

### 2.3 QUADRATIC FORMS

- A **quadratic form**  $Q(x)$  of the  $p$ -vector  $x = (x_1, \dots, x_p)^T$  is defined by

$$Q(x) = \sum_{i=1}^p \sum_{j=1}^p a_{ij} x_i x_j = x^T A x,$$

*scalar*

where  $a_{ij}$  is the  $(i, j)$ th element of a symmetric  $p \times p$  matrix  $A$ .

- If  $A$  is symmetric  $p \times p$  and  $x$  is  $p$ -vector  
 $Q(x) \geq 0$  for all  $x \neq (0, \dots, 0)^T$

then the matrix  $A$  is called **semi positive definite**, which is denoted by  $A \geq 0$ .  
*(二次型大于等于0 半正定)*

- However if the quadratic form satisfies

$$Q(x) > 0 \text{ for all } x \neq (0, \dots, 0)^T$$

*二次型大于0 正定*  
then the matrix  $A$  is called **positive definite**, which is denoted by  $A > 0$ .

- $A > 0$  is equivalent to all the eigenvalues of  $A$  satisfy:

特征值全大于0  $\rightarrow$  正定=满秩, 非奇异

$$\lambda_1 > 0, \dots, \lambda_p > 0.$$

特征值全大于0=正定=满秩=非奇异

Then  $|A| > 0$ ,  $A^{-1}$  exists,  $A$  is of full rank  $p$ ,  $A$  is non singular.

- If  $A \geq 0$  then 不满秩  $p-r$  nonsingular,  $A^{-1}$  不存在

$$\text{rank}(A) = r < p$$

and

•  $p - r$  eigenvalues of  $A$  are equal to zero

• while the other  $r$  are strictly positive.

## 2.4 GEOMETRICAL ASPECTS

For the rest of Chap 2, vectors are columns unless specified otherwise. To get the results for rows, transpose each vector in each expression.

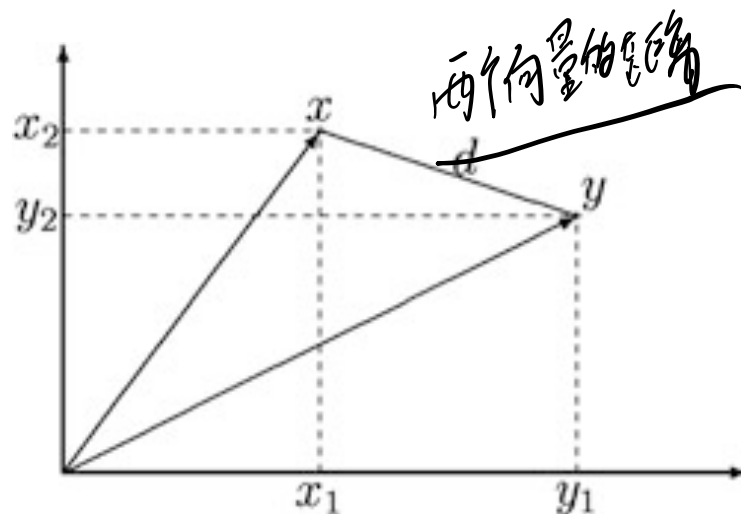
### Distance

欧氏距离:

- The **Euclidian distance**  $d(x, y)$  between  $x, y \in \mathbb{R}^p$  is defined by

$$d(x, y) = \sqrt{\sum_{i=1}^p (x_i - y_i)^2} = \sqrt{(x - y)^T (x - y)}$$

Example in  $\mathbb{R}^2$ , where  $x = (x_1, x_2)^T$  and  $y = (y_1, y_2)^T$ :



加权

- A **weighted version** of this distance can be defined as

$$d(x, y) = \sqrt{\sum_{i=1}^p w_i (x_i - y_i)^2} = \sqrt{(x - y)^T W (x - y)}$$

$\downarrow$  由 designer 决定

where each  $w_i > 0$  and  $W = \text{diag}(w_1, \dots, w_p)$ .

- This can be further generalised into the following distance:

$$d(x, y) = \sqrt{(x - y)^T A (x - y)}$$

$\downarrow$  满秩矩阵

where  $A$  is a  $p \times p$  **positive definite** matrix.

## Norm

- The (Euclidian) **norm** of a vector  $x \in \mathbb{R}^p$  is defined by

$$\|x\| = \sqrt{\sum_{i=1}^p x_i^2} = \sqrt{x^T x}. \quad (\text{vector normalization})$$

- A **unit vector** is a vector of norm 1.
- Multiplication by an orthogonal matrix is ~~norm preserving~~: If  $O$  is a  $p \times p$  orthogonal matrix, then

$$\|Ox\| = \|x\|, \quad \text{模不变.}$$

since

$$\|Ox\|^2 = x^T \underbrace{O^T O}_{\text{identity}} x = x^T x = \|x\|^2.$$

$(Ox)^T(Ox) = x^T O^T O x = x^T x = \|x\|^2$

- Can be generalised into a norm with respect to a **positive definite**  $p \times p$  matrix  $A$ :

