

Introduction to Linear Algebra

- Gilbert Strang

17

Linear Algebra in Probability  
and Statistics



ART NATURE

Note Book

# INDEX

17

Name : ..... **SOORAJ·S.** ..... Subject : .....

*Std. : .....* *Div. : .....* *Roll No. : .....*

School / College : .....

S. No.	Date	Title	Page No.	Teacher's Sign / Remarks
		<p>INTRODUCTION TO LINEAR ALGEBRA</p> <p>- Gilbert Strang, MIT (5<sup>th</sup> edition)</p>		

## □ Covariance Matrices & Joint Probabilities

Linear algebra enters when we run  $M$  different experiments at once. We might measure age, height and weight ( $M=3$  measurements of  $N$  people).

Each expt. has its own mean value.

So we have a vector  $m = (m_1, m_2, m_3)$  containing the  $m$  mean values. These could be sample mean of age & height and weight. Or  $m_1, m_2, m_3$  could be expected values of age, height and weight based on known probabilities.

A matrix becomes involved when we look at variances. Each expt will have a sample variance  $S_i^2$  or an expected  $\sigma_i^2 = E[(x_{ij} - m_i)^2]$  based on the squared distance from its mean. Those  $M$  numbers  $\sigma_1^2, \dots, \sigma_M^2$  will go on the main diagonal of the matrix.

Those  $M$  parallel experiments measure  $M$  different random variables, but the expts are not necessarily independent.

If we measure age, height & weight  $(a, h, w)$  for children, the results will be strongly correlated: older children are generally taller and heavier.

Covariance:

$$\text{Cov}_{ah} = E[(\text{age} - \text{mean age})(\text{height} - \text{mean height})]$$

To compute  $P_{ah}$ , we have to know the joint probability of each pair (age & height). This is because age is connected to height.

$P_{ah}$ : probability that a random child has age = a and height = h : Both at once.

$P_{ij}$ : probability that expt 1 produces  $x_i$  and expt. 2 produces  $y_j$

Suppose, expt 1 (age) has mean  $m_1$ . Expt 2 (height) has mean  $m_2$ . The covariance b/w expts 1 and 2 looks at all pairs of ages  $x_i$ , heights  $y_j$ :

Covariance

$$\sigma_{12} = \sum_{\text{all}} \sum_{i,j} P_{ij} (x_i - m_1)(y_j - m_2)$$

Ex:1. Flip 2 coins separately. With 1 for heads and 0 for tails, the results can be (1,1) or (1,0) or (0,1) or (0,0). Those 4 outcomes all have probability  $P_{11}=P_{10}=P_{01}=P_{00}=\frac{1}{4}$

Independent experiments have,

$$P_{ij} = P(i) P(j)$$

Prob. matrix,  $P = \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix} = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} \end{bmatrix}$

Ex:2. Glue the coins together, facing the same way. The only possibilities are (1,1) and (0,0). Those have probabilities  $\frac{1}{2}$  and  $\frac{1}{2}$ .

$$P_{01} = P_{10} = 0$$

(1,0) & (0,1) won't happen because the coins stick together:

$$P = \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

# Probability matrix,

$$P = \begin{bmatrix} P_{11} & P_{12} \\ P_{21} & P_{22} \end{bmatrix}$$

$P_{11} + P_{12} = P_1$   
 $P_{21} + P_{22} = P_2$

$$P_1 \quad P_2$$

(2<sup>nd</sup> coin) columns sum

Those numbers  $P_1, P_2$  and  $P_1, P_2$  are called the marginals of the matrix  $P$ .

$P_1 = P_{11} + P_{12}$  = chance of heads from coin 1.  
 (coin 2 can be heads or tails)

$P_2 = P_{21} + P_{22}$  = chance of heads from coin 2  
 (coin 1 can be heads or tails).

Zero covariance  $\sigma_{12}$

for independent trials.

$$V = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} = \text{diagonal covariance matrix}$$

\* Independent experiments have  $\sigma_{12} = 0$  because

$$P_{ij} = (P_i)(P_j)$$

$$\sigma_{12} = \sum_i \sum_j P_{ij} (x_i - m_1)(y_j - m_2)$$

$$= \sum_i \sum_j (P_i)(P_j) (x_i - m_1)(y_j - m_2)$$

$$= \left[ \sum_i P_i (x_i - m_1) \right] \left[ \sum_j P_j (y_j - m_2) \right]$$

$$= [0][0] = 0$$

Since sum of deviations from the mean is zero.

The glued coins show perfect correlation.  
Heads on one means heads on the other.

The covariance  $\sigma_{12}$  moves from 0 to  $\sigma_1 \sigma_2 = \frac{1}{4}$  - which is the largest possible value of  $\sigma_{12}$ .

$$\text{means} = \frac{1+1}{2} \cdot (\text{for both})$$

$$\sigma_{12} = \frac{1}{2} \left(1 - \frac{1}{2}\right) \left(1 - \frac{1}{2}\right) + 0 + 0 + \frac{1}{2} \left(0 - \frac{1}{2}\right) \left(0 - \frac{1}{2}\right) = \frac{1}{4}$$

Heads or tails from coin 1 gives complete information about heads or tails from coin 2:

Glued coins give largest possible covariances

Singular covariance matrix : determinant = 0

$$\sqrt{\text{glue}} = \begin{bmatrix} \sigma_1^2 & \sigma_1\sigma_2 \\ \sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

\*  $\sigma_1^2 \sigma_2^2 \geq \sigma_{12}^2 \Rightarrow -\sigma_1\sigma_2 \leq \sigma_{12} \leq \sigma_1\sigma_2$

Proof

Cauchy-Schwarz inequality

$$\left( \sum_i (x_i - m_i)^2 \right) \left( \sum_j (y_j - m_j)^2 \right) \geq \left( \sum_{i,j} (x_i - m_i)(y_j - m_j) \right)^2$$

□ The Covariance Matrix  $\Sigma$  is positive definite

Expected covariance  $\sigma_{12}$  b/w 2 experiments  
1 and 2 (2 coins):

$$\begin{aligned}\sigma_{12} &= \text{expected value of } [(o/p 1-\text{mean}_1) \text{ times } (o/p 2-\text{mean}_2)] \\ &= \sum_{\text{all } i,j} P_{ij} (x_i - m_1)(y_j - m_2)\end{aligned}$$

$P_{ij} \geq 0$  is the probability of seeing  $x_i$ ,  $y_j$  in expt 1  
and  $y_j$  in expt 2.

