

Multivariate  $\rightarrow$  data  $(X_1, X_2 \dots X_p)$ , we consider live in a multivariate space, that is  $P > 2$   
this variables  $(X_1 \dots X_p)$  are linked / dependent

analysis  $\rightarrow$  MA is a collection of methods based on probabilistic models to understand, analyze and interpret the data

$\rightarrow$  data reduction  
discrimination  
sorting  
group  
dependence among variables

different goal

The model depend on the type of variables, the type of the dependence, and the type of the problems.

補充 P<sub>1</sub>: this class  $\rightarrow$  unsupervised learning

include a lot of  $\uparrow$   
classic statistics

you have data  $(X_1 \dots X_p) \in \mathbb{R}^{P \times p}$ ,

you want fulfill some goals.

no label

you have data  $(X_1 \dots X_p) \in \mathbb{R}^{P \times l'}$

and have some response variable

$(Y_1 \dots Y_n) \in \mathbb{R}$  or category

we want do  $Y = f(x)$

have label

usually use to split  
itself  $\rightarrow$

supervised learning

## content P1

Background on multivariate distribution

- ① the data: sample mean, sample covariance
- ② the model: population versions, multivariate families
- ③ densities, cumulative distribution ~~of distributions.~~ functions

## content P2

dependence modeling

- ① How we can describe the dependence between 2 variables
- ② Copulas: separating the marginals from the dependence

## content P3

graphic model

- ① what's the structure of multivariate obs.
- ② graphs as a description of conditional independence relations.
- ③ High-dimensional / data sets and regularization

## content P4

Unsupervised learning

- ① Dimension reduction / data compression.
- ② Visualization of high-dimensional data

- ③ principal components analysis
- ④ clustering method
- ⑤ extend PCA to non-linear transformations  
(Autoencoders)
- ⑥ sample from the distribution of images?

## Variational Autoencoder and GAN

### VAE

Model and data perspective

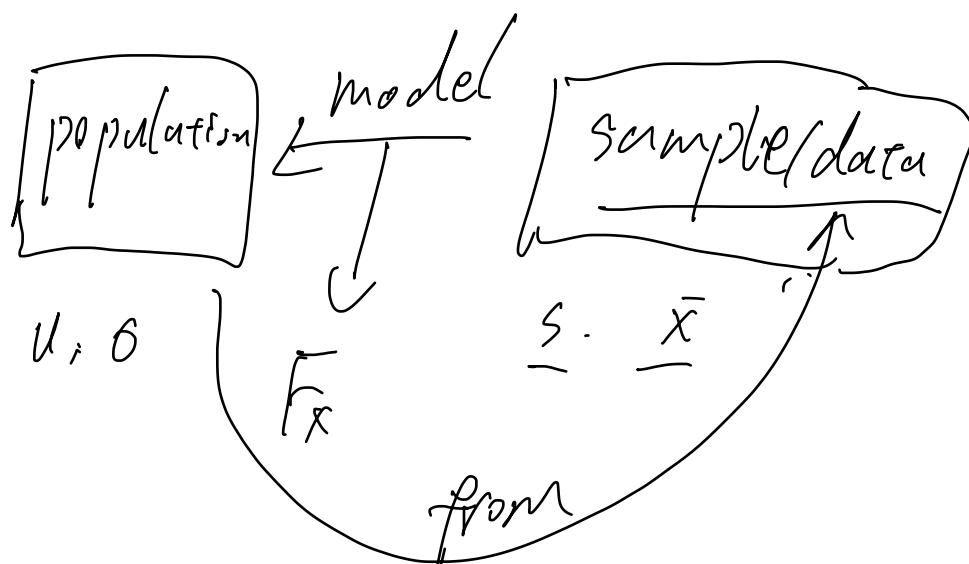
Model perspective:

- ① the model is a random vector  $X = (X_1 \dots X_p)^T \in \mathbb{R}^P$  following a certain distribution  $\underline{F_X}$  describing the entire (infinite) population
- ② we can compute the statistics such as the population mean or population variance for  $X$

Data perspective

- ① the data/sample  $\underline{X_1 \dots X_n} \in \mathbb{R}^P$  are thought of finite number  $n$  of observations from  $X$  this distribution is unknown
- right?  
questions? Range  $n \in [0, P]$   
in feature  $P$  為維度  $\rightarrow n \times p$

② based on data, we can compute sample version of the mean or the variance. they are approximations of the population version.



model {  
 vector  
 distribution  
 infinite population

the model:  $p=2$ , 2-dimensional

Let's  $X = (X_1, X_2)^T$  be a bivariate normal distribution which is  $N_2(\mu, \Sigma)$ , with mean  $(2, 5)$  and covariance matrix:  $\Sigma = \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix}$

$\rightarrow$  model 的角度  
 -  $X$  的 取值.  $P.$   
 $P$

the data:

we simulate  $n=50$  samples  $X_1, \dots, X_n \in \mathbb{R}^2$  from the normal distribution

we get  $\bar{X} = (2.06, 5.07)$   $S = \begin{pmatrix} 0.98 & 0.71 \\ 0.71 & 0.11 \end{pmatrix}$

model  $\xrightarrow{\text{simulation}}$  data

inference, but never know if exactly.

从样本  
 对称分布  
 所得

对称分布  
 所得

Some notation: *data* + *variable*  $\Rightarrow$   
 for a vector  $x \in \mathbb{R}^P$  we write

$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_P \end{pmatrix}$$

$$x^T = (x_1 \cdots x_P)$$

$u = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i / n$  if we  
 mode  $\Rightarrow$  get more  
 sample  $\Rightarrow$  data,  
 $\bar{x} \rightarrow u$

if  $x, y \in \mathbb{R}^P$ ,  $x^T y \in \mathbb{R}$  and  $x y^T \in \mathbb{R}^{P \times P}$

$A \in \mathbb{R}^{m \times m}$  in boldface the diagonal matrix is

$$\text{diag}(a_1 \cdots a_m) = \begin{pmatrix} a_1 & 0 & 0 & \cdots & 0 \\ 0 & a_2 & 0 & \cdots & 0 \\ 0 & 0 & a_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & a_m \end{pmatrix} \in \mathbb{R}^{m \times m}$$

for  $n$  data  $x_1 \cdots x_n \in \mathbb{R}^P$ , we write

$$X = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1P} \\ x_{21} & & & \\ \vdots & & & \\ x_{n1} & \cdots & x_{nP} \end{pmatrix} \in \mathbb{R}^{n \times P}$$

obs  $n$  行数  
feature  $P$  列数

population

Let  $X = (X_1 \dots X_p)^T \sim F$  be a random  $\mathbb{R}^p$ -vector with joint distribution function:

$$F_X(y) = F(y) = F(y_1, y_2 \dots y_p) = P(X_1 \leq y_1, \dots, X_p \leq y_p) \in \mathbb{R}^P$$

the univariate marginal distribution of  $X_j$ ,  $j=1 \dots p$  is

$$F_j(y_j) = P(\infty \dots X_j \leq y_j, \infty)$$

Similarly, we get higher dimensional marginals by

$$\underline{F_{jk}(y_j, y_k) = P(X_j \leq y_j, X_k \leq y_k)}$$

F describes the margins and the dependence structure of the population

if F is differentiable (we assume), the joint density is the derivative:

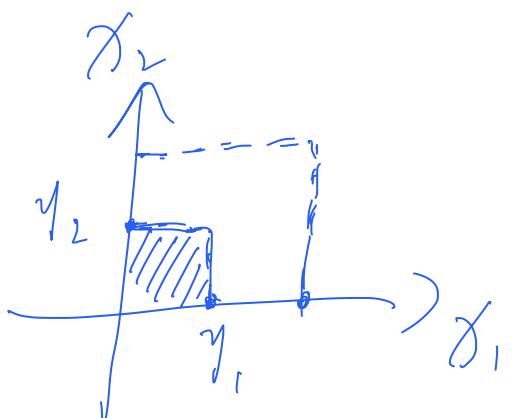
$$f(y) = f(y_1 \dots y_p) = \frac{\partial^p F(y)}{\partial y_1 \dots \partial y_p}$$

or equivalently:

$$F(y) = \int_{-\infty}^{y_1} \dots \int_{-\infty}^{y_p} f(u_1 \dots u_p) du_1 \dots du_p$$

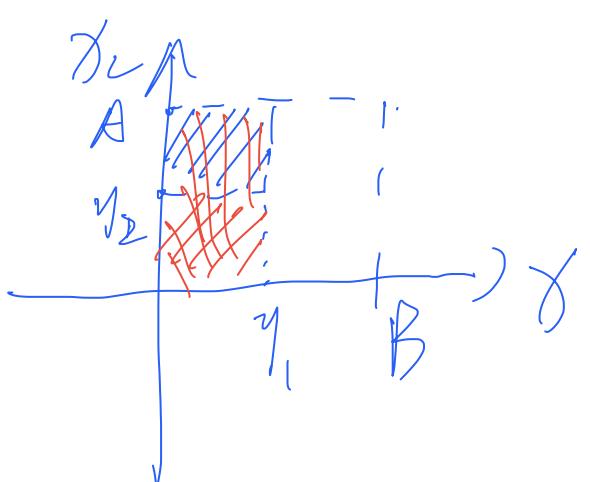
本质上:  $F$  [上页画的] 为 distribution, 可以替换为任意 distribution:

举例: uniform distribution



$$F_X(y_1, y_2) = P(X_1 \leq y_1, X_2 \leq y_2)$$

"/" is the P



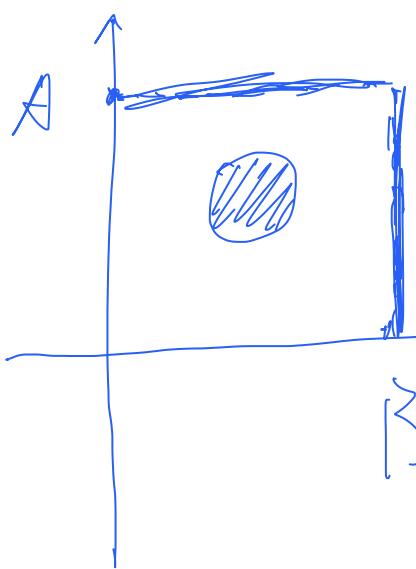
从图可推得

$$P(X_1 \leq y_1, X_2 \leq y_2)$$

A is the range  $[c, A]$

↓ 得 F

$$F(y_1) - F(y_1, y_2)$$



" / " — " / "  
 $\Rightarrow$  if we want to know  
 the  $P(\cdot)$  which is not limit  
 rectangles

(A) one way we try to is  $\#$  rectangles  
 演进.

the joint density exists  $\rightarrow$  the marginal densities

$$f_{\mathbf{y}}(\mathbf{y}) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(u_1, u_2, \dots, u_p) du_1 \cdots du_p$$

the density function integrates to 1

$$F(\infty, \infty) = 1 \quad \text{① PS}$$

Sample quantity: we don't know distribution  
we need to estimate it

Let  $x_1, \dots, x_n \in \mathbb{R}^p$  be the n samples of  $X$   
You should see  $x_1, \dots, x_n$  as random realizations of  $X$   
the empirical distribution function

$$\hat{F}_{\mathbf{y}}(y) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq y, \dots, x_{ip} \leq y_p) \quad \text{② PS}$$

is a non-parametric estimator of  $F$

there is no "empirical" sample-equivalent for the density function without further assumptions. why?

example:

Let  $X = (x_1, x_2)^T$  be a bivariate random variable with the density function:  $f_{\mathbf{y}}(y_1, y_2) = r(y_1^\alpha + y_2^\beta), 0 < y_1, y_2 < 1$ , otherwise  $\alpha, \beta > 0$

$$\int_0^1 \int_0^1 r(y_1^\alpha + y_2^\beta) dy_1 dy_2 = 1 \Rightarrow \int_0^1 \int_0^1 r y_1^\alpha dy_1 dy_2 + \int_0^1 \int_0^1 r y_2^\beta dy_2 dy_1 = 1$$

$$\text{we get: } r = [(\alpha+1)^{-1} + (\beta+1)^{-1}]^{-1} = \int_0^1 r y_1^\alpha dy_1 + \int_0^1 r y_2^\beta dy_2 = \int_0^1 r y_1^{\alpha+1} dy_1 + \int_0^1 r y_2^{\beta+1} dy_2 = \frac{r}{\alpha+1} + \frac{r}{\beta+1} = r$$

$$\alpha=1, \beta=2, r=\frac{1}{3}$$

$$\begin{aligned} F &= \int_0^{y_1} \int_0^{y_2} \frac{6}{5} (u_1 + u_2^2) du_1 du_2 \\ F &= \frac{6}{5} y_1 y_2 [\frac{1}{2} y_1 + \frac{1}{3} y_2^2] \end{aligned}$$

$$F(X_1 < \frac{1}{2}, X_2 > \frac{1}{2}) = \frac{1}{4} \quad \text{Probabilities 的概率}$$

for  $F$  类似的:

$$\hat{P}(X_1 < \frac{1}{2}, X_2 > \frac{1}{2}) = \hat{F}(\frac{1}{2}, 1) - \hat{F}(\frac{1}{2}, \frac{1}{2}) = \frac{1}{n} \sum_{i=1}^n I(X_{i1} \leq \frac{1}{2}, X_{i2} > \frac{1}{2})$$

population quantity:

Let  $X = (x_1, \dots, x_p)^T$  be a random vector

we can define the expectation of  $g_{\mathbf{x}}$ , where  $g: \mathbb{R}^p \rightarrow \mathbb{R}^m$  is some function

$$E(g_{\mathbf{x}}) = \int_{\mathbb{R}^p} g_{\mathbf{x}}(u) dF_{\mathbf{x}}(u) = \int_{\mathbb{R}^p} g_{\mathbf{x}}(u) f_{\mathbf{x}}(u) du \in \mathbb{R}^m$$

for n samples  $x_1, \dots, x_n \in \mathbb{R}^p$  of  $X$ , we define the sample version of  $E(g_{\mathbf{x}})$  as

$$\hat{E}(g_{\mathbf{x}}) = \int_{\mathbb{R}^p} g_{\mathbf{x}}(u) d\hat{F}_{\mathbf{x}}(u) = \frac{1}{n} \sum_{i=1}^n g_{\mathbf{x}}(x_{i1}, \dots, x_{ip})$$

