

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 P₁: this class \rightarrow unsupervised learning

include a lot of \nearrow
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

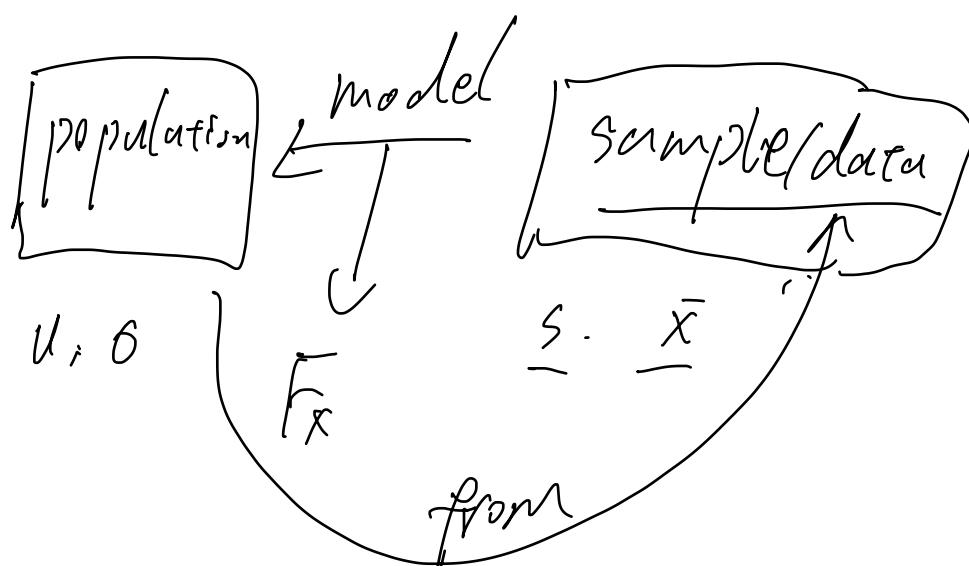
Model perspective:

- 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 $X_1 \dots X_n \in \mathbb{R}^P$ are thought of finite number n of observations from X right?
in feature P 為維度 $\rightarrow n \times p$
- this distribution is unknown

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: $\Sigma = \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix}$

\hookrightarrow model 角度

一个 X 的 取值. $P.$

p

\uparrow

$[2]$

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 樣 本
对 公 署 到 的
所 得.

Some notation: *data + variable*
 for a vector $x \in \mathbb{R}^P$ we write
 $x = \begin{pmatrix} x_1 \\ \vdots \\ x_P \end{pmatrix} \quad x^T = (x_1 \dots x_P)$

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

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

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

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

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ x_{21} & & & \\ \vdots & & & \\ x_{n1} & \dots & x_{nP} \end{pmatrix} \in \mathbb{R}^{n \times P} = \begin{bmatrix} x_1 \\ \vdots \\ x_P \end{bmatrix}^T$$

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:

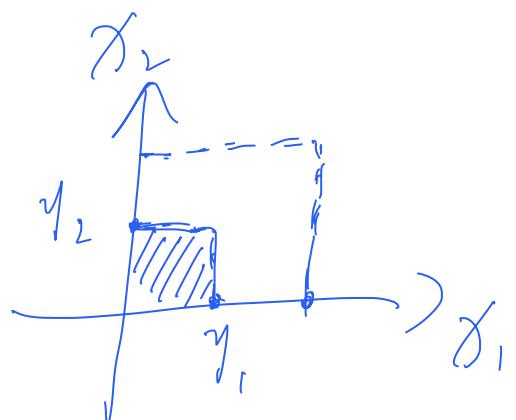
$$f(y) = f(y_1 \dots y_p) = \frac{\partial^p F(y)}{\partial y_1 \dots \partial y_p}$$

or equivalently:

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

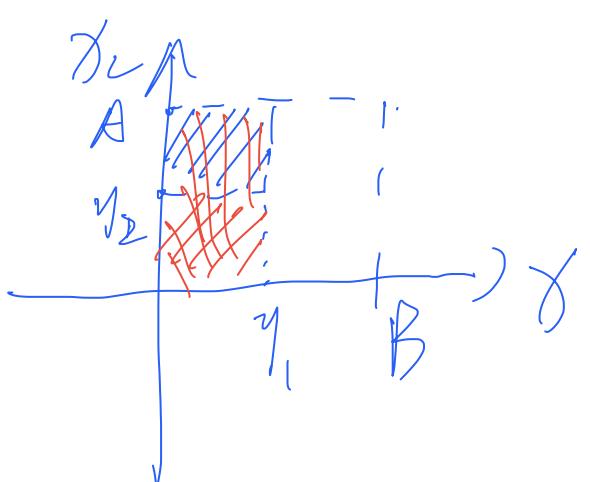
本质上: F [上页画的] 为 distribution, 可以替换为任意 distribution:

举例: uniform distribution



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

"/" is the P



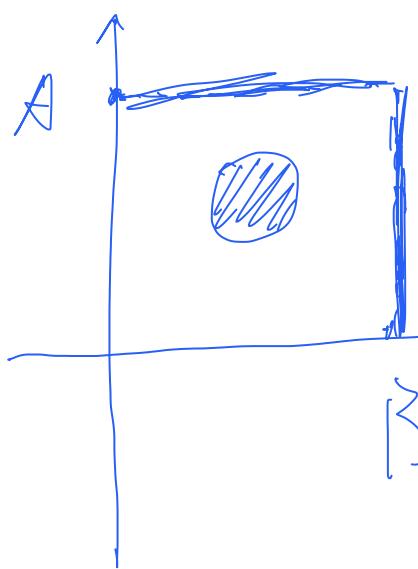
从图可推得

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

A is the range $[0, A]$

↓ 得 F

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



"/" — "/"

\Rightarrow if we want to know
the $P(Y)$ which is not limit
rectangles

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

the joint density exists \Rightarrow the marginal densities
 $f_{ij}(y_j) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(u_1, \dots, u_j, \dots, u_p) du_1 \cdots du_{j-1} du_{j+1} \cdots du_p$
 the density function integrates to 1

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

Sample quantity: we don't know distribution
 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
 $(3) \hat{F}_{nY}(y) = \frac{1}{n} \sum_{i=1}^n I\{x_{ip} \leq y_i, \dots, x_{ip} \leq y_p\} \quad (2) \text{ PS}$

is a non-parametric estimator of F empirical distribution

there is no "empirical" sample-equivalent for the density function without further assumptions. why? (4)

example:

Let $X = (X_1, X_2)^T$ be a bivariate random variable with density function: $f_{XY}(y_1, y_2) = r(y_1^\alpha + y_2^\beta), 0 < y_1, y_2 < 1, \alpha, \beta > 0$

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

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

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

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

$$F(X_1 \leq \frac{1}{2}, X_2 \geq \frac{1}{2}) = \frac{1}{4} \quad (\text{Probabilities的简单表达式})$$

for F 类似地:

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

population quantity:

Let $X = (X_1 \dots X_p)^T$ be a random vector

we can define the expectation of g_{CX_1} , where $g: \mathbb{R}^p \rightarrow \mathbb{R}^m$ is some function
 \Rightarrow g is a density vector

$$E(g_{CX_1}) = \int_{\mathbb{R}^p} g(u) dF(u) = \int_{\mathbb{R}^p} g(u) f(u) du \in \mathbb{R}^m$$

for n samples $x_1 \dots x_n \in \mathbb{R}^p$ of X , we define the sample version of $E(g_{CX_1})$ as recall PS(3)

$$\hat{E}(g_{CX_1}) = \int_{\mathbb{R}^p} g(u) d\hat{F}(u) = \frac{1}{n} \sum_{i=1}^n g(x_{i1} \dots x_{ip}) \Rightarrow \text{empirical expectation}$$

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 $b \in \mathbb{R}^m$

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

$E(Z) = A \cdot E(X) + B \cdot E(Y) + b \Rightarrow$ 线性性质
 the trace is an operator on square matrices that returns the sum of the diagonal elements. 又稱之為迹和.

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 $A \in \mathbb{R}^{p \times p}$

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

how can I prove it? \Rightarrow by linearity $\text{E}[\text{Tr}(c)] = \text{Tr}[\text{E}(c)]$
 population quantity $\Rightarrow X^T A X$ 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

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

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

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

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

$$\text{gamma } (T) = \{P_{kl}\}_{k,l=1}^p = \Delta^{1/2} \Sigma \Delta^{1/2}, \text{ with } \text{diag}(b_1^2, \dots, b_p^2)$$

recall by Cauchy-Schwarz inequality we can get $P_{kl} \in [-1, 1] \Rightarrow |X_k y_l| \leq \langle X_k, X_l \rangle \cdot \langle y_l, y_l \rangle$
 $|\langle X_k, y_l \rangle| \leq \|X_k\| \cdot \|y_l\|$

sample quantity.