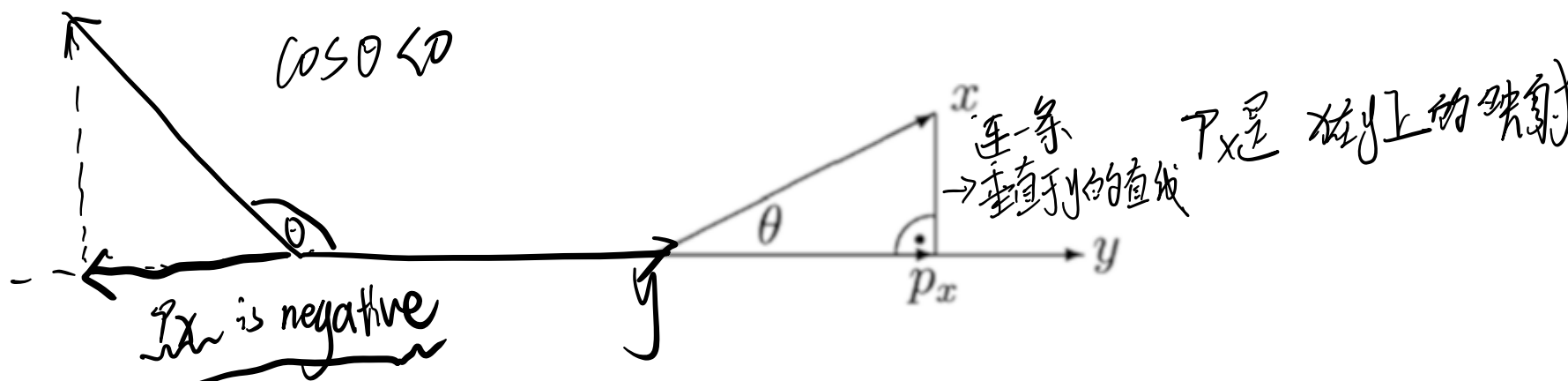
$$\|x\|_A = \sqrt{x^T A x}.$$

## Angle between two vectors

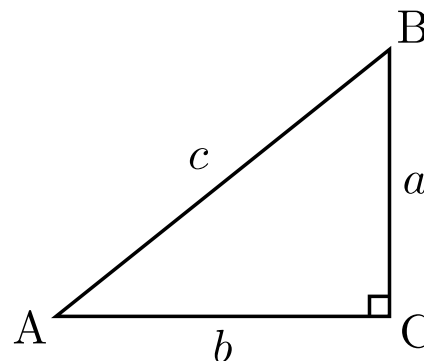
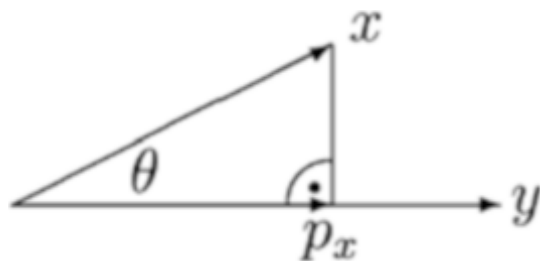
- The angle  $\theta$  between two vectors  $x, y \in \mathbb{R}^p$  is defined through the cosine of  $\theta$  by:

$$\cos(\theta) = \frac{x^T y}{\|x\| \cdot \|y\|}.$$

- Orthogonal **projection**  $p_x \in \mathbb{R}^p$  of  $x \in \mathbb{R}^p$  onto  $y \in \mathbb{R}^p$ ; example in  $\mathbb{R}^2$ :



$p_x$  is projected on the line defined by  $y$ . What is its length  $\|p_x\|$ ?



- From trigono: in right angled triangle ACB with right angle at C,  
 $\cos(\text{angle at A}) = b/c.$

If  $x$  and  $y$  point in same direction ( $x^T y > 0$ ):

$$\cos(\theta) = \|p_x\|/\|x\| \Rightarrow \underbrace{\|p_x\|}_{\text{length of } p_x} = \cos(\theta)\|x\| = x^T y / \|y\|; \quad \|p_x\| = x^T y / \|y\|$$

if point in opposite directions ( $x^T y < 0$ ):  $\|p_x\| = -x^T y / \|y\|.$

In both cases,

$$p_x = \underbrace{\frac{x^T y}{\|y\|}}_{\text{length of } p_x} \cdot \underbrace{\frac{y}{\|y\|}}_{\text{unit vector in the direction of } y}$$

where  $y/\|y\|$  is the unit vector in the direction of  $y$ .

## Rotation

- We often describe a vector in  $\mathbb{R}^p$  through a **system of  $p$  axes** by giving the  $p$  coordinates of the vector in that coordinate system.
- In multivariate statistics it is sometimes useful to **rotate the axes** (all of them at the same time) by an angle  $\theta$ , creating in this way a new  $p$  coordinate system.
- In  $\mathbb{R}^2$ , we can describe a rotation of angle  $\theta$  via the **orthogonal matrix**

$$\Gamma = \begin{pmatrix} \cos(\theta) & \sin(\theta) \\ -\sin(\theta) & \cos(\theta) \end{pmatrix} .$$



逆时针转动  $\Rightarrow$  得到新的向量坐标

If the original axes are **rotated counter clockwise through the origin** by an angle  $\theta$  then the new coordinates  $y$  of a point with coordinates  $x$  in the original system of axes is given by

$$y = \Gamma x.$$

新的  $y = \Gamma x$

If the rotation is **clockwise**, then instead we have

顺时针

$$y = \Gamma^T x.$$

乘一个正交矩阵  $\rightarrow$  对轴转动一个角度.

- More generally, premultiplying a vector  $x$  by an **orthogonal matrix**  $\Gamma$  geometrically corresponds to a rotation of the system of axes.

### 3 MEAN, COVARIANCE, CORRELATION

Sections 3.1, 3.2, 3.3 in Härdle and Simar (2015).

#### 3.1 MEAN

*vector*

- The **mean**  $\mu \in \mathbb{R}^p$  of a random vector  $\mathbf{X} = (X_1, \dots, X_p)^T$  is defined by

$$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix} = \begin{pmatrix} E(X_1) \\ \vdots \\ E(X_p) \end{pmatrix}.$$

- In practice don't usually know  $\mu$  but can estimate it from a sample

$$\mathbf{X}_1 = (X_{11}, \dots, X_{1p})^T, \dots, \mathbf{X}_n = (X_{n1}, \dots, X_{np})^T$$

by the **sample mean**

$$\bar{\mathbf{X}} = \begin{pmatrix} \bar{X}_1 \\ \vdots \\ \bar{X}_p \end{pmatrix},$$

where, for  $j = 1, \dots, p$ , the sample mean

$$\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}$$

of the  $j$ th component  $X_j$  is an estimator of  $\mu_j$ .

- Recall the notation

$$\mathcal{X} = \begin{pmatrix} X_{11} & \dots & X_{1p} \\ X_{21} & \dots & X_{2p} \\ \vdots & & \\ X_{n1} & \dots & X_{np} \end{pmatrix}$$

and  $\mathbf{1}_n = (1, \dots, 1)^T$ , a column vector of length  $n$ .

We can express  $\bar{\mathbf{X}}$  in **matrix notation** as

$$\bar{\mathbf{X}} = n^{-1} \underbrace{\mathcal{X}^T \mathbf{1}_n}_{\text{矩阵乘}} \quad \text{矩阵乘}$$

- Note: in the slides, to avoid too heavy notations, when there is no ambiguity we will not use bold to denote a vector.

### 3.2 COVARIANCE MATRIX

- The **covariance**  $\sigma_{XY}$  between two random variables  $X$  and  $Y$  is a measure of the **linear dependence** between them:

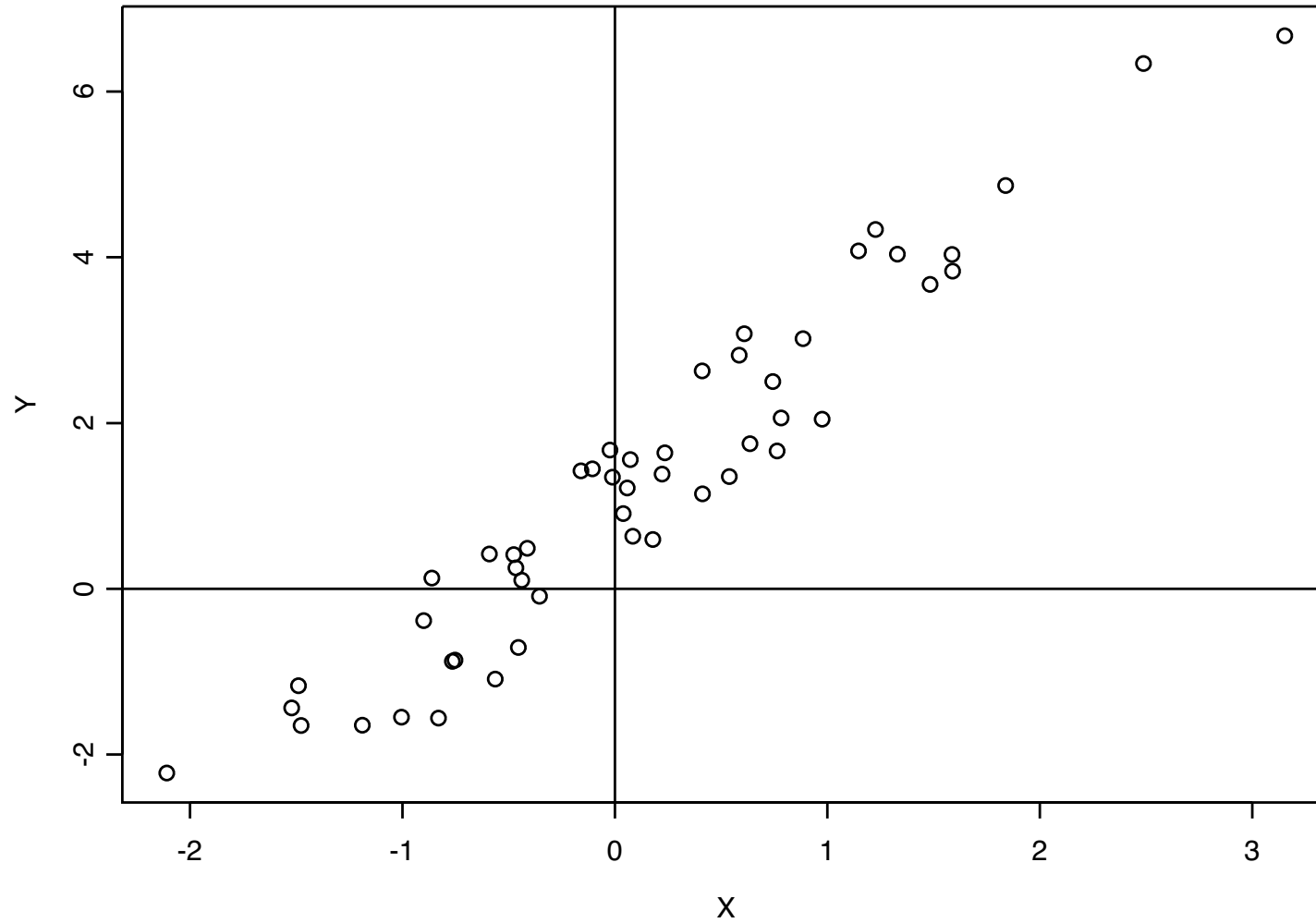
$$\sigma_{XY} = \text{cov}(X, Y) = E(XY) - E(X)E(Y).$$

☞  $\sigma_{XX} = \text{var}(X)$ .

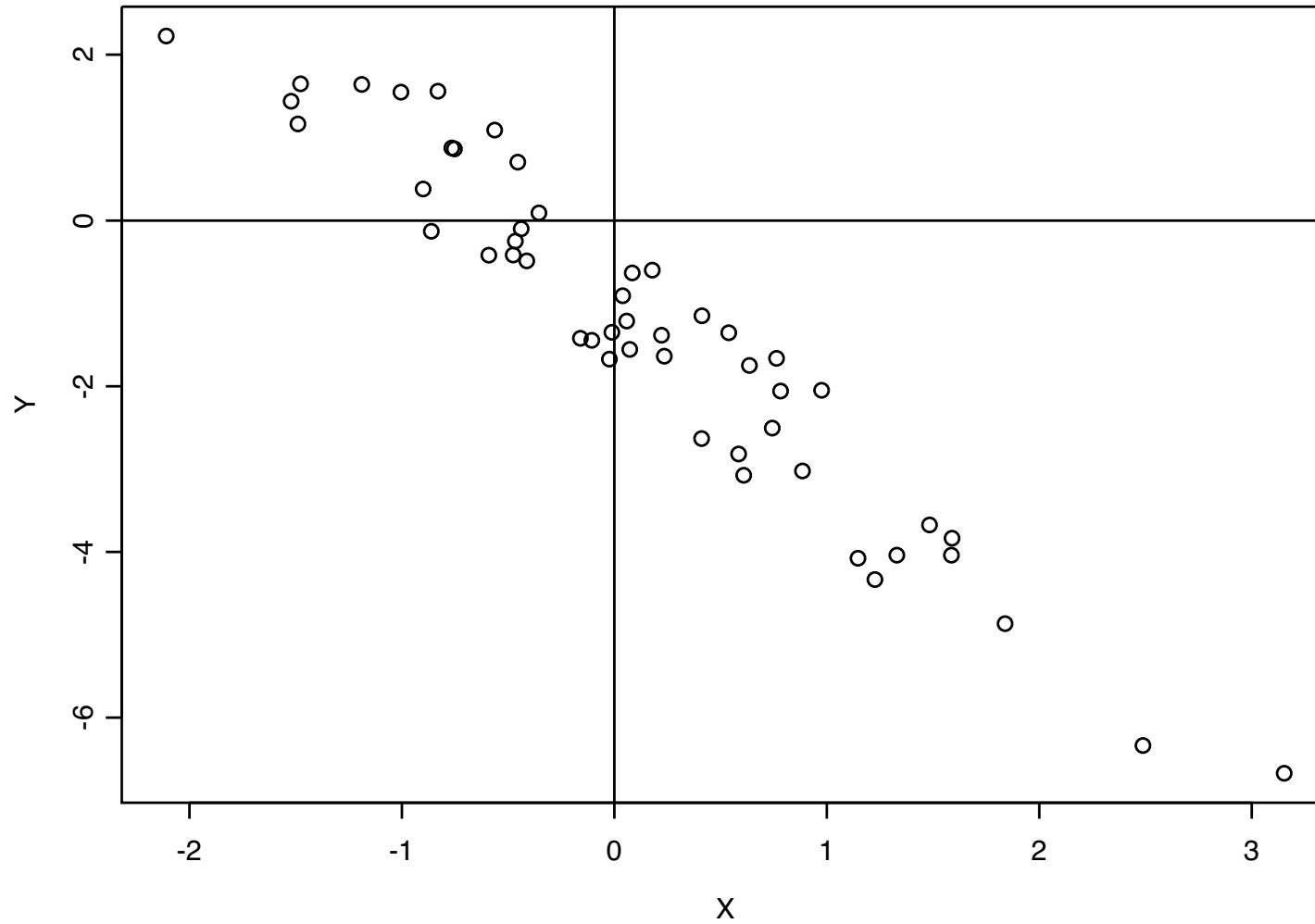
☞ if  $X$  and  $Y$  are independent then  $\sigma_{XY} = 0$ .

