

CHAPTER 2

DESCRIPTIVE STATISTICS

2.1 [Data analysis]

→ model

specification

→ residuals
after fitting
a model

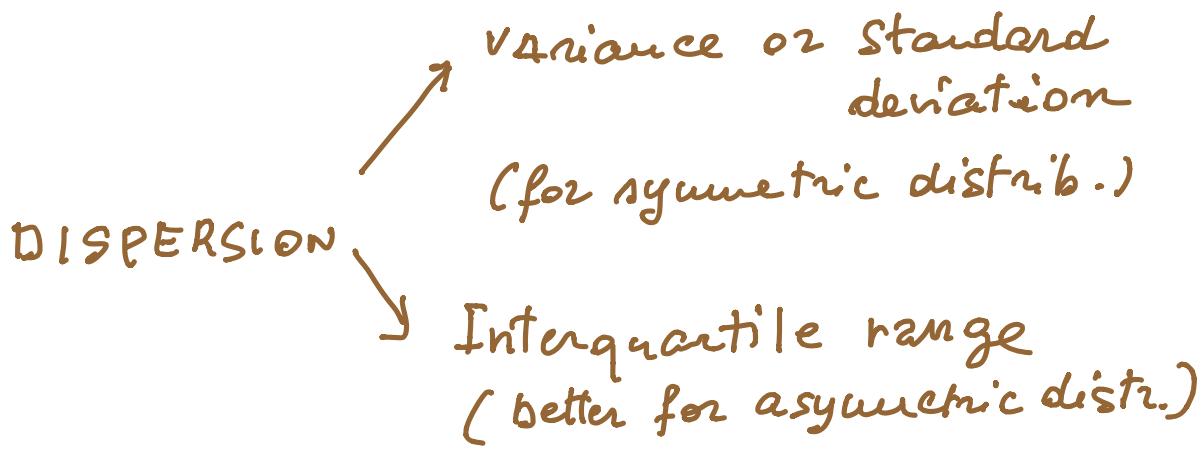
2.2 [Univariate sample]

$\underline{x} = (x_1, \dots, x_n)$ generated by

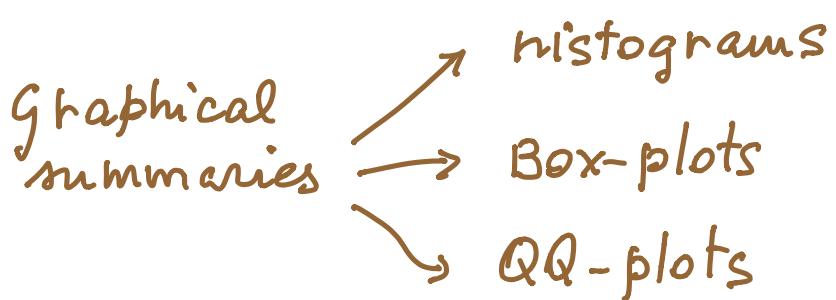
$\underline{X} = (X_1, \dots, X_n)$ i.i.d. F

Important summaries → LOCATION
→ DISPERSION

LOCATION → mean: $\bar{x} = \sum_{i=1}^n x_i / n$
→ median: a middle value
in the sequence of sorted data

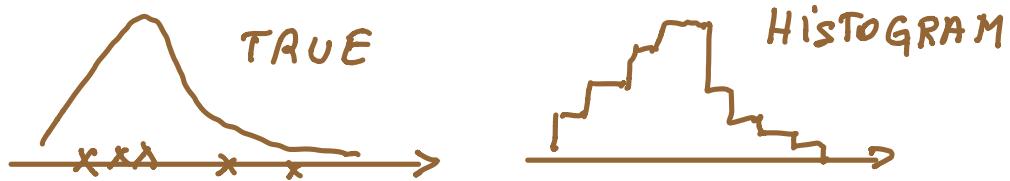


- Sample variance $S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
sample s.d. = S_x
- $IQR = Q_{0.75} - Q_{0.25}$ (percentiles)

