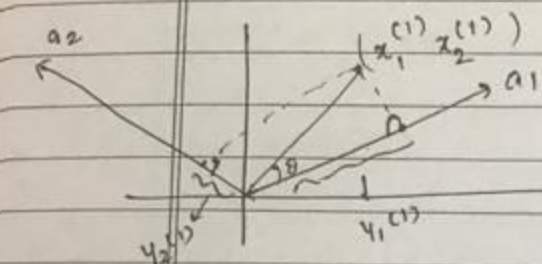
$$a_2 = \begin{bmatrix} a_{21} \\ a_{22} \end{bmatrix}$$

PCA is also called as change of basis

$$\begin{bmatrix} y_1^{(1)} \\ y_2^{(1)} \end{bmatrix} = \underbrace{\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}}_{\text{orthonormal}} \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \end{bmatrix}$$

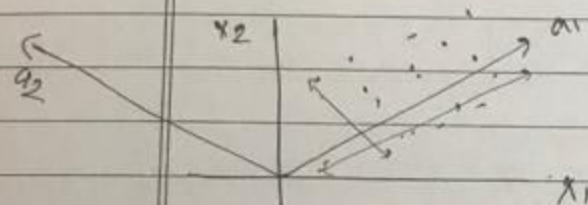
Taken the first example

orthonormal



$\therefore y_2^{(k)}$ can be ignored but not $y_1^{(k)}$

i.e. we are doing dimensional reduction.



a_1 is a measure of x_1 & giving more importance to 1st dimension.

Retaining $y_1^{(k)}$ is not giving much error compared to $y_2^{(k)}$

The ~~var~~ variance is high in the dirⁿ of a_1 because the points are varying high above till a_1 .

= whenever you look at variance
Take the high variance, & give that dimension importance.

Variance of the first row (y) has to be maximum

$$a_1 = \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} \text{ has unit norm}$$

$$\underline{y} = \underline{A} \underline{x}$$

$$\text{var}(\underline{a}_1^T \underline{x}) = \text{var}(y_1)$$

This has to be maximized.

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \leftarrow y \rightarrow \\ \leftarrow y \rightarrow \end{bmatrix} = \begin{bmatrix} \underline{A} \end{bmatrix} \begin{bmatrix} \underline{x} \end{bmatrix}$$

$$\text{var}(\underline{a}_2^T \underline{x}) = \text{var}(y_2)$$

$$\text{var}(\underline{a}_1^T \underline{x}) = \text{var}(y_1)$$

This has to be maximized under the constraint that s.t. $\|\underline{a}_1\|^2 = 1$

PDF.

$$\underline{A} = \text{sq. matrix} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix}$$

Real entries ; Both rows & columns are orthonormal.
; rows column vector have unit norm

$$R_1 \cdot R_2 = 0 \quad C_1 \cdot C_2 = 0$$

If complex entries $\underline{A} = \text{unitary matrix}$

$$\underline{y} = \underline{A} \underline{x}$$

$$\underline{x} = \underline{B} \underline{y} \quad \underline{B} = \underline{A}^{-1} = \underline{A}^T$$

$$N=2$$

$$\begin{bmatrix} y_0 \\ y_1 \end{bmatrix} = \begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \end{bmatrix}$$

$$\begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = y_0 \begin{bmatrix} a_{00} \\ a_{10} \end{bmatrix} + y_1 \begin{bmatrix} a_{01} \\ a_{11} \end{bmatrix}$$

$$a_{00} a_{10} + a_{01} a_{11} = 0$$

$$\boxed{\underline{A}^T \underline{A} = \underline{I} \quad \text{if } \underline{A} \text{ is orthogonal}}$$

$$\boxed{\underline{A}^{-1} = \underline{A}^{*T}}$$

↓
unitary matrix

$$y = Ax$$

$$y = T[x] \quad T \text{ is linear}$$

Analysis : $y = Ax$ $A^{-1} = A^T$
 Synthesis : $x = A^T y$

Why orthogonal transforms

Tend to redistribute energy.