☞ However  $\sigma_{XY} = 0$  does not imply that  $X$  and  $Y$  are independent (there could be a nonlinear dependence).

positive covariance



## negative covariance



- If  $\mathbf{X} = (X_1, \dots, X_p)^T$  is a random vector, we can collect the pairwise covariances between each pair  $X_i$  and  $X_j$  in the  $p \times p$  **covariance matrix**  $\Sigma$ :

$$\Sigma = \begin{pmatrix} \sigma_{X_1 X_1} & \dots & \sigma_{X_1 X_p} \\ & \ddots & \\ \sigma_{X_p X_1} & \dots & \sigma_{X_p X_p} \end{pmatrix} = \begin{pmatrix} \sigma_{11} & \dots & \sigma_{1p} \\ & \ddots & \\ \sigma_{p1} & \dots & \sigma_{pp} \end{pmatrix}$$

- ☞ To highlight that it is the covariance of  $\mathbf{X}$  we can write  $\Sigma_{\mathbf{X}}$ .
- ☞  $\Sigma$  is symmetric:  $\Sigma = \Sigma^T$ .
- ☞  $\Sigma$  is semi positive definite:  $\Sigma \geq 0$ .
- ☞ In matrix notation,

$$\Sigma = \text{var}(\mathbf{X}) = E\{(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T\},$$

where  $\mathbf{X}$  and  $\boldsymbol{\mu}$  are written as column  $p$ -vectors .



- In practice  $\Sigma$  is usually unknown but can be estimated from an iid sample  $\mathbf{X}_1, \dots, \mathbf{X}_n$  by the **sample covariance matrix**

$$S = \begin{pmatrix} s_{X_1 X_1} & \dots & s_{X_1 X_p} \\ & \ddots & \\ s_{X_p X_1} & \dots & s_{X_p X_p} \end{pmatrix} = \begin{pmatrix} s_{11} & \dots & s_{1p} \\ & \ddots & \\ s_{p1} & \dots & s_{pp} \end{pmatrix},$$

where, for  $j, k = 1, \dots, p$ ,

$$s_{X_j X_k} = s_{kj} = \overbrace{\left[ \frac{1}{n-1} \right]}^{\text{unbiased}} \sum_{i=1}^n (X_{ij} - \bar{X}_j)(X_{ik} - \bar{X}_k)$$

is the sample covariance between  $X_j$  and  $X_k$ .

- Again, we may write  $S = S_{\mathbf{X}}$  to highlight the correspondence to  $\mathbf{X}$ .

👉 Like  $\Sigma$ ,  $S$  is symmetric ( $S = S^T$ ) and semipositive definite.

- We can obtain  $S$  by computing

$$S = \frac{1}{n-1} \underset{p \times p}{\mathcal{X}^T \mathcal{X}} - \frac{n}{n-1} \underset{p \times n}{\bar{\mathbf{X}}} \underset{n \times p}{\bar{\mathbf{X}}^T}, \quad \text{where } \bar{\mathbf{X}} = \frac{1}{n} \mathbf{X}^T \mathbf{1}_n \text{ is the empirical mean}$$

where  $\mathcal{X}$  is the  $n \times p$  data matrix and  $\bar{\mathbf{X}}$  is written as column  $p$ -vector.

👉 Hint: always check that matrix dimensions are compatible (i.e. matrices products make sense etc).

### 3.3 CORRELATION MATRIX

- Problem with covariance matrix: it is not unit invariant, i.e. if we change the units, covariances change. 单位不变
- The correlation is a measure of linear dependence which is unit invariant.
- The **correlation matrix**  $P$  of a random vector  $\mathbf{X} = (X_1, \dots, X_p)^T$  is a  $p \times p$  matrix defined by:

$$P = \begin{pmatrix} 1 & \rho_{12} & \dots & \rho_{1p} \\ \rho_{21} & 1 & \dots & \rho_{2p} \\ & \vdots & & \\ \rho_{p1} & \rho_{p2} & \dots & 1 \end{pmatrix}$$

where

$$\rho_{ij} = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}\sigma_{jj}}}$$

(Covariance 的标准化版本)

is the correlation between  $X_i$  and  $X_j$ .

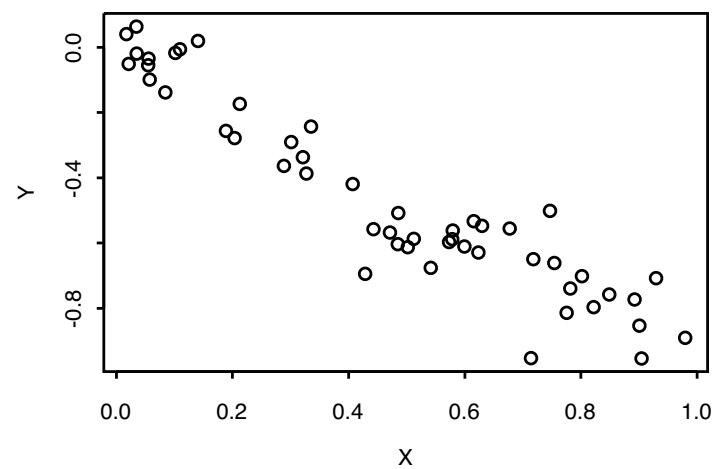
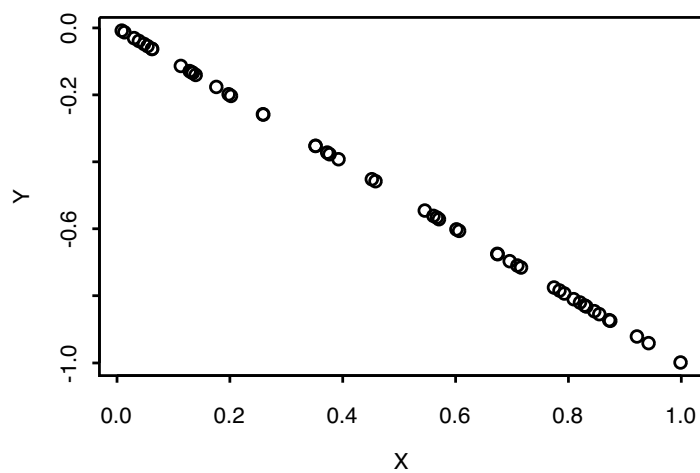
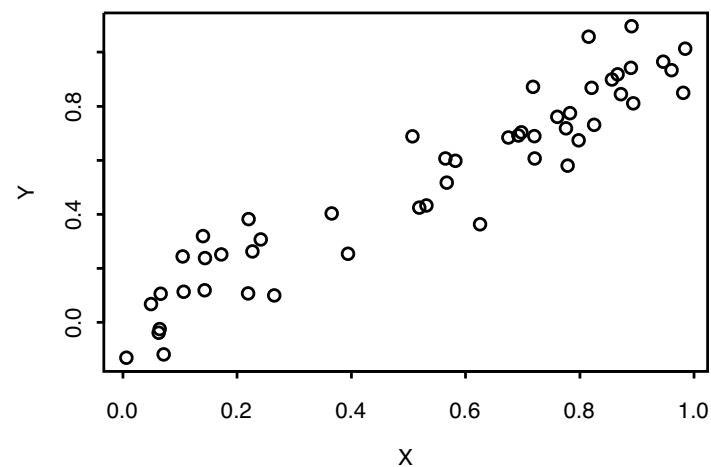
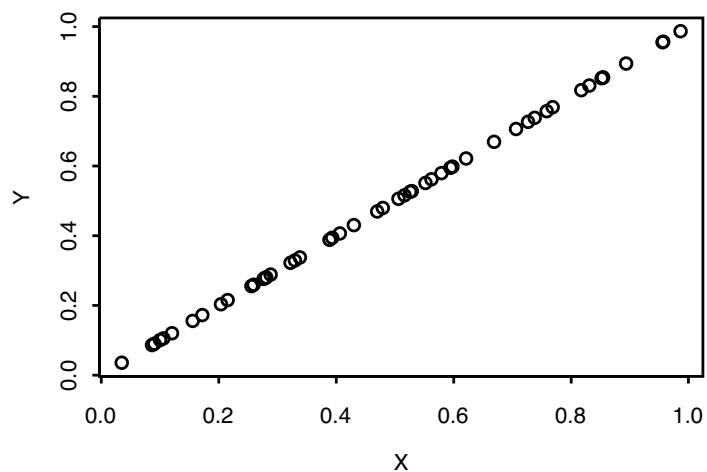
☞ We always have  $-1 \leq \rho_{ij} \leq 1$ .