Total probability (all pairs) is 1 (by definition)

$$\sum_{\text{all } i,j} P_{ij} = 1$$

Fix on one particular o/p  $x_i$  in expt 1.

Allow all o/p  $y_j$  in expt 2.

Add the probabilities of  $(x_i, y_1), (x_i, y_2), \dots, (x_i, y_n)$ :

Row sum  $p_i$  of  $P$ ,

$$\sum_{j=1}^n p_{ij} = \text{probability } P_i \text{ of } x_i \text{ in expt 1}$$

Some  $y_j$  must happen in expt 2.

Whether the 2 coins are completely separate or glued together, we still get  $\frac{1}{2}$  for the probability  $P_H = P_{HH} + P_{HT}$  that coin 1 is heads:

$$(\text{separate}) \quad P_{HH} + P_{HT} = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

$$(\text{glued}) \quad P_{HH} + P_{HT} = \frac{1}{2} + 0 = \frac{1}{2}$$

Covariance matrix,

$$\nabla = \sum \sum \nabla_{ij}$$

$$= \sum_{\text{all}} \sum_{i,j} P_{ij} \begin{bmatrix} (x_i - m_1)^2 & (x_i - m_1)(y_j - m_2) \\ (x_i - m_1)(y_j - m_2) & (y_j - m_2)^2 \end{bmatrix}$$

$$= \sum_{\text{all}} \sum_{i,j} P_{ij} \begin{bmatrix} x_i - m_1 \\ y_j - m_2 \end{bmatrix} \begin{bmatrix} x_i - m_1 & y_j - m_2 \\ (x - \bar{x}) & (x - \bar{x})^\top \end{bmatrix}$$

$$\nabla_{11} = \sum_{\text{all}} \sum_{i,j} P_{ij} (x_i - m_1)^2$$

$$= \sum_{\text{all } i} (\text{probability of } x_i) (x_i - m_1)^2 = \sigma_1^2$$

$V_{ij}$  has diagonal entries  $P_{ij}(x_i - m_1)^2 \geq 0$   
 and  $P_{ij}(y_j - m_2)^2 \geq 0$  and  $\det(V_{ij}) = 0$ .

$$\Rightarrow \text{rank}(V_{ij}) = 1$$

$$\begin{bmatrix} (x_i - m_1)^2 & (x_i - m_1)(y_j - m_2) \\ (x_i - m_1)(y_j - m_2) & (y_j - m_2)^2 \end{bmatrix} = \begin{bmatrix} x_i - m_1 \\ y_j - m_2 \end{bmatrix} \begin{bmatrix} x_i - m_1 & y_j - m_2 \end{bmatrix}$$

$$= UU^T$$

$$x^T(UU^T)x = (U^Tx)^T(U^Tx) = \|U^Tx\|^2 \geq 0$$

ILA⑧

$\Rightarrow$  Every matrix  $UU^T$  is positive semidefinite.

& also, if  $S$  &  $T$  are symmetric positive definite, then so is  $S+T$

$\therefore$  The whole matrix  $V$  (combining these matrices  $UUT^T$  with weights  $P_{ij} \geq 0$ ) is at least semi-definite

- \* The covariance matrix  $V$  is positive definite unless the experiments are dependent.

The o/p from each trial is a vector  $X$  with  $M$  components.

The covariance matrix is now  $M$  by  $M$ .

$\Sigma$  is created from the o/p vectors  $X$  and their average  $\bar{X} = E[X]$ :

Covariance matrix,

$$\Sigma = E[(X - \bar{X})(X - \bar{X})^T]$$

$XX^T$  &  $\bar{X}\bar{X}^T$  = (column)(row) are  $M$  by  $M$  matrices.

Take any linear combination  $c^T X = c_1 X_1 + \dots + c_M X_M$

Linearity  $\Rightarrow E[c^T X] = c^T E[X] = c^T \bar{X}$

$$\begin{aligned} \text{Var}[c^T X] &= E[(c^T X - c^T \bar{X})(c^T X - c^T \bar{X})^T] \\ &= c^T E[(X - \bar{X})(X - \bar{X})^T] c = c^T V c \end{aligned}$$

$$\text{Var}[c^T X] \geq 0 \Rightarrow c^T V c \geq 0$$

∴ The covariance matrix  $V$  is therefore positive semi-definite. By the energy test.  $c^T V c \geq 0$ .

$$V = \sum_{\text{all } x, y, z} P_{x, y, z} U U^\top$$

$$= \sum_{\text{all } x, y, z} P_{x, y, z} \begin{bmatrix} x - \bar{x} \\ y - \bar{y} \\ z - \bar{z} \end{bmatrix} \begin{bmatrix} x - \bar{x} & y - \bar{y} & z - \bar{z} \end{bmatrix}$$

$$= \sum_{\text{all } x, y, z} P_{x, y, z} \begin{bmatrix} (x - \bar{x})^2 & (x - \bar{x})(y - \bar{y}) & (x - \bar{x})(z - \bar{z}) \\ (x - \bar{x})(y - \bar{y}) & (y - \bar{y})^2 & (y - \bar{y})(z - \bar{z}) \\ (x - \bar{x})(z - \bar{z}) & (y - \bar{y})(z - \bar{z}) & (z - \bar{z})^2 \end{bmatrix}$$

$$= E \left[ (X - \bar{X})(X - \bar{X})^\top \right] = E \left[ U U^\top \right]$$

where,

$$U = X - \bar{X} = \begin{bmatrix} x - \bar{x} \\ y - \bar{y} \\ z - \bar{z} \end{bmatrix}$$

The value of the expectation symbol  $E$  is that it also allows pdf's (probability density functions) like  $p(x,y,z)$  for continuous random variables  $x, y$  and  $z$ .

i.e., if we allow all numbers as ages, heights, and weights, instead of age  $i=0, 1, 2, \dots$  then we need  $p(x,y,z)$  instead of  $P_{ijk}$ .