Most of it in few transformed values

Orthogonality also preserves energy

Ex Energy $y^T y$ (final) $x^T x$ (Initial)

$$y = Ax \quad y^T = (Ax)^T = x^T A^T$$

$$y^T y = x^T A^T A x$$

$$= x^T x$$

Correlation : Gradual variations with occasional discontinuities

Correlated data has redundancy i.e. more samples than what is required to represent

Defun correlation : Transformation using orthogonal matrix A results in uncorrelated values of y

Discrete Cosine Transform (DCT)

$$y(k) = C(k) \cdot \sum_{n=0}^{N-1} x(n) \cos \left[\frac{\pi(2n+1)k}{2N} \right]$$

$$x(n) = \sum_{k=0}^{N-1} C(k) y(k) \cos \left[\frac{\pi(2n+1)k}{2N} \right]$$

$$C(k) = \begin{cases} \sqrt{\frac{1}{N}} & k=0 \\ \sqrt{\frac{2}{N}} & \text{for } k \neq 0 \end{cases}$$

N=2.

↓ fixed for N

$$\begin{bmatrix} y(0) \\ y(1) \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} x(0) \\ x(1) \end{bmatrix} \quad y = Ax$$

$$\begin{bmatrix} x(0) \\ x(1) \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} y(0) \\ y(1) \end{bmatrix} \quad a = A^T y$$

DFT (unitary)

$$y(k) = \frac{1}{\sqrt{N}} \sum_{n=0}^{N-1} x(n) e^{-j \frac{2\pi}{N} nk} \quad k = 0, 1, \dots, N-1$$

$$x(n) = \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} y(k) e^{j \frac{2\pi}{N} nk} \quad n = 0, 1, \dots, N-1$$

N=4

$$\begin{bmatrix} y(0) \\ y(1) \\ y(2) \\ y(3) \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -j & -1 & j \\ 1 & -1 & 1 & -1 \\ 1 & j & -1 & -j \end{bmatrix} \begin{bmatrix} x(0) \\ x(1) \\ x(2) \\ x(3) \end{bmatrix}$$

Advantages of DCT over DFT

- Basis are real $\Rightarrow A$ is real \Rightarrow Less computation complexity
- \Rightarrow Better energy compaction
- Less error in DCT, by taking same no of vectors

What is PCA?

- Optimum transform (in MSE sense), better than DFT & DCT

- The transformation matrix is derived from input data (X)
- Rows of transformation matrix correspond to eigen vectors obtained from covariance matrix of input data.

in data on PCA

$$X = (112 \times 92) \times 400 \\ = 10304 \times 400$$

Images = 400 pixels = 10304 = Random variable

Covariance matrix = $C_x = 10304 \times 10304$
10304 eigen values (λ_i) & eigen vectors of
length 10304×1

\Rightarrow eigen vectors are used as rows of transformation matrix A
(first row with highest eigen value)
First few rows represent principal components
eigen vectors are orthonormal

$$A = 10304 \times 10304 \quad X = 10304 \times 400 \quad Y = 10304 \times 400$$

Let us retain 100 eigen vectors

$$A = 100 \times 10304$$

$$A^T = 10304 \times 100 \quad X = 10304 \times 400 \quad Y = 100 \times 400$$

retain only top 100 rows of Y

Here retaining 400 rows gives minimum error (Hit & Trail)

\Rightarrow each image which had 10304 pixels can now be represented using only 400 values

$$m = 400$$

$$n = 10304$$

Minimum the reconstruction error,
the variance will increase

P.C.A.

$$\begin{matrix} x_1 \\ x_2 \end{matrix} \rightarrow \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(m)} \\ x_2^{(1)} & \dots & x_2^{(m)} \end{bmatrix} \xrightarrow{A} \begin{matrix} y_1 \\ y_2 \end{matrix} \rightarrow \begin{bmatrix} y_1^{(1)} & \dots & y_1^{(m)} \\ y_2^{(1)} & \dots & y_2^{(m)} \end{bmatrix}$$

$2 \times m$ 2×2 $2 \times m$

variance decreases as we go down the Y matrix.