☞  $\rho_{ij}$  is a measure of the linear relationship between  $X_i$  and  $X_j$ .

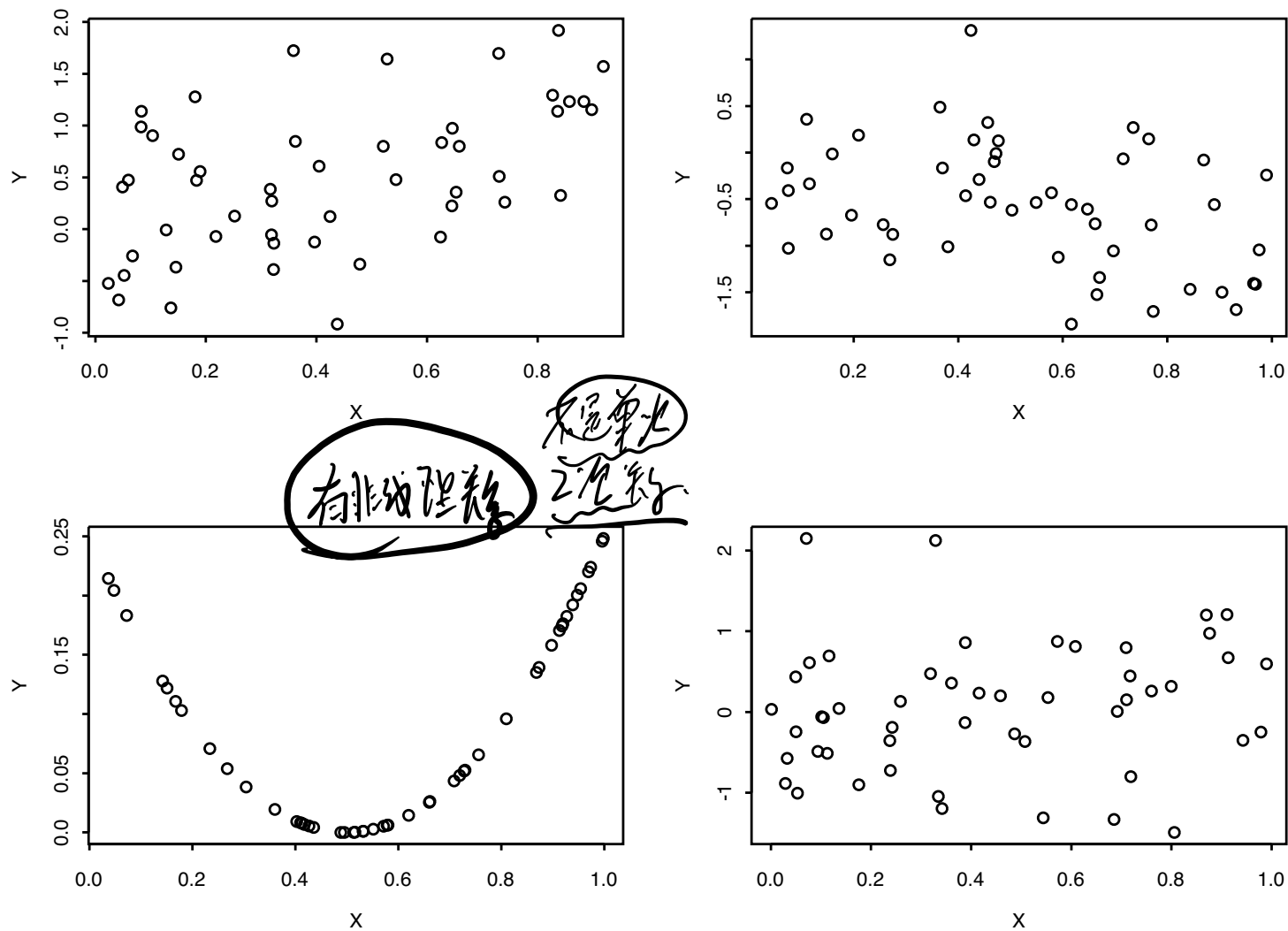
☞  $|\rho_{ij}| = 1$  means perfect linear relationship.

☞  $\rho_{ij} = 0$  means absence of linear relationship, but does not imply independence.  
无线性关系 独立

## Strong positive and negative correlations:



Near zero correlations (does not always imply independence):



- We can obtain  $P$  by computing

$$P = D^{-1/2} \Sigma D^{-1/2}$$

方差-协方差矩阵
对称矩阵

相关系数矩阵
协方差系数矩阵

where  $\Sigma$  is the  $p \times p$  covariance matrix and

$$D = \text{diag}(\sigma_{11}, \dots, \sigma_{pp})$$

is the  $p \times p$  diagonal matrix of variances.

- In practice  $P$  is usually unknown but can be estimated from a iid sample  $\mathbf{X}_1, \dots, \mathbf{X}_n$  by the **sample correlation matrix**

样本的相关系数矩阵

$$R = \begin{pmatrix} r_{11} & \dots & r_{1p} \\ & \ddots & \\ r_{p1} & \dots & r_{pp} \end{pmatrix},$$

where, for  $j, k = 1, \dots, p$ ,

$$r_{jk} = \frac{s_{jk}}{\sqrt{s_{jj}s_{kk}}}$$

is the sample correlation between  $X_j$  and  $X_k$  computed from  $\mathbf{X}_1, \dots, \mathbf{X}_n$ .