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

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

\Rightarrow why not n , because it's

an unbiased estimator of Σ

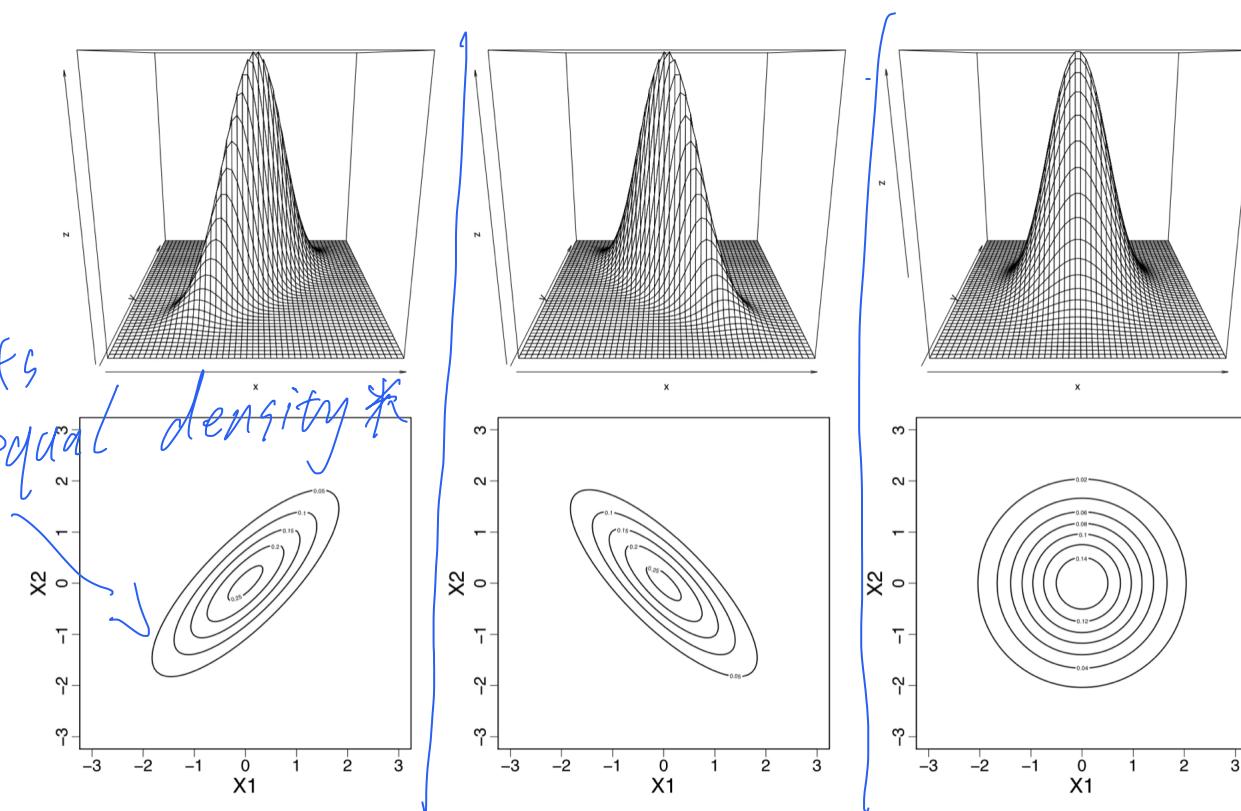
Hessel's correction

$$Q = D^{1/2} S D^{1/2}$$

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

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

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

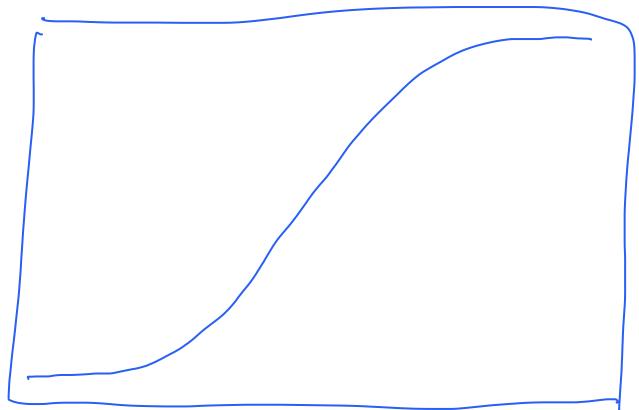
$$\hat{F}(y) = \frac{1}{n} \cdot \sum_{i=1}^n I\{X_{i1} \leq y_1, \dots, X_{ip} \leq y_p\} \quad (3)$$

expectation: [prob x value] \mathcal{M} 估计 Indicator function.