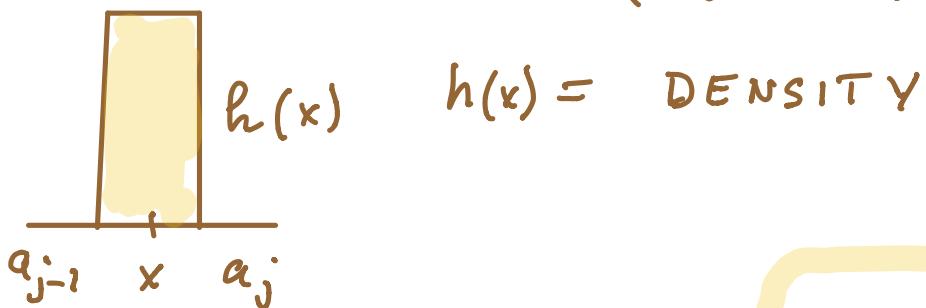


2.3 [Histograms]

To get an idea of the density



$$x \in (a_{j-1}, a_j] \mapsto h(x) = \frac{\text{frequency in a class}}{(a_j - a_{j-1})n}$$



Relative frequency = AREA of rectangle



- Histogram formally

Given (x_1, \dots, x_m)

$$h_m(x) = \frac{1}{a_j - a_{j-1}} \cdot \frac{1}{n} \sum_{i=1}^n 1_{a_{j-1} < x_i \leq a_j}$$

where 1 is an indicator function.

- A histogram is a good density estimator if

- The partition a_0, a_1, \dots, a_m is well chosen
- Sample size is not too small.

Proof. $(X_1, \dots, X_m) \stackrel{iid}{\sim} f(x)$ true density.

Assume that $a_{j-1} < x \leq a_j$ and $f(x) > 0$

$$E h_m(x) = E \frac{1}{n(a_j - a_{j-1})} \sum_{i=1}^n 1_{a_{j-1} < X_i \leq a_j}$$

$$\begin{aligned}
 &= \frac{1}{a_j - a_{j-1}} \cdot \frac{1}{n} \sum_{i=1}^n E(1_{a_{j-1} < X_i \leq a_j}) \\
 &= \frac{1}{a_j - a_{j-1}} \cdot \frac{1}{n} \cdot n E(1_{a_{j-1} < X_1 \leq a_j}) \\
 &= \frac{1}{a_j - a_{j-1}} P(a_{j-1} < X_1 \leq a_j) \\
 &= \frac{\int_{a_{j-1}}^{a_j} f(s) ds}{a_j - a_{j-1}}
 \end{aligned}$$

$$E h_n(x) = \begin{array}{c} \text{---} \\ | \\ \text{---} \\ x \end{array} \approx f(x)$$

By the Law of large numbers.

$$h_n(x) \xrightarrow{P} E h_n(x)$$

for $n \rightarrow \infty$.

See some examples in the notebook.