- In **matrix notation** we can write

样相关矩阵

$$R = D^{-1/2} S D^{-1/2},$$

where  $S$  is the  $p \times p$  sample covariance matrix and, on this occasion,

$$D = \text{diag}(s_{11}, \dots, s_{pp})$$

is the  $p \times p$  diagonal matrix of sample variances.



### 3.4 LINEAR TRANSFORMATIONS

Let  $\mathbf{X} = (X_1, \dots, X_p)^T$  be a  $p$ -vector and let  $\mathbf{Y}$  be  $q$ -vector defined by

$$\mathbf{Y} = A\mathbf{X} + \mathbf{b},$$

where  $A$  is a  $q \times p$  matrix and  $\mathbf{b}$  is a  $q \times 1$  vector. Then we have

$$E(\mathbf{Y}) = A \cdot E(\mathbf{X}) + \mathbf{b}$$

$$\bar{\mathbf{Y}} = A\bar{\mathbf{X}} + \mathbf{b}$$

$$\Sigma_{\mathbf{Y}} = A\Sigma_{\mathbf{X}}A^T$$

$$S_{\mathbf{Y}} = AS_{\mathbf{X}}A^T$$

👉 Hint: to know where to put the transpose, always check that matrix dimensions are compatible.

## 4 MULTIVARIATE DISTRIBUTIONS

### 4.1 DISTRIBUTION AND DENSITY FUNCTION

Sections 4.1, 4.2 in Härdle and Simar (2015).

Let  $\mathbf{X} = (X_1, \dots, X_p)^T$  be a random vector.

- For all  $\mathbf{x} = (x_1, \dots, x_p)^T \in \mathbb{R}^p$ , the cumulative distribution function (cdf), or distribution function, of  $\mathbf{X}$  is defined by

$$\underline{F_{\mathbf{X}}(\mathbf{x}) = P(\mathbf{X} \leq \mathbf{x}) = P(X_1 \leq x_1, \dots, X_p \leq x_p)}$$

- When there is no ambiguity, we can write  $F$  instead of  $F_{\mathbf{X}}$ . Advantage: less heavy notations.

- If  $\mathbf{X}$  is continuous, the probability density function (pdf) or density,  $f_{\mathbf{X}}$ , of  $\mathbf{X}$  is a nonnegative function defined through

$$F_{\mathbf{X}}(\mathbf{x}) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_p} f_{\mathbf{X}}(\mathbf{u}) d\mathbf{u} \equiv \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_p} f_{\mathbf{X}}(u_1, \dots, u_p) du_1 \dots du_p,$$

where  $\mathbf{u} = (u_1, \dots, u_p)$ . n维PDF 与2维类似

- It always satisfies

$$\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_{\mathbf{X}}(\mathbf{u}) d\mathbf{u} = 1.$$

- When there is no ambiguity, we can write  $f$  instead of  $f_{\mathbf{X}}$ .

- The marginal cdf of a subset of  $\mathbf{X}$  is obtained by the marginal of  $\mathbf{X}$  computed at the subset, letting the other values equal to infinity.

➡ For example, the marginal cdf of  $X_1$  is obtained by taking

★ marginal:  $X_1$   
 的分布函数。  
 $X_1$  不动其他

$$\begin{aligned} F_{X_1}(x_1) &= P(X_1 \leq x_1) \\ &= P(X_1 \leq x_1, X_2 \leq \infty, \dots, X_p \leq \infty) \\ &= F_{\mathbf{X}}(x_1, \infty, \dots, \infty) \end{aligned}$$

➡ and the marginal cdf of  $(X_1, X_3)$  is obtained by taking

$$\begin{aligned} F_{X_1, X_3}(x_1, x_3) &= P(X_1 \leq x_1, X_3 \leq x_3) \\ &= P(X_1 \leq x_1, X_2 \leq \infty, X_3 \leq x_3, X_4 \leq \infty, \dots, X_p \leq \infty) \\ &= F_{\mathbf{X}}(x_1, \infty, x_3, \infty, \dots, \infty). \end{aligned}$$

CDF 求导  $\Rightarrow$  PDF

- For a continuous random vector  $\mathbf{X}$ , the **marginal density** of a subset of  $\mathbf{X}$  is obtained from the joint density  $f_{\mathbf{X}}$  of  $\mathbf{X}$  by integrating out the other components.

👉 For example, the marginal density  $X_1$  is obtained by taking

$$f_{X_1}(x_1) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_{\mathbf{X}}(x_1, u_2, \dots, u_p) du_2 \dots du_p$$

👉 and the marginal density of  $(X_1, X_3)$  is obtained by taking

$$f_{X_1, X_3}(x_1, x_3) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_{\mathbf{X}}(x_1, u_2, x_3, u_4, \dots, u_p) du_2 du_4 \dots du_p .$$

- For two continuous random vectors  $X_1$  and  $X_2$ , the **conditional pdf** of  $X_2$  given  $X_1$  is given by

$$f_{X_2|X_1}(x_2|x_1) = f_{X_1,X_2}(x_1, x_2) / f_{X_1}(x_1).$$

It is defined only for values  $x_1$  such that  $f_{X_1}(x_1) > 0$ .

- Two continuous random vectors  $X_1$  and  $X_2$  are **independent** if and only if

$$f_{X_1,X_2}(x_1, x_2) = f_{X_1}(x_1)f_{X_2}(x_2).$$

👉 If  $X_1$  and  $X_2$  are independent then

$$\underbrace{f_{X_2|X_1}(x_2|x_1)}_{\text{pdf of } X_2} = f_{X_1,X_2}(x_1, x_2) / f_{X_1}(x_1) = f_{X_1}(x_1)f_{X_2}(x_2) / f_{X_1}(x_1) = \underbrace{f_{X_2}(x_2)}_{\text{pdf of } X_2}.$$

Thus knowing the value of  $X_1$  does not change probability assessments on  $X_2$  and vice versa.

- The **mean**  $\mu \in \mathbb{R}^p$  of a random vector  $X = (X_1, \dots, X_p)^T$  is defined by

$$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix} = \begin{pmatrix} E(X_1) \\ \vdots \\ E(X_p) \end{pmatrix} = \begin{pmatrix} \int x f_{X_1}(x) dx \\ \vdots \\ \int x f_{X_p}(x) dx \end{pmatrix}.$$

☞ If  $X$  and  $Y$  are two  $p$ -vectors and  $\alpha$  and  $\beta$  are constants then

$$E(\alpha X + \beta Y) = \alpha E(X) + \beta E(Y).$$

☞ If  $X$  is a  $p \times 1$  vector which is independent of the  $q \times 1$  vector  $Y$  then

$$E(XY^T) = E(X)E(Y^T).$$

☞ Hint: Remember to always check that matrix dimensions are compatible.

The **conditional expectation**  $E(X_2|X_1 = x_1)$  is defined by

$$E(X_2|X_1 = x_1) = \int \underbrace{x_2 f_{X_2|X_1}(x_2|x_1)}_{\int x_2 \gamma_1(x_2|x_1) = \int_{x_1, x_2} \gamma_1(x_2|x_1)} dx_2$$

and the **conditional covariance** matrix  $\text{var}(X_2|X_1 = x_1)$  is defined by

$$\text{var}(X_2|X_1 = x_1) = E(X_2 X_2^T | X_1 = x_1) - E(X_2 | X_1 = x_1) E(X_2^T | X_1 = x_1),$$

if  $X_2$  is a column vector.