m : examples $x^{(i)}$

n : subscript(x) no of features

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$Y = AX$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$y_1^{(i)} = a_{11} x_1^{(i)} + a_{12} x_2^{(i)} = a_1^T x^{(i)}$$

$$y_2^{(i)} = a_{21} x_1^{(i)} + a_{22} x_2^{(i)} = a_2^T x^{(i)}$$

$$a_1 = \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} \quad x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \end{bmatrix} \quad a_2 = \begin{bmatrix} a_{21} \\ a_{22} \end{bmatrix}$$

$$\text{var}(y_1) = \text{var}(a_1^T x)$$

$$\begin{bmatrix} y_1^{(1)} & \dots & y_1^{(m)} \\ y_2^{(1)} & \dots & y_2^{(m)} \end{bmatrix} = \begin{bmatrix} A \\ A \end{bmatrix} \begin{bmatrix} x \\ x \end{bmatrix}$$

$2 \times m$

$$\begin{pmatrix} y_1^{(1)} \\ y_2^{(1)} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x_1^{(1)} \\ x_2^{(1)} \end{pmatrix}$$

$$a_1 = \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix}$$

$$a_1^T = [a_{11} \quad a_{12}]$$

working on 1 example

$$y = Bx \quad \text{where } B = A^T$$

$$Y = AX$$

$\text{var}(a_1^T x)$ can be written in terms of covariance of x

Covariance matrix of \underline{x}

when $n=2$

$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ x_2 & x_1 now vector of m values

$$C_x = \text{cov}(\underline{x}) = \begin{bmatrix} E(x_1 - \overset{\sigma_{x_1}^2}{m x_1})^2 & \text{cov}(x_1, x_2) \\ \text{cov}(x_1, x_2) & E(x_2 - m x_2)^2 \end{bmatrix}_{2 \times 2}$$

$$m x_1 = \frac{x_1^{(1)} \dots x_1^{(m)}}{m}$$

$x_1 - m x_1$ = from each value you subtract $m x_1$

Sq. each of them & divide by m = sample variance

$$\text{cov}(x_1, x_2) = \text{cov}(x_2, x_1) = E \left[\underbrace{(x_1 - m x_1)}_{\substack{\downarrow \\ \text{element wise multiply}}} \underbrace{(x_2 - m x_2)}_{\substack{\downarrow \\ \text{element wise multiply}}} \right]$$

$$\left[\begin{array}{l} C_x = \text{semi +ve definite} \\ \lambda_i \geq 0 \\ \text{Symmetric} \end{array} \right.$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$E[y_1] = a_{11} m x_1 + a_{12} m x_2$$

$$E[y_2] = a_{21} m x_1 + a_{22} m x_2$$

$$E \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = A \underline{m x}$$

$$m x = \begin{bmatrix} m x_1 \\ m x_2 \end{bmatrix}$$

$$\underline{C}_Y = \text{Cov. matrix of } \underline{Y}$$

$$= E \left[\underset{2 \times 1}{(\underline{Y} - \underline{m}_Y)} \underset{1 \times 2}{(\underline{Y} - \underline{m}_Y)^T} \right] \quad = 2 \times 2 \text{ matrix}$$

$$\underline{Y} = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} \quad \underline{m}_Y = \begin{bmatrix} m_{Y1} \\ m_{Y2} \end{bmatrix}$$

Y_1 & Y_2 are R.V.

$$E \underline{C}_X = E \left[\underset{2 \times 1}{\begin{pmatrix} (X_1 - m_{X1}) \\ (X_2 - m_{X2}) \end{pmatrix}} \underset{1 \times 2}{\begin{pmatrix} (X_1 - m_{X1}) & (X_2 - m_{X2}) \end{pmatrix}} \right] \quad = 2 \times 2$$

$$\begin{aligned} C_Y &= E(Y Y^T) \\ &= E[(Y - m_Y)(Y - m_Y)^T] \\ &= E[(A \underline{X} - A \underline{m}_X)(A \underline{X} - A \underline{m}_X)^T] \\ &= E[A(\underline{X} - \underline{m}_X)(\underline{X} - \underline{m}_X)^T A^T] \end{aligned}$$

