

Es 11)

$$X, Y \sim \text{Exp}(1), \quad \frac{X}{X+Y}, \quad X+Y$$

$$\begin{cases} u = \frac{x}{x+y} \\ v = x+y \end{cases} \rightarrow \begin{cases} x = uv \\ y = v - uv \end{cases}$$

$$f_{u,v}(u,v) = f_{x,y}(uv, v(1-u)) |\det(J)| =^*$$

$$J = \begin{pmatrix} v & u \\ -v & 1-u \end{pmatrix} \rightarrow |\det(J)| = v(1-u) + uv = v$$

$$= e^{-uv} e^{-(1-u)v} |v|$$

$$= |v| e^{-v} \cdot \mathbb{1}_{(0,1)}(u)$$

$\hookrightarrow u$ can't miss from this function

it's a product between a gamma and an uniform distribution

\Downarrow

u, v are independent

$X_n \sim \text{Bin}(n, \theta) \rightarrow$ Sum of Bernoulli (a.i.) random variables

a) $X_n/n \xrightarrow{P} \theta \rightarrow$ Average of the Bernoulli r.v.

\rightarrow linearity of expectation

$$E(X_n/n) = \frac{1}{n} E(X_n) = \frac{1}{n} n\theta = \theta$$

$$\text{Var}(X_n/n) = \frac{1}{n^2} \text{Var}(X_n) = \frac{1}{n^2} n\theta(1-\theta) = \frac{\theta(1-\theta)}{n}$$

\hookrightarrow for $n \rightarrow \infty$, $\text{Var}(X_n/n) \rightarrow 0$

let $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|X_n/n - \theta| \geq \epsilon) \leq \lim_{n \rightarrow \infty} \frac{\theta(1-\theta)}{n\epsilon^2} \rightarrow 0$$

\hookrightarrow Proved that $\frac{X_n}{n} \rightarrow \theta$

b) $1 - \frac{X_n}{n} \xrightarrow{P} 1 - \theta$

$$E(1 - \frac{X_n}{n}) = 1 - E(\frac{X_n}{n}) = 1 - \theta$$

$$\text{Var}(1 - \frac{X_n}{n}) = -\text{Var}(\frac{X_n}{n}) = -\frac{\theta(1-\theta)}{n} \rightarrow 0$$

\hookrightarrow Proved that converges to 0

c) $(X_n/n)(1 - X_n/n) \xrightarrow{P} \theta(1-\theta)$

\hookrightarrow True due to algebraic laws of operations of r.v.

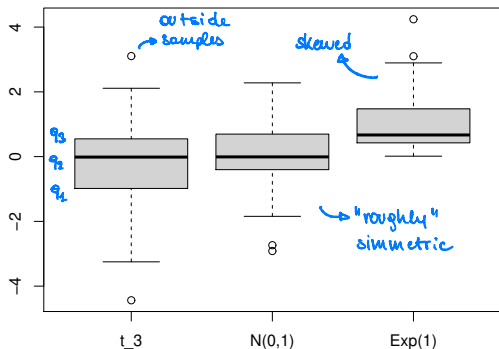
Boxplot (Box-and-whiskers plot)

It's a 5-summary statistics description of a (typically observed) sample.

It provides informations about the location, spread and the shape of the distribution of the sample.

In the vertical orientation:

- the middle line represents q_2 , and vertical edges of the box represent q_1 and q_3 , resp.
- the upper whisker is the largest $x_i \leq q_3 + 1.5 \cdot \text{iqr}$
- the lower whisker is the smallest $x_i \geq q_1 - 1.5 \cdot \text{iqr}$
- observations outside the whiskers are typically marked by a “*”



↓
box is wider → higher variability

When we see whiskers that are almost equal or we see that the median is roughly in the middle



The distribution is (roughly) symmetric

Quantile-Quantile plot

The QQ plot is useful for checking if an observed sample is compatible with a population with continuous F .

It works by comparing a list of observed sample quantiles with the corresponding quantiles of a distribution F .

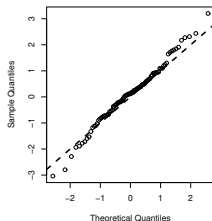
It consists in plotting the pairs $(x_{(i)}, F^{-1}(i/(n+1)))$ and looking for a linear pattern.

observed quantiles *quantile function (true)*

The QQ plot with F the normal distribution is the most widely used. In practice F involves unknown parameters which have to be estimated beforehand.