👉 Hint: In doubt check the dimension of the resulting matrices to see if you get them right.



- As seen earlier, the **covariance** matrix  $\Sigma$  of a vector  $X$  of mean  $\mu$  is defined by

$$\Sigma = \text{var}(X) = E\{(X - \mu)(X - \mu)^T\}.$$

We write

$$X \sim (\mu, \Sigma)$$

to denote a vector  $X$  with mean  $\mu$  and covariance matrix  $\Sigma$ .

- We can also define a **covariance matrix** between a  $p \times 1$  vector  $X$  of mean  $\mu$  and a  $q \times 1$  vector  $Y$  of mean  $\nu$  by

$$\Sigma_{X,Y} = \text{cov}(X, Y) = E\{(X - \mu)(Y - \nu)^T\} = E(XY^T) - E(X)E(Y^T).$$



The elements of this matrix are the pairwise covariances between the components of  $X$  and those of  $Y$ .

👉 For  $p \times 1$  vectors  $X$  and  $Y$  and a  $q \times 1$  vector  $Z$ , we have

$$\text{cov}(X + Y, Z) = \text{cov}(X, Z) + \text{cov}(Y, Z)$$

👉 For  $p \times 1$  vectors  $X$  and  $Y$  we have  $X$  不是向量是必须矩阵

$$\text{var}(X + Y) = \text{var}(X) + \boxed{\text{cov}(X, Y) + \text{cov}(Y, X)} + \text{var}(Y)$$

$E(XY) - E(X)E(Y)$

👉 For matrices  $A$  and  $B$  and random vectors  $X$  and  $Y$  of dimensions such that the below quantities are well defined we have

$$\text{cov}(AX, BY) = A \text{cov}(X, Y) B^T.$$

## 4.2 MULTINORMAL DISTRIBUTION

Sections 4.4, 4.5, 5.1 in Härdle and Simar (2015).

A very useful and commonly encountered distribution is the **multinormal** distribution, also simply called **normal** distribution.

☞ Recall that in the univariate case, the density of a  $N(\mu, \sigma^2)$  is given by

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ - (x - \mu)^2 / (2\sigma^2) \right\} .$$

☞ In the multivariate case, need to deal with vectors and matrices.

👉 The density of a normal vector  $X = (X_1, \dots, X_p)^T$  with mean  $\mu = (\mu_1, \dots, \mu_p)^T$  and positive definite covariance matrix

$$\Sigma = \begin{pmatrix} \sigma_{11} & \dots & \sigma_{1p} \\ & \ddots & \\ \sigma_{p1} & \dots & \sigma_{pp} \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \dots & \sigma_{1p} \\ & \ddots & \\ \sigma_{p1} & \dots & \sigma_p^2 \end{pmatrix},$$

where  $\sigma_j^2 = \text{var}(X_j)$ , is given by

$$f(x) = |2\pi\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu) \right\}. \quad (1)$$

Handwritten note:  $\Sigma$  is a  $p \times p$  matrix

👉 If the  $p$ -vector  $X$  is normal with mean  $\mu$  and cov matrix  $\Sigma$  we write

$$X \sim N_p(\mu, \Sigma).$$

If the  $X_i$ 's are independent, then

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \dots & 0 \\ & \ddots & \\ 0 & \dots & \sigma_p^2 \end{pmatrix} = \text{diag}(\sigma_1^2, \dots, \sigma_p^2).$$

Thus

$$|2\pi\Sigma| = |\text{diag}(2\pi\sigma_1^2, \dots, 2\pi\sigma_p^2)| = (2\pi)^p \sigma_1^2 \dots \sigma_p^2$$

对称矩阵行列式的算法 (不用元素连乘)

and

$$\Sigma^{-1} = \text{diag}(\sigma_1^{-2}, \dots, \sigma_p^{-2})$$

so that

$$\begin{aligned} f(x) &= \frac{1}{\sqrt{(2\pi)^p} \prod_{j=1}^p \sigma_j} \exp \left\{ -\frac{1}{2} \sum_{j=1}^p (x_j - \mu_j)^2 / \sigma_j^2 \right\} \\ &= \frac{1}{\sqrt{(2\pi)^p} \prod_{j=1}^p \sigma_j} \prod_{j=1}^p \exp \left\{ -\frac{1}{2} (x_j - \mu_j)^2 / \sigma_j^2 \right\} \\ &= \prod_{j=1}^p \left[ \frac{1}{\sqrt{2\pi} \sigma_j} \exp \left\{ - (x_j - \mu_j)^2 / (2\sigma_j^2) \right\} \right]. \end{aligned}$$

is the product of densities of  $p$  univariate  $N(\mu_j, \sigma_j^2)$ .

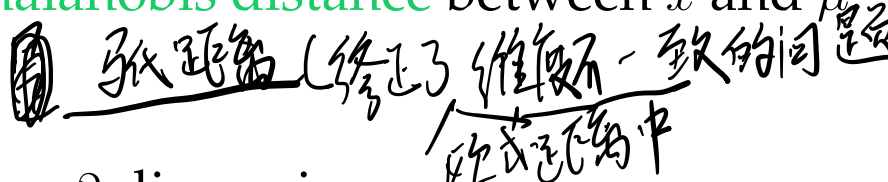
We see from (1) that  $f(x)$  takes the same value for all  $x \in \mathbb{R}^p$  such

