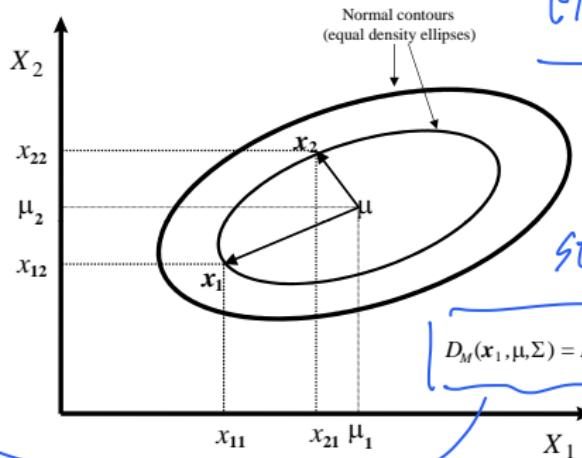


Contours: ellipse 橢圓

大綱 The density f of X is constant on ellipsoids of the form

$$x \in \mathbb{R}^p : (x - \mu)^T \Sigma^{-1}(x - \mu) = c.$$

The Mahalanobis distance between $x, y \in \mathbb{R}^p$ \Rightarrow mahalanobis distance is fixed by c
takes into account the covariance structure of X . All points on the same ellipsoid have the same Mahalanobis distance from the center μ .



elliptical distribution

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

$$D_M(x_1, \mu, \Sigma) = D_M(x_2, \mu, \Sigma)$$

Quadratic forms

- If $Y \sim N_p(0_p, I_p)$, then

$$Y^T Y = \sum_{j=1}^p Y_j^2 \sim \chi_p^2$$

where χ_p^2 denotes the chi-squared distribution with p degrees of freedom which has density

$$f(x) = \frac{1}{2^{p/2}\Gamma(p/2)} x^{p/2-1} \exp(-x/2).$$

- Consequently, for $X \sim N_p(\mu, \Sigma)$, the square of the Mahalanobis distance is distributed as a χ_p^2 ,

$$D^2 = (X - \mu)^T \Sigma^{-1} (X - \mu) = Y^T Y \sim \chi_p^2.$$

$$\frac{(X-\mu)^T \Sigma^{-1} (\Sigma^{-1})^T (X-\mu)}{N(\mu, I_p)^T N(\mu, I_p)}$$

Conditional distribution

6. A multivariate normal variable $X \sim N_p(\mu, \Sigma)$ can be partitioned into $(X_{(1)}, X_{(2)}) \in \mathbb{R}^{p_1+p_2}$, $\mu = (\mu_{(1)}, \mu_{(2)})$ and

$$\Sigma = \begin{bmatrix} \Sigma_{(11)} & \Sigma_{(12)} \\ \Sigma_{(21)} & \Sigma_{(22)} \end{bmatrix}.$$

The conditional distribution of $X_{(2)}$ given $X_{(1)} = x_{(1)}$ can be computed explicitly by the formula

$$f_{X_{(2)}|X_{(1)}}(x_{(2)} | X_{(1)} = x_{(1)}) = \frac{f_{(X_{(1)}, X_{(2)})}(x_{(1)}, x_{(2)})}{f_{X_{(1)}}(x_{(1)})}.$$



It turns out that this is the density of p_2 -variate normal distribution

$$N_{p_2} \left(\mu_{(2)} + \Sigma_{(21)} \Sigma_{(11)}^{-1} (x_{(1)} - \mu_{(1)}), \Sigma_{(22)} - \Sigma_{(21)} \Sigma_{(11)}^{-1} \Sigma_{(12)} \right)$$



Central Limit Theorem

The CLT is the reason why the multivariate normal distribution is omnipresent.

- Let X_i , $i = 1, \dots, n$ be independent and identically distributed random p -vectors with mean μ and variance Σ , then, as $n \rightarrow \infty$,

From many distribution

$$\sqrt{n}(\bar{X} - \mu) \xrightarrow{D} N_p(0, \Sigma)$$
$$\bar{X} \sim N\left(\mu, \frac{\Sigma}{n}\right)$$

where $\bar{X} = \frac{1}{n} \sum_i X_i$ and \xrightarrow{D} denotes convergence in distribution.

Recall: One denotes $X_n \xrightarrow{D} X$ if $\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$ for all $x \in \mathbb{R}^p$ such that F is continuous at x .

Central Limit Theorem

The CLT is the reason why the multivariate normal distribution is omnipresent.

- Let $X_i, i = 1, \dots, n$ be independent and identically distributed random p -vectors with mean μ and variance Σ , then, as $n \rightarrow \infty$,

$$\sqrt{n} (\bar{X} - \mu) \xrightarrow{D} N_p(0, \Sigma)$$

where $\bar{X} = \frac{1}{n} \sum_i X_i$ and \xrightarrow{D} denotes convergence in distribution.

Recall: One denotes $X_n \xrightarrow{D} X$ if $\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$ for all $x \in \mathbb{R}^p$ such that F is continuous at x .

- This means, that for large n , $\sqrt{n} (\bar{X} - \mu)$ is approximately normal, sometimes denoted by

$$\bar{X} \sim N_p \left(\mu, \frac{1}{n} \Sigma \right)$$

 The distribution of $\sqrt{n} (\bar{X} - \mu)$ is exactly normal when the X_i 's are themselves normally distributed.

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

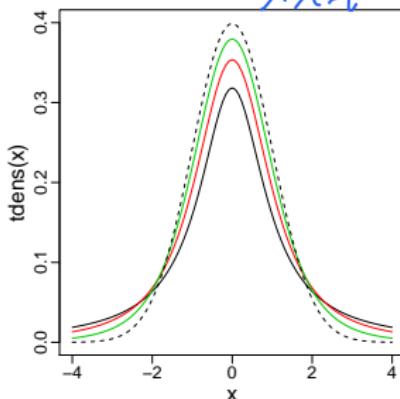
Student t -distribution

- Recall the univariate Student t -distribution $X \sim t(\nu, \mu, \sigma^2)$ has density

$$f(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\sigma^2 \nu \pi} \Gamma(\frac{\nu}{2})} \left(1 + \frac{(x - \mu)^2}{\sigma^2 \nu}\right)^{-\frac{\nu+1}{2}}, \quad x \in \mathbb{R},$$

with $E(X) = \mu$, $\nu > 1$, and $\text{Var}(X) = \frac{\sigma^2 \nu}{\nu - 2}$, $\nu > 2$.

- The tail of the t -distribution is polynomial and thus heavier than the normal.
- When the degrees of freedom ν tends to ∞ , the t -distribution converges to the normal distribution (dashed line).



$$\frac{\bar{x} - \mu}{\sqrt{\frac{s^2}{n}}} \sim t(n)$$

$s < \hat{\sigma}$

(13)

多元如何推广到多元

Multivariate Student t -distribution

- Similarly as for the multivariate normal we can apply scaling and rotation transformation to the density

$$f(x) = c_p \left(1 + \frac{(x - \mu)^T (x - \mu)}{\nu} \right)^{-\frac{\nu+p}{2}}.$$

- If X has a multivariate Student t -distribution, denoted $X \sim t_p(\nu, \mu, \Sigma)$, then its density is

$$f(x) = c_p \det(\Sigma)^{-1/2} \left(1 + \frac{(x - \mu)^T \Sigma^{-1} (x - \mu)}{\nu} \right)^{-\frac{\nu+p}{2}}$$

where ν is the degrees of freedom and

$$c_p = \frac{\Gamma((\nu + p)/2)}{(\pi\nu)^{p/2} \Gamma(\nu/2)}$$

- We have $E(X) = \mu$, $\nu > 1$, and $\text{cov}(X) = \frac{\nu}{\nu-2} \Sigma$, $\nu > 2$.

Properties

1. The marginals of the multivariate t -distribution have univariate t -distributions, and thus heavy tails.
2. If Σ is diagonal the components of X are uncorrelated but **not** independent, unlike for the multivariate normal distribution.
3. When $\nu = 1$, it is called the multivariate Cauchy distribution.
4. Like the multivariate normal, the multivariate Student is an **elliptical distribution**: the density function only depends on the quadratic form $(x - \mu)^T \Sigma^{-1} (x - \mu)$



Wishart distribution

This is a random matrix distribution.

- Let X_1, \dots, X_n a random sample (i.i.d.) with $X_i \sim N_p(\mu, \Sigma)$, then the distribution of

sample covariance matrix $\sum_{i=1}^n (X_i - \mu)(X_i - \mu)^T$ ERPKP

if we change $\mu \rightarrow A\mu$
nothing will be changed

has the Wishart distribution $W_p(n, \Sigma)$, n is the degrees of freedom.

- If $M \sim W_p(n, \Sigma)$, the density is

R PXP^T

$$f(M) = \frac{|M|^{(n-p-1)/2} \exp\left(-\frac{1}{2} \text{tr}(\Sigma^{-1} M)\right)}{2^{np/2} |\Sigma|^{n/2} \Gamma_p(n/2)}$$

it's independent from μ .

just keep in mind
where $\Gamma_p(\cdot)$ is the multivariate gamma function
not very important

example: $X_i \sim N_p(\mu, \Sigma)$

$$\Gamma_p(n/2) = \pi^{p(p-1)/4} \prod_{j=1}^p \Gamma\left(n/2 + (1-j)/2\right)$$

$$\sum_{i=1}^n X_i X_i^T \sim W_p(n, \Sigma)$$

and $\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$ is the gamma function.

Properties, remarks and cases

- The Wishart distribution is the multivariate analogue of the chi-square distribution.
- By Cochran's theorem (cf., Mardia et al., p. 68), the scaled sample covariance is Wishart distributed

$$\left[\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T = (n-1)\mathbf{S} \right] \sim W_p(n-1, \Sigma).$$

and independent of the sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N_p(\mu, n^{-1}\Sigma). \quad \text{(14) 3.}$$

(1b)

- If $p = 1$, then $(n-1)\mathbf{S}$ is $\sigma^2 \chi_{n-1}^2$. ~~(14) 2.~~
- If \mathbf{W} is $W_p(m, \Sigma)$, then \mathbf{BWB}^T is $W_q(m, \mathbf{B}\Sigma\mathbf{B}^T)$ for any $q \times p$ matrix \mathbf{B} of rank $q \leq p$.



$$\left\{ \begin{array}{l} W \sim W_p(m, \Sigma) \\ W = \sum_{i=1}^m X_i^T X_i, \quad X_i \sim N_p(0, \Sigma) \\ BWB^T = B \left(\sum_{i=1}^m X_i^T X_i \right) B^T \end{array} \right.$$

$BWB^T = \sum_{i=1}^m (B^T X_i) \cdot (B X_i)^T$
 it's linear
 by defn:
 $B^T X_i \sim N_q(0, B \Sigma B^T)$
 $\therefore BWB^T \sim W_q(m, B \Sigma B^T)$

Hotelling's T^2 distribution

This is a univariate distribution, an extension of the Student's t -distribution used in the context of multivariate tests.

- Let $Z \sim N(0, I_p)$ and $M \sim W_p(m, I_p)$ be independent then, by definition, the distribution of

$$mZ^T M^{-1} Z \stackrel{R}{=} mZ^T M^{-1} Z \sim T^2(p, m)$$

is Hotelling $T^2(p, m)$ (stochastic representation).

- It is equivalent to the F distribution:

$$T^2(p, m) = \{mp/(m - p + 1)\} F_{p, m-p+1}$$

Recall: The F_{d_1, d_2} distribution is the distribution of the ratio of two independent, scaled chi-squared distribution $U_1 \sim \chi^2_{d_1}$ and $U_2 \sim \chi^2_{d_2}$:

$$\frac{U_1/d_1}{U_2/d_2} \sim F_{d_1, d_2}.$$

\hookrightarrow degree of freedom

Hotelling's T^2 distribution: properties

- If $X \sim N_p(\mu, \Sigma)$ and $M \sim W_p(m, \Sigma)$ are independent, then

$$m(X - \mu)^T M^{-1}(X - \mu) \sim T^2(p, m)$$

$$x_i \sim N_p(0, \Sigma)$$

$$m \sim W_p(m, \Sigma)$$

$$M = \sum_{i=1}^m x_i x_i^T = \Sigma^{-\frac{1}{2}} \left(\sum_{i=1}^m x_i x_i^T \right) \Sigma^{-\frac{1}{2}}^T$$

$$\tilde{x}_i = \Sigma^{\frac{1}{2}} x_i \sim N_p(0, I_p)$$

$$\tilde{M}^{-1} = (\Sigma^{\frac{1}{2}})^T \cdot M^{-1} \cdot \Sigma^{\frac{1}{2}}$$

$$m \cdot X^T \tilde{M}^{-1} X = m \cdot \tilde{M} \cdot (\Sigma^{\frac{1}{2}} X)^T$$

by definition $\tilde{M}^{-1} = 1 \cdot \tilde{M}^{-1} \cdot (\Sigma^{\frac{1}{2}} X)$

$$= T^2(p, m)$$

$$N_p(0, I_p)$$

Hotelling's T^2 distribution: properties

- If $X \sim N_p(\mu, \Sigma)$ and $\mathbf{M} \sim W_p(m, \Sigma)$ are independent, then

$$m(X - \mu)^T \mathbf{M}^{-1} (X - \mu) \sim T^2(p, m)$$

(15)

- Let X_1, \dots, X_n i.i.d. $N_p(\mu, \Sigma)$, then

$$T^2 = n(\bar{X} - \mu)^T \mathbf{S}^{-1} (\bar{X} - \mu) \sim T^2(p, n-1) = \{(n-1)p/(n-p)\} F_{p, n-p},$$

Proof: Write

$$T^2 = (n-1)Z^T \mathbf{M}^{-1} Z, \quad Z \sim N(0, \Sigma)$$

with $Z = \sqrt{n}(\bar{X} - \mu)$ and $\mathbf{M} = (n-1)\mathbf{S}$.

$M \sim W_p(n-1, \Sigma)$

$$T^2 = \{(n-1) \cdot [\sqrt{n}(\bar{X} - \mu)]^T \frac{\mathbf{S}^{-1}}{(n-1)} \cdot \sqrt{n}(\bar{X} - \mu)\}$$

$$T^2 = n \cdot (X - \mu)^T \mathbf{S}^{-1} (\bar{X} - \mu) \quad \bar{X} \perp\!\!\!\perp S \Rightarrow \text{independent}$$

(15)

Hotelling's T^2 distribution: properties

$$T^2(p, n-1)$$

- If $X \sim N_p(\mu, \Sigma)$ and $\mathbf{M} \sim W_p(m, \Sigma)$ are independent, then

$$m(X - \mu)^T \mathbf{M}^{-1} (X - \mu) \sim T^2(p, m)$$

- Let X_1, \dots, X_n i.i.d. $N_p(\mu, \Sigma)$, then

$$T^2 = n(\bar{X} - \mu)^T \mathbf{S}^{-1} (\bar{X} - \mu) \sim T^2(p, n-1) = \{(n-1)p/(n-p)\} F_{p, n-p},$$

Proof: Write

$$T^2 = (n-1)Z^T \mathbf{M}^{-1} Z,$$

with $Z = \sqrt{n}(\bar{X} - \mu)$ and $\mathbf{M} = (n-1)\mathbf{S}$.

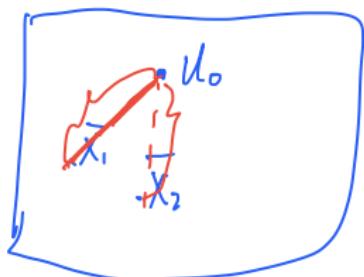
- This statistics (i.e. the estimated Mahalanobis distance) is used for hypothesis testing of population means (multivariate extensions of the t-tests).

Hypothesis testing in p dimension for normal distribution. Only for normal distribution.

$$H_0: \mu = \mu_0 \in \mathbb{R}^p [i-p \geq \mu_0 \text{ is a value.}] \quad \text{known } \Sigma$$

you observe the sample mean \bar{X} , we want to test the true mean μ_0 if it's equal \bar{X}
 $\bar{X} \cong \mu_0$.

use mahalanobis distance.



$$(x - \mu_0)^T I_p^{-1} (x - \mu_0) = c$$

if c is large, $\bar{x} \neq \mu_0$

if c is small, $\bar{x} \cong \mu_0$

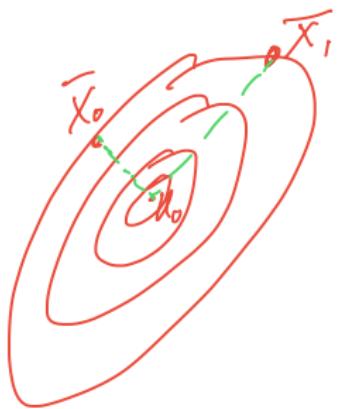
but we don't

\Rightarrow use it.

because of

M -distance

use the Σ



for m-distance \bar{x}_1 , and \bar{x}_0 is equal, but we think \bar{x}_0 is better. So we use:

$$T^2 = n \cdot (\bar{x} - u_0)^T \Sigma^{-1} (\bar{x} - u_0)$$

$$\begin{aligned} T^2 &= \sqrt{n} \cdot (\bar{x} - u_0)^T \bar{\Sigma}^{-\frac{1}{2}} \cdot \bar{\Sigma}^{-\frac{1}{2}} (\bar{x} - u_0) \\ &= \underbrace{\sqrt{n} \cdot \left(\sum \bar{\Sigma}^{-\frac{1}{2}} \cdot (\bar{x} - u_0) \right)^T}_{S} \underbrace{\sqrt{n} \cdot \left[\bar{\Sigma}^{-\frac{1}{2}} \cdot (\bar{x} - u_1) \right]}_{\text{known } \bar{\Sigma}} \\ &= N(0, I_p) \quad (16) \quad (n-1)S \sim \chi_p^2 \\ &= Y^T Y = \sum_{i=1}^p Y_i^2 \sim \chi_p^2 \text{ when } p=1 \end{aligned}$$

$$\alpha = P(T^2 \geq k | H_0) = 1 - \chi_p^2(k) \quad \text{if } T^2 \geq k, \text{ reject } H_0$$

if unknown Σ . \rightarrow positive semidefinite.

$$T^2 = n \cdot (\bar{x} - \mu_0)^T S^{-1} (\bar{x} - \mu_0) \text{ by (1)}$$

$$\sim T^2(p, n-1)$$

if $T^2 \geq k$, where $\alpha = 1 - T(p, n-1)_{(k)} \Rightarrow$ reject H_0 .

summary, if T^2 is large, reject.