For a random vector  $X \in \mathbb{R}^p$  and  $Y \in \mathbb{R}^q$  matrix  $A \in \mathbb{R}^{m \times p}$   $B \in \mathbb{R}^{m \times q}$  and constant vector  $b \in \mathbb{R}^m$

$$Z = AX + BY + b \in \mathbb{R}^m$$

$$E(Z) = A \cdot E(X) + B \cdot E(Y) + b.$$

the trace is an operator on square matrices that returns the sum of the diagonal elements.  $\text{对角元素和}$

Let  $A = \{a_{kl}\}_{k,l=1}^m$  be  $m \times m$  matrix, then

$$\text{Tr}(A) = \sum_{i=1}^m a_{ii}$$

the trace is linear and circular (with compatible matrix)

$$\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B), \text{Tr}(ABC) = \text{Tr}(CAB)$$

For a random vector  $X \in \mathbb{R}^p$  and a deterministic matrix  $A \in \mathbb{R}^{p \times p}$

$$E(X^TAX) = \mathbb{E}[\text{Tr}(X^TAX)] = \mathbb{E}[\text{Tr}(AXX^T)] = \text{Tr}[A \cdot E(XX^T)]$$

how can I prove it?  $\hookrightarrow$  by linearity  $\mathbb{E}[\text{Tr}(c)] = \text{Tr}[\mathbb{E}(c)]$   
population quantity  $\hookrightarrow X^TAX$  is a scalar,  $\text{Tr}(\text{scalar}) = \text{scalar}$ .

the variance-covariance matrix of  $X$  is the matrix of variances and covariances between the p components

$$\Sigma = \text{cov}(X) = \{b_{kl}\}_{k,l=1}^p = \mathbb{E}[(X - E(X))(X - E(X))^T]$$

$$b_{kl} = \mathbb{E}[(X_k - E(X_k)) \cdot (X_l - E(X_l))], b_{kk} = b_{kk}$$

the correlation matrix is the matrix of correlations obtained through normalization.

$$p_{kl} = \frac{b_{kl}}{\sqrt{b_{kk} b_{ll}}}, p_{kk} = 1$$

$$T = \text{cor}(X) = \{p_{kl}\}_{k,l=1}^p = D^{-\frac{1}{2}} \Sigma D^{-\frac{1}{2}}, \text{ with } \text{diag}(D) = (b_{11}, \dots, b_{pp})$$

recall by Cauchy-Schwarz inequality we can get  $p_{kl} \in [-1, 1] \hookrightarrow |(X_k, Y_l)| \leq \|(X_k)\| \cdot \|(Y_l)\|$

sample quantity.

$$S = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T \in \mathbb{R}^{p \times p}$$

$$s_{kl} = \frac{s_{kl}}{\sqrt{s_{kk} s_{ll}}}$$

$\hookrightarrow$  why not  $n$ , because it's

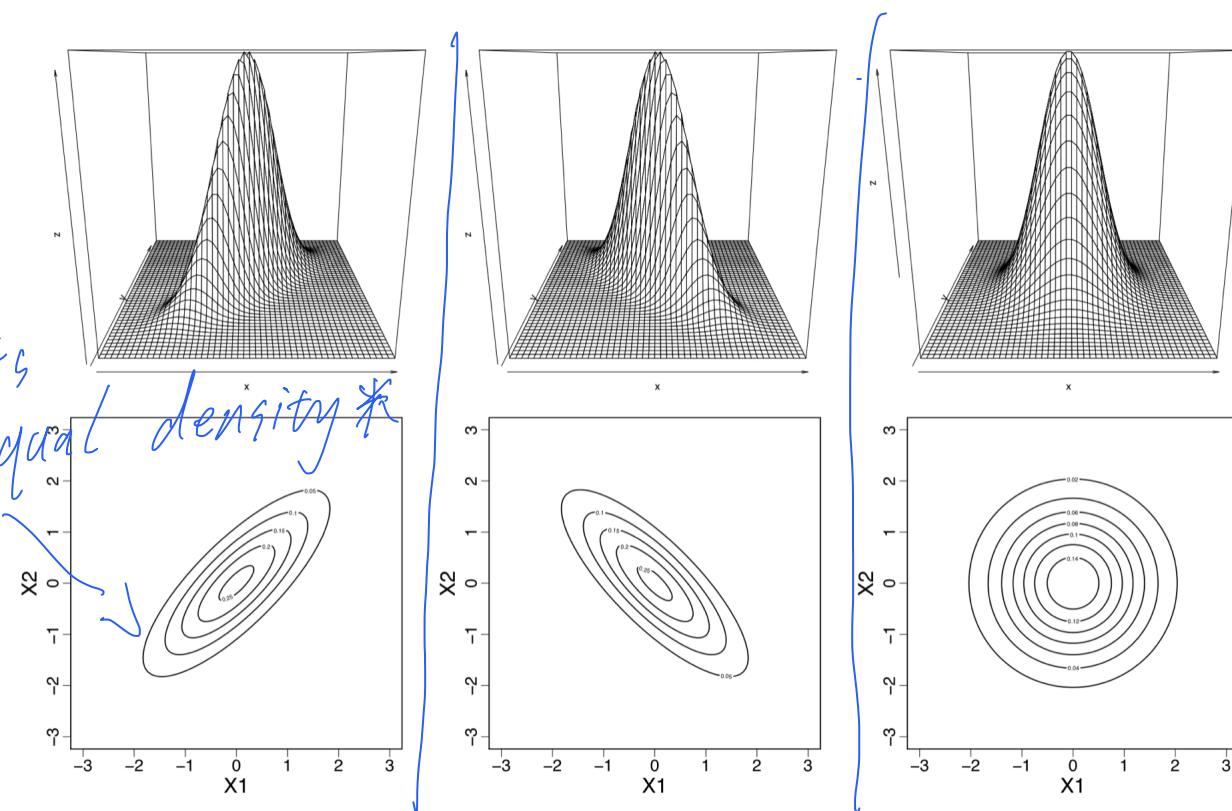
an unbiased estimator of  $\Sigma$   
Hessel's correction

$$Q = D^{-\frac{1}{2}} S D^{-\frac{1}{2}}$$

$$\text{with } D = \text{diag}(s_1^2, \dots, s_p^2)$$

PS:①

bivariate normal distribution



這是 different density function

上面的圖可以理解為:  $Z = f_{x,y}(X, Y)$   
Z 是 probability

下圖為其在 x-y 平面上的投影

其  $\Sigma = \begin{pmatrix} a & b \\ b & c \end{pmatrix}$  變動的是下圖的位置

$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix}$   $\sigma_{11}, \sigma_{12}$  變動在 x or y 的跨度  
 $\sigma_{21}, \sigma_{12}$  是相关系数