$$(x - \mu)^T \Sigma^{-1} (x - \mu) = c$$

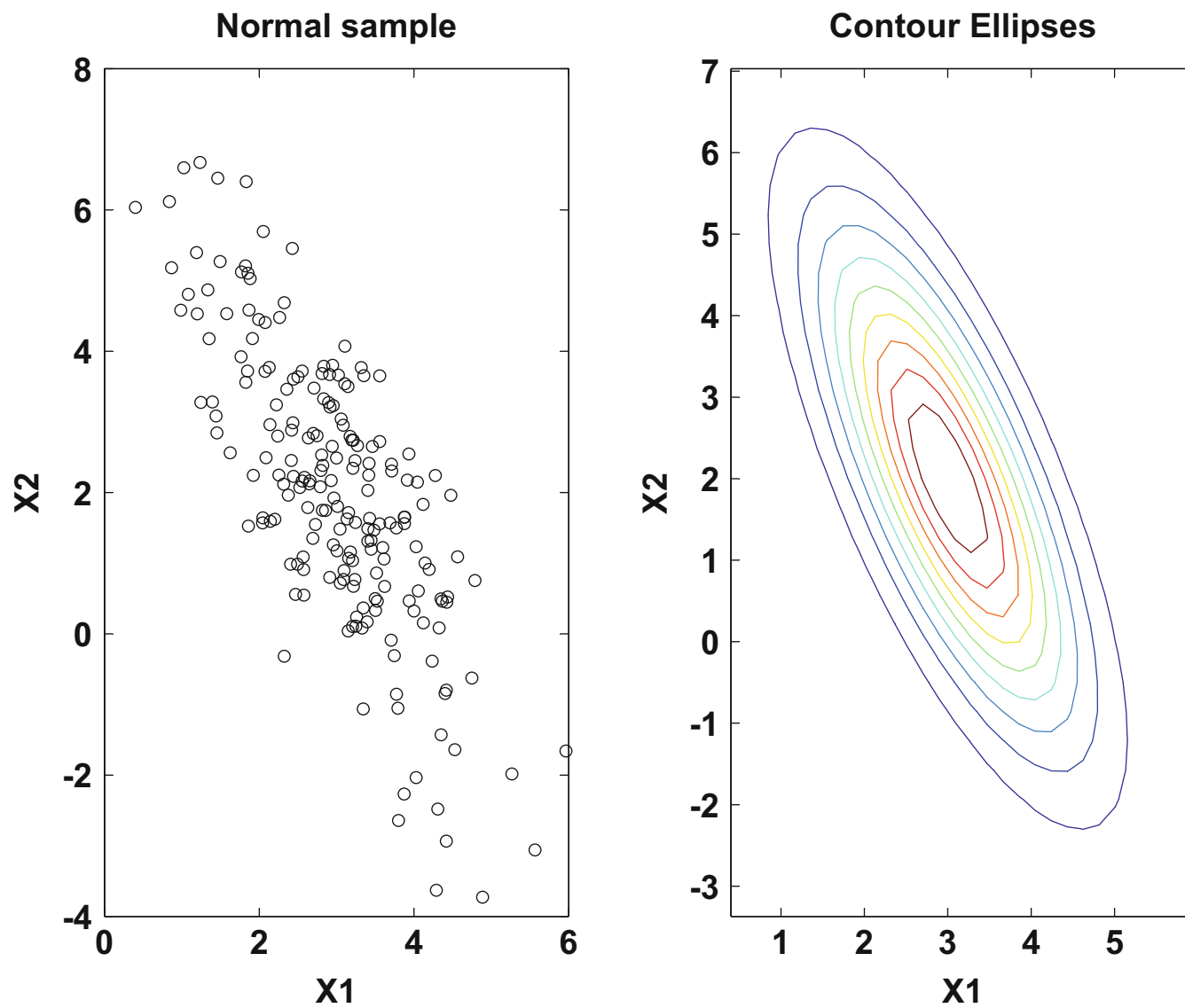
where  $c$  is a positive constant. For each  $c > 0$ , these  $x$ -values correspond to an **ellipsoid** (a different one for each  $c > 0$ ; they are called contour ellipsoids).

The quantity

$$\sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

is called the **Mahalanobis distance** between  $x$  and  $\mu$ .  


For example in  $p = 2$  dimensions:



**Fig. 4.3** Scatterplot of a normal sample and contour ellipses for  $\mu = \begin{pmatrix} 3 \\ 2 \end{pmatrix}$  and  $\Sigma = \begin{pmatrix} 1 & -1.5 \\ -1.5 & 4 \end{pmatrix}$  

☞ Let  $X \sim N_p(\mu, \Sigma)$ ,  $A$  a  $q \times p$  matrix and  $b$  a  $q \times 1$  vector. Then

$$Y = AX + b \sim N_q(A\mu + b, A\Sigma A^T).$$

☞ Let  $X = (X_1^T, X_2^T)^T \sim N_p(\mu, \Sigma)$  where  $X_1$  and  $X_2$  are two column vectors. Then

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

where

$$\Sigma_{11} = \text{var}(X_1), \quad \Sigma_{22} = \text{var}(X_2), \quad \Sigma_{12} = \text{cov}(X_1, X_2) \quad \Sigma_{21} = \text{cov}(X_2, X_1).$$

Then one can prove (not us)

$$\Sigma_{12} = 0 \iff X_1 \text{ and } X_2 \text{ are independent.}$$



• If  $X \sim N_p(\mu, \Sigma)$  and  $A$  and  $B$  are matrices with  $p$  columns, then

$$\underline{AX \text{ and } BX \text{ are independent} \iff A\Sigma B^T = 0.} \quad (2)$$

充要条件.

$A\Sigma B^T = 0 \Rightarrow AX$  与  $BX$  独立

• If  $X \sim N_p(\mu, \Sigma)$  and  $\Sigma$  is invertible, then

$$Y = (X - \mu)^T \Sigma^{-1} (X - \mu) \sim \chi_p^2 \quad (\text{standardise the variable}) \quad (3)$$

(chi square with  $p$  degrees of freedom).

• If  $X_1, \dots, X_n$  are i.i.d.  $\sim N_p(\mu, \Sigma)$ , then

covariance matrix.

$$\bar{X} \sim N_p(\mu, \Sigma/n)$$

(4)

proof:

$$y = (x - \mu)^T \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}} (x - \mu)$$

$$= z^T z$$

$$\text{where } z = \Sigma^{-\frac{1}{2}} (x - \mu) \sim N_p(0, I_p)$$

$$\Sigma^{-1} = \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}}$$

$$\sum z_i^2 \sim \chi_p^2(0)$$

### 4.3 WISHART DISTRIBUTION

extension of  $\chi^2$

$\frac{\lambda}{\sqrt{b}} \left( \frac{\sum}{\sum} \right)$

- The **Wishart distribution** is a generalisation to multiple dimensions of the **chi square** distribution.

It depends on 3 parameters:  $p$ , a  $p \times p$  scale matrix  $\Sigma$  and the number of degrees of freedom  $n$ :

$$W_p(\Sigma, n).$$

- Recall that if  $Z_1, \dots, Z_n$  are independent  $N(0, 1)$  then

$$X = \sum_{k=1}^n Z_k^2 \sim \chi_n^2$$

is a chi square with  $n$  degrees of freedom.

- If  $M$  is an  $p \times n$  matrix whose columns are independent and all have a  $N_p(0, \Sigma)$  distribution, then the matrix

$M$ : 各列独立

$$MM^T \sim W_p(\Sigma, n),$$

columns are normal distribution  $p \times p$

i.e.  $MM^T$  has a Wishart distribution with parameters  $p, \Sigma$  and  $n$ .

行数  
列数  
自由度  
形状

$$W_1(b^2, n)$$

$$y = MM^T = (\pi_1, \dots, \pi_n) \begin{pmatrix} \pi_1 \\ \pi_2 \\ \vdots \\ \pi_n \end{pmatrix} = \pi_1^2 + \dots + \pi_n^2$$

where  $\pi_j$ 's are independent  $N(0, b^2)$

Thus  $\pi_j = bz_j$  where  $z_j \sim N(0, 1)$  and the  $z_j$  are independent.

$$\text{Thus } M = b^2 \sum_{j=1}^n z_j^2$$

/43