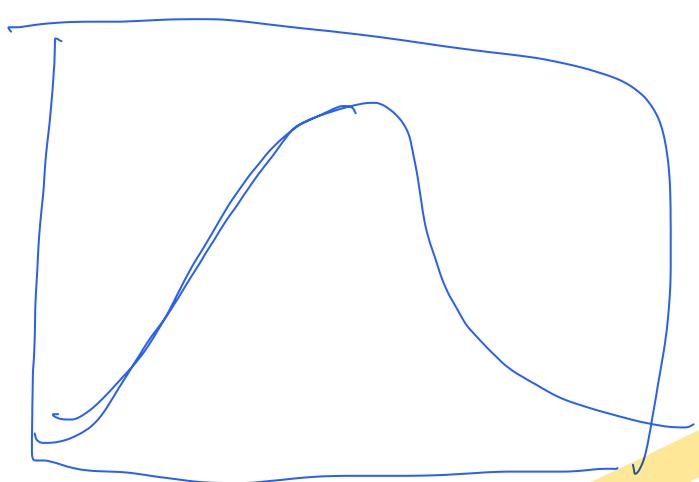
PS: it's non-parametric estimator

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

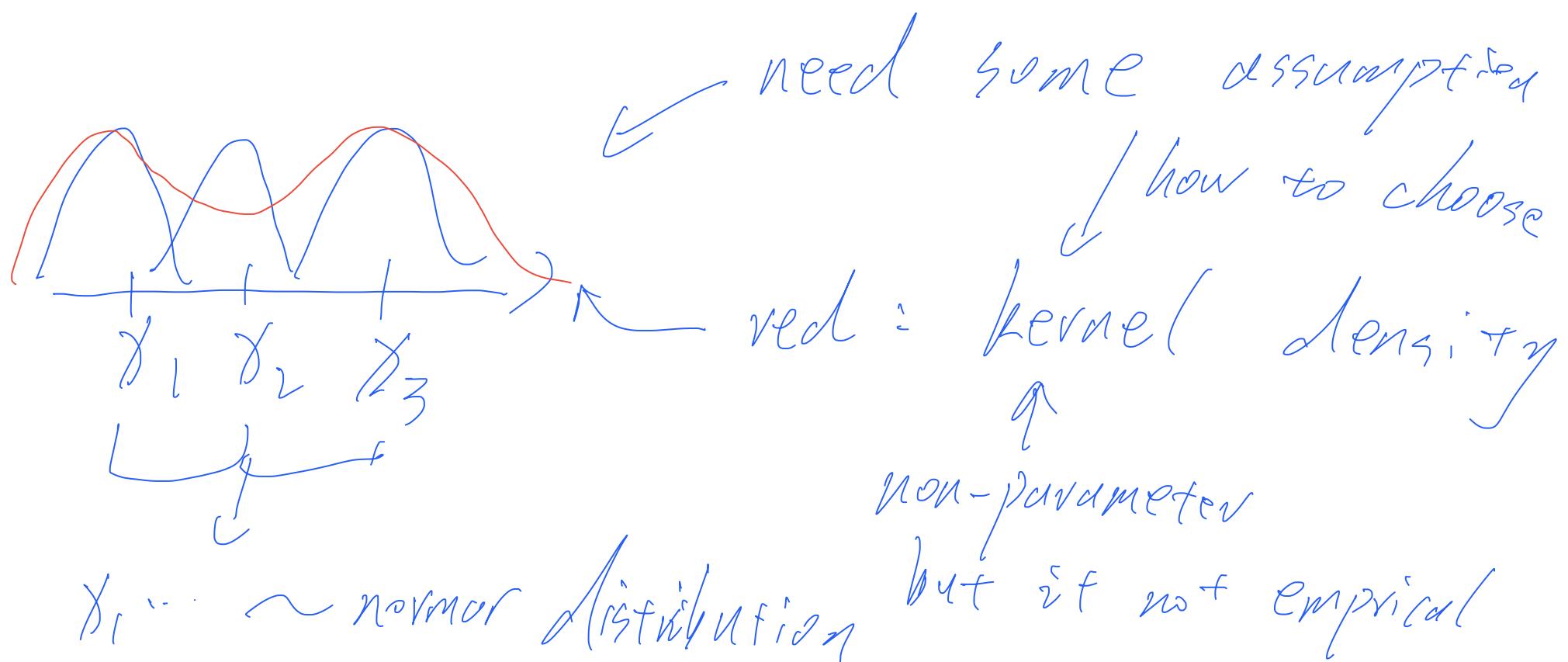
CDF of normal distribution



density function of normal distribution



這裡只是估計了 CDF
 沒有估計 density
 如何应对? way 1: 連續下, 微分.
 Way 2: kernel function?
 estimate



下文(5): 本復 $\xrightarrow{m \rightarrow p} \text{cov}(b + AX) \rightarrow A\Sigma A^T$
是转化了 dimension

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

$$X^T \Sigma X \geq 0, \quad \forall \text{ for all } X \in \mathbb{R}^p$$

if the above holds strictly with > 0 , then Σ is positive definite.

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

$$\Sigma^{1/2} (\Sigma^{1/2})^T = \Sigma \quad \text{Choleske-decomposition}$$

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

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

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 $\Sigma_{xy} = \Sigma_{yx}$

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:

$$\frac{X - \mu}{\sigma} \stackrel{\text{def}}{=} Y = \Delta^{-1/2}(X - E(X)) \Rightarrow E(Y) = 0 \quad \text{cov}(Y) = \Gamma \quad (i)$$

For fair standardization, consider:

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

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

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

$$\Delta = \text{diag}(\sigma_1^2, \dots, \sigma_p^2) \Rightarrow (\Delta^{1/2}) = (\Delta^{1/2})^T$$

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

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

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

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

If two random variables are independent, then the $E(xy) = E(x) \cdot E(y)$

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

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:

$$f(\gamma_1, \gamma_2) = \frac{1}{2\pi\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)}(\gamma_1^2 - 2\rho\gamma_1\gamma_2 + \gamma_2^2)}$$

原形式: (补充)

$$f(x, y) = \frac{1}{2\pi\sqrt{6x\delta_y}} e^{-\frac{1}{2(1-\rho^2)}\left(\frac{(x-\mu_x)}{\delta_x}\right)^2 - 2\rho \cdot \frac{x-\mu_x}{\delta_x} \cdot \frac{y-\mu_y}{\delta_y} + \left(\frac{y-\mu_y}{\delta_y}\right)^2}$$

$$\mu = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \Sigma = \begin{pmatrix} \delta_x^2 & \delta_x\delta_y & \rho \\ \delta_x\delta_y & \delta_y^2 & \rho \end{pmatrix}$$

$$f_1(\gamma_1) = \int_{-\infty}^{+\infty} f(\gamma_1, \gamma_2) d\gamma_2 = \frac{1}{\sqrt{2\pi}} e^{-\frac{\gamma_1^2}{2}} \text{ same to } \gamma_2$$

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

$$\therefore f_1(\gamma_1) f_2(\gamma_2) = \frac{1}{2\pi} e^{-\frac{1}{2}(\gamma_1^2 + \gamma_2^2)} = 1.$$

$$\text{if } T = f(\gamma_1, \gamma_2) \Rightarrow P = 0.$$

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 $Y=y$

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

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

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

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

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

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

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

more generally: $E(g_{(S)}|Y=y) = \int_S g_{(S)} \cdot f_{X|Y=y}(x|y) dx$

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

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

similarly:

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

if X and Y independent.