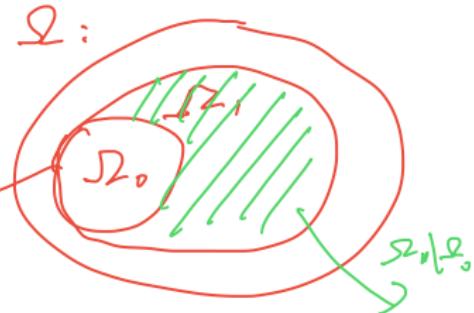
Another procedure: the likelihood ratio test – LRT

A more general framework for testing is based on the log-likelihood of the data x_1, \dots, x_n that are supposed to have density $f(\cdot; \theta)$ for some parameter $\theta \in \Omega$:

$$\ell(\theta) = \sum_{i=1}^n \log f(x_i; \theta) = \log L(\theta) = \log \prod_{i=1}^n f(x_i; \theta).$$

you can think
it's covariance
for ~~example~~
example
它可以是任意变量.

we want to
test if $S_2 \not\subset$
有 θ , 即 $\theta \in S_2$
subspace



Another procedure: the likelihood ratio test – LRT

A more general framework for testing is based on the log-likelihood of the data x_1, \dots, x_n that are supposed to have density $f(\cdot; \theta)$ for some parameter $\theta \in \Omega$:

$$\ell(\theta) = \sum_{i=1}^n \log f(x_i; \theta) = \log L(\theta) = \log \prod_{i=1}^n f(x_i; \theta).$$

Likelihood ratio test: Suppose we want to test

$$H_0 : \theta \in \Omega_0, \text{ versus } H_1 : \theta \in \Omega_1 \setminus \Omega_0,$$

where $\Omega_0 \subset \Omega_1 \subset \Omega$, $\dim(\Omega_0) = r$ and $\dim(\Omega_1) = q$. Define the (log)-likelihood ratio

$$LRT = -2 \log(\hat{\ell}_0 / \hat{\ell}_1) = 2(\hat{\ell}_1 - \hat{\ell}_0),$$

where $\hat{\ell}_0 = \max_{\theta \in \Omega_0} \ell(\theta)$ and $\hat{\ell}_1 = \max_{\theta \in \Omega_1} \ell(\theta)$. Under mild conditions, for every true $\theta \in \Omega_0$,

$$LRT \xrightarrow{n \rightarrow \infty} D \chi^2_{q-r},$$

sometimes written as $LRT \sim \chi^2_{q-r}$.

if $LRT \geq k$,
 $\alpha = 1 - \chi^2_{q-r}(k)$
then we reject H_0

$$\dim(\Omega_1) = R^P$$

$$H_0: u = u_0 \in R^P$$

$$\Omega_0 = \{u_0\}$$

$$\dim(\Omega_0) = 0$$

Another procedure: the likelihood ratio test – LRT

$$\begin{array}{c} x_2 \\ \vdots \\ x_1 \end{array} \Rightarrow R^2 \quad \dim(R^2) = 2$$

A more general framework for testing is based on the log-likelihood of the data x_1, \dots, x_n that are supposed to have density $f(\cdot; \theta)$ for some parameter $\theta \in \Omega$:

$$\ell(\theta) = \sum_{i=1}^n \log f(x_i; \theta) = \log L(\theta) = \log \prod_{i=1}^n f(x_i; \theta).$$

Likelihood ratio test: Suppose we want to test

bnt one point in
dimension will be
consider as 0

$$H_0 : \theta \in \Omega_0, \text{ versus } H_1 : \theta \in \Omega_1 \setminus \Omega_0,$$

where $\Omega_0 \subset \Omega_1 \subset \Omega$, $\dim(\Omega_0) = r$ and $\dim(\Omega_1) = q$. Define the (log)-likelihood ratio

$$LRT = -2 \log(\hat{L}_0 / \hat{L}_1) = 2(\hat{\ell}_1 - \hat{\ell}_0),$$

where $\hat{\ell}_0 = \max_{\theta \in \Omega_0} \ell(\theta)$ and $\hat{\ell}_1 = \max_{\theta \in \Omega_1} \ell(\theta)$. Under mild conditions, for every true $\theta \in \Omega_0$:

$$LRT \xrightarrow{n \rightarrow \infty} D \chi^2_{q-r},$$

sometimes written as $LRT \sim \chi^2_{q-r}$.

- ▶ This provides an asymptotically motivated test that does not require normally distributed samples.

example:

$$X_1, \dots, X_n \sim \text{Exp}(\lambda)$$

$$\text{test : } H_0: \lambda = \lambda_0$$

$$\Omega_0 = \{\lambda_0\} \quad \Omega_1 = \mathbb{R}.$$

$$q=1 \quad r=0$$

$$LR T \sim \chi^2_n$$

$$f(x; \lambda) = \lambda \cdot \exp(-\lambda x), x > 0$$

$$\log f(x; \lambda) = \log(\lambda) - \lambda x$$

$$l(\lambda) = \sum_{i=1}^n \log(f(x_i; \lambda)) = n \cdot \log \lambda - n \cdot \bar{x} \cdot \lambda$$

$$\Omega_0: \hat{h}(\lambda_0) = n \cdot \log \lambda_0 - n \cdot \bar{x} \cdot \lambda_0 \leq \max_{\lambda \in \Omega_0} h(\lambda)$$

$$\Omega_1: \frac{d h(\lambda)}{d \lambda} = 0 \Rightarrow \frac{n}{\lambda} - n \cdot \bar{x} = 0 \quad \left(= \max_{\lambda \in \Omega_1} h(\lambda)\right)$$

$\lambda = \frac{1}{\bar{x}}$

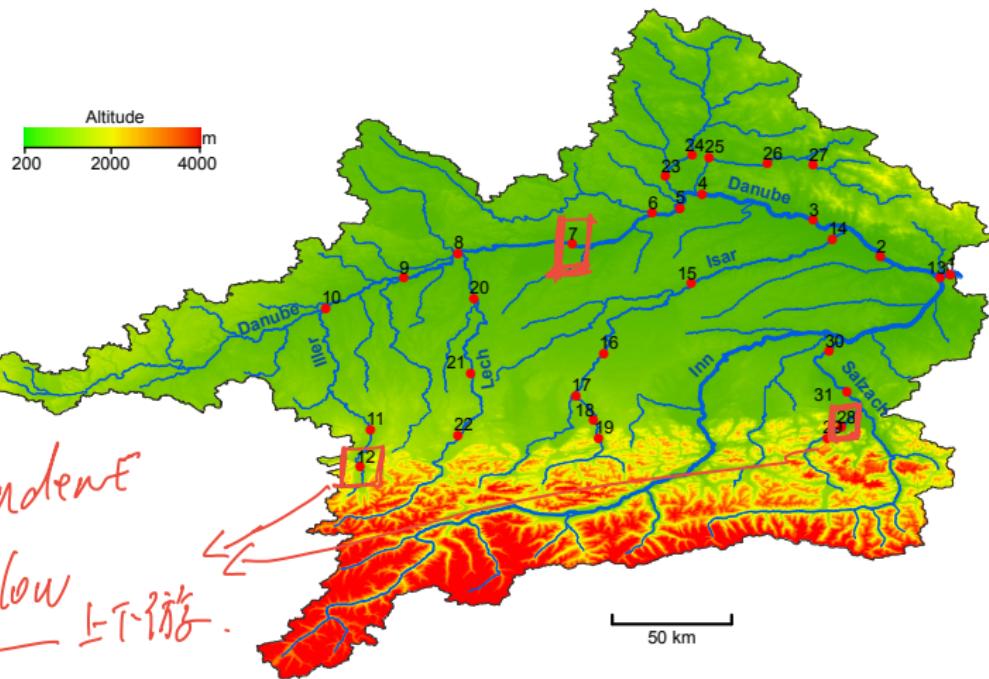
$$\begin{aligned}\hat{h}(\lambda_1) &= n \cdot \log \left(\frac{1}{\bar{x}} \right) - n \cdot \frac{1}{\bar{x}} \cdot \bar{x} \\ &= -n \log \bar{x} - n\end{aligned}$$

↑ close to 0, reject H_0
 n is small

$$\begin{aligned}LRT &= 2 \left(\hat{h}(\lambda_0) - \hat{h}(\lambda_1) \right) = 2n \cdot \left[\lambda_0 \bar{x} - 1 - \log \left(\frac{\lambda_0}{\bar{x}} \right) \right] \\ &\approx \chi^2_1, \text{ for } n \rightarrow \infty\end{aligned}$$

River discharge data

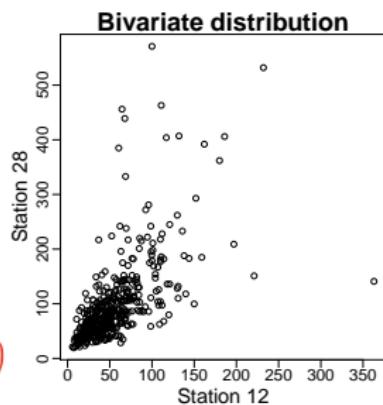
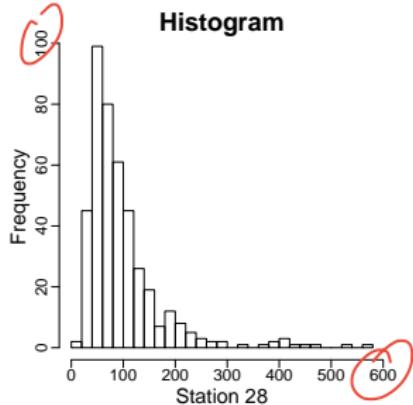
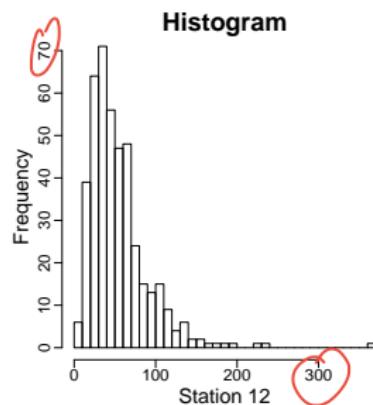
排放



① dependent
waterflow 上下游

rain 下雨 范围包含 {12, 28}

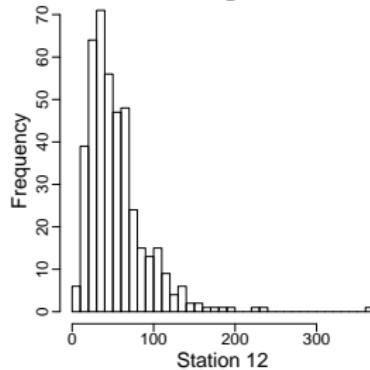
River discharge data



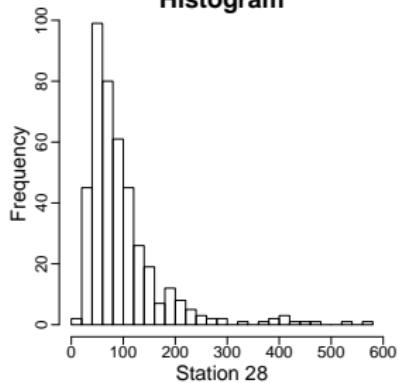
X it shows strong dependence
but hard to say!!

River discharge data

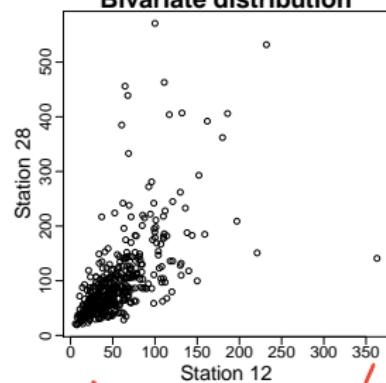
Histogram



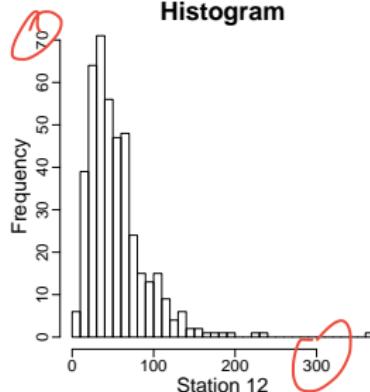
Histogram



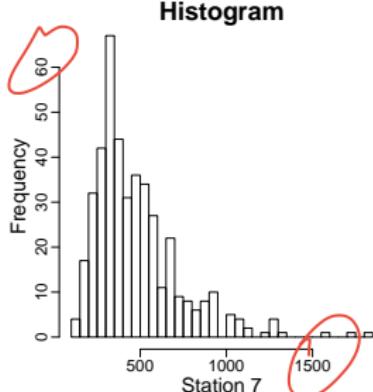
Bivariate distribution



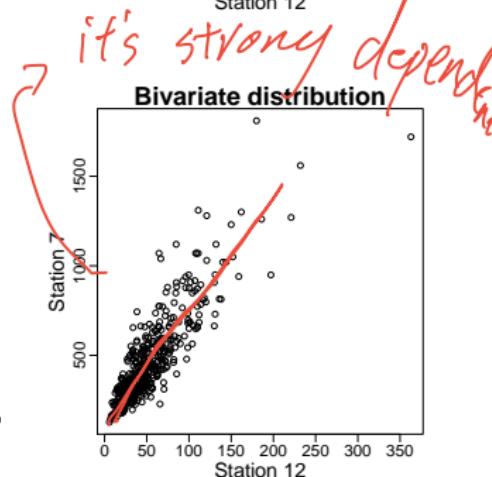
Histogram



Histogram



Bivariate distribution

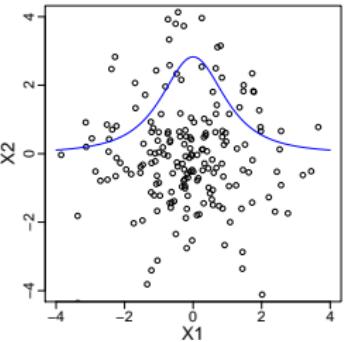
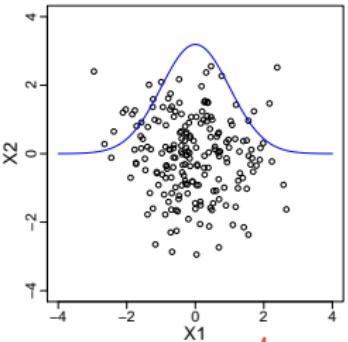
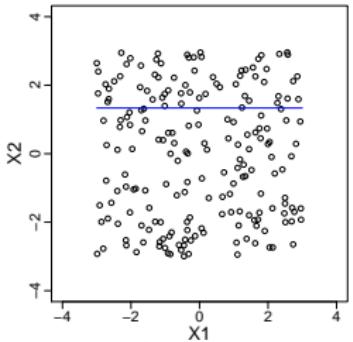
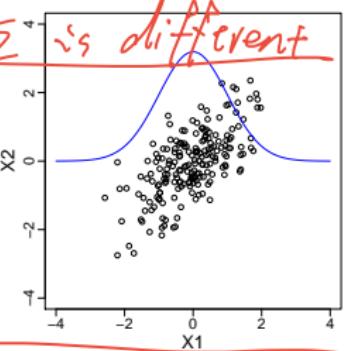
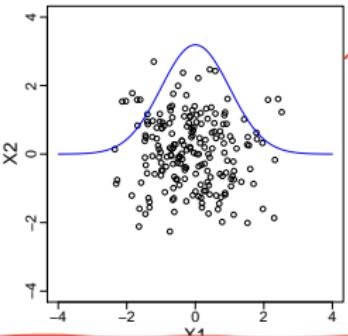
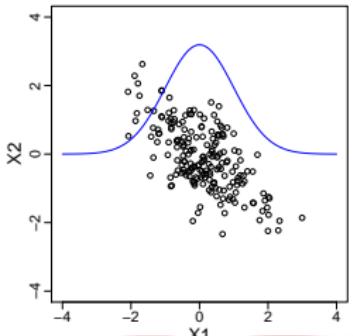


Simulated data

點圖 X₁X₂

they are same mean.

the marginal distribution
is normal distribution



marginal \rightarrow uniform
but they are independent

normal

T

Multivariate distributions

上面就想說有相同的 marginal distribution
有 dependent, but 不相同的 marginal distribution
却也可以有 independent.

A multivariate random vector $(X_1, \dots, X_p) \sim F$ has a distribution function

$$F(y_1, \dots, y_p) = P(X_1 \leq y_1, \dots, X_p \leq y_p), \text{ joint distribution}$$

and marginal distributions

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

The idea of copulas is to separate the marginals from the dependence:

$$F = \text{marginal cdfs } F_1, \dots, F_p + \text{ dependence structure } C.$$

idea.

The normalized dependence structure is called the copula C .

Advantages

- We can understand and model the dependence independently of the margins.
- Model building and fitting is easier on the standardized level.
- There might be different amounts of information about margins and dependence.

Preliminaries: the quantile transformation

F can be not invertible

Let F be a univariate distribution function, the generalized inverse (also called quantile function) is

$$F^{-\leftarrow}(p) = \inf\{x \in \mathbb{R} : F(x) \geq p\}.$$

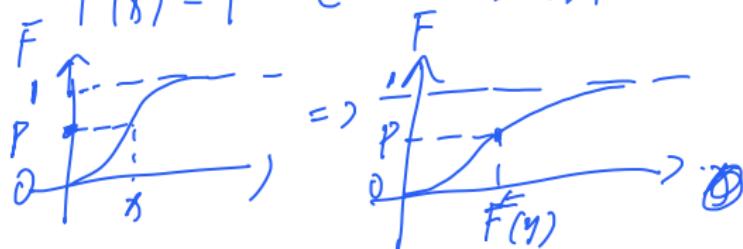
Note that if F is continuous then $F(F^{-\leftarrow}(p)) = p$ and $F^{-\leftarrow}(F(x)) = x$, for all $p \in [0, 1]$ and $x \in \mathbb{R}$.

$$F : \mathbb{R} \rightarrow [0, 1]$$

$$F : [0, 1] \rightarrow \mathbb{R}$$

exponential distribution

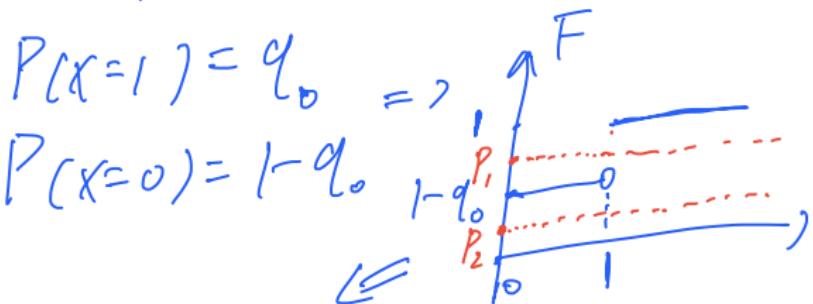
$$F(x) = 1 - e^{-x} \Rightarrow \text{CDF}$$



$$F^-(y) = -\log(1-p)$$

bernouilli distribution:

$$\text{Ber}(q_0) \quad q_0 \in [0, 1]$$



$$F(p_1) = \inf \{x \in \mathbb{R} : F(x) \geq p_1\}$$

all x_1 的 $F(x_1) \geq p_1$.

$$F(p_2) = 0 \Leftarrow \text{only } 0 \text{ 滿足.}$$

$$F(p_2) = \begin{cases} 0, & \text{if } p \in [0, 1 - q_0] \\ 1, & \text{if } p \in [1 - q_0, 1] \end{cases}$$

Preliminaries: the quantile transformation

Let F be a univariate distribution function, the generalized inverse (also called quantile function) is

$$F^{-1}(p) = \inf\{x \in \mathbb{R} : F(x) \geq p\}.$$

Note that if F is continuous then $F(F^{-1}(p)) = p$ and $F^{-1}(F(x)) = x$, for all $p \in [0, 1]$ and $x \in \mathbb{R}$.

Lemma

- Let $X \sim F$, F continuous, then $F(X) \sim U[0, 1]$.
- If $U \sim U[0, 1]$, then $F^{-1}(U) \sim F$.

$$\begin{aligned} P(F(x) \leq x) &= P(x \leq F^{-1}(x)) \\ &= F(F^{-1}(x)) = x \sim U[0, 1] \end{aligned}$$

Probability and quantile transformations are crucial to all applications involving copulas.

if $X \sim Ber(p_0)$, $Y = F(x)$, then $y = \begin{cases} 1 & x \geq p_0 \\ 0 & x < p_0 \end{cases}$

Examples:

- If $X \sim Ber(p_0)$, then $F(x) = (1 - p_0)\mathbf{1}\{x \geq 0\} + p_0\mathbf{1}\{x \geq 1\}$, and $F^{-1}(U) = \mathbf{1}\{U > 1 - p_0\} \sim Ber(p_0)$.
- If $X \sim \exp(\lambda)$ then $F(x) = 1 - \exp(-\lambda x)$ and $F(X) \sim U(0, 1)$.

if $Y = F(u) \Rightarrow P(Y \leq x) = P(F^{-1}(u) \leq x) = P(u \leq F_x) = F_x$

Definition of a copula

Summary: $X \sim F \xrightarrow{F} U \sim U[0, 1]$
all sition

Definition

A p -dimensional copula C is a multivariate distribution function $C : [0, 1]^p \rightarrow [0, 1]$ with standard uniform margins $U(0, 1)$.