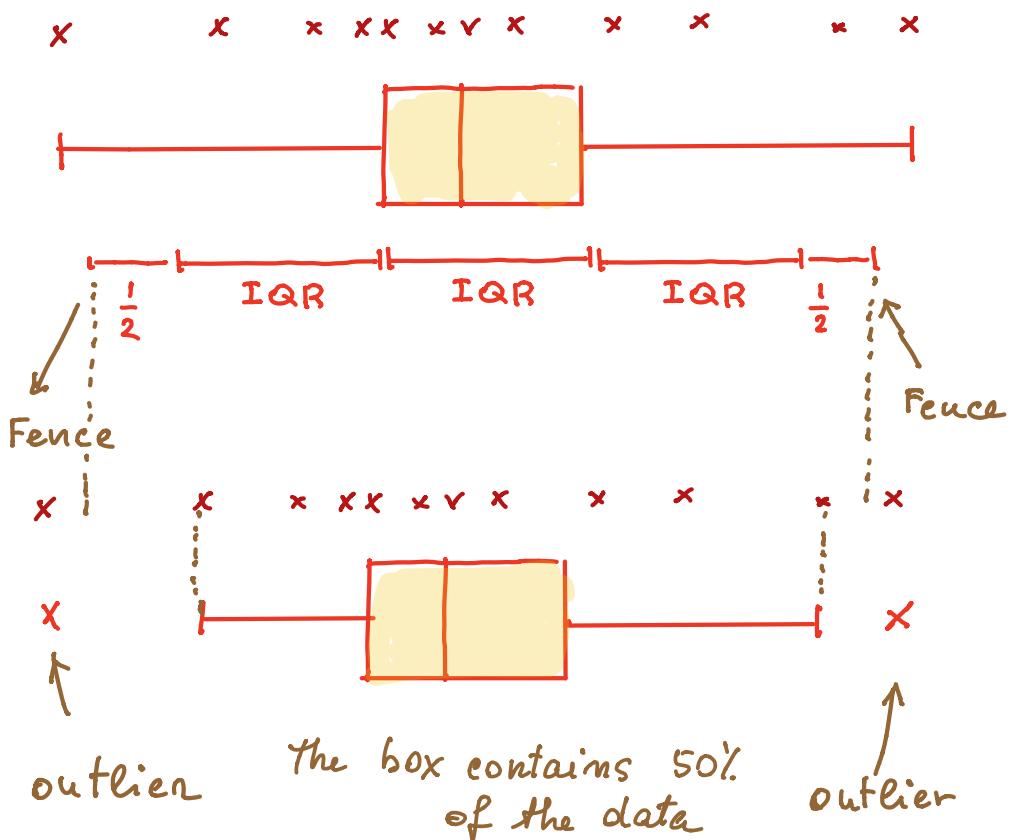
Box-plots

They can be vertical or horizontal

$$\text{Data} = (x_1, \dots, x_n)$$

↓
5-number
SUMMARY

min $Q_{0.25}$ Q_{50} $Q_{0.75}$ max



Example 2.5

Differences between Normal
and Cauchy.

Standard Cauchy density

$$f(x) = \frac{1}{\pi(1+x^2)} \quad x \in \mathbb{R}$$

The expected value $E(x)$ does not
exist since $E(|x|)$ is not finite.

LOCATION-SCALE FAMILY

• Preliminaries

If $X \sim F_X$ (distribution function)

Then $Y = a + bX$ has distribution

$$F_Y(y) = F_X\left(\frac{y-a}{b}\right) \quad b > 0$$

Proof:

$$P(a+bX \leq y) = P\left(X \leq \frac{y-a}{b}\right)$$

$F_y(y) = F_{a,b}(y)$ is a family
depending on $a \in \mathbb{R}$, $b > 0$.

called location-scale family

generated by X .

• Density.

if F_x has a density f_x

then $F_{a,b}$ has a density $f_{a,b}$

$$f_{a,b}(y) = \frac{d}{dy} F_x\left(\frac{y-a}{b}\right) = \frac{1}{b} f_x\left(\frac{y-a}{b}\right)$$

Note

if X is standardized i.e.

$$EX = 0 \quad \text{var } X = 1$$

Then $a = E(Y)$
 $b^2 = \text{Var}(Y)$

• Examples 1) $X \sim N(0, 1)$

$$a + bX \sim N(a, b^2)$$

2) $X \sim \text{Cauchy}(0, 1)$ $f_x(x) = \frac{1}{\pi(1+x^2)}$

$$a + bX \rightarrow f_y(y) = \frac{1}{b\pi \left[1 + \left(\frac{y-a}{b} \right)^2 \right]}$$

• NOTE

$X \sim \chi^2_1$ chi-square with 1 degree of freedom

$Y \sim \chi^2_2$ 2 degrees of freedom.