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

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

$$F_{X|Y}(x|y) = F_X(x), \quad \text{both is}$$

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

PDF

對應 (但只有在 normal distribution)

$$f_{X,Y}(x_1, y_1)$$

$$= f_X(x) \cdot f_Y(y)$$

實在理解解就

$$E(X) = \int_{-\infty}^{+\infty} x f(x) dx$$

$$F(x) = \int_{-\infty}^x F(t) dt \Rightarrow \text{CDF}$$

Module 02:

multivariate normal and other distributions

the univariate uniform distribution: 均匀分布(沿着剖面) 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:

$$\underbrace{P(F(X) \leq u)}_{\text{we can thus generate the samples of } X \text{ using a standard uniform random variable } U \sim U(0, 1)} = P(X \leq F^{-1}(u)) = F(F^{-1}(u)) = u, \forall u \in (0, 1)$$

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

$$\underbrace{P(X \leq x)}_{x \sim F} = \underbrace{P(F^{-1}(U) \leq x)}_{P(U \leq F(x))} = F(x)$$

PS: ⑥ $P(X \leq x) = F(x)$ 解釋了 第二个文字內容
 F 是 CDF

因為是 uniform distribution $\therefore F(x) = x$ 解釋三个文字
 F^{-1} 是反函數 \Rightarrow 解釋一个文字

⑦: 問題在求 inverse function. $X \sim F$ 是知道的
but we don't know 是否可以求反

for example: $F_\lambda(x) = 1 - \exp(-\lambda \cdot x) = y$ (\Leftarrow 若 CDF 是这个

$$\text{反} \rightarrow F_\lambda^{-1}(y) = \frac{\log(1-y)}{-\lambda} \quad \stackrel{=?}{=} \text{exponential function}$$

若令: $x = \frac{\log(1-y)}{-\lambda} \sim \exp(\lambda) = F_\lambda(v)$ \Leftarrow 定義
 $v \sim \text{uniform distribution}$

$$P(X \leq x) = P(F_\lambda(v) \leq x) = P(v \leq F_\lambda(x)) = F_\lambda(x)$$

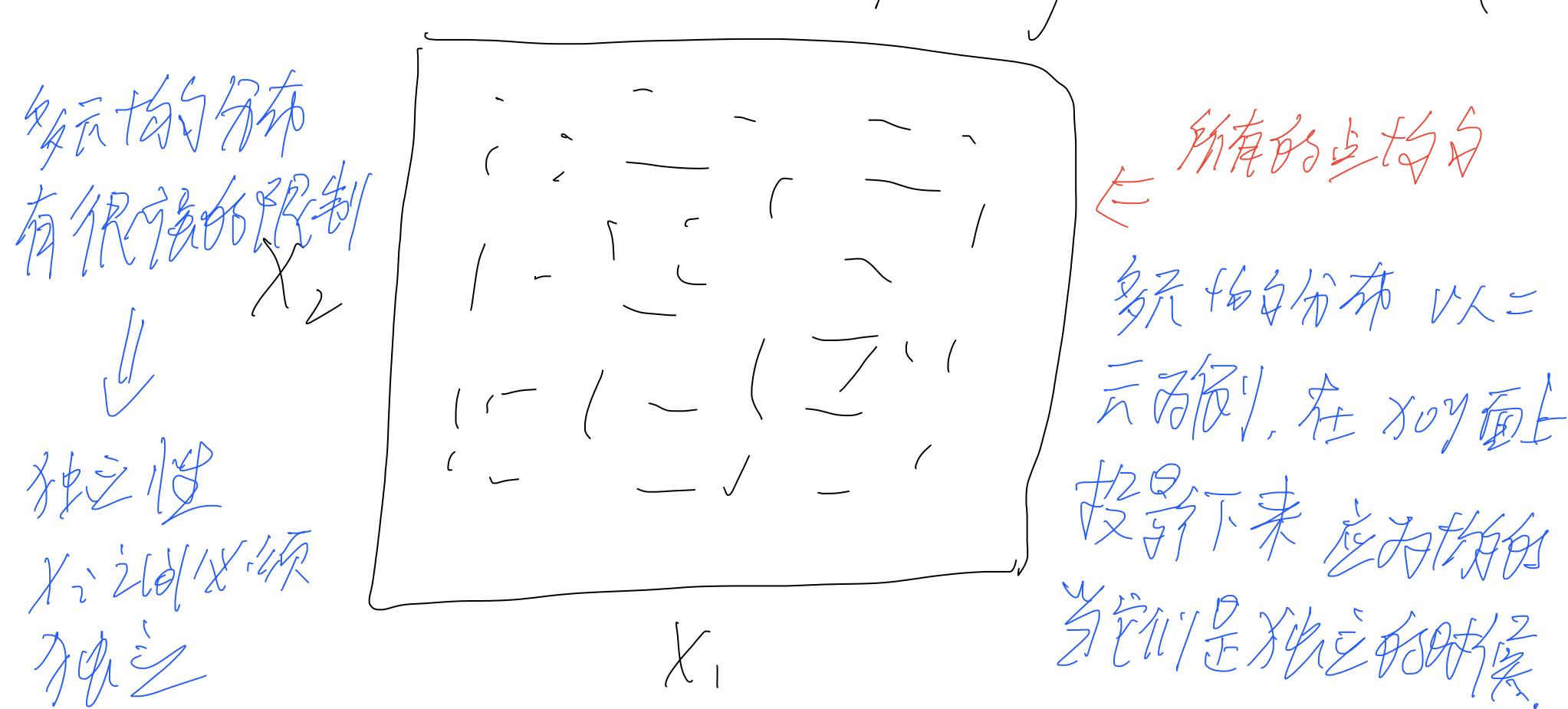
In general, we can construct a multivariate distribution for $X = (X_1 \cdots X_p)^T$ by letting $X_1 \cdots X_p$ be independent univariate random variables. The density is:

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

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

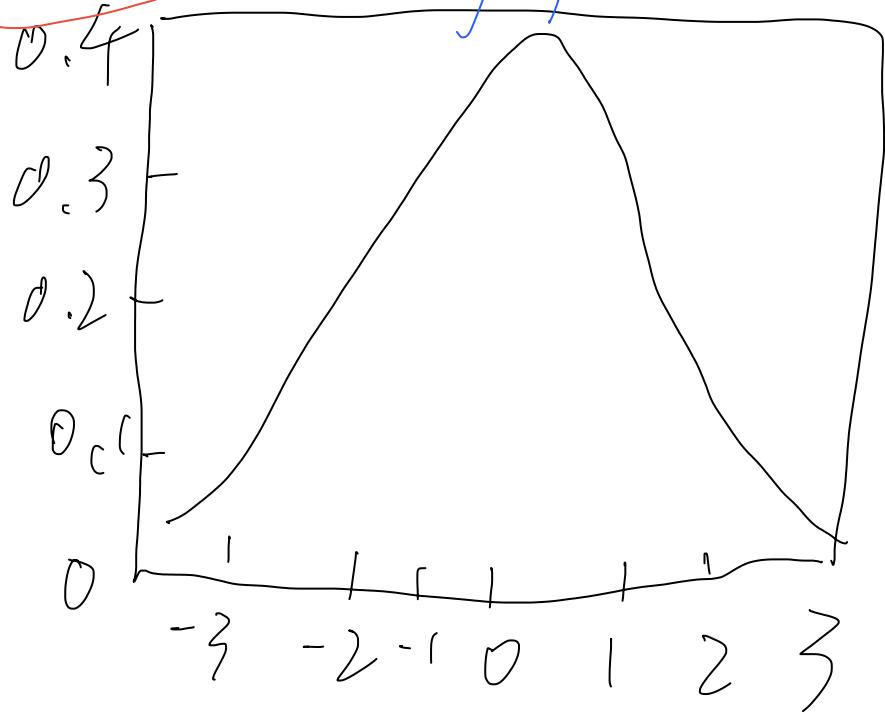
recall the $X \sim N(\mu, \sigma^2)$ has a univariate normal distribution if its density is

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

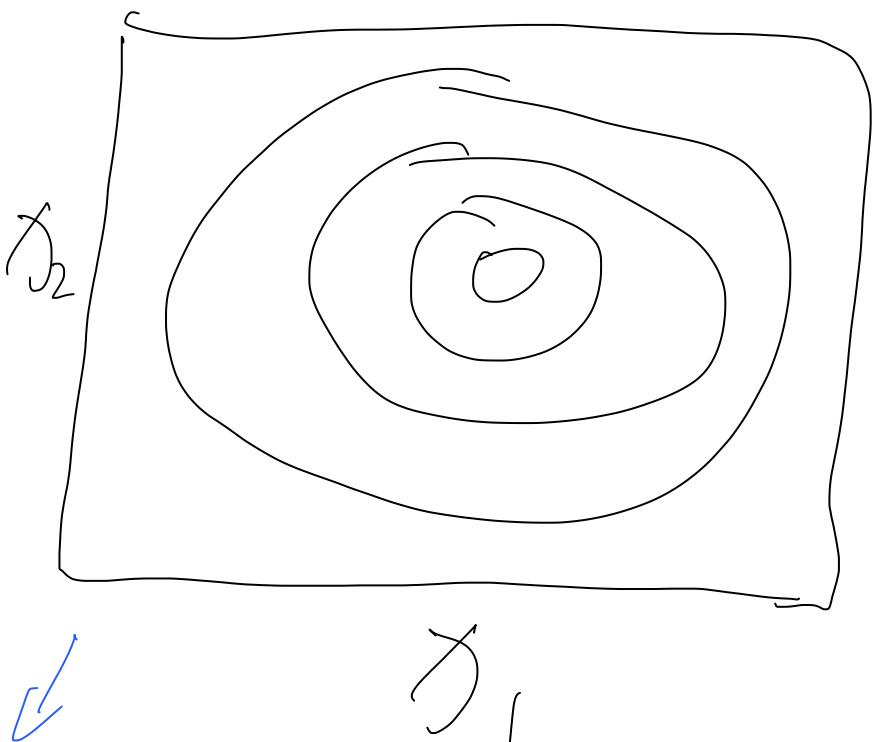
the standard multivariate normal distribution follows by the same "independent construction" as before, namely let $X_1 \cdots X_p$ be standard normal $N(0, 1)$, then $X = (X_1 \cdots X_p)^T$ has density