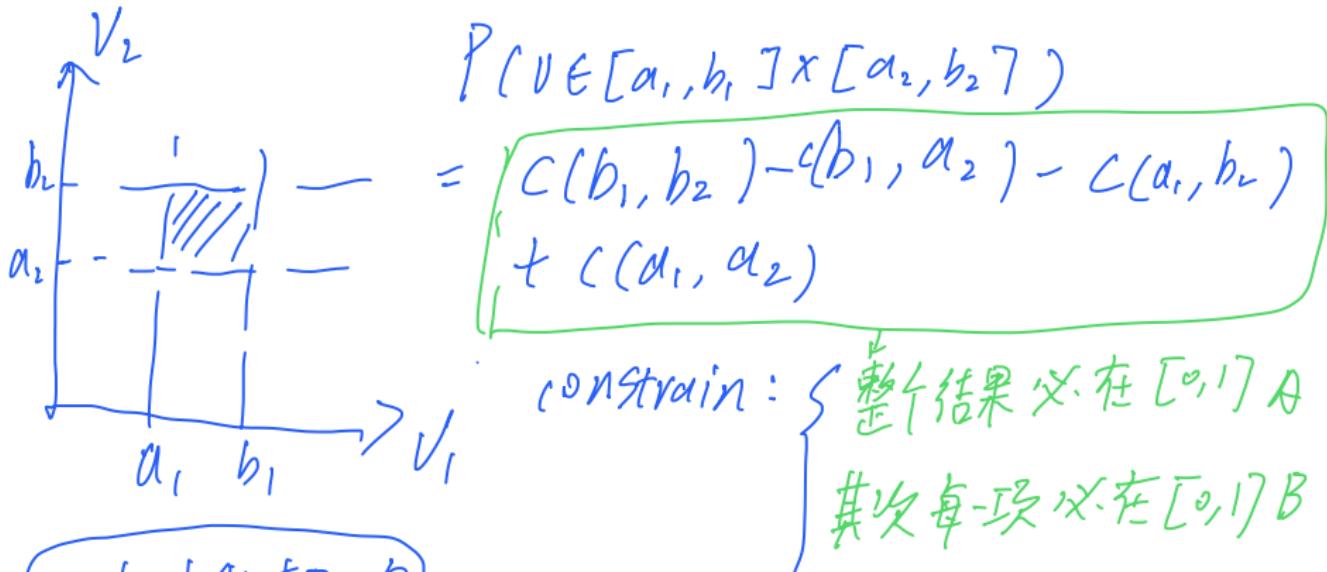
uniform margin of C

U_i is random

$$C(1, 1, 1, x_i, \dots) = x_i \quad C \text{ is distribution}$$

$$U = (U_1, \dots, U_p) \sim C$$

$$P(U_1 \leq 1, U_2 \leq 1, U_3 \leq x_3, \dots, U_p \leq 1) = x_i \in [0, 1]$$



copula 本质上把一个
 多元标准化的边缘
 把边缘用均匀替换.
 dependent 结构可用平均
 分布的联合表子, 即 copula.

A: P 必在 $[0,1]$

B: 首先将 $C(b_1, b_2)$ 表子边缘分布被
转化为 \dots 的分布

我们要求这些分布在 $[0,1]$

Definition of a copula

Definition

A p -dimensional copula C is a multivariate distribution function $C: [0, 1]^p \rightarrow [0, 1]$ with standard uniform margins $U(0, 1)$.

- A stochastic representation is $U \sim C$ with random variable.

$$P(U \leq u) = P(U_1 \leq u_1, \dots, U_p \leq u_p) = C(u_1, \dots, u_p),$$

- The margins of U are uniformly distributed

$$P(U_j \leq u_j) = C_j(u_j) = C(1, \dots, u_j, \dots, 1) = u_j.$$

- Any k -dimensional ($k < p$) marginal-subset is also a copula.

Sklar's theorem

溝通前文分析

Sklar's theorem

- Let F be a joint distribution function with margins F_1, \dots, F_p . There exists a copula C such that, for all $x = (x_1, \dots, x_p) \in \mathbb{R}^p$,
$$F(x_1, \dots, x_p) = C(F_1(x_1), \dots, F_p(x_p)).$$
- C is unique on $\text{ran}(F_1) \times \dots \times \text{ran}(F_p)$, where $\text{ran}(F_j) = \{F_j(x) : x \in \mathbb{R}\}$ is the range of F_j . In particular, C is unique if all the margins F_1, \dots, F_p are continuous.

For $F_1 \dots F_p$,
continuous only.
let $X \sim F$, $X \in \mathbb{R}^p$, and define
因为 $V_i \in [0, 1]$, 用 V_i 构建 C
所以由 $X \sim F$ 可取得 $V_i = F_i(x_i)$

proof of sklar's theorem

$\forall x_i \in \mathbb{R} \rightarrow F_i(x_i) \in [0, 1]$
 \therefore satisfy $[0, 1]$ 条件
 $V = (V_1, V_2, \dots, V_p)$ where
 $V_i = F_i(x_i) \in [0, 1]$

$$\begin{aligned} F(x_1, x_2, \dots, x_p) &= P(X_j \leq x_j, \forall j) = P(F_j(x_j) \leq \underline{F_j}(x_j), \forall j) \\ &= P(U_j \leq \underline{F_j}(x_j), \forall j) \end{aligned}$$

we define $U \sim C$ thus $\rightarrow G(F_j(x_j), F_i(x_i), \dots)$

Sklar's theorem

Sklar's theorem

1. Let F be a joint distribution function with margins F_1, \dots, F_p . There exists a copula C such that, for all $x = (x_1, \dots, x_p) \in \mathbb{R}^p$

$$F(x_1, \dots, x_p) = C(F_1(x_1), \dots, F_p(x_p)).$$

C is unique on $\text{ran}(F_1) \times \dots \times \text{ran}(F_p)$, where $\text{ran}(F_j) = \{F_j(x) : x \in \mathbb{R}\}$ is the range of F_j . In particular, C is unique if all the margins F_1, \dots, F_p are continuous.

2. ~~Conversely~~, if C is a copula and F_1, \dots, F_p are univariate distribution functions, then F defined as

$$F(x_1, \dots, x_p) = C(F_1(x_1), \dots, F_p(x_p))$$

is a multivariate distribution function with margins F_1, \dots, F_p .

and copula C

Sklar's theorem

$\rightarrow [0,1]$ $\rightarrow \mathbb{R}^p$

If the margins of $X \sim F$, F_1, \dots, F_p are all continuous, then the (unique) copula of X is

$$C(u_1, \dots, u_p) = F(F_1^{-1}(u_1), \dots, F_p^{-1}(u_p)).$$

Thus, a distribution function F is equivalent to a copula C and its margins F_j . The copula fully characterizes the dependence between the components and we have

\textcircled{X} has copula $C \Leftrightarrow (F_1(X_1), \dots, F_p(X_p)) \sim C.$

$\hookrightarrow X \sim F$

Sklar's theorem

If the margins of $X \sim F$, F_1, \dots, F_p are all continuous, then the (unique) copula of X is

$$C(u_1, \dots, u_p) = F(F_1^{-1}(u_1), \dots, F_p^{-1}(u_p)). \quad (16)$$

Thus, a distribution function F is equivalent to a copula C and its margins F_j . The copula fully characterizes the dependence between the components and we have

$$X \text{ has copula } C \Leftrightarrow (F_1(X_1), \dots, F_p(X_p)) \sim C.$$

不变性

Invariance principle

Let $X \sim F$ with continuous margins F_1, \dots, F_p and copula C . Then for any strictly monotonically increasing T_j on $\text{ran}(X_j)$, we have that $(T_1(X_1), \dots, T_p(X_p))$ has also copula C .

we need the direction unchanged. $\xrightarrow{\text{单调}} \xleftarrow{\text{单偶}}$ $(X_1, \dots, X_p) \sim N(0, \Sigma)$ $(X_1^2, \dots, X_p^2) \sim \chi_i^2$ has same C because

Examples

这也说明了边缘分布与同分布无关.

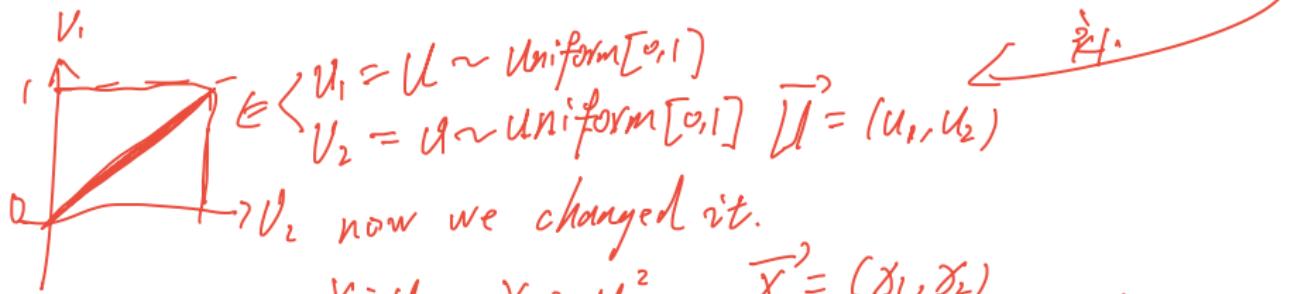
if $F(x_j) = \lambda_j x_j$, $\lambda_j > 0$ 但可能有 same copulas
 $\lambda_j \leq 0$, C is unchanged
 C is changed

Bivariate Bernoulli distribution

Let (X_1, X_2) follow a bivariate Bernoulli distribution with

$$P(X_1 = k, X_2 = l) = 1/4, k, l \in \{0, 1\},$$

then $X_1, X_2 \sim Ber(1/2)$. Any copula with $C(1/2, 1/2) = 1/4$ satisfies the equation in point 1. of Sklar's theorem on $\{0, 1/2, 1\} \times \{0, 1/2, 1\}$, e.g., $C(u_1, u_2) = u_1 u_2$.



$$X_1 = U \quad X_2 = U^2 \quad \vec{X} = (X_1, X_2)$$

$$F_{X_1}(x_1) \sim \text{uniform}[0,1] \quad F_{X_2}(x_2) = P(X_2 \leq x_2) \leftarrow \text{marginal/} \\ \left. \begin{aligned} &= P(U^2 \leq x_2) \\ &= P(U \leq \sqrt{x_2}) = \int_{-\infty}^{\sqrt{x_2}} f_U(u) du \end{aligned} \right\} \because U \sim \text{uniform}[0,1] \text{ 且 } U \text{ 为随机变量}$$

joint distribution $(X_1, X_2) \neq (U_1, U_2)$ 因为 marginal different

$$C_{\vec{U}}(u_1, u_2) = P(U_1 \leq u_1, U_2 \leq u_2)$$

why? 因为 $U_1 = u$
 $U_2 = u$

$$B = \underbrace{P[U \leq u_1, U \leq u_2]}_{= P[U \leq \min(u_1, u_2)]} = \min(u_1, u_2)$$

$$C_{\vec{X}}(x_1, x_2) = P(X_1 \leq F_{X_1}^{-1}(u_1), X_2 \leq F_{X_2}^{-1}(u_2))$$

PS ⑯的定义

$$= P[U \leq u_1, V \leq u_2] \quad \begin{matrix} \leftarrow \\ \text{因为 } X_1 = U \\ X_2 = V \end{matrix}$$

$$A = P[U < u_1, V \leq u_2] \quad \begin{matrix} \leftarrow \\ \text{因为单调上升且} > 0 \end{matrix}$$

A与B一样。

has the same Copulas.

直接去。

$$\begin{aligned} F_{X_1}(u_1) &= u_1 \\ F_{X_2}(u_2) &= u_2^2 \end{aligned}$$

两边
从 $F_X(x)$ 可得。

summary $U \begin{cases} U_1 = u \\ U_2 = u \end{cases} \rightarrow X \begin{cases} X_1 = u \\ X_2 = u^2 \end{cases} \Leftarrow \begin{cases} T_1(X_1) = x_1 \\ T_2(X_2) = x_2^2 \end{cases}$

let me another example:

$$\vec{X} \begin{cases} X_1 = u \\ X_2 = -u \end{cases} \subset \vec{Y}_1, Y_2) = P(X_1 \in F_{Y_1}(u_1), X_2 \in F_{Y_2}(u_2))$$

$$= P(u \in u_1, -u \in u_2^-)$$

$$= P(u \in u_1, u \geq 1-u_2)$$

$$= P(1-u_2 \leq u \leq u_1) = P(u_1) - P(1-u_2)$$

$$F_{Y_1}(Y_1) = u.$$

$$F_{Y_2}(Y_2) : \text{let } y = 1 + Y_2 \Rightarrow Y_2 = y - 1 = u_1 - (1-u_2)$$

$$= u_1 + u_2 - 1$$

$F_{Y_2}(Y_2) = 1 - Y_2$

condition: $u_1 > 1 - u_2$

$$\text{true} = P(Y_2 \leq Y_1)$$

$$= P(-u \leq Y_1)$$

$$= P(u \geq -Y_1)$$

$$= P(1+Y_1 \in (0, 1))$$

$$= 1 - P(-Y_1) = 1 - (-Y_1) = (+Y_1) \begin{cases} 1+Y_1 \in (0, 1) \\ Y_1 \in (-1, 0) \end{cases}$$

Examples

$$\rightarrow F(F_{u_p}^{-1}(U_1) \cdots F_p^{-1}(U_p))$$

Bivariate Bernoulli distribution

Let (X_1, X_2) follow a bivariate Bernoulli distribution with

$$P(X_1 = k, X_2 = l) = 1/4, k, l \in \{0, 1\},$$

then $X_1, X_2 \sim Ber(1/2)$. Any copula with $C(1/2, 1/2) = 1/4$ satisfies the equation in point 1. of Sklar's theorem on $\{0, 1/2, 1\} \times \{0, 1/2, 1\}$, e.g., $C(u_1, u_2) = u_1 u_2$.

Gaussian copula $\rightarrow P=2, R = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ When $P=2$, Gaussian copula only has one parameter.

Let $R \in \mathbb{R}^{P \times P}$ be a correlation matrix and Φ the univariate normal cdf, and Φ_R the p -dimensional cdf of a multivariate normal with mean zero and covariance matrix R , then

$$C_R(u_1, \dots, u_p) = \Phi_R(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_p))$$

is called the Gaussian copula with parameter matrix R .

$$X_1 \sim N(\mu, \Sigma_1)$$

$$X_2 \sim N(\mu, \Sigma_2)$$

have the same Copulas

$$\left| \begin{array}{l} \tilde{X}_1 = \bar{\Sigma}^{-1/2} X_1 \in T_{X_1} \\ \tilde{X}_2 = \bar{\Sigma}^{-1/2} X_2 \in T_{X_2} \end{array} \right| \Rightarrow \tilde{X}_1 = \bar{\Sigma}^{-1/2} X_1 \in T_{X_1} \Rightarrow \text{monotonic/increasing so no changed}$$

Examples

Independence copula

The X_1, \dots, X_p are independent if and ~~only if~~ their copula is $C(u) = \prod_{j=1}^p u_j$.

Examples

Independence copula

The X_1, \dots, X_p are independent if and only if their copula is $C(u) = \prod_{j=1}^p u_j$.

Countermonotonicity copula *反單調性 copula*.

The countermonotonicity copula is the copula of $(U, 1 - U)$ where $U \sim U[0, 1]$. It can be shown that (X_1, X_2) has this copula if and only if they are perfectly negatively dependent, i.e., X_2 is a strictly decreasing function of X_1 .

Examples

Independence copula

The X_1, \dots, X_p are independent if and only if their copula is $C(u) = \prod_{j=1}^p u_j$.

Countermonotonicity copula

The countermonotonicity copula is the copula of $(U, 1 - U)$ where $U \sim U[0, 1]$. It can be shown that (X_1, X_2) has this copula if and only if they are perfectly negatively dependent, i.e., X_2 is a strictly decreasing function of X_1 .

Comonotonicity copula

The comonotonicity copula is the copula of (U, \dots, U) where $U \sim U[0, 1]$. It can be shown that (X_1, \dots, X_p) has this copula if and only if the components are perfectly positively dependent, i.e., X_2, \dots, X_p are strictly increasing functions of X_1 .

Examples

Independence copula

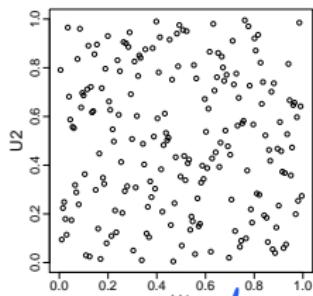
The X_1, \dots, X_p are independent if and only if their copula is $C(u) = \prod_{j=1}^p u_j$.

Countermonotonicity copula

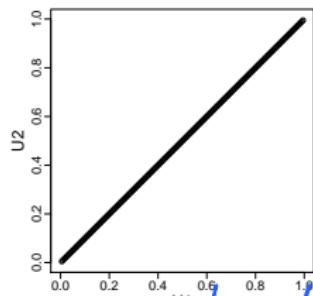
The countermonotonicity copula is the copula of $(U, 1 - U)$ where $U \sim U[0, 1]$. It can be shown that (X_1, X_2) has this copula if and only if they are perfectly negatively dependent, i.e., X_2 is a strictly decreasing function of X_1 .

Comonotonicity copula 共同性 copula.