\*如何理解:  $Z = 0 = f_{(X,Y)}(X, Y) = \frac{1}{(2\pi)^{d/2} \cdot |\Sigma|^{1/2}} e^{-\frac{1}{2} (X-\mu)^\top \Sigma^{-1} (X-\mu)}$

其  $d = \text{len}(X)$

$X = (X_1, X_2, X_3, \dots)$  covariance = density value

$$\text{OPS: } F(y) = P(X_1 \leq y_1, \dots, X_p \leq y_p)$$

$$\hat{F} = E[\text{Indicator function}]$$

$$= E[I(X_1 \leq y_1, \dots, X_p \leq y_p)]$$

core:  
 prob is expectation of indicator function

$$I\{\text{FA}\} = \begin{cases} 1 & \text{if it's true} \\ 0 & \text{if it's false} \end{cases}$$

$\therefore$  if we want estimate the  $F(y)$ ,  
 we need to estimate the expectation of  
 the indicator function.

So how to estimate the expectation of  
 the indicator function?  $\Rightarrow$  use monte carlo

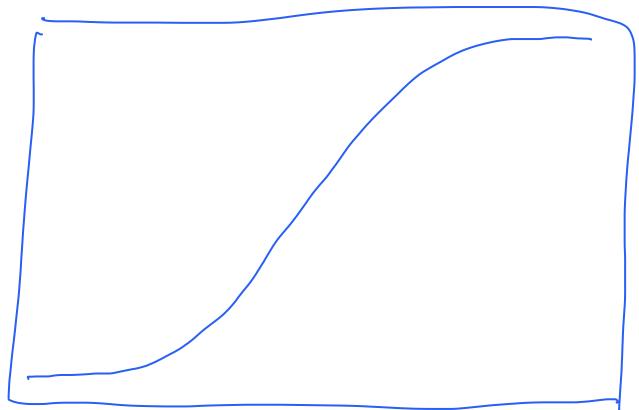
$$\hat{F}(y) = \frac{1}{n} \cdot \sum_{i=1}^n I\{X_{i1} \leq y_1, \dots, X_{ip} \leq y_p\} \quad (3)$$

expectation: [prob x value]  $\nearrow$  估计 Indicator function.

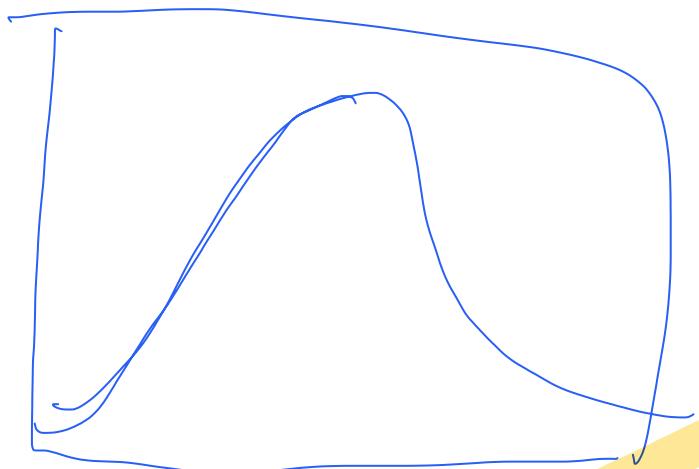
PS: it's non-parametric estimator

good estimator:  
 $\hat{F}_n(y) \xrightarrow{\text{more data}} F_n(y)$   
 it means it  
 converge to the  
 true CDF

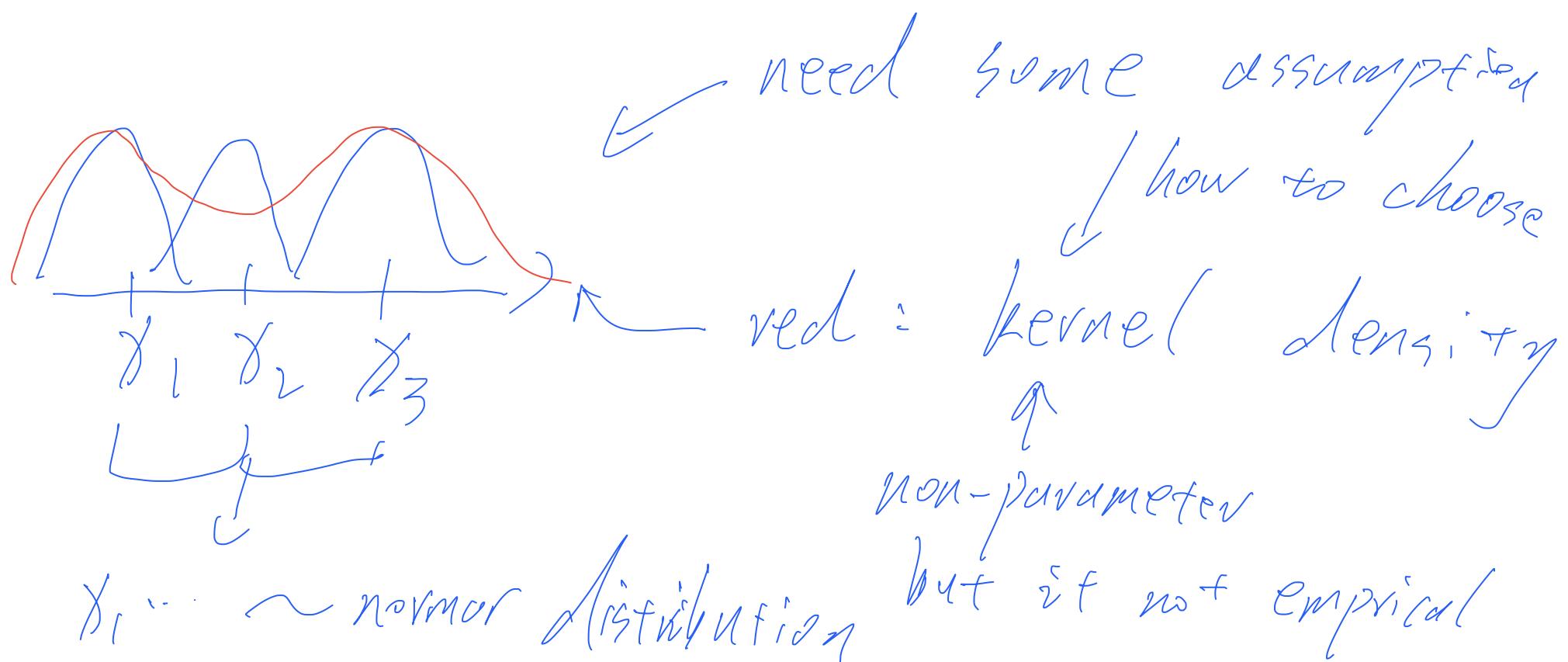
CDF of normal distribution



density function of normal distribution



這裡只是估計了 CDF  
 沒有估計 density  
 如何应对? way1: 連續下, 微分.  
 Way 2: kernel function?  
 estimate



a covariance matrix  $\Sigma = \text{cov}(X)$  is symmetric and positive-semidefinite, that is,

$$x^T \Sigma x \geq 0, \quad \forall x \in \mathbb{R}^p$$

if the above holds strictly with  $> 0$ , then  $\Sigma$  is positive definite.