Covariance matrix,

$$V = E[UV^T]$$

$$= \iiint p(x,y,z) UV^T dx dy dz$$

$$\text{with } U = X - \bar{X} =$$

$$\begin{bmatrix} x - \bar{x} \\ y - \bar{y} \\ z - \bar{z} \end{bmatrix}$$

$$= (S, V, C)$$

that if  $\mathbb{E}$  belongs to the set of center of  
various planes (fixed) along walls due to  
various material conditions of  $(x,y)$  off  
( $x$  from  $y$ )

it fixed, says no random no walls are off  
so  $\dots$   $\mathbb{E}$  lies off the boundary, always has  
to lie to boundary ( $x,y$ )  $\mathbb{E}$  can not

where,  $\iiint p(x,y,z) dx dy dz = 1$

~~extreme - improved~~

Independent variables  $x, y, z$   $[U U] \mathbb{E} = V$

$$p(x,y,z) = P_1(x) P_2(y) P_3(z)$$

$$\therefore \text{prob. } [U U(x,y,z)] =$$

Dependent variables  $\begin{bmatrix} F \\ E-B \\ X-X \end{bmatrix} \quad x, y, z \quad U \quad \text{Ans}$

2.  $p(x,y,z) = 0$  except when  $cx + dy + ez = 0$

## □ Diagonalization of the Covariance Matrix

Covariance matrices  $\Sigma$  opens up the link b/w probability and linear algebra.

For problems with many simultaneous random variables, put them into vectors:

$$\vec{x} = \begin{bmatrix} R \\ S \end{bmatrix}, \vec{y} = \begin{bmatrix} T \\ U \\ V \end{bmatrix}$$

the covariance matrix,

$$\text{Cov}(\vec{x}, \vec{y}) = E[(\vec{x} - E(\vec{x}))(\vec{y} - E(\vec{y}))^\top]$$

$$= \begin{bmatrix} \text{cov}(R, T) & \text{cov}(R, U) & \text{cov}(R, V) \\ \text{cov}(S, T) & \text{cov}(S, U) & \text{cov}(S, V) \end{bmatrix}$$

Often there is one vector with all the variables:

$$\vec{X} = \begin{bmatrix} R \\ S \\ T \end{bmatrix}$$

$$\text{Cov}(\vec{X}) = \text{Cov}(\vec{X}, \vec{X})$$

$$= E[(\vec{X} - E(\vec{X}))(\vec{X} - E(\vec{X}))^T]$$

$$= \begin{bmatrix} \text{cov}(R, R) & \text{cov}(R, S) & \text{cov}(R, T) \\ \text{cov}(S, R) & \text{cov}(S, S) & \text{cov}(S, T) \\ \text{cov}(T, R) & \text{cov}(T, S) & \text{cov}(T, T) \end{bmatrix}$$

$$= \begin{bmatrix} \text{var}(R) & \text{cov}(R, S) & \text{cov}(R, T) \\ \text{cov}(S, R) & \text{var}(S) & \text{cov}(S, T) \\ \text{cov}(T, R) & \text{cov}(T, S) & \text{var}(T) \end{bmatrix}$$

$$\text{Cov}(\vec{x}, \vec{y}) = \text{Cov}(\vec{y}, \vec{x})^T$$

$$\text{Cov}(A\vec{x} + B, \vec{y}) = A \text{Cov}(\vec{x}, \vec{y})$$

$$\text{Cov}(\vec{x}, C\vec{y} + D) = \text{Cov}(\vec{x}, \vec{y}) C^T$$

$$\text{Cov}(A\vec{x} + B) = A \text{Cov}(\vec{x}) B^T$$

$$\text{Cov}(\vec{x}_1 + \vec{x}_2, \vec{y}) = \text{Cov}(\vec{x}_1, \vec{y}) + \text{Cov}(\vec{x}_2, \vec{y})$$

$$\text{Cov}(\vec{x}, \vec{y}_1 + \vec{y}_2) = \text{Cov}(\vec{x}, \vec{y}_1) + \text{Cov}(\vec{x}, \vec{y}_2)$$

$A, C$  : constant matrices

~~$A, C$~~

$\vec{B}, \vec{D}$  : constant vectors

2

the first few principal components  
will capture most of the variance.  
 $\text{and } \mathbf{Q}^T \mathbf{A} \mathbf{Q} = \mathbf{V}$  represents the  
two vectors as the first two.

A key question in analyzing data analysis is to figure out which variables are responsible for most of the variation in the data.

These may not be the variables that are measured, but may instead be some linear combination of the measured variables.

→ Mathematically, this corresponds to diagonalizing the covariance (or correlation) matrix.

ie,

The diagonalization of covariance matrices is used to transform correlated random variables into uncorrelated random variables.

- \* Covariance matrix  $\Sigma$  is real-symmetric & +ve definite (or atleast semidefinite), so diagonalization  $\Sigma = Q \Lambda Q^T$  finds real, +ve eigenvalues  $\lambda_k = \sigma_k^2 \geq 0$  and an orthonormal basis of eigenvectors  $q_1$  to  $q_n$ .
- \* The eigenvectors form a coordinate system in which the ~~covariance~~ matrix  $\Sigma$  becomes diagonal.  
i.e., coordinates in which the variables are uncorrelated.
- \* The diagonal entries in this coordinate system, the eigenvalues, are the variances of these uncorrelated components.
- \* The process of diagonalizing the covariance matrix is called principal component analysis (PCA).