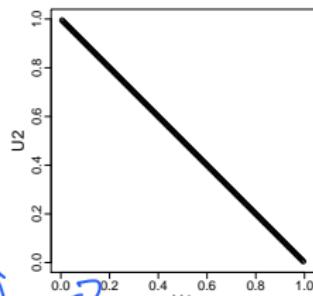
The comonotonicity copula is the copula of (U, \dots, U) where $U \sim U[0, 1]$. It can be shown that (X_1, \dots, X_p) has this copula if and only if the components are perfectly positively dependent, i.e., X_2, \dots, X_p are a strictly increasing functions of X_1 .



independent



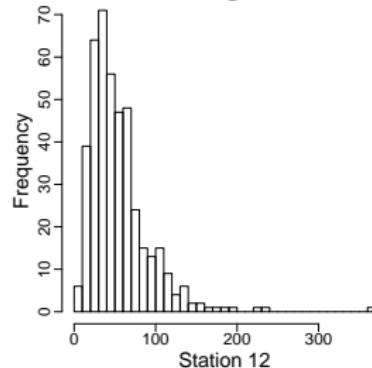
strong dependent
comonotonicity



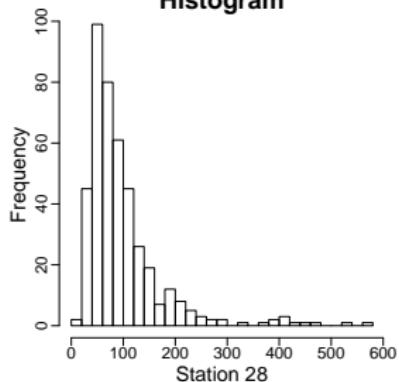
countermonotonicity

River discharge data

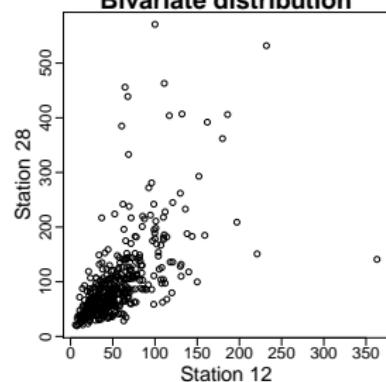
Histogram



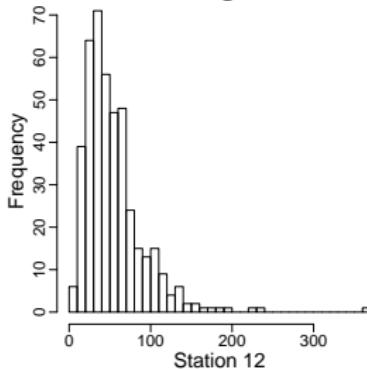
Histogram



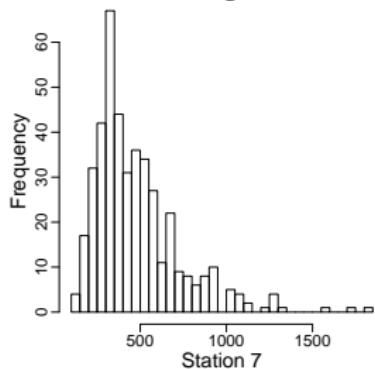
Bivariate distribution



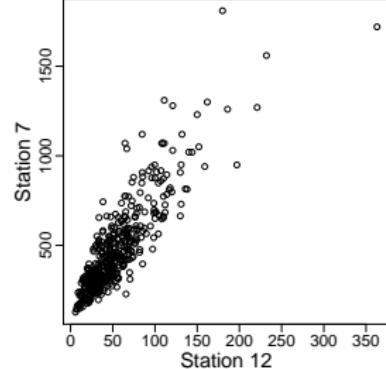
Histogram



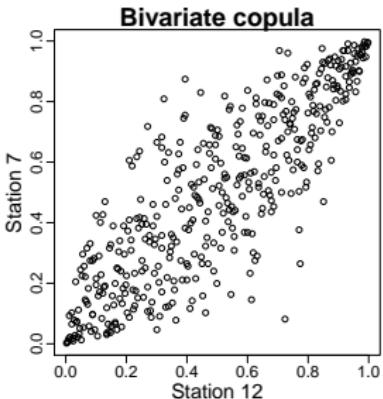
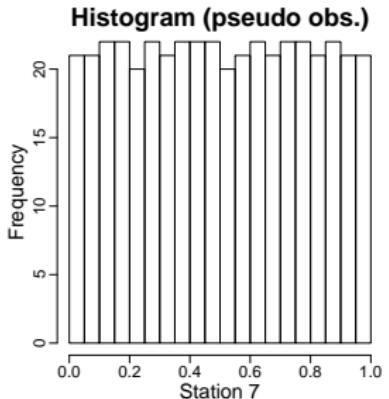
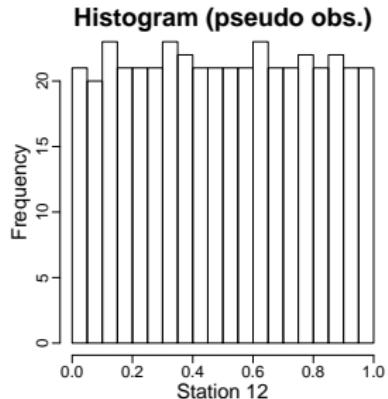
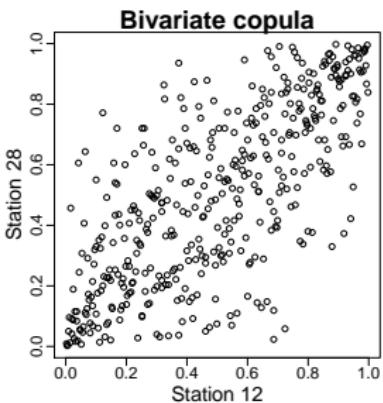
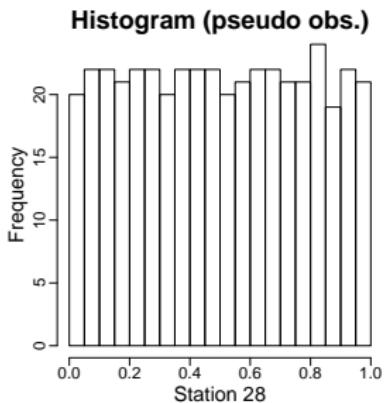
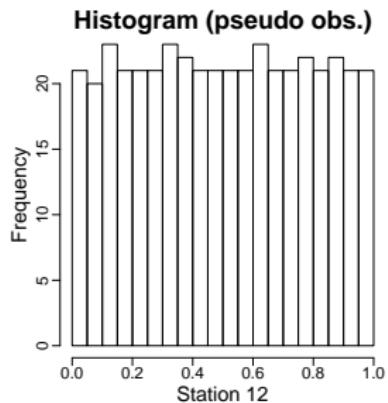
Histogram



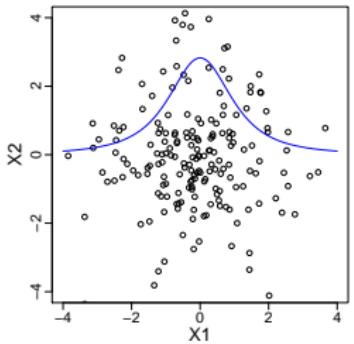
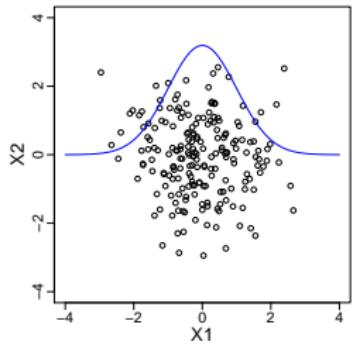
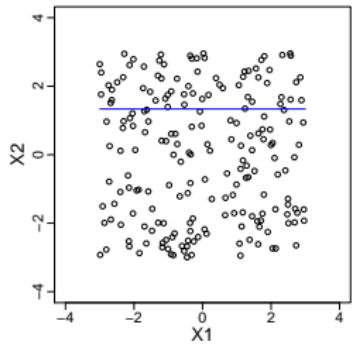
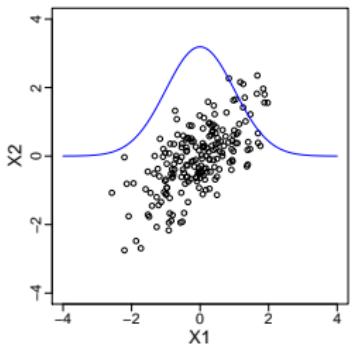
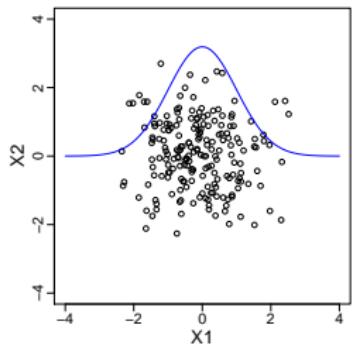
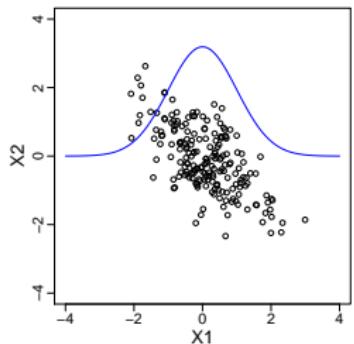
Bivariate distribution



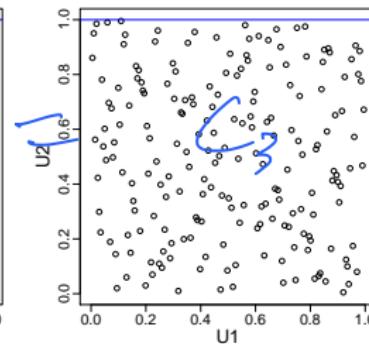
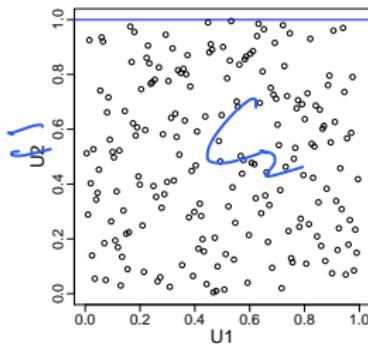
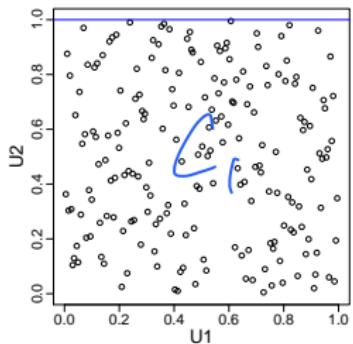
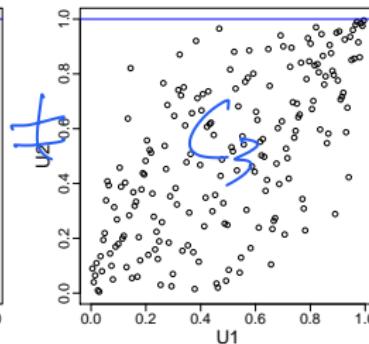
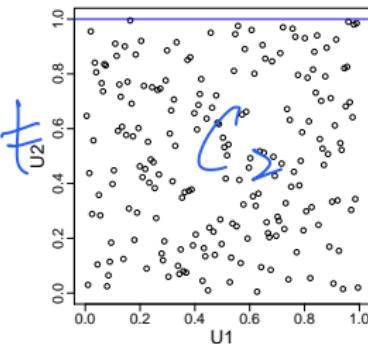
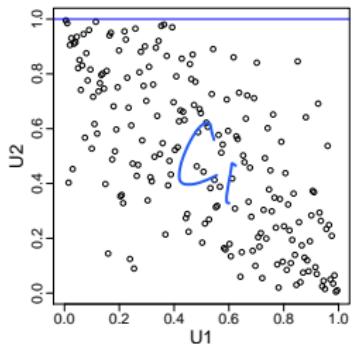
River discharge data



Simulated data



Simulated data



Fréchet–Höffding bounds

Theorem (Fréchet–Höffding bounds)

Let $W(u) = \max\{\sum_{j=1}^p u_j - p + 1, 0\}$ and $M(u) = \min_{j=1,\dots,p} u_j$.

(1) For any p -dimensional copula C

$$W(u) \leq C(u) \leq M(u), \quad u \in [0, 1]^p.$$

(2) W is a copula if and only if $p = 2$.

(3) M is a copula for all $p \geq 2$

Note:

It is easy to see that for $U \sim U[0, 1]$

► $(U, \dots, U) \sim M$ for any p .

► $(U, 1 - U) \sim W$ for $p = 2$.

Proof of (2) :

$$\begin{aligned} W(1, \dots, u_j - 1) &= \max\{p-1+u_j-1-p, 0\} \\ &= \max\{u_j, 0\} \end{aligned}$$

M we also don't show but the key problem is that we can construct it. but we can construct it. W is only function We not show it as joint CDT

$u_j \in [0, 1]$
all marginal distribution uniform.

Densities of copulas

- Let (X_1, \dots, X_p) have differentiable cdf F with joint and marginal densities given by f and f_1, \dots, f_p , respectively. Then the corresponding copula has density

$$c(u_1, \dots, u_p) = \frac{\partial^p C(u_1, \dots, u_p)}{\partial u_1 \dots \partial u_p}$$

$$C(u) = F(F_1^{-1}(u_1) \dots F_p^{-1}(u_p)) = \frac{f(F_1^{-1}(u_1), \dots, F_p^{-1}(u_p))}{\prod_{j=1}^p f_j(F_j^{-1}(u_j))} \quad \text{chain rule}$$

Example:

The Gaussian copula has density

$$\frac{F_j(u_j)}{f_j(F_j(u_j))} = \frac{1}{f_j(F_j(u_j))}$$

$$c_{\mathbf{R}}(u_1, \dots, u_p) = \frac{1}{\sqrt{\det \mathbf{R}}} \exp\left(-\frac{1}{2} x^T (\mathbf{R}^{-1} - I_p) x\right),$$

where $x = (\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_p))$.

shortcoming: it's for ellipse

$$X \sim N(0, R)$$

$$f_X(x) = \frac{1}{\sqrt{2\pi} \cdot \sqrt{\det(R)^T}} \cdot \exp\left(-\frac{1}{2} x^T R^{-1} x\right)$$

$$f_1(x_1) \cdots f_p(x_p) = \frac{1}{(2\pi)^{p/2}} \cdot \exp\left(-\frac{1}{2} X^T X\right)$$

分子

$$= \frac{1}{(2\pi)^{p/2}} \exp\left(-\frac{1}{2} X^T I_p X\right)$$

Archimedean copulas



why ψ cannot use function in here $\Rightarrow C: [0, 1]^P \rightarrow [0, 1]$

$\Rightarrow \psi^{-1}: [0, 1] \rightarrow [0, \infty)$

$$C(u_1, \dots, u_p) = \psi(\psi^{-1}(u_1) + \dots + \psi^{-1}(u_p)) = \psi(\psi(u_1) + \dots + \psi(u_p)) = u_i$$

Archimedean copulas are copulas of the form

$$C(u) = \psi(\psi^{-1}(u_1) + \dots + \psi^{-1}(u_p)), \quad G[0, 1]$$

where ψ is called the (Archimedean) generator, where $\psi: \mathbb{R} \rightarrow \mathbb{R}$

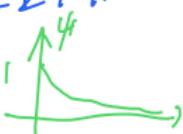
A $\psi: [0, \infty) \rightarrow [0, 1]$ is strictly decreasing on its support and it satisfies $\psi(0) = 1$ and $\psi(\infty) = \lim_{t \rightarrow \infty} \psi(t) = 0$.

There are further conditions to be satisfied in order that ψ generates a proper copula (e.g., complete monotonicity).

PS: Copulas $\xrightarrow{\text{joint distribution (i) 且证明中未提及}}$ marginal distribution are uniform (ii)

这里不满足 (iii)

有兴趣自己看
太 technical 了



我们要求 $u_i \in [0, 1]$

所以有了 A

Archimedean copulas: Gumbel and Clayton

Taking particular (admissible) functions ψ we can construct particular examples of Archimedean copulas.

Gumbel copula

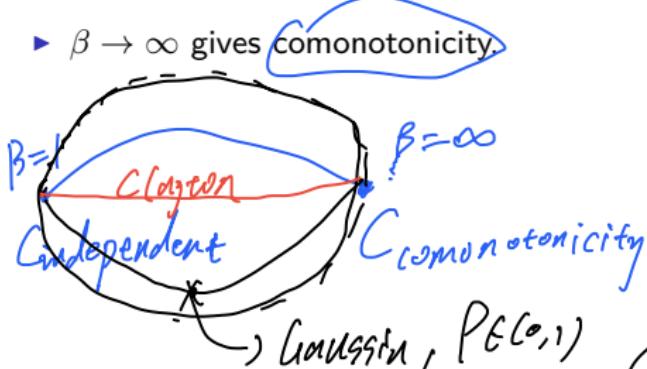
- Let $\psi(t) = \exp(-t^{1/\beta})$, for $\beta \geq 1$, then

$$C_{\beta}^{Gu}(u_1, \dots, u_p) = \exp \left(- \left((-\log u_1)^{\beta} + \dots + (-\log u_p)^{\beta} \right)^{1/\beta} \right).$$

- $\beta = 1$ gives independence since

$$C_{\beta}^{Gu}(u_1, \dots, u_p) = \exp(\log u_1 \cdots + \log u_p) = \prod_{j=1}^p u_j$$

- $\beta \rightarrow \infty$ gives comonotonicity



but we cannot get all the copulas. from independence to dependence !!!
may all copulas is like O , we only get blue by Gumbel

Archimedean copulas: Gumbel and Clayton

Taking particular (admissible) functions ψ we can construct particular examples of Archimedean copulas.

Gumbel copula

- Let $\psi(t) = \exp(-t^{1/\beta})$, for $\beta \geq 1$, then

$$C_{\beta}^{Gu}(u_1, \dots, u_p) = \exp \left(- \left((-\log u_1)^{\beta} + \dots + (-\log u_p)^{\beta} \right)^{1/\beta} \right).$$

- $\beta = 1$ gives independence since

$$C_{\beta}^{Gu}(u_1, \dots, u_p) = \exp(\log u_1 + \dots + \log u_p) = \prod_{j=1}^p u_j$$

- $\beta \rightarrow \infty$ gives comonotonicity.

Clayton copula

- Let $\psi(t) = (1+t)^{-1/\beta}$, for $\beta \in (0, \infty)$, then

$$C_{\beta}^{Cl}(u_1, \dots, u_p) = \left(u_1^{-\beta} + \dots + u_p^{-\beta} - p + 1 \right)^{-1/\beta}.$$

- $\beta = 0$ gives independence and $\beta \rightarrow \infty$ gives comonotonicity.

Archimedean copulas: properties

Advantages

- ▶ Typically they are explicit (if ψ^{-1} is available). ✓
- ▶ Properties can be expressed in terms of the generator ψ .
- ▶ Densities are available in many examples.
- ▶ Simulation is often simple.
- ▶ Flexible: not restricted to radial symmetry.

Drawbacks

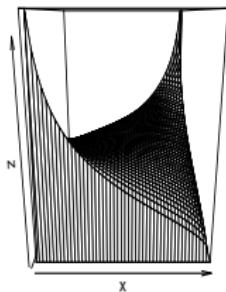
- ▶ All d -dimensional margins, $d < p$, have the same distribution.
- ▶ Often used only with few parameters, much less than $p(p - 1)/2$ as in the Gaussian case.

边缘都一样, 但高维时
以次不一样

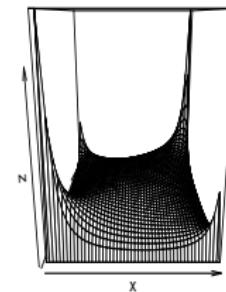
参数的
数量

density of different copulas

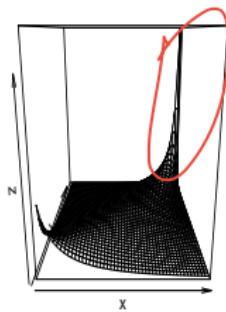
Gaussian, $r = 0.3$



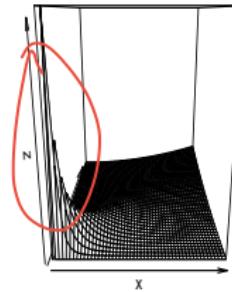
Student, $r = 0.3$, df = 2



Gumbel, beta = 2

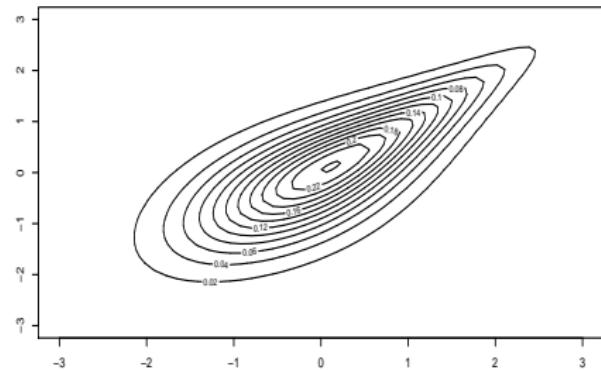


Clayton, beta = 2

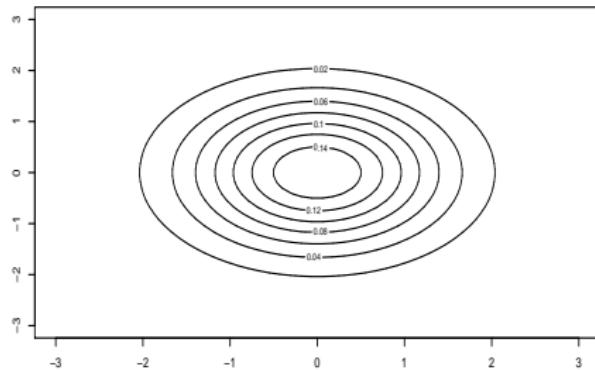


X Y 技巧

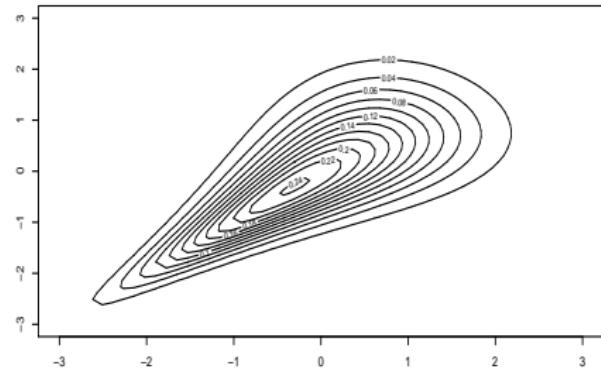
Gumbel on Normal marginal, beta = 2



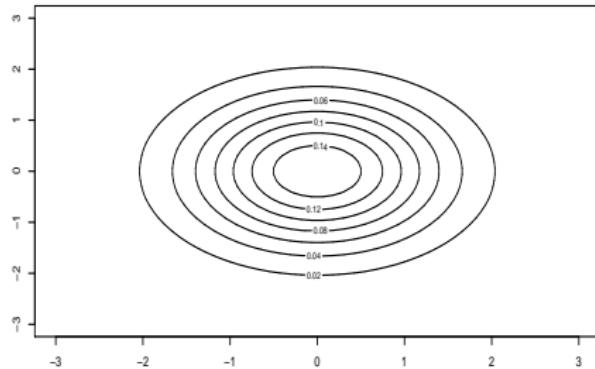
Gumbel on Normal marginal, beta = 1



Clayton on Normal marginal, beta = 2



Clayton on Normal marginal, beta = 0.00001



Statistics with copulas

- ▶ Suppose we have data $X_1, \dots, X_n \in \mathbb{R}^p$ from a distribution F with copula C .
- ▶ In order to apply the copula framework, we first need to model the univariate marginal distributions with estimators \hat{F}_j of F_j , $j = 1, \dots, p$.
- ▶ We can then introduce pseudo-observations

$$\hat{U}_i = (\hat{U}_{i1}, \dots, \hat{U}_{ip}) = (\hat{F}_1(X_{i1}), \dots, \hat{F}_p(X_{ip})) \quad i = 1, \dots, n.$$

- ▶ Note that the \hat{U}_{ij} are only approximately $U[0, 1]$ distributed.