If  $\Sigma$  is positive definite and  $\Sigma^{-1}$  with  $\{\text{invertible}\}$  we can have a decomposition  $\Sigma^{1/2}$  satisfying

$$\Sigma^{1/2} (\Sigma^{1/2})^T = \Sigma$$

for a matrix  $A \in \mathbb{R}^{n \times p}$ , a constant vector  $b \in \mathbb{R}^m$ , and a random vector  $X \in \mathbb{R}^p$

$$\text{cov}(b + Ax) = E[(a \cdot (X - E(X)))(X - E(X))^T A^T] = A \Sigma A^T \in \mathbb{R}^{n \times n}$$

The covariance notion also generalizes to two vectors. Let  $X \in \mathbb{R}^p$  and  $Y \in \mathbb{R}^q$  be random vectors, then their covariance matrix is  $p \times q$ .

$$\text{cov}(X; Y) = E[(X - E(X))(Y - E(Y))^T] = \{\text{cov}(X_i; Y_j)\}_{i,j=1}^{p,q}$$

Note that  $\Sigma_{xy} = \Sigma_{yx}^T$

population quantity

Let  $X \in \mathbb{R}^p$  be a random vector with  $\bar{\Sigma} = \text{cov}(X)$  and  $\Delta = \text{diag}(\sigma_1^2, \dots, \sigma_p^2)$  and correlation matrix:

$$\text{cor}(X) = \Gamma = \Delta^{-1/2} \bar{\Sigma} \Delta^{-1/2}$$

then the random vector  $Y$  satisfies:

$$Y = \Delta^{1/2}(X - E(X)) \Rightarrow E(Y) = 0 \quad (i)$$

for fair standardization, consider:

$$Y = \bar{\Sigma}^{1/2}(X - E(X)) \quad (ii)$$

then  $E(Y) = 0$  and  $\text{cov}(Y) = I_p = \text{diag}(1, \dots, 1)$

$$\begin{aligned} \text{for (i)} \quad \text{cov}(Y) &= \text{cov}(\Delta^{1/2}X - \Delta^{1/2}E(X)) = \text{cov}(\Delta^{1/2}X) \\ &= E[\Delta^{1/2}(X - E(X))(X - E(X))^T (\Delta^{1/2})^T] \\ &= \Delta^{1/2} \cdot \bar{\Sigma} \cdot (\Delta^{1/2})^T = \Delta^{1/2} \cdot \bar{\Sigma} \cdot \Delta^{1/2} = \Gamma \end{aligned}$$

$$\begin{aligned} \Delta &= \text{diag}(\sigma_1^2, \dots, \sigma_p^2) \\ &= \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_p^2 \end{bmatrix} \Rightarrow (\Delta^{1/2}) = (\Delta^{-1})^T \end{aligned}$$

$$\begin{aligned} \text{for (ii)} : \quad \text{cov}(Y) &= \text{cov}(\bar{\Sigma}^{1/2}X - \bar{\Sigma}^{1/2}E(X)) = \text{cov}(\bar{\Sigma}^{1/2}X) \\ &= E[\bar{\Sigma}^{1/2} \cdot (X - E(X))(X - E(X))^T (\bar{\Sigma}^{1/2})^T] \\ &= \bar{\Sigma}^{1/2} \cdot \bar{\Sigma} \cdot (\bar{\Sigma}^{1/2})^T = \bar{\Sigma}^{-1} \cdot \bar{\Sigma} = I \end{aligned}$$

sample quantity

the analogous sample standardization for the data  $X_1, \dots, X_n$  are

$$y_i = D^{1/2}(x_i - \bar{x}) \quad y_c = S^{1/2}(x_i - \bar{x})$$

$S$  is the sample covariance

$D$  is the diag  $(S_{11}^{-2}, \dots, S_{pp}^{-2})$

Independence:

By definite: the  $p$  components  $X_1, \dots, X_p$  of  $X \sim F$  are mutually independent

$$F(y) = \prod_{j=1}^p F_j(y_j), \quad \forall y \in \mathbb{R}^p$$

$$f(y) = \prod_{j=1}^p f_j(y_j), \quad \forall y \in \mathbb{R}^p$$

if two random variables are independent, then the  $E(xy)$

$$E(xy) = E(x) \cdot E(y)$$

consequently, if  $X = (X_1, \dots, X_p)^T$  mutually independent,

$$\begin{aligned} \rho_{pk} &= \text{cov}(X_k; X_l) = E[(X_k - E(X_k))(X_l - E(X_l))] \\ &= 0 \end{aligned}$$

$$(\text{cov}(X) = \text{diag}(\sigma_1^2, \dots, \sigma_p^2))$$

Let  $X = (X_1, X_2)^T$  be a bivariate normal random variables with mean = 0 and covariance  $\Sigma = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$

the density function:

$$f(y_1, y_2) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)}(y_1^2 - 2\rho y_1 y_2 + y_2^2)}$$

原形式: (补充)

$$f(x, y) = \frac{1}{2\pi\sqrt{\sigma_x\sigma_y}} e^{-\frac{1}{2(\sigma_x^2+\sigma_y^2)}\left(\frac{(x-\mu_x)^2}{\sigma_x^2} - 2\rho \cdot \frac{x-\mu_x}{\sigma_x} \cdot \frac{y-\mu_y}{\sigma_y} + \frac{(y-\mu_y)^2}{\sigma_y^2}\right)}$$

$$\mu = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \Sigma = \begin{pmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{pmatrix}$$

$$f_1(y_1) = \int_{-\infty}^{+\infty} f(x, y_1) dx = \frac{1}{\sqrt{2\pi}} e^{-\frac{y_1^2}{2}} \text{ same to } y_2$$

if  $\rho = 0 \Rightarrow X_1, X_2$  independent?

$$\therefore f_1(y_1) f_2(y_2) = \frac{1}{2\pi} e^{-\frac{1}{2}(y_1^2 + y_2^2)} = 1.$$

$$\text{if } T = f(y_1, y_2) \Rightarrow P = 0.$$

summary:

$$\begin{cases} X_1, X_2 \sim N \\ \rho = 0 \end{cases} \Rightarrow X_1, X_2 \text{ independent}$$

$$X_1, X_2 \text{ indep} \Rightarrow \rho = 0$$

conditional distribution:

Let  $X \in \mathbb{R}^p$ ,  $Y \in \mathbb{R}^q$  be random vectors, and let  $(X, Y)^T$  be the joint vector with density  $f(x, y)$ . Then the conditional distribution of  $X$  given the  $Y=y$

$$f_{X|Y=y} = \frac{f_{X,Y}(x, y)}{f_Y(y)}$$

if  $X, Y$  independent,  $f_X(x) = f_Y(y)$

for the bivariate normal distribution our numerical example with  $X = X_1$ ,  $Y = X_2$ .