$$E C_Y = A C_X A^T$$

diagonal

first element in $C_Y = \sum Y^{(1)2}$

$$\max \text{var}(a_1^T \underline{X}) = \text{var}(Y_1)$$

$$\max \begin{bmatrix} C_Y \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} C_X \end{bmatrix} \begin{bmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{bmatrix}$$

$$\max (a_1^T C_X a_1) \text{ s.t. } \|a_1\|^2 = 1$$

$$L(a_1, \lambda) = a_1^T C_X a_1 - \lambda [\|a_1\|^2 - 1]$$

$$\downarrow$$

diff w.r.t a_1 & equate to 0

$$= 2 C_X a_1 - 2 \lambda a_1$$

$$C_X a_1 = \lambda a_1$$

$$C_y = A C_x A^T$$

$$= A (A^T \lambda A) A^T$$

$$= A A^T \lambda A A^T$$

$$A^T = A^{-1}$$

$$C_y = \lambda$$

$$[\text{new data } Y]_{k \times n} = (\text{top } k)_{k \times m} (\text{orig})_{m \times n}$$

13/11/19 PCA continued.

$$Y = A X$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

↓

Output R.V.

↘ Input R.V.

$$\begin{array}{l} x_1 \text{ takes } \begin{pmatrix} x_1^{(1)} \\ x_2^{(1)} \end{pmatrix} \dots \begin{pmatrix} x_1^{(m)} \\ x_2^{(m)} \end{pmatrix} \\ x_2 \text{ takes } \end{array}$$

example

$$\begin{pmatrix} y_1^{(1)} \\ y_2^{(1)} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x_1^{(1)} \\ x_2^{(1)} \end{pmatrix}$$

If you apply each data vector we get corresponding y .
(for a 2D case)

$$\text{Given data set } \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(m)} \\ x_2^{(1)} & \dots & x_2^{(m)} \end{bmatrix}$$

$$\text{we arrive at } A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \cdot \text{It can be } nD.$$

The problem can be solved as follows

- use constraint optimization to A .

$$a_1 = \begin{bmatrix} a_{11} \\ a_{12} \end{bmatrix} \quad a_1^T = [a_{11} \quad a_{12}]$$

$$a_2 = \begin{bmatrix} a_{21} \\ a_{22} \end{bmatrix} \quad a_2^T = [a_{21} \quad a_{22}]$$

Now Variance of $a_1^T x$ in order to get y_1 's variance
 $\text{var}(a_1^T x) = \text{var}(y_1)$

$$y = Ax$$

$$C_y = A C_x A^T$$

$$\begin{bmatrix} \text{var} y_1 & \text{cov}(y_1, y_2) \\ \text{cov}(y_1, y_2) & \text{var} y_2 \end{bmatrix}$$

; I am interested only in diagonal elements for variance of y_1 & y_2 .

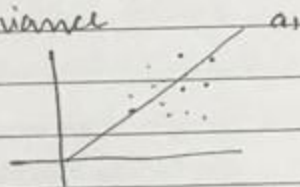
$$A C_x A^T = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} \sigma_{x_1}^2 & \text{cov}(x_1, x_2) \\ \text{cov}(x_1, x_2) & \sigma_{x_2}^2 \end{bmatrix} \begin{bmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{bmatrix}$$

$\text{var}(y_1) = a_1^T C_x a_1$ has to be maximized.

so that y_1 has maximum variance

$$\text{s.t. } \|a_1\|^2 = 1$$

Above is solved using Lagrangian



$$L(a_1, \lambda) = a_1^T C_x a_1 - \lambda (\|a_1\|^2 - 1)$$

diff. w.r.t a_{11} & a_{12}

$$L(a_1, \lambda) = \begin{pmatrix} a_{11} & a_{12} \end{pmatrix} \begin{pmatrix} \sigma_{x_1}^2 & \text{cov}(x_1, x_2) \\ \text{cov}(x_1, x_2) & \sigma_{x_2}^2 \end{pmatrix} \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} - \lambda (a_{11}^2 + a_{12}^2 - 1)$$