因为是估计所以是 approximately.
所以问题转化成如何估计 \hat{F}_j from data.

Statistics with copulas

- ▶ Suppose we have data $X_1, \dots, X_n \in \mathbb{R}^p$ from a distribution F with copula C .
- ▶ In order to apply the copula framework, we first need to model the univariate marginal distributions with estimators \hat{F}_j of F_j , $j = 1, \dots, p$.
- ▶ We can then introduce **pseudo-observations**

$$\hat{U}_i = (\hat{U}_{i1}, \dots, \hat{U}_{ip}) = (\hat{F}_1(X_{i1}), \dots, \hat{F}_p(X_{ip})) \quad i = 1, \dots, n.$$

- ▶ Note that the \hat{U}_{ij} are only approximately $U[0, 1]$ distributed.

Possible choices for \hat{F}_j :

- ▶ Parametric estimator using a parametric family, e.g. Gaussian or t -distribution.

Statistics with copulas

- ▶ Suppose we have data $X_1, \dots, X_n \in \mathbb{R}^p$ from a distribution F with copula C .
- ▶ In order to apply the copula framework, we first need to model the univariate marginal distributions with estimators \hat{F}_j of F_j , $j = 1, \dots, p$.
- ▶ We can then introduce **pseudo-observations**

$$\hat{U}_i = (\hat{U}_{i1}, \dots, \hat{U}_{ip}) = (\hat{F}_1(X_{i1}), \dots, \hat{F}_p(X_{ip})) \quad i = 1, \dots, n.$$

- ▶ Note that the \hat{U}_{ij} are only approximately $U[0, 1]$ distributed.

Possible choices for \hat{F}_j :

- ▶ Parametric estimator using a parametric family, e.g. Gaussian or t -distribution.
- ▶ Semi-parametric estimators, e.g. based on extreme value theory where the bulk of the distribution is modeled empirically and the tail is approximated by a so-called generalized Pareto distribution.

Statistics with copulas

- ▶ Suppose we have data $X_1, \dots, X_n \in \mathbb{R}^p$ from a distribution F with copula C .
- ▶ In order to apply the copula framework, we first need to model the univariate marginal distributions with estimators \hat{F}_j of F_j , $j = 1, \dots, p$.
- ▶ We can then introduce **pseudo-observations**

$$\hat{U}_i = (\hat{U}_{i1}, \dots, \hat{U}_{ip}) = (\hat{F}_1(X_{i1}), \dots, \hat{F}_p(X_{ip}))$$

$$i = 1, 2, \dots, n$$

- ▶ Note that the \hat{U}_{ij} are only approximately $U[0, 1]$ distributed.

$$\hat{F}_j(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{X_{ij} \leq x\}}$$

Possible choices for \hat{F}_j :

- ▶ Parametric estimator using a parametric family, e.g. Gaussian or t -distribution.
- ▶ Semi-parametric estimators, e.g. based on extreme value theory where the bulk of the distribution is modeled empirically and the tail is approximated by a so-called generalized Pareto distribution.
- ▶ Non-parametric estimators using the empirical cdf

$$\hat{U}_{ij} = \frac{n}{n+1} \hat{F}_{n,j}(X_{ij}) = \frac{R_{ij}}{n+1},$$

where R_{ij} is the rank of X_{ij} among all X_{1j}, \dots, X_{nj} .

要排序的！排序后的①, ②, ③进行④

MC

蒙特卡洛仿真

Statistics with copulas

Estimation

- The log-likelihood of the data $X_1, \dots, X_n \in \mathbb{R}^p$ from a distribution F_θ with

copula C_{θ_C} is

log likelihood \rightarrow for the dependence

$$\ell(\theta; X_1, \dots, X_n) = \sum_{i=1}^n \ell(\theta, X_i) = \sum_{i=1}^n \log[f(x_i; \theta)]$$

但是 $X_1 \cdots X_n$, 不是 $\ell(\theta)$

$$= \sum_{i=1}^n \ell_C(\theta_C; F_1(X_{i1}, \theta_1), \dots, F_p(X_{ip}, \theta_p)) + \sum_{i=1}^n \sum_{j=1}^p \ell_j(\theta_j; X_{ij}),$$

for Copulas

where

$$\ell_C(\theta_C; u_1, \dots, u_p) = \log c(u_1, \dots, u_p; \theta_C)$$

$$\ell_j(\theta_j; x) = \log f_j(x; \theta_j), \quad j = 1, \dots, p.$$

and the θ_C are the copula parameters and the θ_j the marginal parameters.

$$\theta = (\theta_1, \theta_p, \theta_C)$$

not one just

θ contain Dependence

and marginal distributions

$$\sum_{i=1}^n \sum_{j=1}^p \ell_j(\theta_j; X_{ij}),$$

for marginal

Statistics with copulas

Estimation

- The log-likelihood of the data $X_1, \dots, X_n \in \mathbb{R}^p$ from a distribution F_θ with copula C_{θ_C} is

$$\begin{aligned}\ell(\theta; X_1, \dots, X_n) &= \sum_{i=1}^n \ell(\theta, X_i) \\ &= \sum_{i=1}^n \ell_C(\theta_C; F_1(X_{i1}, \theta_1), \dots, F_p(X_{ip}, \theta_p)) + \sum_{i=1}^n \sum_{j=1}^p \ell_j(\theta_j; X_{ij}),\end{aligned}$$

where

$$\begin{aligned}\ell_C(\theta_C; u_1, \dots, u_p) &= \log c(u_1, \dots, u_p; \theta_C) \\ \ell_j(\theta_j; x) &= \log f_j(x; \theta_j), \quad j = 1, \dots, p.\end{aligned}$$

and the θ_C are the copula parameters and the θ_j the marginal parameters.

This suggests a **two-step estimation** procedure:

- First estimate the marginal parameters θ_j separately (or non-parameterically);
- and then use the pseudo-observations to estimate θ_C by maximizing

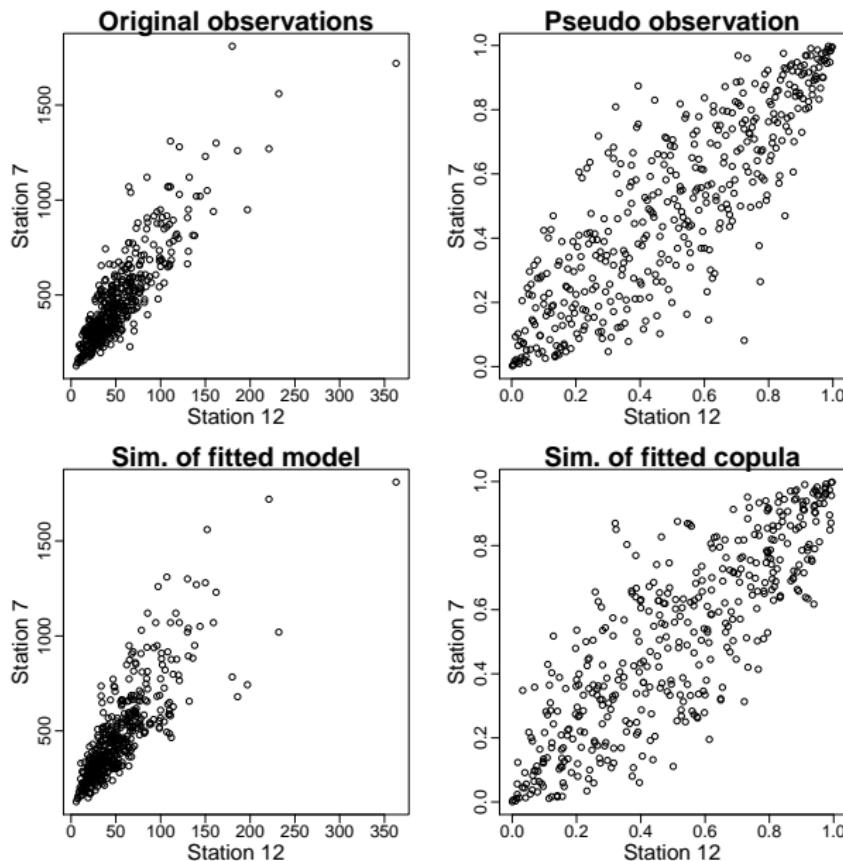
$$\sum_{i=1}^n \ell_C(\theta_C; \widehat{U}_{i1}, \dots, \widehat{U}_{ip})$$

River discharge data: fitting in R

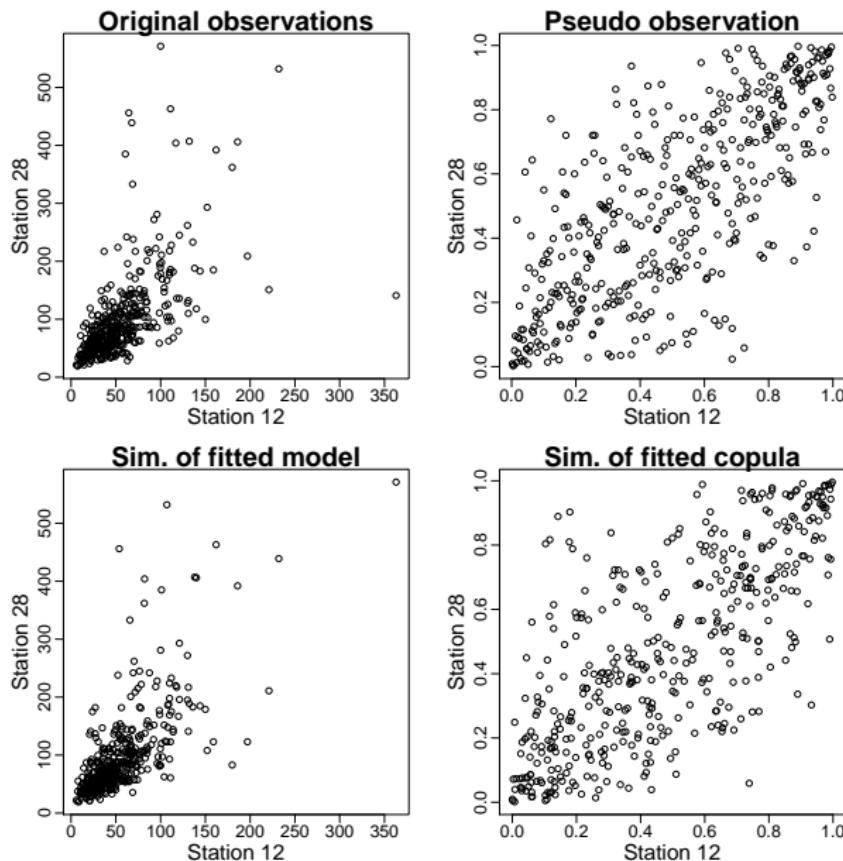
We transform the data non-parameterically and then fit a Gaussian copula. We can simulate from the copula and on the original scale.

```
> library(copula)
> dim(X)
[1] 428 31      n= 428    P = 31
> U <- pobs(X)
> i <- 12
> j <- 7
> plot(X[,c(i,j)])
> plot(U[,c(i,j)])
> cop.fit <- fitCopula(normalCopula(), U[,c(i,j)], method="ml")
> cop.fit
fitCopula() estimation based on 'maximum likelihood'
and a sample of size 428.
      Estimate Std. Error z value Pr(>|z|)
rho.1  0.863887   0.009297  92.92  <2e-16 ***
> Usim <- rCopula(n = nrow(U), copula = cop.fit@copula)
> Xsim <- cbind(sort(X[,i])[rank(Usim[,1])], sort(X[,j])[rank(Usim[,2])])
> plot(Xsim)
> plot(Usim)
```

River discharge data



River discharge data



Measures of dependence: linear correlation

- ▶ Measures of dependence summarize a complicated dependence structure in a single number (in the bivariate case) or a matrix.
- ▶ For two random variables X_1 and X_2 with $E(X_j^2) < \infty$ the (Pearson) correlation coefficient is defined as

$$\rho(X_1, X_2) = \frac{\text{cov}(X_1, X_2)}{\sqrt{\text{Var}X_1}\sqrt{\text{Var}X_2}}.$$

disadv : only linear

adv : put data in original scale

ρ is covariance for distribu

Properties of Pearson correlation

Some properties:

- ▶ $|\rho(X_1, X_2)| \leq 1$. Moreover, $|\rho(X_1, X_2)| = 1$ if and only if $X_2 = aX_1 + b$ almost surely for constants $a \neq 0$, $b \in \mathbb{R}$. This discards other strong functional dependence such as $X_2 = X_1^2$ for instance. 棄却
3.36M
- ▶ X_1, X_2 independent $\Rightarrow \rho(X_1, X_2) = 0$. But $\rho(X_1, X_2) = 0 \not\Rightarrow X_1, X_2$ independent except for the normal distribution.
- ▶ “Marginal distributions A and correlation B determine the joint distribution”C: This is true for the class of multivariate normal distributions.
- ▶ $\rho(X_1, X_2)$ is invariant under strictly increasing linear transformations, but not under strictly increasing functions in general.
- ▶ For a class of elliptical distributions (e.g., normal or t), the correlation matrix together with the marginal distributions determine the whole distribution.

Alternative measures based on the ranks: rank correlations.

Spearman's rho

- ▶ Let $X = (X_1, X_2)$ be a random vector with distribution function F and copula C .
- ▶ The Spearman's rho between X_1 and X_2 is

$$\rho_S(X_1, X_2) = \rho(F_1(X_1), F_2(X_2)),$$

where ρ is the Pearson's correlation. *用 copular scale 代替*

Spearman's rho

- Let $X = (X_1, X_2)$ be a random vector with distribution function F and copula C .
- The Spearman's rho between X_1 and X_2 is

$$\rho_S(X_1, X_2) = \rho(F_1(X_1), F_2(X_2)),$$

where ρ is the Pearson's correlation.

- One can show that this is equivalent to

exercice I you can see it

$$\rho_S(X_1, X_2) = 12 \int_0^1 \int_0^1 \{C(u_1, u_2) - u_1 u_2\} du_1 du_2.$$

- An estimator for ρ_S is given by the sample correlation on the ranks of the data.

Properties of Spearman's rho

- ▶ ρ_S is invariant under any strictly increasing transformations of the random variables.
- ▶ $\rho_S(X_1, X_2) = 1 \iff X_1, X_2$ comonotonic. 同
- ▶ $\rho_S(X_1, X_2) = -1 \iff X_1, X_2$ countermonotonic. 反

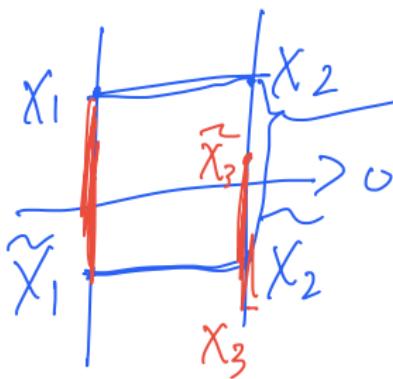
Kendall's tau \rightarrow 2 random variables.

- Take an independent copy of (X_1, X_2) denoted $(\tilde{X}_1, \tilde{X}_2)$.

- The Kendall's tau is defined as the probability of concordance

$\xrightarrow{\text{Eq 1.8}}$

$$\rho_\tau(X_1, X_2) = 2P\left((X_1 - \tilde{X}_1)(X_2 - \tilde{X}_2) > 0\right) - 1 \in [-1, 1]$$



$$(X_1 - \tilde{X}_1)(X_2 - \tilde{X}_2) > 0$$

$$(X_1 - \tilde{X}_1)(X_3 - \tilde{X}_3) < 0$$

> 0

< 0

Kendall's tau

- ▶ Take an independent copy of (X_1, X_2) denoted $(\tilde{X}_1, \tilde{X}_2)$.
- ▶ The Kendall's tau is defined as the probability of concordance

$$\rho_\tau(X_1, X_2) = 2P\left((X_1 - \tilde{X}_1)(X_2 - \tilde{X}_2) > 0\right) - 1 \in [-1, 1]$$

- ▶ Once can show that this is equivalent to

$$\rho_\tau(X_1, X_2) = 4 \int_0^1 \int_0^1 C(u_1, u_2) dC(u_1, u_2) - 1.$$

- ▶ The properties of Kendall's tau are similar to those of Spearman's rho.

Kendall's tau

- ▶ Take an independent copy of (X_1, X_2) denoted $(\tilde{X}_1, \tilde{X}_2)$.
- ▶ The Kendall's tau is defined as the probability of concordance

$$\rho_\tau(X_1, X_2) = 2P\left((X_1 - \tilde{X}_1)(X_2 - \tilde{X}_2) > 0\right) - 1 \in [-1, 1]$$

- ▶ Once can show that this is equivalent to

$$\rho_\tau(X_1, X_2) = 4 \int_0^1 \int_0^1 C(u_1, u_2) dC(u_1, u_2) - 1.$$

- ▶ The properties of Kendall's tau are similar to those of Spearman's rho.
- ▶ A sample version of ρ_τ is given by

$$\hat{\rho}_\tau(X_1, X_2) = \frac{1}{\binom{n}{2}} \sum_{1 \leq i < j \leq n} \text{sign} [(X_{i,1} - X_{j,1})(X_{i,2} - X_{j,2})].$$

Kendall's tau

- ▶ Take an independent copy of (X_1, X_2) denoted $(\tilde{X}_1, \tilde{X}_2)$.
- ▶ The Kendall's tau is defined as the probability of concordance

$$\rho_\tau(X_1, X_2) = 2P\left((X_1 - \tilde{X}_1)(X_2 - \tilde{X}_2) > 0\right) - 1 \in [-1, 1]$$

- ▶ One can show that this is equivalent to

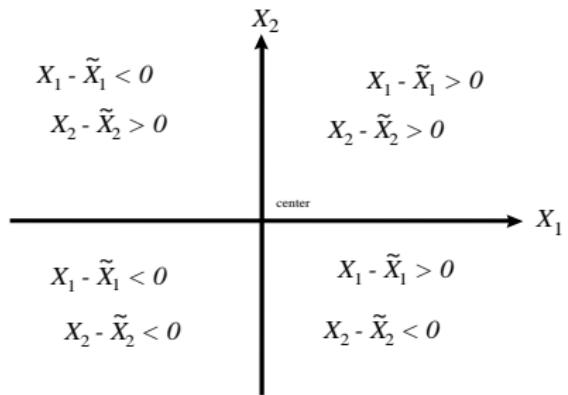
$$\rho_\tau(X_1, X_2) = 4 \int_0^1 \int_0^1 C(u_1, u_2) dC(u_1, u_2) - 1.$$

- ▶ The properties of Kendall's tau are similar to those of Spearman's rho.
- ▶ A sample version of ρ_τ is given by

$$\hat{\rho}_\tau(X_1, X_2) = \frac{1}{\binom{n}{2}} \sum_{1 \leq i < j \leq n} \text{sign} [(X_{i,1} - X_{j,1})(X_{i,2} - X_{j,2})].$$

- ▶ All dependence measures ρ , ρ_S and ρ_τ can be used to estimate the copula parameters θ_C of a parametric family of copulas by comparing the empirical to the model versions.