Ex:- (Decorrelation of random vectors)

A random vector  $X = (X_1, X_2, X_3)^T$  has covariance matrix

$$K_X = \begin{bmatrix} 2 & -1 & 1 \\ -1 & 2 & 0 \\ 1 & 0 & 2 \end{bmatrix}$$

Design a non-trivial transformation that will generate from  $X$  a new random vector  $Y$  whose components are uncorrelated.

Ans:  $|K_X - \lambda I| = 0 = \begin{vmatrix} 2-\lambda & -1 & 1 \\ -1 & 2-\lambda & 0 \\ 1 & 0 & 2-\lambda \end{vmatrix}$

$$(2-\lambda)^3 + 1(\lambda-2) + (\lambda-2) = 0.$$

$$\Rightarrow (2-\lambda)[\lambda^2 - 4\lambda + 4 - 1] = (2-\lambda)(\lambda^2 - 4\lambda + 3) = 0$$

$$\Delta = 16 - 48 = 0$$

$$\lambda_1 = 2, \lambda_2 = 2 + \sqrt{2}, \lambda_3 = 2 - \sqrt{2}$$

$$\frac{4 \pm 2\sqrt{2}}{2}$$

$$V_1 = \begin{bmatrix} 0 \\ Y_2 \\ Y_{\sqrt{2}} \end{bmatrix}, V_2 = \begin{bmatrix} Y_2 \\ -Y_2 \\ Y_2 \end{bmatrix}, V_3 = \begin{bmatrix} Y_{\sqrt{2}} \\ Y_2 \\ -Y_2 \end{bmatrix}$$

$$V = Q \Lambda Q^T$$

$$A = Q^{-1} = Q^T = \begin{bmatrix} v_1 & v_2 & v_3 \end{bmatrix}^T$$

$$= \begin{bmatrix} 0 & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{\sqrt{2}} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix}$$

$$= \begin{bmatrix} 0 & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{\sqrt{2}} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix}$$

$$\begin{aligned}
 Y &= \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \end{bmatrix} = AX = \begin{bmatrix} 0 & \sqrt{2} & \sqrt{2} \\ \sqrt{2} & -\sqrt{2} & \sqrt{2} \\ \sqrt{2} & \sqrt{2} & -\sqrt{2} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \\
 &= \begin{bmatrix} \frac{1}{\sqrt{2}}X_2 + \frac{1}{\sqrt{2}}X_3 \\ \frac{1}{\sqrt{2}}X_1 - \frac{1}{2}X_2 + \frac{1}{2}X_3 \\ \frac{1}{\sqrt{2}}X_1 + \frac{1}{2}X_2 - \frac{1}{2}X_3 \end{bmatrix}
 \end{aligned}$$

The covariance of  $Y$  is given by,

$$\begin{aligned}
 K_Y &= Q^T K_X Q = \frac{1}{4} \begin{bmatrix} 0 & \sqrt{2} & \sqrt{2} \\ \sqrt{2} & -1 & 1 \\ \sqrt{2} & 1 & -1 \end{bmatrix} \begin{bmatrix} 2 & -1 & 1 \\ -1 & 2 & 0 \\ 1 & 0 & 2 \end{bmatrix} \begin{bmatrix} 0 & \sqrt{2} & \sqrt{2} \\ \sqrt{2} & -1 & 1 \\ \sqrt{2} & 1 & -1 \end{bmatrix} \\
 &= \frac{1}{4} \begin{bmatrix} 0 & \sqrt{2} & \sqrt{2} \\ \sqrt{2} & -1 & 1 \\ \sqrt{2} & 1 & -1 \end{bmatrix} \begin{bmatrix} 0 & 2\sqrt{2}+2 & 2\sqrt{2}-2 \\ 2\sqrt{2} & -\sqrt{2}-2 & -\sqrt{2}+2 \\ 2\sqrt{2} & \sqrt{2}+2 & \sqrt{2}-2 \end{bmatrix} \\
 &= \frac{1}{4} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 8+4\sqrt{2} & 0 \\ 0 & 0 & 8-4\sqrt{2} \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2+\sqrt{2} & 0 \\ 0 & 0 & 2-\sqrt{2} \end{bmatrix}
 \end{aligned}$$

## Mean & Variance of $Z = x + y$

### sample mean

We have  $N$  samples of  $x$ .  
Their mean (= average) is  $m_x$ .

We also have  $N$  samples of  $y$  and  
their mean is  $m_y$ .

The sample mean of  $Z = x + y$  is  $m_z = m_x + m_y$

Mean of sums = Sum of means

$$\frac{1}{N} \sum_{i=1}^N (x_i + y_i) = \frac{1}{N} \sum_{i=1}^N x_i + \frac{1}{N} \sum_{i=1}^N y_i$$

Expected mean

$P_{ij}$  : joint probability of the pair  $(x_i, y_j)$

its value depends on whether the expt  
are independent, which we don't know.

$$E[x+y] = \sum_i \sum_j P_{ij} (x_i + y_j)$$

$$= \sum_{i=1}^N \sum_{j=1}^N P_{ij} x_i + \sum_{i=1}^N \sum_{j=1}^N P_{ij} y_j$$

$$= \sum_{i=1}^N (P_{i1} + P_{i2} + \dots + P_{iN}) x_i + \sum_{j=1}^N (P_{1j} + P_{2j} + \dots + P_{Nj}) y_j$$

$$= \sum_{i=1}^N P_i x_i + \sum_{j=1}^N P_j y_j$$

$$= E[x] + E[y]$$

$$E[x+y] = E[x] + E[y]$$

Variance

$$\begin{aligned}
 \sigma_z^2 &= \sum_i \sum_j P_{ij} (x_i + y_j - m_x - m_y)^2 \\
 &= \sum_i \sum_j P_{ij} (x_i - m_x)^2 + \sum_i \sum_j P_{ij} (y_j - m_y)^2 \\
 &\quad + 2 \sum_i \sum_j P_{ij} (x_i - m_x)(y_j - m_y) \\
 &= \sum_i P_i (x_i - m_x)^2 + \sum_j P_j (y_j - m_y)^2 \\
 &\quad + 2 \sum_i \sum_j P_{ij} (x_i - m_x)(y_j - m_y) \\
 &= \sigma_x^2 + \sigma_y^2 + 2\sigma_{xy}
 \end{aligned}$$