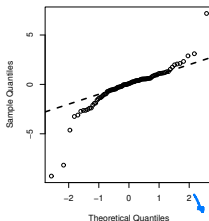
Good compatibility

(a) observed vs $N(0,1)$: ok



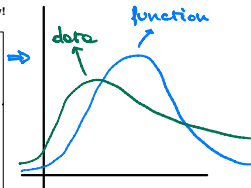
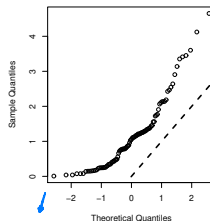
Tails of our funct.
are too rapid

(b) observed vs $N(0,1)$: tails!

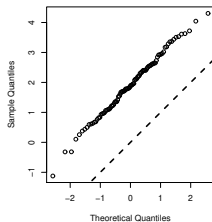


The chosen a symmetric
function while it should be skewed

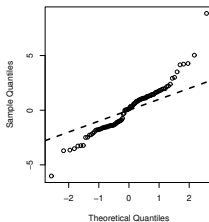
(c) observed vs $N(0,1)$: symmetry!



(d) observed vs $N(0,1)$: location!

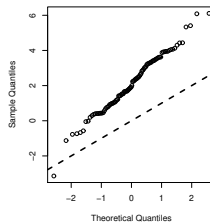


(e) observed vs $N(0,1)$: scale!



Must find new distribution

(f) observed vs $N(0,1)$: location and scale!



Un-matched
sample space

My function has too
low variability

Multivariate data

In realistic applications we may collect observation for several variables, thus a typical dataset looks like

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}$$

where n is the number of observations and p is the number of variables.

The j th column represents the overall sample for the j th variable, and the i th row is the i th sample point for all variables.

For example, the columns could be pollutants s.t. $\text{PM}_{2.5}$, PM_{10} , CO_2 , etc. and the rows may be values measured hourly.

Summaries for multivariate data

A common query is if the p variables are related to each other.

A first approach could be to plot pairs of variables and inspect the graph for possible associations.

For pairs of variables the sample covariance and the sample Pearson's correlation are widely used measures of association.

In particular, for a pair of variables x, y , the sample covariance is

$$s_{xy} = (n - 1)^{-1} \sum_i (x_i - \bar{x})(y_i - \bar{y}),$$

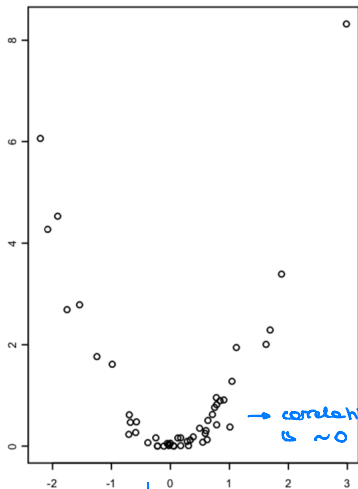
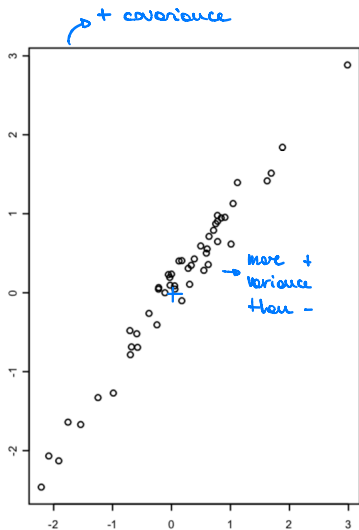
and the sample Pearson correlation is

$$r_{xy} = \frac{s_{xy}}{s_x s_y} \rightarrow \in [-1, 1]$$

↗ strong positive association
↘ strong negative association

s_{xy} targets its population version σ_{XY} , whereas r_{xy} targets its population version ρ_{XY} .

Caution: lack of correlation \nRightarrow lack of association!



are associated but not (linearly) correlated

it's a more general term

\Rightarrow can't tell that there's no association

Statistical models

Let X_1, \dots, X_n be random sample with $X \sim F_\theta$. If, in addition, X_i are also independent we call it iid random sample.

Depends on a parameter that I don't know

↳ identically independently distributed
The joint pdf of the sample is $f(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$.

By a statistical model we mean the set

$$\{f(x_1, \dots, x_n; \theta) : \theta \in \Theta, x_i \in \mathcal{X}\},$$

where Θ is the parameter space, i.e. the set of all possible values for θ .

Typically, $\Theta \subseteq \mathbb{R}^d$ for some integer $d > 0$ and X_i could be a rv or a rve of any dimension.