$$= \begin{pmatrix} a_{11} & a_{12} \end{pmatrix} \begin{pmatrix} \sigma_{x_1}^2 a_{11} + \text{cov}(x_1, x_2) a_{12} \\ \text{cov}(x_1, x_2) a_{11} + \sigma_{x_2}^2 a_{12} \end{pmatrix} - \lambda (a_{11}^2 + a_{12}^2 - 1)$$

$$\lambda (a_{11}^2 + a_{12}^2 - 1)$$

$$L(a_1, \lambda) = \sigma_{x_1}^2 a_{11}^2 + 2 a_{11} a_{12} \text{cov}(x_1, x_2) + a_{12}^2 \sigma_{x_2}^2 - \lambda (a_{11}^2 + a_{12}^2 - 1)$$

Differentiate w.r.t a_{11} & equate to 0.

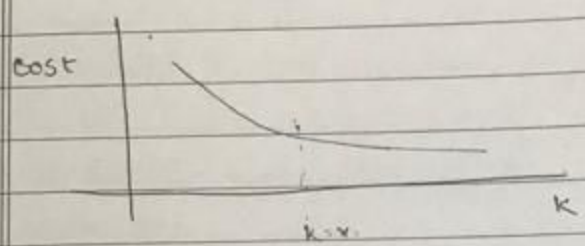
$$2 a_{11} \sigma_{x_1}^2 + 2 a_{12} \text{cov}(x_1, x_2) - 2 \lambda a_{11} = 0 \quad (1)$$

- Application

(1) Market Segmentation

(2) Social networking - Group of people using only email
Facebook
WhatsApp

Elbow Method : To get k



Algorithm:

Input k (number of clusters)

Input training set $\{x^{(1)} \dots x^{(m)}\}$
 $x^{(i)} \in \mathbb{R}^n$

Cluster assignment { Randomly initialize k cluster centroids
 $\mu_1, \mu_2 \dots \mu_k \in \mathbb{R}^n$

Step. \rightarrow Repeat {

for $i = 1$ to m

$c^{(i)} = \text{index (from 1 to } k) \text{ of cluster centroid closest to } x^{(i)}$

$$\min \|x^{(i)} - \mu_k\|^2 \quad \forall k$$

}

Move centroid

for $k = 1$ to K

$\mu_k = \text{mean of points assigned to cluster } k$

}

Andrew NG.

papergrid

Date: / /

Independent Component Analysis (ICA)

(different from PCA)

: Andrew NG Notes

PCA - dimensionality reduction

ICA is also called Blind Source Separation.

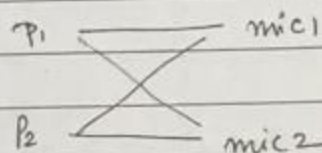
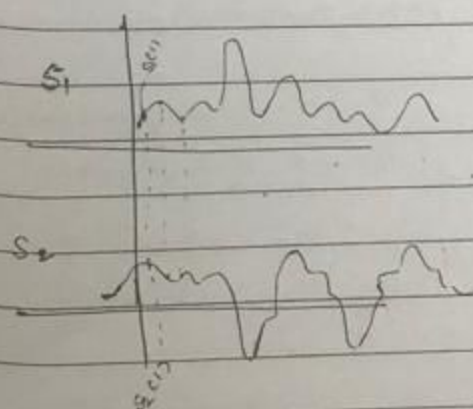
Independence is a statistical component

ICA works on separation

Independence ~~means~~ refers to the statistical component of the final (separated) output.

$$\underline{X} = \underline{A} \underline{S}$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$$



The sampled values of s_1 take up first row of S

I will sample the s (@ m time instances)

$$\begin{bmatrix} s_1^{(1)} & \dots & s_1^{(m)} \\ s_2^{(1)} & \dots & s_2^{(m)} \end{bmatrix} \quad 2 \times m$$

you will get a X of $2 \times m$

In ICA you are given X ; your objective is to find S & also A

$$X = \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(m)} \\ x_2^{(1)} & \dots & x_2^{(m)} \end{bmatrix}$$