Illustration



Tail dependence



copulas dependence

\Rightarrow now we only want estimate
part of this

- Another important measure is given by the tail dependence or extremal dependence (bivariate case).
- When the limit exists, the coefficient of upper tail dependence is

if $\lambda_u(X_1, X_2)$ is positive

$$\lambda_u(X_1, X_2) = \lim_{q \rightarrow 1} P(X_2 > F_2^{-1}(q) | X_1 > F_1^{-1}(q)),$$

we get always some points

- Analogously the coefficient of lower tail dependence is

$$\lambda_l(X_1, X_2) = \lim_{q \rightarrow 0} P(X_2 \leq F_2^{-1}(q) | X_1 \leq F_1^{-1}(q)).$$

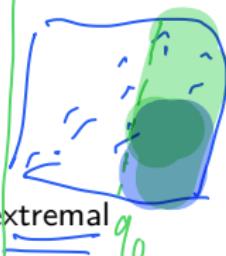
$$X_2 > q_0$$

only look this

$$F_1(X_1) > q_0$$

but we only
need green part

$$so \cdot a$$



Tail dependence

- ▶ Another important measure is given by the tail dependence or extremal dependence (bivariate case).
- ▶ When the limit exists, the coefficient of upper tail dependence is

$$\lambda_u(X_1, X_2) = \lim_{q \rightarrow 1} P(X_2 > F_2^{-1}(q) \mid X_1 > F_1^{-1}(q)),$$

- ▶ Analogously the coefficient of lower tail dependence is

$$\lambda_\ell(X_1, X_2) = \lim_{q \rightarrow 0} P(X_2 \leq F_2^{-1}(q) \mid X_1 \leq F_1^{-1}(q)).$$

- ▶ When $\lambda_u > 0$ or $\lambda_\ell > 0$ we call (X_1, X_2) upper or lower tail dependent, respectively.
- ▶ If $\lambda_u = 0$ or $\lambda_\ell = 0$ we call (X_1, X_2) upper or lower tail independent, respectively.

Copulas and tail dependence

- ▶ For a bivariate copula $U \sim C$ let the survival copula be $\bar{C}(u, v) = P(U_1 > u, U_2 > v)$.
- ▶ Using the corresponding copulas, we have

$$\begin{aligned}\lambda_u &= \lim_{q \rightarrow 1} \frac{\bar{C}(q, q)}{1 - q} = \lim_{q \rightarrow 1} \frac{1 - 2q + C(q, q)}{1 - q}, \\ \lambda_\ell &= \lim_{q \rightarrow 0} \frac{C(q, q)}{q}.\end{aligned}$$

- ▶ Note that $\lambda_u \in [0, 1]$ and $\lambda_\ell \in [0, 1]$.
- ▶ Basically, we need to look at the slope of the copula when approaching $(0, 0)$ or $(1, 1)$ along the diagonal.

Example: Clayton copula

- We have

$$C_{\beta}^{CI}(q, q) = (2q^{-\beta} - 1)^{-1/\beta}.$$

$$\begin{aligned}\lambda_{\ell} &= \lim_{q \rightarrow 0} \frac{C(q, q)}{q} = \lim_{q \rightarrow 0} \frac{(2q^{-\beta} - 1)^{-1/\beta}}{q} \\ &= \lim_{q \rightarrow 0} (q^{\beta})^{-1/\beta} (2q^{-\beta} - 1)^{-1/\beta} \\ &= \lim_{q \rightarrow 0} (2 - q^{\beta})^{-1/\beta} = 2^{-1/\beta}\end{aligned}$$

- The Clayton copula is lower tail dependent for $\beta > 0$.

Further properties

- For elliptical copulas $\lambda_u = \lambda_\ell =: \lambda$.
- The t copula is tail dependent when $\rho > -1$.

$$\lambda = 2\bar{t}_{\nu+1} \left(\sqrt{\nu+1} \sqrt{1-\rho} / \sqrt{1+\rho} \right).$$

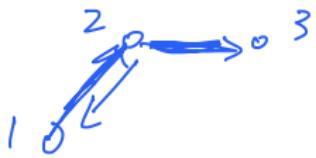
- The Gumbel copula is upper tail dependent for $\beta > 1$.

$$\lambda_u = 2 - 2^{1/\beta}.$$

- The Gaussian copula is upper and lower tail independent, that is, $\lambda_u = \lambda_\ell = 0$.

(1) in off no different model influence different

Graphs: notation and terminology



$$V = \{1, 2, 3\}$$

$$E = \{(1, 2), (2, 3)\}$$

- ▶ A graph $G = (V, E)$ is a pair with

- ▶ a set of nodes/vertices $V = \{1, \dots, p\}$;
- ▶ a set of edges $E \subset V \times V$ between the nodes.

~~节点~~ P nodes
~~连线~~ connected

- ▶ Edges with $(u, v) \in E$ and $(v, u) \in E$ are called undirected and we write $u \sim v$.

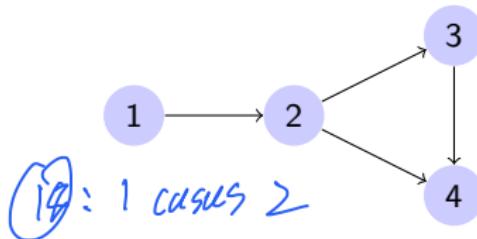


Graphs: notation and terminology

- ▶ A graph $G = (V, E)$ is a pair with
 - ▶ a set of nodes/vertices $V = \{1, \dots, p\}$;
 - ▶ a set of edges $E \subset V \times V$ between the nodes.
- ▶ Edges with $(u, v) \in E$ and $(v, u) \notin E$ are called undirected and we write $u \sim v$.
- ▶ Edges with $(u, v) \in E$ but $(v, u) \notin E$ are called directed and we write $u \rightarrow v$.
- ▶ Directed graphs only have directed edges. They are important in the field of causal inference; e.g., Peters et al. (Elements of Causal Inference, 2017) (18)

Example for directed graph: $V = \{1, \dots, 4\}$, $E = ??$.

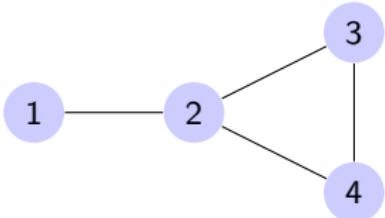
$$E = \{(1, 2), (2, 3), (3, 4), (2, 4)\}$$



Undirected graphs

- In this course we will only consider **undirected graphs**, that is, all edges are undirected.
- To simplify notation, for $u \sim v$ we write either $(u, v) \in E$ or $(v, u) \in E$, but always assume that both directions are in the edge set.
- Two nodes $u \sim v$ are called adjacent or neighbors.
- See Lauritzen (*Graphical Models*, 1996) for details.

Example for undirected graph: $V = \{1, \dots, 4\}$, $E = ??$.



Terminology

- ▶ A subset $A \subset V$ of vertices is called **complete**, if all vertices are adjacent.

Terminology

vertices $\{1-2\}$ $\{3\}$ $A=1, 2$, complete
 $\{1-2\}$ $\{3\}$ $\{4\}$ $A=1, 2, 3$, ~~complete~~

- ▶ A subset $A \subset V$ of vertices is called **complete**, if all vertices are adjacent.
- ▶ A complete subset that is maximal is called a **clique**.

$A = \{3, 4\}$ complete
but not clique

$A = \{2, 3, 4\}$ clique
complete

we have

$$C_1 \cap C_2 = \emptyset$$

(1-2)

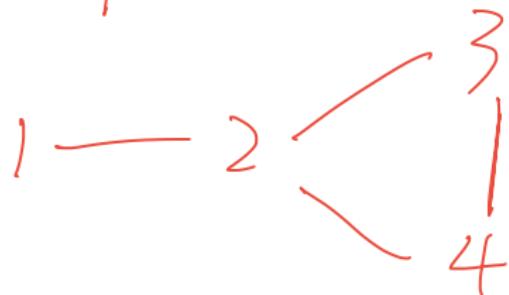
2 clique
(2-3-4)

Terminology

- ▶ A subset $A \subset V$ of vertices is called **complete**, if all vertices are adjacent.
- ▶ A complete subset that is maximal is called a **clique**.
- ▶ A **path** from node u to node v of length n is a sequence $u = w_0, w_1 \dots, w_n = v$, where w_1, \dots, w_n are distinct nodes and $(w_{i-1}, w_i) \in E$ for all $i = 1, \dots, n$.

path from 3 to 1

{ length 2 w_3 w_2 w_1
| length 3 w_3 w_4 w_2 w_1



Terminology

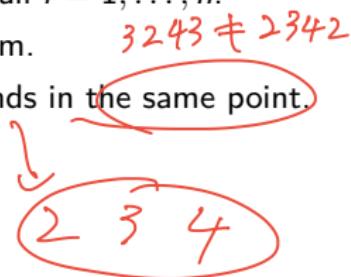
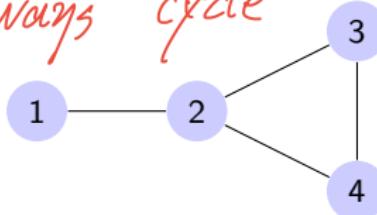
- ▶ A subset $A \subset V$ of vertices is called **complete**, if all vertices are adjacent.
- ▶ A complete subset that is maximal is called a **clique**.
- ▶ A **path** from node u to node v of length n is a sequence $u = w_0, w_1 \dots, w_n = v$, where w_1, \dots, w_n are distinct nodes and $(w_{i-1}, w_i) \in E$ for all $i = 1, \dots, n$.
- ▶ Two nodes are connected if there exists a path between them.

Terminology

- ▶ A subset $A \subset V$ of vertices is called **complete**, if all vertices are adjacent.
- ▶ A complete subset that is maximal is called a **clique**.
- ▶ A **path** from node u to node v of length n is a sequence $u = w_0, w_1 \dots, w_n = v$, where w_1, \dots, w_n are distinct nodes and $(w_{i-1}, w_i) \in E$ for all $i = 1, \dots, n$.
- ▶ Two nodes are connected if there exists a path between them.
- ▶ A cycle is a path with $u = v$, i.e., a path that begins and ends in the same point.

Example: Find the complete subsets, cliques and cycles.

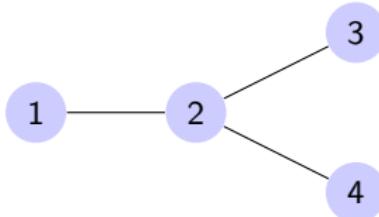
clique is not always cycle



Examples for graphs

- ▶ A **tree** is a connected, undirected graph without cycles, or equivalently, there exists a unique path between any two nodes.

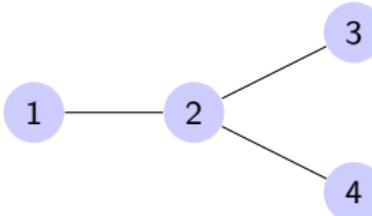
Example for a tree



Examples for graphs

- A **tree** is a connected, undirected graph without cycles, or equivalently, there exists a unique path between any two nodes.

Example for a tree



- A **decomposable/chordal/triangulated** graph is a graph where each cycle of length $n \geq 4$ has a chord, i.e., two non-consecutive nodes that are adjacent.

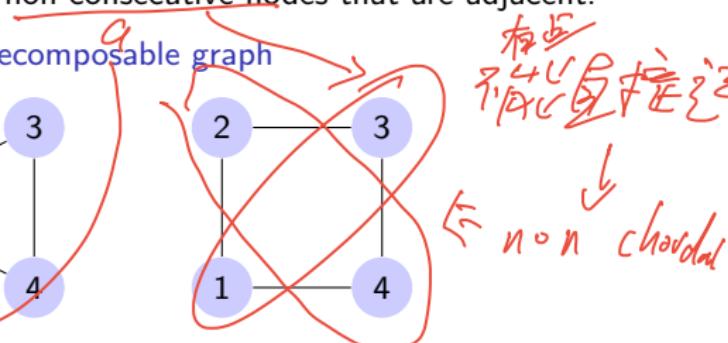
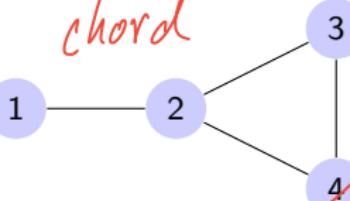
Example for a decomposable and non-decomposable graph

1. 每到图

2. 找 length ≥ 4 的 cycle

3. 找 Q

chord



有
Q

无
Q

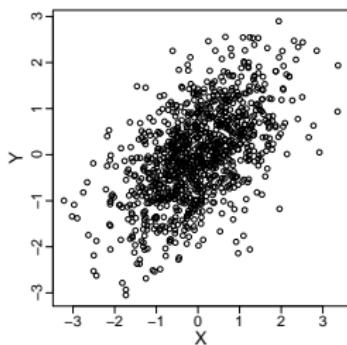
non chordal

Conditional independence

- ▶ Recall: Two random vectors with joint density $f(x, y)$ are independent if $f(x, y) = f(x)f(y)$.
- ▶ This is a very strong assumption, and many variables cannot be assumed to be independent.

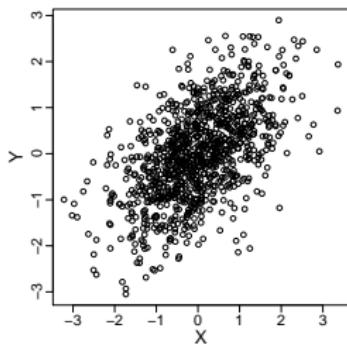
Conditional independence

- ▶ Recall: Two random vectors with joint density $f(x, y)$ are independent if $f(x, y) = f(x)f(y)$.
- ▶ This is a very strong assumption, and many variables cannot be assumed to be independent.
- ▶ Example 1: X = "Ice cream sales" and Y = "Violent crime"



Conditional independence

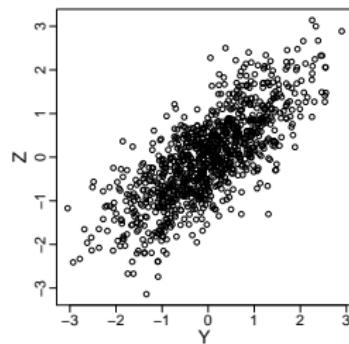
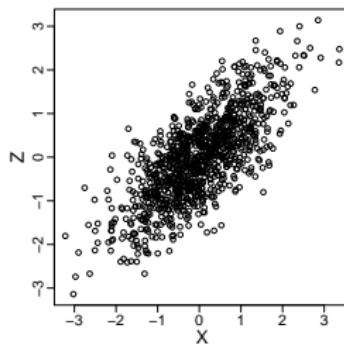
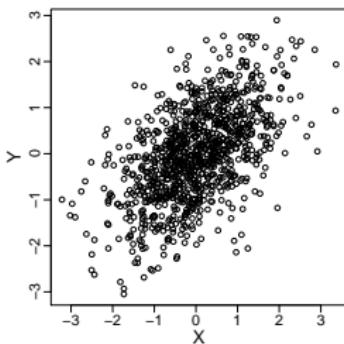
- ▶ Recall: Two random vectors with joint density $f(x, y)$ are independent if $f(x, y) = f(x)f(y)$.
- ▶ This is a very strong assumption, and many variables cannot be assumed to be independent.
- ▶ Example 1: X = "Ice cream sales" and Y = "Violent crime"
- ▶ Example 2: X = "Student 1 late" and Y = "Student 2 late"



Conditional independence

- ▶ Recall: Two random vectors with joint density $f(x, y)$ are independent if $f(x, y) = f(x)f(y)$.
- ▶ This is a very strong assumption, and many variables cannot be assumed to be independent.
- ▶ Example 1: X = "Ice cream sales" and Y = "Violent crime"
- ▶ Example 2: X = "Student 1 late" and Y = "Student 2 late"
- ▶ Conditional independence of X and Y given a third variable Z means:

"If I already know Z , then information on Y is irrelevant for X ".



Conditional independence

- Let X, Y and Z be random vectors with joint density $f_{X,Y,Z}$. We say that X is **conditionally independent** of Y given Z , if

$$f_{\text{_____}}(x \mid z, y) = f_{X|Z}(x \mid z).$$

In this case we write $\underbrace{X \perp\!\!\!\perp Y \mid Z}$.

Conditional independence

- Let X, Y and Z be random vectors with joint density $f_{X,Y,Z}$. We say that X is **conditionally independent** of Y given Z , if

$$f_{X,Y|Z}(x, y | z) = f_{X|Z}(x | z) \cdot f_{Y|Z}(y | z).$$

In this case we write $X \perp\!\!\!\perp Y | Z$.

- There are many equivalent definitions, e.g.,

$$f_{X,Y|Z}(x, y | z) = f_{X|Z}(x | z) f_{Y|Z}(y | z).$$

or

$$f_{X,Y,Z}(x, y, z) = h(x, z) k(y, z),$$

for suitable functions h and k .

$$f_{X|Z}(x | z) = \frac{f(x, z, y)}{f(z, y)} =$$

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

conditional density

$$= f_{(X|Z)}(x | z)$$

Conditional independence

- Let X, Y and Z be random vectors with joint density $f_{X,Y,Z}$. We say that X is **conditionally independent** of Y given Z , if

$$f_{X,Y|Z}(x \mid z, y) = f_{X|Z}(x \mid z).$$

In this case we write $X \perp\!\!\!\perp Y \mid Z$.

- There are many equivalent definitions, e.g.,

$$f_{X,Y|Z}(x, y \mid z) = f_{X|Z}(x \mid z)f_{Y|Z}(y \mid z).$$

or

$$f_{X,Y,Z}(x, y, z) = h(x, z)k(y, z),$$

for suitable functions h and k .

- This clearly generalizes the classical notion of independence, in which case we would have that Z is a trivial, empty vector.

The pairwise Markov property

- ▶ How do we link conditional independence to graph structures?
- ▶ We consider an undirected graph $G = (V, E)$ and a collection of random variables $X = (X_u)_{u \in V}$ indexed by the nodes of the graph.

The pairwise Markov property

- ▶ How do we link conditional independence to graph structures?
- ▶ We consider an undirected graph $G = (V, E)$ and a collection of random variables $X = (X_u)_{u \in V}$ indexed by the nodes of the graph.
- ▶ We say that X satisfies the **pairwise Markov property** if for any $u, v \in V$ with $(u, v) \notin E$ it holds that

$$X_u \perp\!\!\!\perp X_v \mid X_{V \setminus \{u, v\}}.$$

- ▶ In this case we also say that X is a graphical model on G .

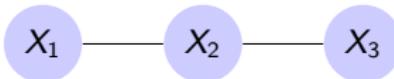
The pairwise Markov property

- ▶ How do we link conditional independence to graph structures?
- ▶ We consider an undirected graph $G = (V, E)$ and a collection of random variables $X = (X_u)_{u \in V}$ indexed by the nodes of the graph.
- ▶ We say that X satisfies the **pairwise Markov property** if for any $u, v \in V$ with $(u, v) \notin E$ it holds that

$$X_u \perp\!\!\!\perp X_v \mid X_{V \setminus \{u, v\}}.$$

- ▶ In this case we also say that X is a graphical model on G .