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

$d\text{norm}(x) \Rightarrow$ density function



某个截面 X



拉下來是圖！

當 X_1, X_2 独立。

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)
transformation (8)

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
反函数 無倒推解

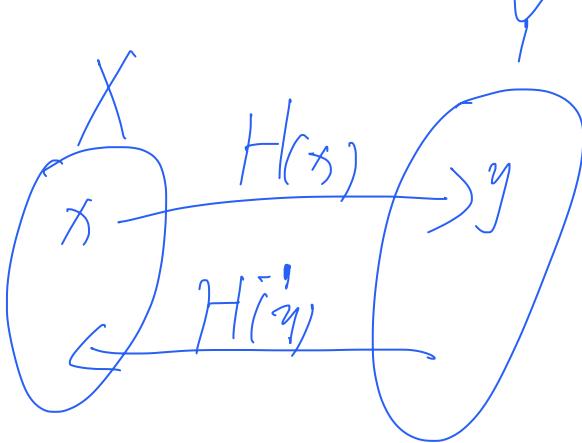
where for a mapping $T = (T_1 \cdots T_P): \mathbb{R}^P \rightarrow \mathbb{R}^P$

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

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

the matrix J_T is the best linear approximation of T at y how to prove (9)

(8)



(9)

你认为T是一个未知的函数，比如二次函数，然后你想让计算机在某一点上，然后你试着在这一点上找到最佳超平面，它在这一点上最接近函数，所以它是在点Y上函数T的最佳线性近似

know $X \sim F$

to get $Y \sim ?$

$$P(Y \leq y) = P(H(X) \leq y)$$

$$F_Y(y) = P(X \leq H^{-1}(y))$$

$$F_Y(y) = F_X(H^{-1}(y)) \in \text{CDF}$$

$$f_Y(y) = f_X(H^{-1}(y)) \left(\frac{d(H^{-1}(y))}{dy} \right) \in \text{chain rules}$$

这是一维的，高维如次。

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

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

$$\Rightarrow J_H(\gamma) = A \quad \text{该导数, 具体?} \quad (10)$$

scaling and rotation: 改变形状

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 $V \Lambda^{\frac{1}{2}} x$ (11)

$$x \in \mathbb{R}^P \quad f_{V \Lambda^{\frac{1}{2}} x}(x) = \frac{1}{\sqrt{(2\pi)^P \det(\Lambda)}} \exp \left\{ -\frac{1}{2} x^T \Lambda^{-1} x \right\} \quad (12)$$

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 $V \Lambda^{\frac{1}{2}} x$, (13)

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

$$(1) \quad A^{\frac{1}{2}}x = y$$

$$H(x) = A^{\frac{1}{2}}x \Rightarrow A^{\frac{1}{2}}A^{\frac{1}{2}}x = A^{\frac{1}{2}}y$$

$$H^{-1}(y) = \leftarrow x = A^{-\frac{1}{2}}y$$

$$\det(J_{H^{-1}(y)}) = \prod_{i=1}^p \lambda_i^{-\frac{1}{2}} \Rightarrow \sqrt{\det(A)} = ? \quad (12)$$

$$-\sum \lambda_i^{-1} x^T A^{-1} x = (x \cdot A^{-\frac{1}{2}})^T \cdot x \cdot A^{-\frac{1}{2}}$$

$$-\sum \lambda_i^{-1} x^T A^{-1} x = x^T A^{-\frac{1}{2}} A^{-\frac{1}{2}} \cdot x$$