$$\sigma_{xy}^2 = \sigma_x^2 + \sigma_y^2 + 2\sigma_{xy}$$

□ Covariance Matrix for  $Z = AX$

Let  $X = \begin{bmatrix} x \\ y \end{bmatrix}$  and  $A = \begin{bmatrix} 1 & 1 \end{bmatrix}$

$$AX = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x+y \\ y \end{bmatrix} = \begin{bmatrix} z \\ y \end{bmatrix}$$

$$\sigma_z^2 = \sigma_{x+y}^2 = \sigma_x^2 + \sigma_y^2 + 2\sigma_{xy}$$



$$\sigma_z^2 = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \iff \sigma_z^2 = A V A^T$$

The covariance matrix of  $Z = AX$  is:

$$\Sigma_Z = A \Sigma_X A^T$$

$X$ : vector

$$E[AX] = A E[X]$$

$$\text{Cov}[AX] = A \text{Cov}[X] A^T$$

$$E[X] = E[X_1] + \dots + E[X_n]$$

Proof

$$\begin{aligned} E[AX] &= E\left(\begin{bmatrix} \sum_{j=1}^n a_{1j} X_j \\ \vdots \\ \sum_{j=1}^n a_{mj} X_j \end{bmatrix}\right) \\ &= E\left[\sum_{j=1}^n a_{1j} X_j\right] + \dots + E\left[\sum_{j=1}^n a_{mj} X_j\right] \\ &= \sum_{j=1}^n a_{1j} E[X_j] + \dots + \sum_{j=1}^n a_{mj} E[X_j] \\ &= A E[X] \end{aligned}$$

$$\text{Cov}[AX] = E \left[ (AX - E[AX])(AX - E[AX])^T \right]$$

$$= E \left[ A(x - E[x])(x - E[x])^T A^T \right]$$

$$= A E \left[ (x - E[x])(x - E[x])^T \right] A^T$$

$$= A \text{cov}(x) A^T$$

## □ The Correlation $\rho$

Rescaling or standardizing the random variables  $x$  and  $y$ :

$$X = \frac{x}{\sigma_x} \quad \text{and} \quad Y = \frac{y}{\sigma_y}$$

$$\sigma_x^2 = 1 \quad \text{and} \quad \sigma_y^2 = 1 \quad \leftarrow \text{variances of the new variables } X \text{ and } Y.$$

This is just like dividing a vector ' $v$ ' by its length to produce a unit vector  $\frac{v}{\|v\|}$  of length 1.

The correlation of  $\alpha$  &  $\beta$  is the covariance of  $X$  and  $Y$ .

Correlation,

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \text{covariance of } X = \frac{x}{\sigma_x} \text{ & } Y = \frac{y}{\sigma_y}$$

$$-1 \leq \rho_{xy} \leq 1$$

$$\sigma_x^2 \sigma_y^2 \geq \sigma_{xy}^2$$

$$-\sigma_x \sigma_y \leq \sigma_{xy} \leq \sigma_x \sigma_y$$

Zero covariance gives zero correlation.

Independent random variables produce  $\rho_{xy} = 0$ .

nce always  $\sigma_{xy}^2 \leq \sigma_x^2 \sigma_y^2 \Rightarrow$  the covariance matrix  $V$  is at least positive semidefinite.

$$\rightarrow \text{circle } r_{xy}^2 \leq 1$$

$\frac{\sigma_x}{\sigma_y}$

Correlation near  $r = +1$  means strong dependence in the same direction.

i.e.,

$y$  tends to be above its mean when  $x$  is also above its mean and vice versa.  
 $\Rightarrow$  often voting the same.

$r = 0$ .  
-ve correlation means that  $y$  tends to be below its mean when  $x$  is above its mean  
 $\Rightarrow$  Voting in opposite directions.

Correlation matrix,  $R = \begin{bmatrix} 1 & \rho_{xy} \\ \rho_{xy} & 1 \end{bmatrix}$

- \*  $R$  always has 1's on the diagonal because we normalized to  $\sigma_x = \sigma_y = 1$ .  
i.e.,

$R$  is the correlation matrix for  $x$  and  $y$ , and the covariance matrix for  $X = \frac{x}{\sigma_x}$  and  $Y = \frac{y}{\sigma_y}$ .

The number  $\rho_{xy}$  is also called the Pearson coefficient.

- \* The correlation matrix  $R$  comes from the covariance matrix  $V$ , when we rescale every row and every column. Divide each row  $i$  and column  $i$  by the  $i^{\text{th}}$  standard deviation  $\sigma_i$ .

②  $R = DVD$  for the diagonal matrix  $D$ ,

where  $D = \text{diag} [\gamma_1, \dots, \gamma_M]$

③ If covariance  $V$  is positive definite, correlation  $R = DVD$  is also positive definite.

$$V = \begin{bmatrix} \sigma_{xx}^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_{yy}^2 \end{bmatrix} ; R = \begin{bmatrix} 1 & \rho_{xy} \\ \rho_{xy} & 1 \end{bmatrix}$$

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

Ex:3. Suppose that  $y$  is just  $-\alpha$ . A coin flip has  $\text{op } \alpha=0$  or  $1$ . The same flip has  $\text{op } y=0$  or  $-1$ .

Ans: The mean,  $m_x = \frac{1}{2}$

$$m_y = -\frac{1}{2}$$

The covariance is,  $\sigma_{xy} = -\sigma_x \sigma_y$

$$\Rightarrow \rho_{xy} = -1$$

$$R = \begin{bmatrix} 1 & \rho_{xy} \\ \rho_{xy} & 1 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \text{ when } y = -\alpha$$

The correlation matrix has determinant zero (singular & only semidefinite).