$$f_{X|Y=y} = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{(x-y)^2}{2(1-\rho^2)}}$$

therefore  $(X|Y=y) \sim N(\rho y, 1-\rho^2)$

we also can define the  $E(X|Y=y)$ :

$$E(X|Y=y) = \int_X x \cdot f_{X|Y=y}(x) dx$$

more generally:  $E(g_{(S)}|Y=y) = \int_S g_{(S)} \cdot f_{S|Y=y}(s) ds$

the conditional expectation can be seen as a random variable:  $E(X|Y)$  with the important property

$$E(X) = E[E(X|Y)]$$

similarly:

$$\text{cov}(X|Y=y) = E[(X - E(X|Y))(X - E(X|Y))^T]$$

if  $X$  and  $Y$  independent.

$$E(X|Y) = E(X)$$

$$\text{cov}(X|Y) = \text{cov}(X)$$

$$E_{X|Y}(X|Y) = E_X(X), \quad \text{both is}$$

$$f_{X|Y}(X|Y) = f_X(x) \quad \Rightarrow \text{sufficient condition}$$

## Module 02:

multivariate normal and other distributions

the univariate uniform distribution: 均匀分布(别看那个 distribution 就是正态分布)

Let  $S = (a, b)$  be interval, then the uniform distribution  $U \sim U(a, b)$  is a uniformly distributed random variable if  $f(u) = \frac{1}{b-a}$ , for  $u \in (a, b)$

prob density function

a important example:  $S = (0, 1) \rightarrow$  standard uniform distribution

recall that any random variable  $X$  with continuous cumulative distribution function  $F$  satisfies:

$F(x) \sim U(0, 1)$ , since:

$$P(F(x) \leq u) = P(X \leq F^{-1}(u)) = F(F^{-1}(u)) = u, \forall u \in (0, 1)$$

we can thus generate the samples of  $X$  using a standard uniform random variable  $U \sim U(0, 1)$ , by setting  $x = F^{-1}(u)$

$$P(X \leq x) = P(F^{-1}(u) \leq x) = F(x)$$

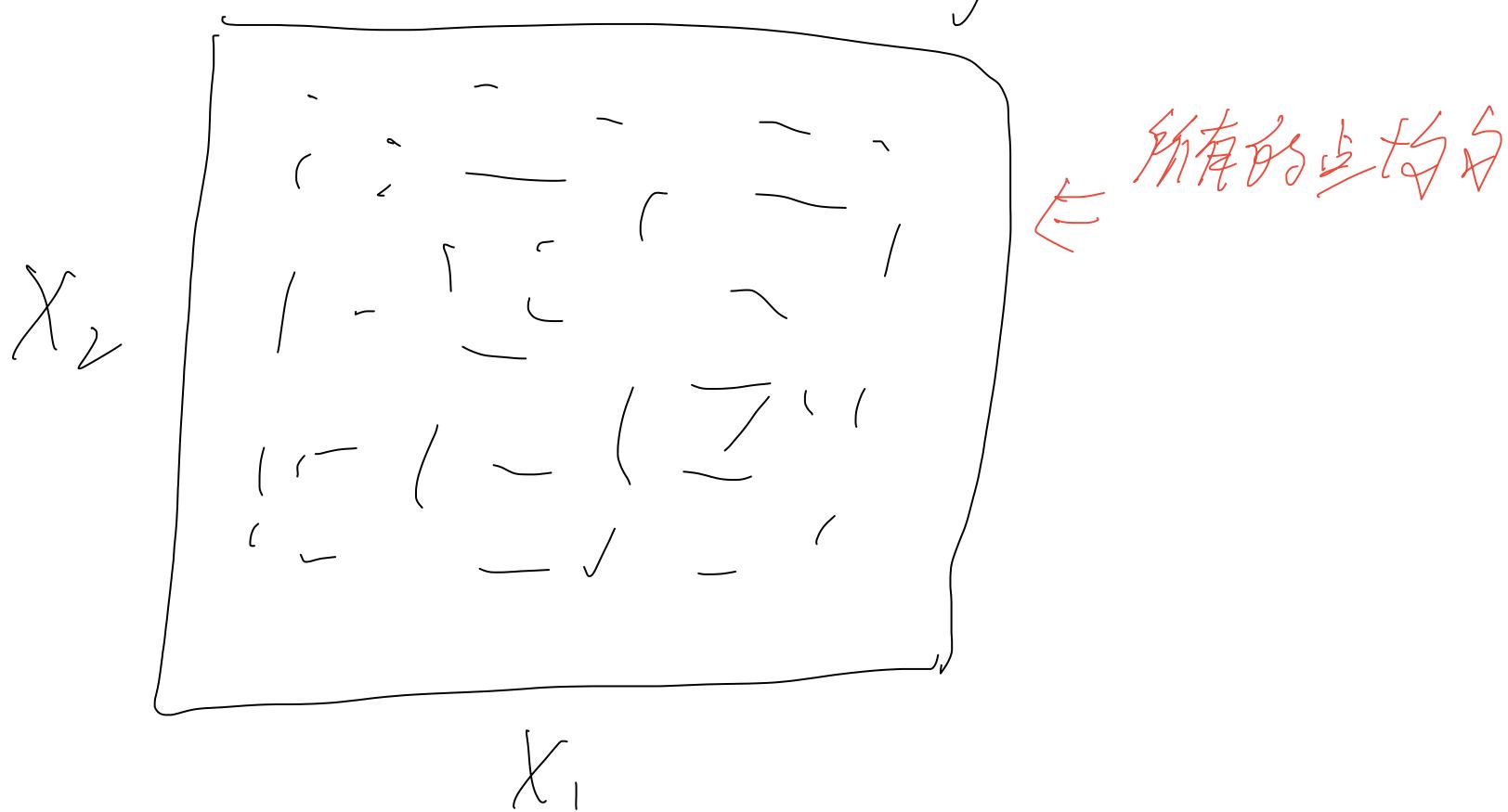
In general, we can construct a multivariate distribution for  $X = (X_1 \cdots X_p)^T$  by letting  $X_1 \cdots X_p$  be independent univariate random variables. The density is:

$$f(x_1, \dots, x_p) = f_{x_1}(x_1) \cdot f_{x_2}(x_2) \cdots f_{x_p}(x_p)$$

Let  $X_j, j=1 \cdots p$  be random variables uniformly distributed on interval ( $S_j = (a_j, b_j)$ ,  $j=1 \cdots p$ ). Then  $X = (X_1 \cdots X_p)^T$  has a multivariate uniform distribution with density

$$f(x_1, \dots, x_p) = \frac{1}{\prod_{j=1}^p (b_j - a_j)}, \quad x_j \in S_j.$$

by definition the marginals  $X_j$  are mutually independent and all uniformly distributed