Example:



$$X_1 \perp\!\!\!\perp X_3 \mid X_2$$

$$f(x_1, x_2, x_3) = f(x_1 \mid x_2, x_3)f(x_2, x_3) = f(x_1 \mid x_2)f(x_2 \mid x_3)f(x_3) = f(x_1, x_2)f(x_2, x_3)/f(x_2)$$

The global Markov property

- ▶ Pairwise Markov property: for any $u, v \in V$ with $(u, v) \notin E$ it holds that

$$X_u \perp\!\!\!\perp X_v \mid X_{V \setminus \{u, v\}}.$$

The global Markov property

- ▶ Pairwise Markov property: for any $u, v \in V$ with $(u, v) \notin E$ it holds that

$$X_u \perp\!\!\!\perp X_v \mid X_{V \setminus \{u, v\}}.$$

- ▶ For disjoint subsets $A, B, S \subset V$, S is said to separate A and B if all paths from a node $u \in A$ and $v \in B$ intersect S .

The global Markov property

- ▶ Pairwise Markov property: for any $u, v \in V$ with $(u, v) \notin E$ it holds that

$$X_u \perp\!\!\!\perp X_v \mid X_{V \setminus \{u, v\}}.$$

- ▶ For disjoint subsets $A, B, S \subset V$, S is said to separate A and B if all paths from a node $u \in A$ and $v \in B$ intersect S .
- ▶ We say that X satisfies the **global Markov property** on G if for any disjoint subsets $A, B, S \subset V$ such that S separates A and B in G , we have

$$X_A \perp\!\!\!\perp X_B \mid X_S.$$

The global Markov property

- ▶ Pairwise Markov property: for any $u, v \in V$ with $(u, v) \notin E$ it holds that

$$X_u \perp\!\!\!\perp X_v \mid X_{V \setminus \{u, v\}}.$$

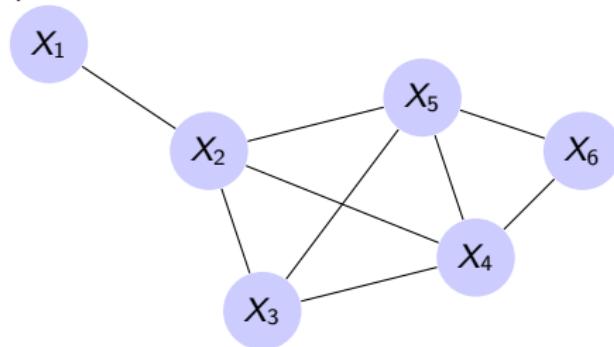
- ▶ For disjoint subsets $A, B, S \subset V$, S is said to separate A and B if all paths from a node $u \in A$ and $v \in B$ intersect S .
- ▶ We say that X satisfies the **global Markov property** on G if for any disjoint subsets $A, B, S \subset V$ such that S separates A and B in G , we have

$$X_A \perp\!\!\!\perp X_B \mid X_S.$$

- ▶ It turns out, that if $X = (X_u)_{u \in V}$ has a positive and continuous density, then pairwise and global Markov properties are **equivalent**.

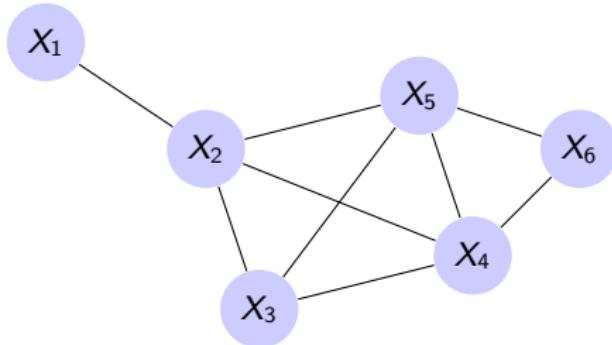
Examples

- ▶ List all conditional independencies and simplify the joint density $f(x_1, x_2, x_3, x_4, x_5, x_6)$:



Examples

- ▶ List all conditional independencies and simplify the joint density $f(x_1, x_2, x_3, x_4, x_5, x_6)$:



- ▶ Draw the corresponding graph:

$$f(x_1, x_2, x_3, x_4, x_5) = f(x_3)f(x_5 \mid x_3)f(x_1, x_4 \mid x_5)f(x_2)$$

The Hammersley–Clifford theorem

Theorem

Let $X = (X_u)_{u \in V}$ be a random vector with positive and continuous density f on \mathcal{X} .

Then the vector X satisfies the **pairwise Markov property** on the graph $G = (V, E)$ if and only if its **density factorizes** as

$$f(x) = \prod_{C \in \mathcal{C}} \psi_C(x_C), \quad x \in \mathcal{X},$$

for suitable functions ψ_C with arguments $x_C = (x_u)_{u \in C}$, where \mathcal{C} is the set of all cliques of G .

The Hammersley–Clifford theorem

Theorem

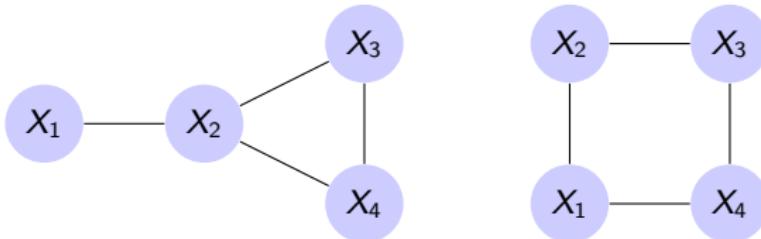
Let $X = (X_u)_{u \in V}$ be a random vector with positive and continuous density f on \mathcal{X} .

Then the vector X satisfies the **pairwise Markov property** on the graph $G = (V, E)$ if and only if its **density factorizes** as

$$f(x) = \prod_{C \in \mathcal{C}} \psi_C(x_C), \quad x \in \mathcal{X},$$

for suitable functions ψ_C with arguments $x_C = (x_u)_{u \in C}$, where \mathcal{C} is the set of all cliques of G .

Examples:



Hammersley–Clifford for decomposable graphs

Corollary

Let $X = (X_u)_{u \in V}$ be a random vector with positive and continuous density f on \mathcal{X} . If X satisfies the pairwise Markov property on the **decomposable** graph $G = (V, E)$ then

$$f(x) = \frac{\prod_{C \in \mathcal{C}} f_C(x_C)}{\prod_{D \in \mathcal{D}} f_D(x_D)}, \quad x \in \mathcal{X},$$

where \mathcal{D} is a multiset containing intersections between the maximal cliques called separator sets.

Hammersley–Clifford for decomposable graphs

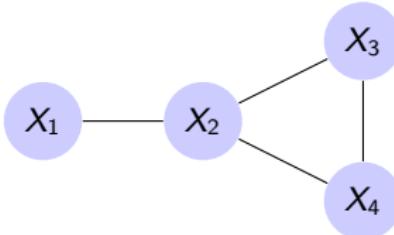
Corollary

Let $X = (X_u)_{u \in V}$ be a random vector with positive and continuous density f on \mathcal{X} . If X satisfies the pairwise Markov property on the **decomposable** graph $G = (V, E)$ then

$$f(x) = \frac{\prod_{C \in \mathcal{C}} f_C(x_C)}{\prod_{D \in \mathcal{D}} f_D(x_D)}, \quad x \in \mathcal{X},$$

where \mathcal{D} is a multiset containing intersections between the maximal cliques called separator sets.

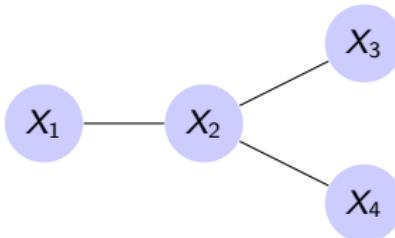
Example:



Tree graphical models

- If X is a graphical model on the tree $T = (V, E)$, then its density has the simple form

$$f(x) = \prod_{(u,v) \in E} \frac{f_{uv}(x_u, x_v)}{f_u(x_u)f_v(x_v)} \prod_{u \in V} f_u(x_u), \quad x \in \mathcal{X}.$$



Recap: formulas for Gaussian distributions

- If $X \sim N_p(\mu, \Sigma)$ can be partitioned into $(X_{(1)}, X_{(2)}) \in \mathbb{R}^{p_1+p_2}$, $\mu = (\mu_{(1)}, \mu_{(2)})$ and

$$\Sigma = \begin{bmatrix} \Sigma_{(11)} & \Sigma_{(12)} \\ \Sigma_{(21)} & \Sigma_{(22)} \end{bmatrix},$$

then the conditional distribution of $X_{(2)}$ given $X_{(1)} = x_{(1)}$ is the density of p_2 -variate normal distribution

$$N_{p_2} \left(\mu_{(2)} + \Sigma_{(21)} \Sigma_{(11)}^{-1} (x_{(1)} - \mu_{(1)}), \Sigma_{(22)} - \Sigma_{(21)} \Sigma_{(11)}^{-1} \Sigma_{(12)} \right)$$

Gaussian graphical models

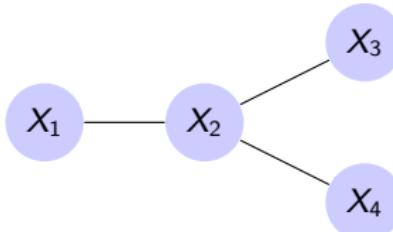
Proposition

Let Σ be a positive-definite $p \times p$ -matrix and $\mu \in \mathbb{R}^p$, and $X \sim N_p(\mu, \Sigma)$. For any $i, j \in \{1, \dots, p\}$ with $i \neq j$ we have

$$X_i \perp\!\!\!\perp X_j \mid X_{\setminus\{i,j\}} \Leftrightarrow K_{ij} = 0,$$

where $K = \{K_{ij}\}_{i,j=1,\dots,p} = \Sigma^{-1}$ is called the concentration/precision matrix.

Example:



$$K = \begin{bmatrix} * & * & 0 & 0 \\ * & * & * & * \\ 0 & * & * & 0 \\ 0 & * & 0 & * \end{bmatrix}$$

Structure learning

- ▶ In most applications, the conditional independencies and the underlying graph structure are **not known**.

Structure learning

- ▶ In most applications, the conditional independencies and the underlying graph structure are **not known**.
- ▶ Therefore, the best graphical structure should be learned from observations.

Maximum likelihood trees

Recall the density for a tree $T = (V, E)$:

$$f(x) = \prod_{(u,v) \in E} \frac{f_{uv}(x_u, x_v)}{f_u(x_u)f_v(x_v)} \prod_{u \in V} f_u(x_u), \quad x \in \mathcal{X}.$$

- ▶ Suppose we have observations $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$ and it is known that the concentration matrix $K = \Sigma^{-1}$ has a **tree structure**, but the specific tree structure T and the parameters are unknown.

Maximum likelihood trees

Recall the density for a tree $T = (V, E)$:

$$f(x) = \prod_{(u,v) \in E} \frac{f_{uv}(x_u, x_v)}{f_u(x_u)f_v(x_v)} \prod_{u \in V} f_u(x_u), \quad x \in \mathcal{X}.$$

- ▶ Suppose we have observations $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$ and it is known that the concentration matrix $K = \Sigma^{-1}$ has a **tree structure**, but the specific tree structure T and the parameters are unknown.
- ▶ For a given tree T , let $\mathcal{S}^+(T)$ denote the space of all positive-definite concentration matrices that factorize according to the tree T .

Maximum likelihood trees

Recall the density for a tree $T = (V, E)$:

$$f(x) = \prod_{(u,v) \in E} \frac{f_{uv}(x_u, x_v)}{f_u(x_u)f_v(x_v)} \prod_{u \in V} f_u(x_u), \quad x \in \mathcal{X}.$$

- ▶ Suppose we have observations $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$ and it is known that the concentration matrix $K = \Sigma^{-1}$ has a **tree structure**, but the specific tree structure T and the parameters are unknown.
- ▶ For a given tree T , let $\mathcal{S}^+(T)$ denote the space of all positive-definite concentration matrices that factorize according to the tree T .
- ▶ For the fixed tree T , we first maximize the likelihood over all parameters of the concentration matrix for this given tree:

$$\widehat{\ell}(T) = \widehat{\ell}(T \mid X_1, \dots, X_n) = \max_{K \in \mathcal{S}^+(T)} \widehat{\ell}(K \mid X_1, \dots, X_n)$$

Maximum likelihood trees

- ▶ Thus, the maximized likelihood is

$$\hat{\ell}(T) = \prod_{(u,v) \in E} \frac{\hat{f}_{uv}(x_u, x_v)}{\hat{f}_u(x_u)\hat{f}_v(x_v)} \prod_{u \in V} \hat{f}_u(x_u),$$

where \hat{f}_A denotes the maximized likelihood based only on data X_{ij} , $j \in A \subset V$, $i = 1, \dots, n$.

Maximum likelihood trees

- ▶ Thus, the maximized likelihood is

$$\hat{\ell}(T) = \prod_{(u,v) \in E} \frac{\hat{f}_{uv}(x_u, x_v)}{\hat{f}_u(x_u)\hat{f}_v(x_v)} \prod_{u \in V} \hat{f}_u(x_u),$$

where \hat{f}_A denotes the maximized likelihood based only on data X_{ij} , $j \in A \subset V$, $i = 1, \dots, n$.

- ▶ We compute that:

$$\log \frac{\hat{f}_{uv}(x_u, x_v)}{\hat{f}_u(x_u)\hat{f}_v(x_v)} = -\frac{n}{2} \log(1 - r_{uv}^2),$$

where r_{uv} is the **empirical correlation coefficient**

$$r_{uv} = w_{uv} / \sqrt{w_{uu} w_{vv}},$$

and $\mathbf{W} = \{w_{uv}, u, v \in V\} = (n-1)\mathbf{S}$ is the Wishart matrix of sums of products.

Maximum likelihood trees

- ▶ Define the **empirical correlation weight** w_{uv} of the edge $(u, v) \in E$ as

$$w_{uv} = -\frac{n}{2} \log(1 - r_{uv}^2),$$

and let the total empirical weight of the tree T be

$$w(T) = \sum_{(u,v) \in E} w_{uv}.$$

Maximum likelihood trees

- ▶ Define the **empirical correlation weight** w_{uv} of the edge $(u, v) \in E$ as

$$w_{uv} = -\frac{n}{2} \log(1 - r_{uv}^2),$$

and let the total empirical weight of the tree T be

$$w(T) = \sum_{(u,v) \in E} w_{uv}.$$

- ▶ The maximized likelihood of all nodes independent, i.e., with no edges, is

$$\widehat{\ell}(\emptyset) = \prod_{u \in V} \widehat{f}_u(x_u).$$

Maximum likelihood trees

- ▶ Define the **empirical correlation weight** w_{uv} of the edge $(u, v) \in E$ as

$$w_{uv} = -\frac{n}{2} \log(1 - r_{uv}^2),$$

and let the total empirical weight of the tree T be

$$w(T) = \sum_{(u,v) \in E} w_{uv}.$$

- ▶ The maximized likelihood of all nodes independent, i.e., with no edges, is

$$\widehat{\ell}(\emptyset) = \prod_{u \in V} \widehat{f}_u(x_u).$$

- ▶ With this notation we have

$$\log \widehat{\ell}(T) - \log \widehat{\ell}(\emptyset) = w(T) = \sum_{(u,v) \in E} w_{uv}.$$

Maximum likelihood trees

- ▶ The **maximum likelihood tree** is therefore given by the maximizer

$$\widehat{T} = \arg \max_{T=(V,E)} \log \widehat{\ell}(T) - \log \widehat{\ell}(\emptyset) = \arg \max_{T=(V,E)} \sum_{(u,v) \in E} w_{uv}.$$

Maximum spanning trees

Maximum spanning tree

Suppose that a positive real number $w_{uv} > 0$ is attached to each pair of nodes (u, v) , $u, v \in E$, with $w_{uu} = 0$ and $w_{uv} = w_{vu}$. The **maximum spanning tree** is the tree $T_{\max} = (V, E_{\max})$, $E \subset V \times V$ that maximizes the sum of weights on that tree, i.e.,

$$T_{\max} = \arg \max_{T=(V,E)} \sum_{(u,v) \in E} w_{uv}. \quad (1)$$

Maximum spanning trees

Maximum spanning tree

Suppose that a positive real number $w_{uv} > 0$ is attached to each pair of nodes (u, v) , $u, v \in E$, with $w_{uu} = 0$ and $w_{uv} = w_{vu}$. The **maximum spanning tree** is the tree $T_{\max} = (V, E_{\max})$, $E \subset V \times V$ that maximizes the sum of weights on that tree, i.e.,

$$T_{\max} = \arg \max_{T=(V,E)} \sum_{(u,v) \in E} w_{uv}. \quad (1)$$

- Given a set of weights $(w_{uv})_{u,v \in V}$, there are very efficient greedy algorithms to solve (1).

Gaussian trees

Theorem

Let $X \sim N_p(\mu, \Sigma)$ and suppose it satisfies the pairwise Markov property on the tree $T = (V, E)$. Let

$$w_{uv} = -\frac{n}{2} \log(1 - r_{uv}^2),$$

based on observations $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$, then the maximum spanning tree T_{\max} with these weights is equal to the maximum likelihood tree and it satisfies

$$\mathbb{P}(T_{\max} = T) \rightarrow 1, \quad n \rightarrow \infty,$$

that is, the true tree structure is recovered almost surely.

Gaussian trees

Theorem

Let $X \sim N_p(\mu, \Sigma)$ and suppose it satisfies the pairwise Markov property on the tree $T = (V, E)$. Let

$$w_{uv} = -\frac{n}{2} \log(1 - r_{uv}^2),$$

based on observations $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$, then the maximum spanning tree T_{\max} with these weights is equal to the maximum likelihood tree and it satisfies

$$\mathbb{P}(T_{\max} = T) \rightarrow 1, \quad n \rightarrow \infty,$$

that is, the true tree structure is recovered almost surely. Moreover, it follows that

$$T = \arg \max_{T=(V,E)} \sum_{(u,v) \in E} -\frac{n}{2} \log(1 - \rho_{uv}^2),$$

where ρ_{uv} is the **correlation coefficient** between X_u and X_v , $u, v \in V$.

Kruskal's algorithm

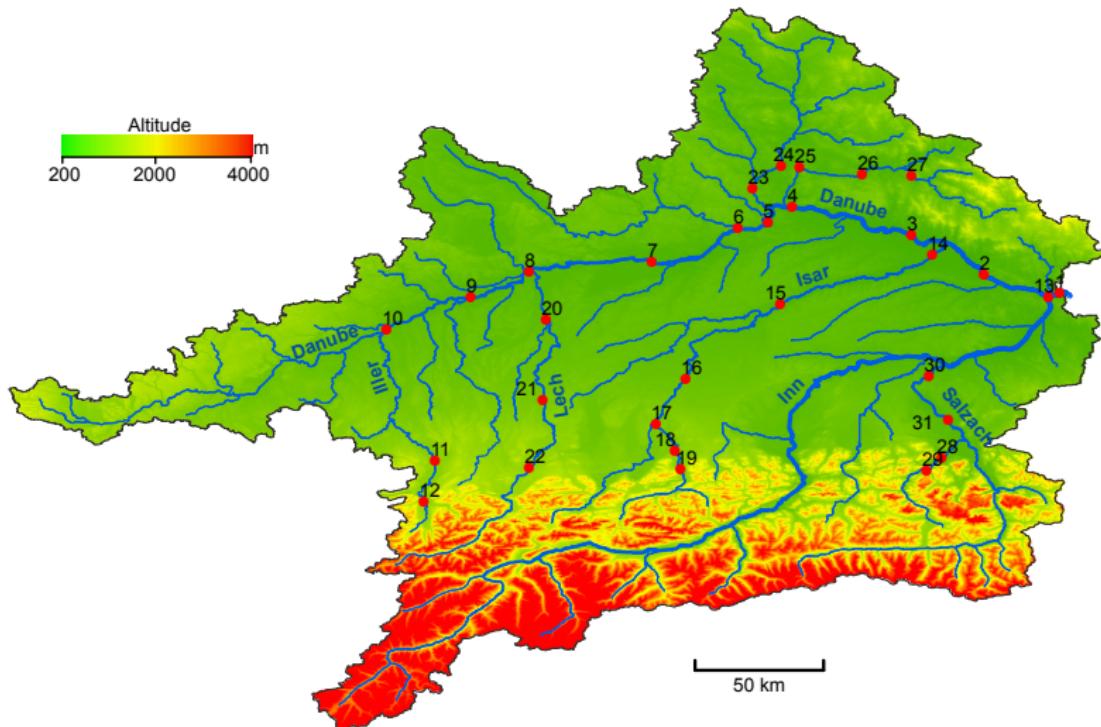
Kruskal's algorithm

Under the notation from above perform the following steps to obtain the maximum spanning tree T_{\max} .