Ex:4 Suppose the random variables  $x, y, z$  are independent. What matrix is  $R$ ?

Dns:  $\rho_{xx} = \rho_{yy} = \rho_{zz} = 1$

All 3 cross correlations,  $\rho_{xy} = \rho_{yz} = \rho_{xz} = 0$

$$R = \begin{bmatrix} \rho_{xx} & \rho_{xy} & \rho_{xz} \\ \rho_{yx} & \rho_{yy} & \rho_{yz} \\ \rho_{zx} & \rho_{yz} & \rho_{zz} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

is the identity matrix.

(12.2(A)) Suppose  $x$  &  $y$  are independent random variables with mean 0 and variance 1. Then the covariance matrix  $V_x$  for  $X = (x, y)$  is  $2 \times 2$  identity matrix. What are the mean  $m_z$  & the covariance matrix  $V_z$  for the 3-comp. vector  $Z = (x, y, ax+by)$ ?

Ans:

$$Z = \begin{bmatrix} x \\ y \\ ax+by \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ a & b \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = AX$$

$$m_z = Am_x$$

mean of  $ax+by$  is  $am_x + bm_y$ .

$$V_z = A V_x A^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ a & b \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & b \\ 0 & b & a^2+b^2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & a \\ 0 & 1 & b \\ a & b & a^2+b^2 \end{bmatrix}$$

$$\sigma_{ax+by}^2 = \sigma_{ax}^2 + \sigma_{by}^2 + 2\sigma_{ax,by} = a^2 + b^2 + 0 = a^2 + b^2$$

$$\det(V_2) = a^2 + b^2 - b^2 - a^2 = 0$$

The 3<sup>rd</sup> component  $z = ax + by$  is completely dependent on  $x$  and  $y$ . The rank of  $V_2$  is only 2.

□ Multivariate Gaussian & Weighted  
Least Squares

The normal probability density  $p(x)$  [the Gaussian] depends on only 2 numbers:

— Mean  $m$  & Variance  $\sigma^2$

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}}$$

The graph of  $p(x)$  is a bell-shaped curve centered at  $x=m$ . The const. variable  $x$  can be anywhere b/w  $-\infty$  and  $+\infty$ . With probability close to  $\frac{2}{3}$ , the random  $x$  will lie b/w  $m-\sigma$  and  $m+\sigma$ .

$$\int_{-\infty}^{+\infty} p(x) dx = 1 \quad \times \int_{m-\sigma}^{m+\sigma} p(x) dx = \frac{1}{\sqrt{2\pi}} \int_{-1}^{+1} e^{-\frac{x^2}{2}} dx \approx \frac{2}{3}$$

~~$$X = \frac{x-m}{\sigma} \Rightarrow dx = \frac{dX}{\sigma}$$~~

Every Gaussian terms into a standard Gaussian  $p(x)$  with mean  $m=0$  and variance  $\sigma^2=1$ .

The standard normal distribution  $N(0, 1)$

has  $p(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

Integrating  $p(x)$  from  $-\infty$  to  $x$  gives the cumulative distribution  $F(x)$ : the prob. that a random variable sample is below  $x$ .

$$F = \frac{1}{2} \text{ at } x=0 \text{ (mean).}$$

## Two-dimensional Gaussians

We have  $M=2$  Gaussian random variables  $x$  and  $y$ . They have means  $m_1$  &  $m_2$ . They have variances  $\sigma_1^2$  and  $\sigma_2^2$ .

### Independent $x$ and $y$

If they are independent, then their probability density  $p(x,y)$  is just  $p_1(x)$  times  $p_2(y)$ . Multiply probabilities when variables are independent:

$$p_1(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x-m_1)^2}{2\sigma_1^2}}$$
$$p_2(y) = \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-\frac{(y-m_2)^2}{2\sigma_2^2}}$$

Independent:

$$x \& y \sim N(m-\underline{x}, V(m-\underline{x}))$$

$$p(x,y) = p_1(x)p_2(y) = \frac{1}{2\pi\sigma_1\sigma_2} e^{-\frac{(x-m_1)^2}{2\sigma_1^2}} \times e^{-\frac{(y-m_2)^2}{2\sigma_2^2}}$$

The covariance of  $x$  &  $y$  will be  $\sigma_{12} = 0$ .

The covariance matrix will be diagonal.

The variances  $\sigma_1^2$  and  $\sigma_2^2$  are always on the main diagonal of  $V$ .

The exponent in  $p(x|y)$  is just the sum of the  $x$ -exponent & the  $y$ -exponent.

Good to notice that the  $x$  exponents can be combined into  $-\frac{1}{2}(x-m)^T V^{-1}(x-m)$ :

$$\begin{aligned} -\frac{(x-m_1)^2}{2\sigma_1^2} - \frac{(y-m_2)^2}{2\sigma_2^2} &= -\frac{1}{2} \begin{bmatrix} x-m_1 & y-m_2 \end{bmatrix} \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}^{-1} \begin{bmatrix} x-m_1 \\ y-m_2 \end{bmatrix} \\ &= -\frac{1}{2} (x-m)^T V^{-1} (x-m) \end{aligned}$$

## Non-independent $x$ and $y$

When  $M=2$ , the first variable  $x$  may give partial information about the second variable  $y$  (and vice versa). Maybe part of  $y$  is decided by  $x$  and part is truly independent. It is the  $M \times M$  covariance matrix  $V$  that accounts for dependencies b/w the  $M$  variables  
 $\mathbf{x} = x_1, x_2, \dots, x_M$ .

Multivariate Gaussian probability distribution :

$$P(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^M \sqrt{\det V}} e^{-\frac{-(\mathbf{x}-\mathbf{m})^T V^{-1} (\mathbf{x}-\mathbf{m})}{2}}$$

$$\mathbf{x} = (x_1, \dots, x_M) \quad \& \quad \mathbf{m} = (m_1, \dots, m_M)$$

The  $M$  square roots of  $\omega^2$  and the determinant of  $V$  are included to make the total probability equal to 1.