the multivariate normal distribution

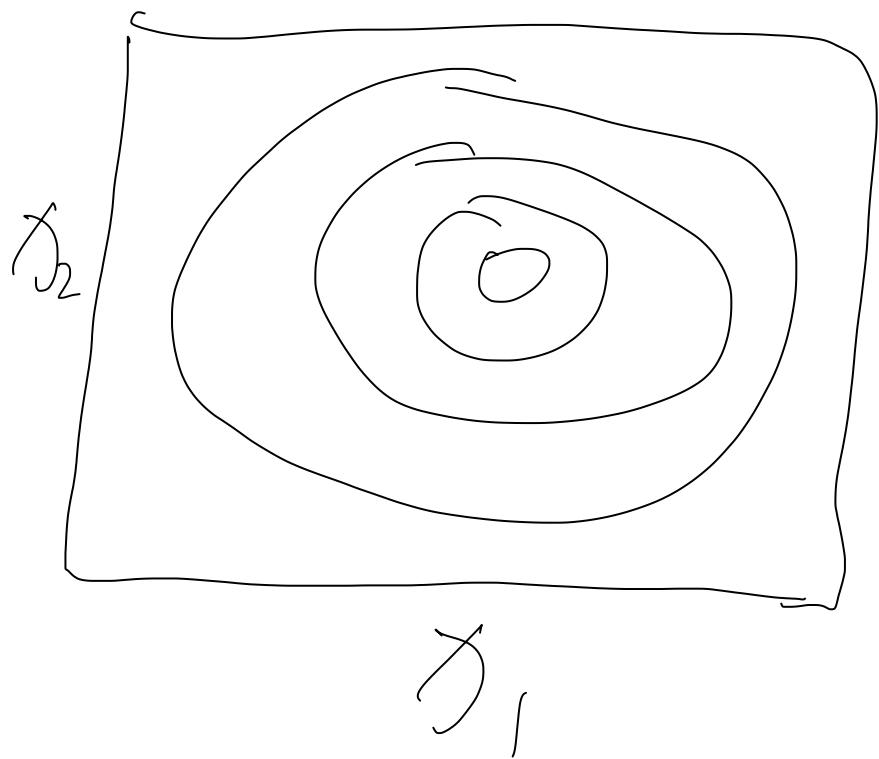
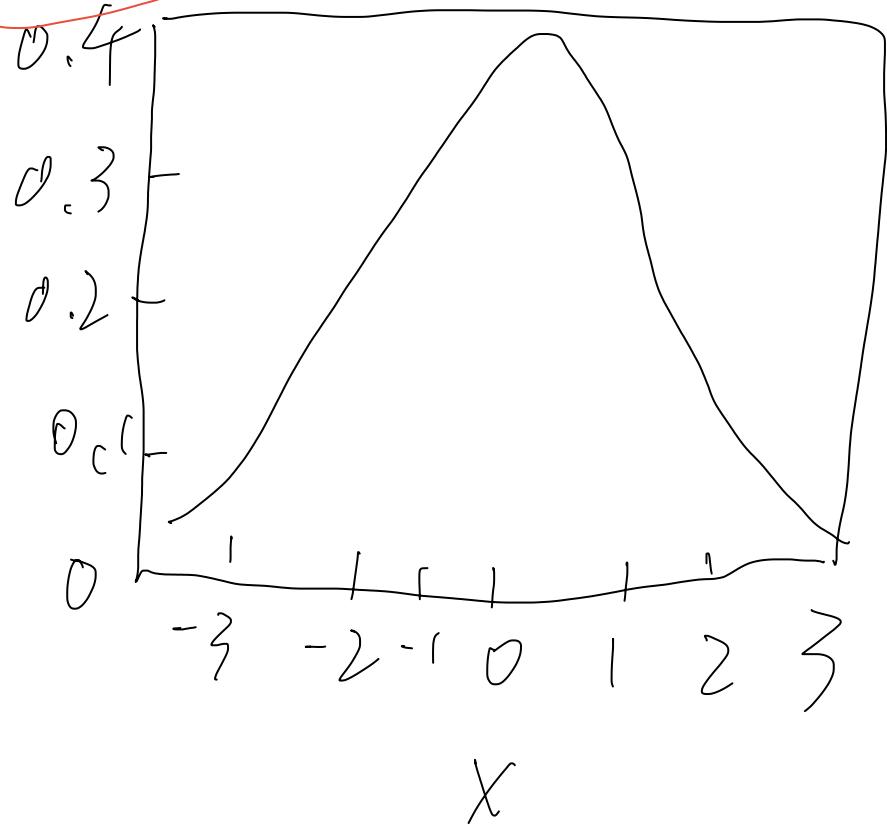
recall the  $X \sim N(\mu, \sigma^2)$  has a univariate normal distribution if its density is

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} \quad \mu = E(X) \\ \sigma^2 = \text{Var}(X)$$

the standard multivariate normal distribution follows by the same "independent construction" as before, namely let  $X_1 \cdots X_p$  be standard normal  $N(0, 1)$ , then  $X = (X_1 \cdots X_p)^T$  has density

$$f(x) = f_{X_1}(x_1) \cdots f_{X_p}(x_p) = \frac{1}{\sqrt{(2\pi)^p}} \exp\left\{-\frac{1}{2} x^T x\right\}, \quad x \in \mathbb{R}^p$$

dnorm(x)



Density of transformed random vectors.

Let  $X$  be a  $P$  random vector and  $\tilde{Y} = H(X)$   
for  $H: \mathbb{R}^P \rightarrow \mathbb{R}^P$  a bijection (i.e.: invertible)  
双射的  $\rightarrow$  (可逆的)  
transformation

if  $f_X$  is the density of  $X$ , then  $Y$  has the

density  $f_Y(y) = f_X(H^{-1}(y)) / \det(J_{H^{-1}}(y))$  要搞懂

where for a mapping  $T = (T_1 \cdots T_P): \mathbb{R}^P \rightarrow \mathbb{R}^P$

$J_T(y)$  is the Jacobian matrix at point  
 $y \in \mathbb{R}^P$

$$J_T(y) = \begin{pmatrix} \frac{\partial T_1}{\partial y_1} & \cdots & \frac{\partial T_1}{\partial y_P} \\ \vdots & \ddots & \vdots \\ \frac{\partial T_P}{\partial y_1} & \cdots & \frac{\partial T_P}{\partial y_P} \end{pmatrix} \in \mathbb{R}^{P \times P}$$

the matrix  $J_T$  is the best linear  
approximation of  $T$  at  $y$  how to prove

example: matrix  $A \in \mathbb{R}^{P \times P}$ ,  $b \in \mathbb{R}^P$

the mapping  $H(x) = Ax + b$  is affine

$$\Rightarrow J_H(\gamma) = A$$

(?)