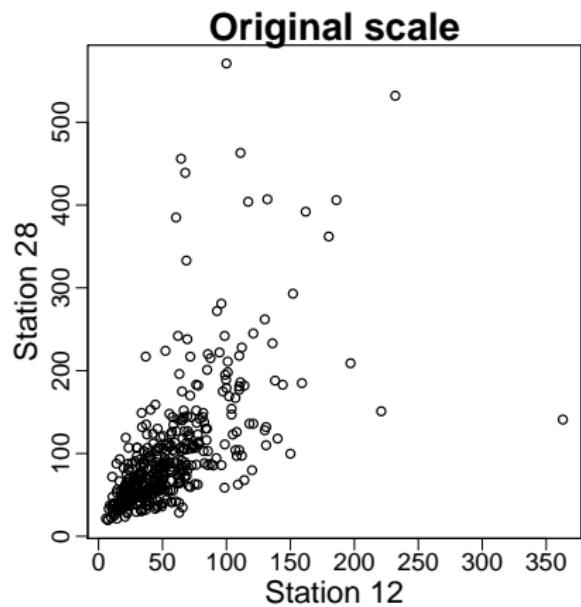
1. Initiate $E_{\max} = \emptyset$.
2. Among all edges in (u, v) not yet chosen, choose the one with the shortest length that does not form any loops with the edges already chosen, and add it to E_{\max} .
3. Repeat 2. until no more edges can be added to E_{\max} without creating a loop.

River discharge data

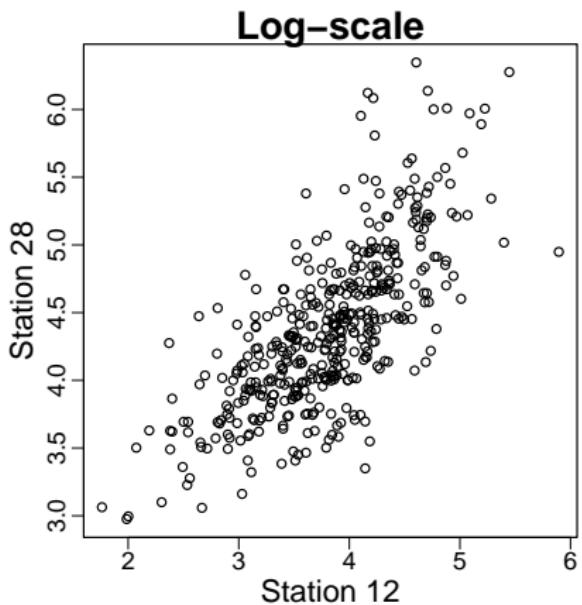
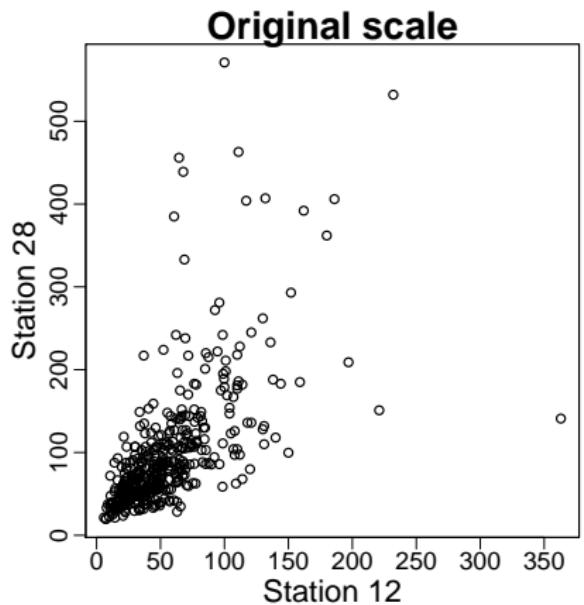
We have $n = 428$ independent observations of $X \in \mathbb{R}^p$, where $p = 31$ is the number of stations on the river network.



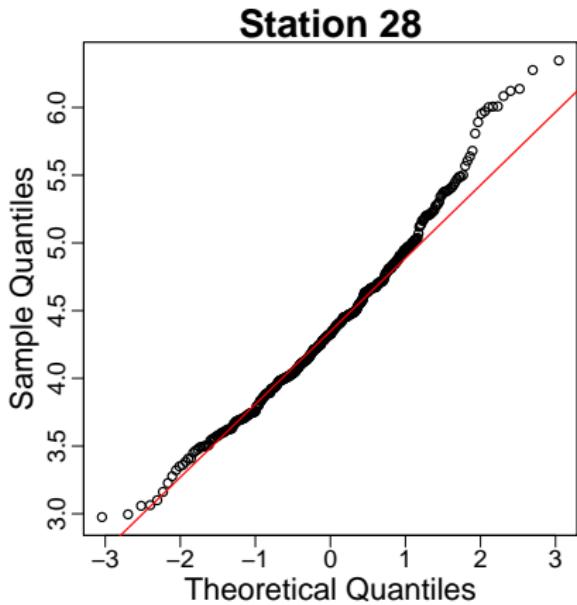
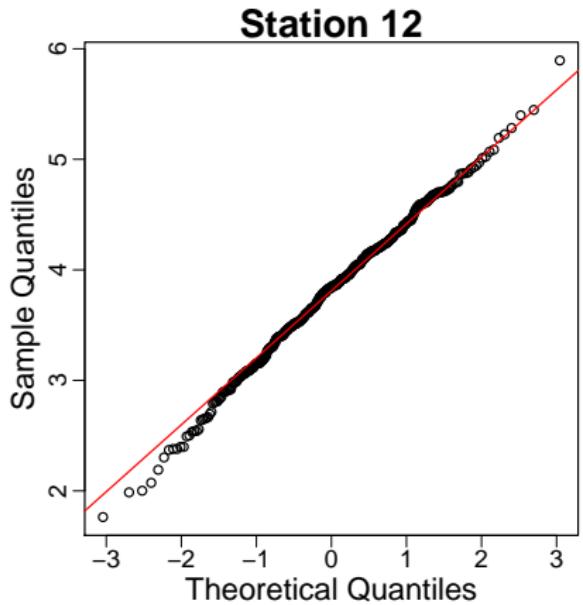
River discharge data



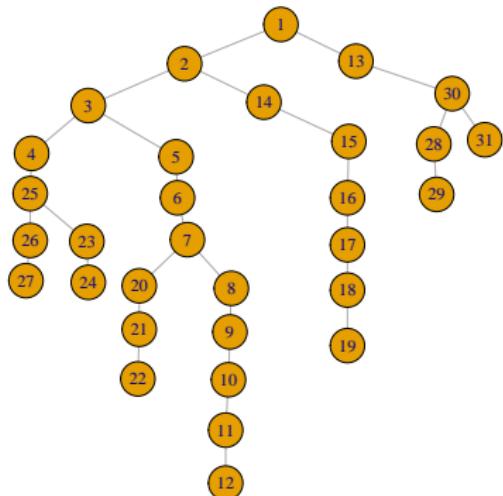
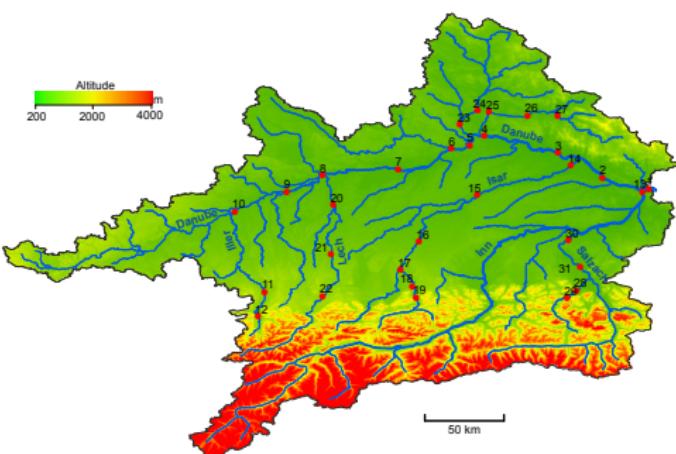
River discharge data



River discharge data



River discharge data: maximum spanning tree



Structure learning for general Gaussian graphical models

Recall the Gaussian log-density for data $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$:

$$\ell(\mu, \Sigma) \propto -\frac{n}{2} \log \det(\Sigma) - \frac{n}{2} \text{tr}(\widehat{\Sigma} \Sigma^{-1}) - \frac{n}{2} (\bar{X} - \mu)^\top \Sigma^{-1} (\bar{X} - \mu).$$

Structure learning for general Gaussian graphical models

Recall the Gaussian log-density for data $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$:

$$\ell(\mu, \Sigma) \propto -\frac{n}{2} \log \det(\Sigma) - \frac{n}{2} \text{tr}(\widehat{\Sigma} \Sigma^{-1}) - \frac{n}{2} (\bar{X} - \mu)^\top \Sigma^{-1} (\bar{X} - \mu).$$

Maximizing this partially wrt μ gives

$$\ell(\Sigma) \propto -\frac{n}{2} \log \det(\Sigma) - \frac{n}{2} \text{tr}(\widehat{\Sigma} \Sigma^{-1}).$$

Graphical lasso

Instead of solving $\arg \max_{\Sigma \succeq 0} \ell(\Sigma)$ directly (which would give $\widehat{\Sigma}$) we reformulate this problem in terms of the precision matrix K , that is

$$\widehat{K} = \arg \max_{K \succeq 0} \log \det(K) - \text{tr}(\widehat{\Sigma} K)$$

Structure learning for general Gaussian graphical models

Recall the Gaussian log-density for data $X_1, \dots, X_n \sim N_p(\mu, \Sigma)$:

$$\ell(\mu, \Sigma) \propto -\frac{n}{2} \log \det(\Sigma) - \frac{n}{2} \text{tr}(\widehat{\Sigma} \Sigma^{-1}) - \frac{n}{2} (\bar{X} - \mu)^\top \Sigma^{-1} (\bar{X} - \mu).$$

Maximizing this partially wrt μ gives

$$\ell(\Sigma) \propto -\frac{n}{2} \log \det(\Sigma) - \frac{n}{2} \text{tr}(\widehat{\Sigma} \Sigma^{-1}).$$

Graphical lasso

Instead of solving $\arg \max_{\Sigma \succeq 0} \ell(\Sigma)$ directly (which would give $\widehat{\Sigma}$) we reformulate this problem in terms of the precision matrix K , that is

$$\widehat{K}_\lambda = \arg \max_{K \succeq 0} \log \det(K) - \text{tr}(\widehat{\Sigma} K) - \lambda |K|_1$$

where $|K|_1 = \sum_{ij} |K_{ij}|$ and $\lambda \geq 0$ is a tuning parameter. This is a convex optimization problem that can be efficiently solved.

Graphical lasso

$$\widehat{K}_\lambda = \arg \max_{K \succeq 0} \log \det(K) - \text{tr}(\widehat{\Sigma} K) - \lambda |K|_1$$

- ▶ The ℓ_1 -penalization term $\lambda|K|_1$ penalizes the **complexity** of the concentration matrix.

Graphical lasso

$$\widehat{K}_\lambda = \arg \max_{K \succeq 0} \log \det(K) - \text{tr}(\widehat{\Sigma} K) - \lambda |K|_1$$

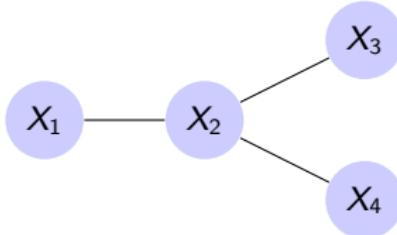
- ▶ The ℓ_1 -penalization term $\lambda|K|_1$ penalizes the **complexity** of the concentration matrix.
- ▶ For increasing λ , one obtains simpler models and for decreasing λ models closer to the MLE $K = \widehat{\Sigma}^{-1}$.

Graphical lasso

$$\widehat{K}_\lambda = \arg \max_{K \succeq 0} \log \det(K) - \text{tr}(\widehat{\Sigma} K) - \lambda |K|_1$$

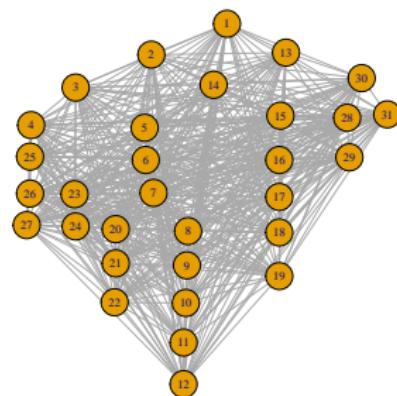
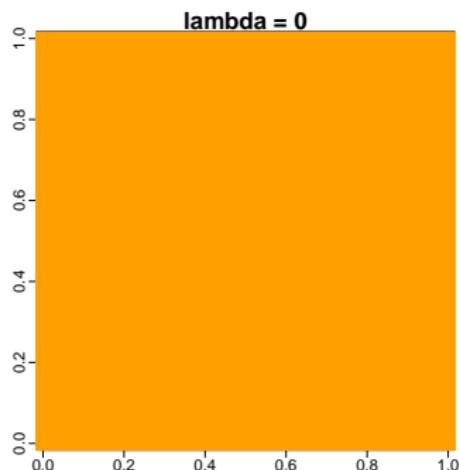
- ▶ The ℓ_1 -penalization term $\lambda|K|_1$ penalizes the **complexity** of the concentration matrix.
- ▶ For increasing λ , one obtains simpler models and for decreasing λ models closer to the MLE $K = \widehat{\Sigma}^{-1}$.
- ▶ The geometry of the ℓ_1 -penalization forces some of the entries of K to be **exactly zero**, and therefore one obtains a graphical model!

Example:

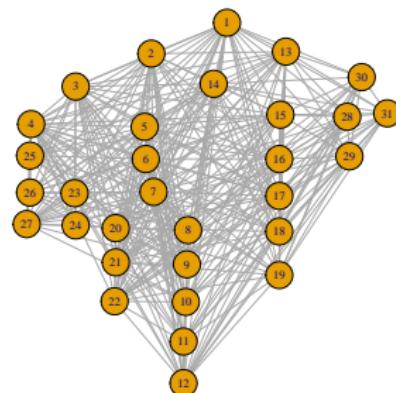
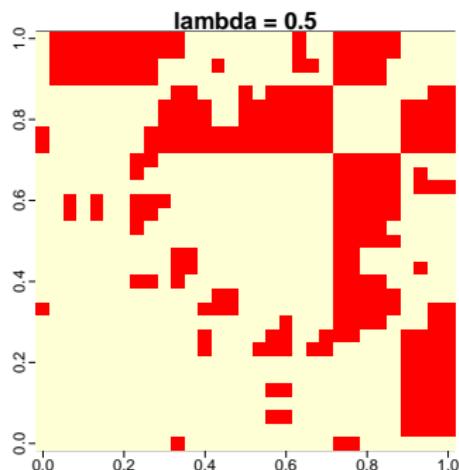


$$\widehat{K}_\lambda = \begin{bmatrix} * & * & 0 & 0 \\ * & * & * & * \\ 0 & * & * & 0 \\ 0 & * & 0 & * \end{bmatrix}$$

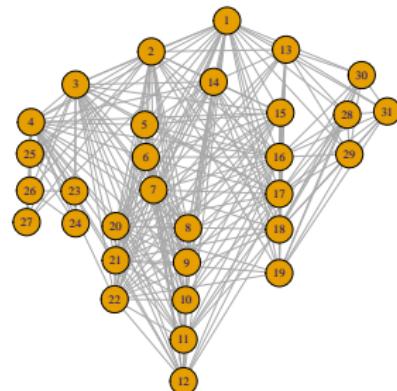
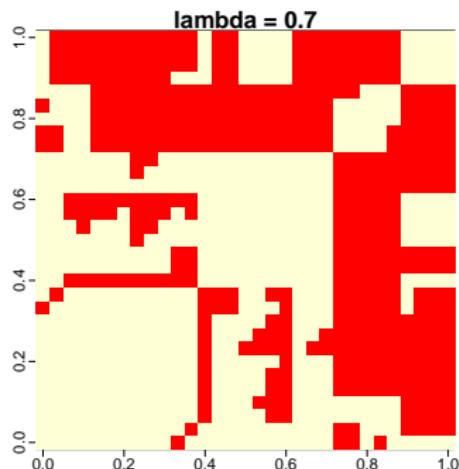
River discharge data: graphical lasso



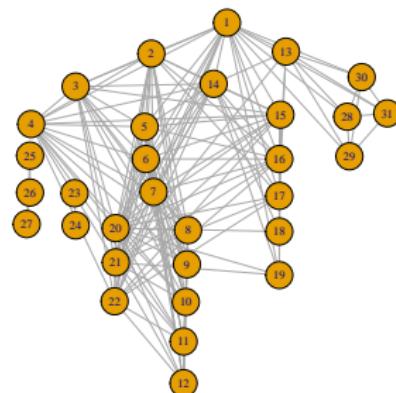
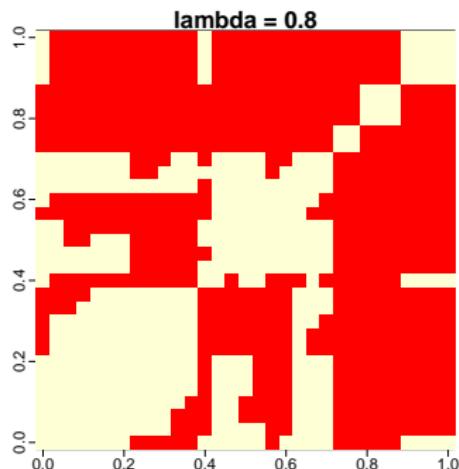
River discharge data: graphical lasso



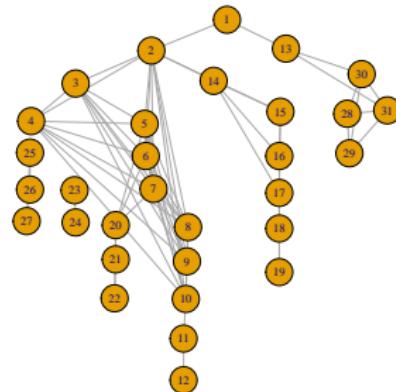
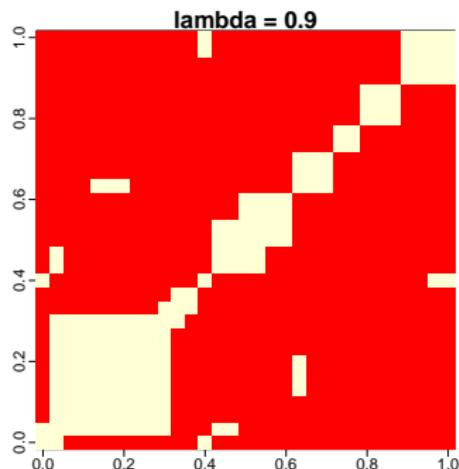
River discharge data: graphical lasso



River discharge data: graphical lasso



River discharge data: graphical lasso



River discharge data: graphical lasso