~~Conjugate~~ Vectors

Vector

We'll use the eigenvalues  $\lambda$  and orthonormal eigenvectors  $Q$  of the symmetric matrix  $V = Q \Lambda Q^T$ .

$$\therefore V^{-1} = Q \Lambda^{-1} Q^T$$

$$X = x - m$$

minant

$$(x - m)^T V^{-1} (x - m) = X^T V^{-1} X = X^T Q \Lambda^{-1} Q^T X \\ = Y^T \Lambda^{-1} Y$$

where,  $Y = Q^T X = Q^T (x - m)$  are statistically independent & their covariance matrix  $\Lambda$  is diagonal.

→ This step of diagonalizing  $V$  by its eigenvector matrix  $Q$  is the same as "uncorrelating" the random variables. Covariances are zero for the new variables  $X_1, \dots, X_m$ . ?

- This is the point where linear algebra helps calculus to compute multidimensional integrals.

$$-\frac{1}{2}(\mathbf{x}-\mathbf{m})^T \nabla^T (\mathbf{x}-\mathbf{m}) = \mathbf{x}^T Q \Lambda^{-1} Q^T \mathbf{x} = \mathbf{y}^T \Lambda^{-1} \mathbf{y}$$

$$= -\frac{1}{2} \begin{bmatrix} y_1 & \dots & y_M \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_M \end{bmatrix}^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}$$

$$= -\frac{y_1^2}{2\lambda_1} - \frac{y_2^2}{2\lambda_2} - \dots - \frac{y_M^2}{2\lambda_M}$$

$$\int \dots \int e^{-\mathbf{y}^T \Lambda^{-1} \mathbf{y}/2} d\mathbf{y} = \int_{-\infty}^{+\infty} e^{-y_1^2/2\lambda_1} dy_1 \dots \int_{-\infty}^{+\infty} e^{-y_M^2/2\lambda_M} dy_M$$

$$\text{Det } \Lambda = \left( \sqrt{2\pi\lambda_1} \right) \dots \dots \dots \left( \sqrt{2\pi\lambda_M} \right)$$

$$= \left( \sqrt{2\pi} \right)^M \sqrt{\det V}$$

Since,  $\int_{-\infty}^{+\infty} e^{-at^2} dt = \sqrt{\pi/a}$

Vector  $m$  of means:

$$\int \dots \int \alpha P(\alpha) d\alpha = (m_1, m_2, \dots) = m$$

Covariance matrix  $V$ :

$$\int \dots \int (\alpha - m)(\alpha - m)^T P(\alpha) d\alpha = V$$

## Weighted Least Squares

Least squares started from an unsolvable system  $A\alpha = b$ . We chose  $\hat{\alpha}$  to minimize the error

$\|b - A\hat{\alpha}\|^2$ . This lead to the least square equation  $A^T A \hat{\alpha} = A^T b$ . The best  $A \hat{\alpha}$  is the projection of  $b$  onto the  $C(A)$ . But, is this squared distance  $E = \|b - A\hat{\alpha}\|^2$  the right error measure to minimize?

If the measurement errors in  $b$  are independent random variables, with mean = 0 and variance  $\sigma^2 = 1$  and a normal distribution.

→ use least squares.

If the errors are not independent or their variances are not equal

→ use weighted least squares.

Suppose, there is an  $m \times n$  matrix  $A'$  and a vector  $\alpha \in \mathbb{R}^n$  such that

$$Y = A\alpha + \epsilon$$

Here  $A$  and  $\alpha$  are non-random but  $\epsilon$  is normally distributed random vector, with mean  $0$  and covariance matrix  $V$ .

$Y$  is normally distributed with mean  $A\alpha$  and covariance matrix  $V$ .

Let  $b$  be the observed value of  $Y$ .

Our goal is to estimate  $\hat{\alpha}$ , given  $A$  &  $b$ .

Let  $f_y$  be the probability density function  
for a normally distributed random ~~vector~~  
vector with mean  $A\alpha$ , ~~covar~~ and  
covariance matrix  $V$ .

\* The argument of the maxima (arg max) are the points, or elements, of the domain of some function at which the function values are maximized.

In contrast to global maxima, which refers to the largest o/p of a function, arg max refers to the inputs, at which the function o/p are as large as possible.

Given an arbitrary set  $X$ , a totally ordered set  $Y$  and a function  $f: X \rightarrow Y$ , the arg max over some subset  $S$  of  $X$  is defined by

$$\text{arg max}_S f := \arg \max_{x \in S} f(x) := \left\{ x \in S : f(s) \leq f(x) \text{ for all } s \in S \right\}$$

$::=$  - defined to be equal to

A natural estimate of  $\alpha$  is the max.  
likelihood estimate

$$\hat{\alpha} = \arg \max_{\alpha} f_y(b)$$

$$-\frac{1}{2}(b - A\alpha)^T V^{-1}(b - A\alpha)$$

$$= \arg \max_{\alpha} \frac{1}{(2\pi)^m \det(V)} e^{-\frac{1}{2}(b - A\alpha)^T V^{-1}(b - A\alpha)}$$

maximizing  $f_y(b)$  is equivalent to minimizing  
 $\ln [f_y(\alpha)]$ .

i.e., minimizing,  $E(\alpha) = (b - A\alpha)^T V^{-1}(b - A\alpha)$

\* Let  $c_1, c_2, \dots, c_m > 0$  be a set of positive numbers. The corresp. weighted inner product and weighted norm on  $\mathbb{R}^m$  are defined by

$$\langle v|w \rangle = \sum_{i=1}^m c_i v_i w_i$$

$$\|v\| = \sqrt{\langle v|v \rangle} = \sqrt{\sum_{i=1}^m c_i v_i^2}$$

The numbers  $c_i$  are the weights.