$$(2) \quad V \cdot V^T = I, \quad y = V A^{\frac{1}{2}} x_i.$$

$$V^T = V^{-1} \quad z = H(y) = V \cdot y$$

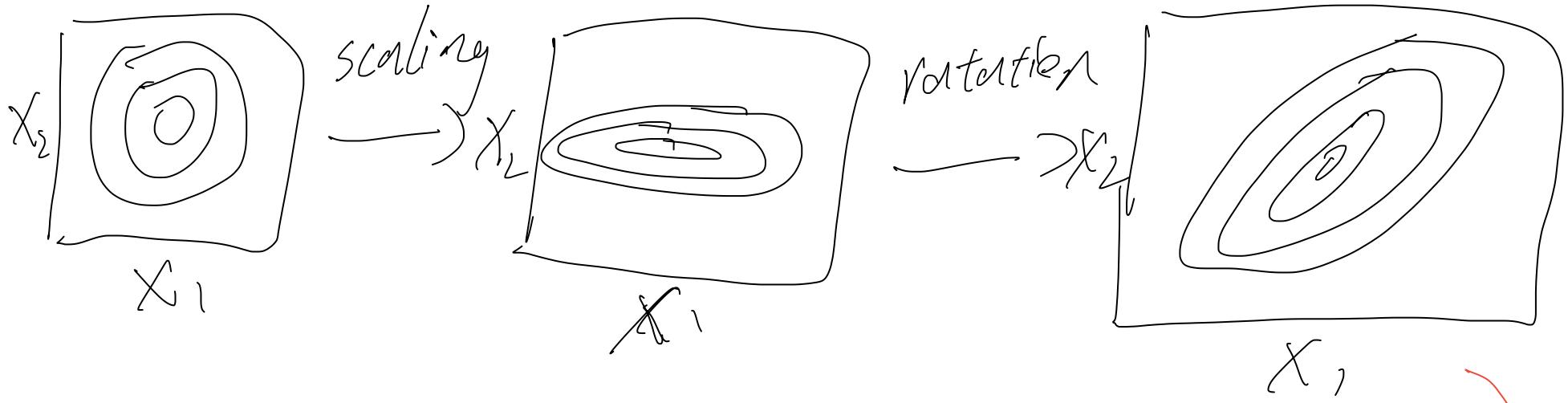
$$\rightarrow H^{-1}(z) = V^{-1} \cdot z = V^T z \Rightarrow J_{H^{-1}(z)} = V^T$$

$$\det(VAV^T) = \det(V) \cdot \det(A) \cdot \det(V^T)$$

$$\det(A) \cdot \det(B) = \det(A \cdot B)$$

$$\begin{aligned} \text{cov}(VA^{\frac{1}{2}}x) &= V \cdot A^{\frac{1}{2}} \cdot \text{cov}(x) \cdot (V \cdot A^{\frac{1}{2}})^T \\ &= V \cdot A^{\frac{1}{2}} \underline{I} \cdot (V \cdot A^{\frac{1}{2}})^T \\ &= VA^{\frac{1}{2}}V^T \end{aligned}$$

$$x \sim \mathcal{N}(0, I)$$



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

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

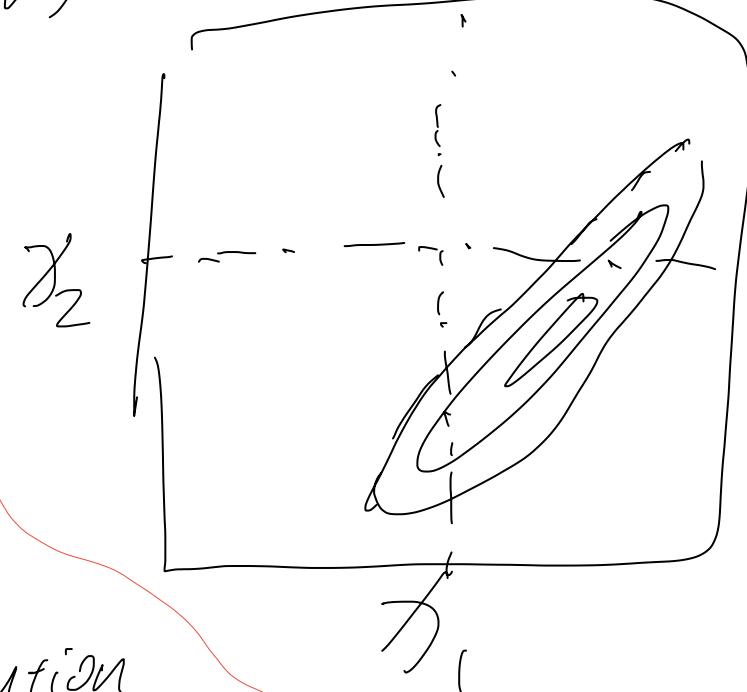
let Σ be a positive
 -definite $P \times P$ -matrix

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 Σ has eigendecomposition

$\Sigma = U\Lambda U^T$ where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$, λ_i is the eigenvalues of Σ , and U is a rotation.

matrix. Therefore, we get this representation

$$X \sim u + \underbrace{\Sigma^{1/2} N_p(0, I_p)}_{\text{该步}} V$$

by the rules for expectation and covariance

this gives: $E[X] = u$, $\text{cov}(X) = \Sigma$

Let $X \sim N_p(u, \Sigma)$

① the marginal (X_j) has a univariate normal distribution $N(u_j, \sigma_j^2)$

② Let A be a $m \times p$ matrix, $b \in R^m$, then $b + Ax$ is multivariate normal $N_m(b + Au, A\Sigma A^T)$

in particular

$$\Sigma^{1/2}(X - u) \sim N_p(0, I_p)$$

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

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