scaling and rotation:

further construction principles for multivariate distributions are

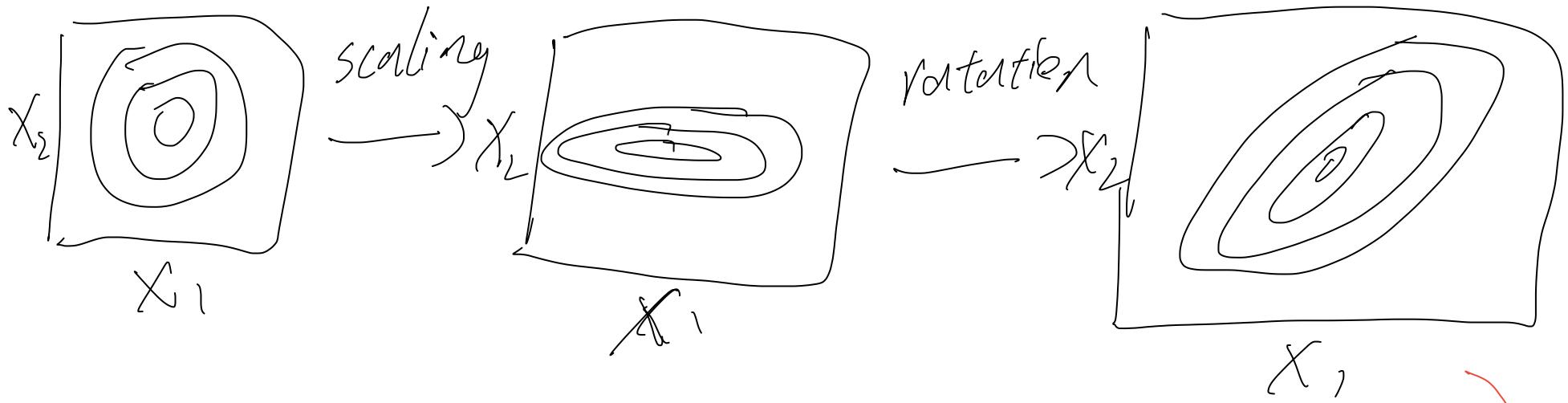
scaling: Let  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_P)$  be a scaling matrix and define a new random vector  $V \Lambda^{\frac{1}{2}} x$

$$x \in \mathbb{R}^P \quad f_{V\Lambda^{\frac{1}{2}}x}(x) = \frac{1}{\sqrt{(2\pi)^P \det(\Lambda)}} \exp \left\{ -\frac{1}{2} x^\top \Lambda^{-1} x \right\}$$

PS  $y = H(x) = \Lambda^{\frac{1}{2}} x$  why?  
 $f_y(y) = f_x(H^{-1}(y)) |\det(J_{H^{-1}}(y))|$

rotation: let  $V$  be a rotation matrix and define a new random vector  $V \Lambda^{\frac{1}{2}} x$ ,

$$f_{V\Lambda^{\frac{1}{2}}x}(x) = \frac{1}{\sqrt{(2\pi)^P \det(V\Lambda V^\top)}} \cdot \exp \left\{ -\frac{1}{2} x^\top (V\Lambda V^\top)^{-1} x \right\}$$



location shifting: for some  $u \in \mathbb{R}^P$ , let  $H(x) = x + u$

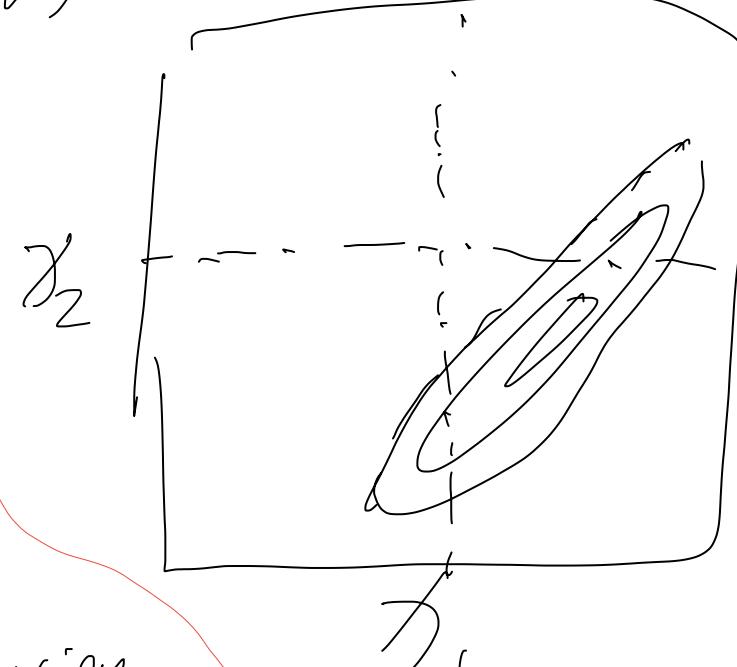
$$f_{VN^2x}(x) = \frac{1}{\sqrt{(2\pi)^P \det(VN^T)}} \exp \left\{ -\frac{1}{2}(x-u)^T (VN^T)^{-1} (x-u) \right\}$$

let  $\Sigma$  be a positive-definite  $P \times P$ -matrix and  $u \in \mathbb{R}^P$ . Then

$x = (x_1 \cdots x_P)^T$  has a multivariate normal distribution

$N_p(u, \Sigma)$  if its density satisfies

$$f(x) = \frac{1}{\sqrt{(2\pi)^P \det(\Sigma)}} \exp \left\{ -\frac{1}{2}(x-u)^T \Sigma^{-1} (x-u) \right\}, x \in \mathbb{R}^P$$



every positive-definite matrix  $\bar{\Sigma}$  has eigendecomposition

$\bar{\Sigma} = U\Lambda U^T$  where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ ,  $\lambda_i$  is the eigenvalues of  $\bar{\Sigma}$ , and  $U$  is a rotation matrix. Therefore, we get this representation

$$X \sim \mu + \sqrt{\Lambda}^{\frac{1}{2}} N_p(0, I_p)$$

by the rules for expectation and covariance this gives:  $E[X] = \mu$ ,  $\text{cov}(X) = \bar{\Sigma}$

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

- ① the marginal  $(X_j)$  has a univariate normal distribution  $N(\mu_j, \sigma_j^2)$
- ② Let  $A$  be a  $m \times p$  matrix,  $b \in R^m$ , then  $b + Ax$  is multivariate normal  $N_m(b + Au, A\bar{\Sigma}A^T)$

in particular  $\bar{\Sigma}^{\frac{1}{2}}(X - \mu) \sim N_p(0, I_p)$

where  $\bar{\Sigma}^{\frac{1}{2}}$  is the choleski decomposition of  $\bar{\Sigma}^{-1}$  satisfying  $\bar{\Sigma}^{\frac{1}{2}} \cdot \bar{\Sigma}^{\frac{1}{2}} = \bar{\Sigma}^{-1}$

3)  $X_j$  and  $X_k$  are independent if their covariance  $\sigma_{jk} = 0$ . That is the components of  $X$  are mutually independent if and only if

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

this equivalence between "independent" and "uncorrelated" is specific to the normal distribution.

The normal distribution is completely characterized by its mean vector and covariance matrix.