The larger the weights  $c_i$ , the more the  $i^{\text{th}}$  coordinate of  $v$  contributes to the norm.

$$\langle v|w \rangle = v^T C w,$$

where  $C = \begin{bmatrix} c_1 & 0 & 0 & \cdots & 0 \\ 0 & c_2 & 0 & \cdots & 0 \\ 0 & 0 & c_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & c_m \end{bmatrix}$

$C$ : diagonal weight matrix

\* Weighted norms are particularly relevant in statistics and data fitting, where one wants to emphasise certain quantities and deemphasize others; this is done by assigning weights to the different components of the data vector  $v$ .

$$\begin{bmatrix} w_1 & w_2 & \dots & w_n \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} = \langle w, v \rangle$$

The weight matrix  $W$  is necessarily positive definite.

$W = B^T B$  for a matrix  $B$  with independent columns.

ILA(8)

$$B = QR$$

|  $Q$  has orthonormal columns

$$B^T B = (QR)^T (QR) = R^T Q^T Q R = R^T R$$

Cholesky factorization

ILA(6)

For any vector  $\alpha$ , the weighted norm satisfies:

$$\alpha^T W \alpha = \alpha^T R^T R \alpha = (\alpha^T R)(R \alpha)$$

$$\min_{\alpha \in \mathbb{R}^m} E(\alpha) = \min_{\alpha \in \mathbb{R}^m} (b - A\alpha)^T W (b - A\alpha)$$

$$= \min_{\alpha \in \mathbb{R}^m} (b - A\alpha)^T R^T R (b - A\alpha)$$

$$= \min_{\alpha \in \mathbb{R}^m} [R(b - A\alpha)]^T [R(b - A\alpha)]$$

$$= \min_{\alpha \in \mathbb{R}^m} [Rb - RA\alpha]^T [Rb - RA\alpha]$$

i.e.,

the solution that we are looking for  
is the least square solution to the equation

$$RA\alpha = Rb$$

Using the ordinary least squares, the solution satisfies

$$(RA)^T(RA)\hat{x} = (RA)^T(Rb)$$

$$A^T(R^T R)A\hat{x} = A^T(R^T R)b$$

$$A^TWA\hat{x} = A^Tw^Tb$$

$$W = \begin{bmatrix} w_1 & 0 & 0 & \dots & 0 \\ 0 & w_2 & 0 & \dots & 0 \\ 0 & 0 & w_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & w_m \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & 0 & \dots & 0 \\ 0 & \frac{1}{\sigma_2^2} & 0 & \dots & 0 \\ 0 & 0 & \frac{1}{\sigma_3^2} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \frac{1}{\sigma_m^2} \end{bmatrix} = V^{-1}$$

- Since each weight is inversely proportional to the error variance, it reflects the information in that observation. So an observation with small error variance has a large weight since it contains relatively more information than an observation with large error variance (small weight).

Weighted least squares

$$A^T V^{-1} A \hat{x} = A^T V^{-1} b$$

$$\hat{x} = (A^T V^{-1} A)^{-1} A^T V^{-1} b$$

## The Variance in the Estimated $\hat{a}$

Often the important question is not the best  $\hat{a}$  for one particular set of measurements b. This is only one sample.

The real goal is to know the reliability of the whole experiment. That is measured (as reliability always is) by the variance in the estimate  $\hat{a}$ .

$$E[AX] = A E[X]$$

$$\text{Cov}[AX] = A \text{Cov}[X] A^T$$

If  $b$  has covariance matrix  $V$ , then  
 $\hat{a} = Lb$  has covariance matrix  $LVL^T$ .

where,  $L = (A^T V^{-1} A)^{-1} A^T V^{-1}$

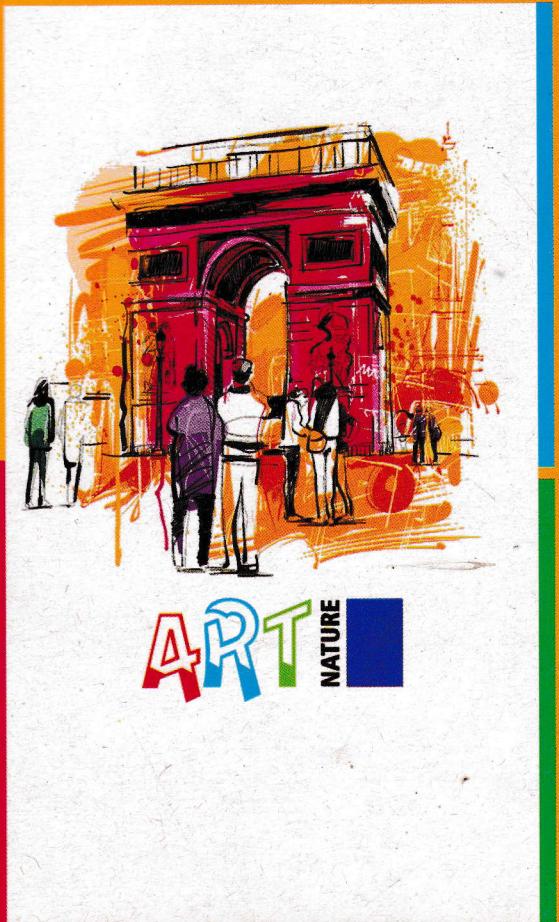
$$\text{Cov}[\hat{\alpha}] = \text{Cov}[L b] = E[(\hat{\alpha} - \alpha)(\hat{\alpha} - \alpha)^T]$$

$$= L \text{Cov}(b) L^T$$

$$= L V L^T$$

$$= (A^T V^{-1} A)^{-1} A^T V^{-1} V V^T A (A^T V^{-1} A)^{-1}$$

$$= (A^T V^{-1} A)^{-1}$$



9446390009  
**KUNNAMKULAM STATIONERY**  
PATHANAMTHITTA

Pages  
with cover 160  
Price ₹ 50.00  
incl of GST

Size : 18.4 X 24.8 cms

Mfrs: Ramsons Enterprises, Madurai, Ph : 0452-2671408