They are not in the same location-scale family

but they have the same name

Verify. $X \sim \text{Gamma} \left(\frac{1}{2}, \frac{1}{2} \right)$

$Y \sim \text{Gamma} \left(1, \frac{1}{2} \right)$

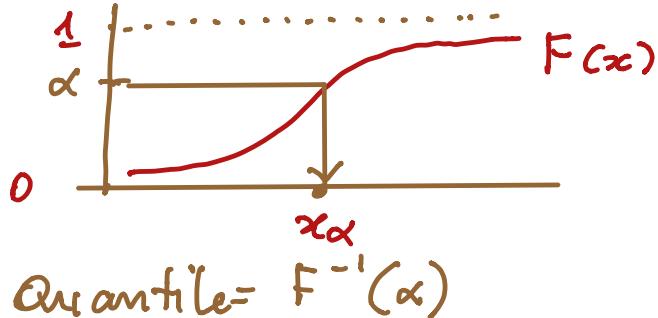
Different shape, same rate

$a=0; bX \sim \text{Gamma} \left(\frac{1}{2}, \frac{1}{2b} \right)$

Then bX cannot change shape.

Quantile function

- Idea



$$\text{Quantile} = F^{-1}(\alpha)$$

x_α is an α -quantile of F

if exist a single value $\alpha \in (0, 1)$
such that

$$F(x_\alpha) = \alpha$$

The α -quantile is denoted by $\tilde{F}^{-1}(\alpha)$

[Note] $F^{-1}(0.25) = 25$ percentile
 $F^{-1}(0.75) = 75$ percentile
 $F^{-1}(0.5) = \text{median}$

The function $\alpha \mapsto F^{-1}(\alpha)$
is the quantile function provided
that F^{-1} is well-defined.

If F is strictly increasing and continuous

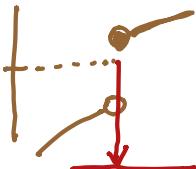
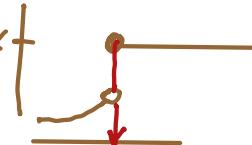
$$F(F^{-1}(\alpha)) = \alpha \quad \text{for all } \alpha \in (0,1)$$
$$F^{-1}(F(x)) = x \quad \text{for all } x \in \mathbb{R}$$

Example - Exponential (λ)

$$F(x) = 1 - e^{-\lambda x} \quad x \geq 0$$

$$F^{-1}(x) = -\log(1-\alpha)/\lambda$$

In general the equation $F(x) = \alpha$ may have

- one solution 
- no solution 
- ∞ solutions 

To define well the quantile function in all cases let :

$$\bar{F}(\alpha) = \inf \{ x : F(x) \geq \alpha \}$$

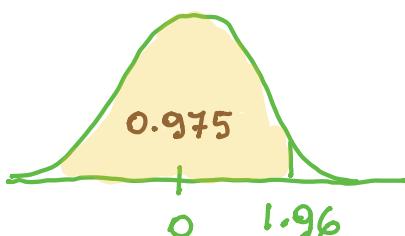
→

In location-scale family.

$$F_{a,b}^{-1}(\alpha) = \alpha + \beta F^{-1}(\alpha)$$

Classical example

$$X \sim N(0,1) \quad x_{0.975} = 1.96$$



$$Y \sim N(\mu, \sigma^2) \quad y_{0.975} = \mu + 1.96\sigma$$

In general, the points with coordinates

$$(\bar{F}(\alpha), \bar{F}_{a,b}^{-1}(\alpha))$$

are on the line

$$y = a + b x$$

this is the base of the **QQ plot**

—

Order statistics

Given a sample (X_1, \dots, X_n)

take the **sorted sample**

$$(X_{(1)}, X_{(2)}, \dots, X_{(n)})$$

placed in increasing order

$X_{(1)}, \dots, X_{(n)}$ are called order statistics.

It can be proved that

$$E(F(X_{(i)})) = \frac{i}{n+1}$$

It is expected that the points

$$\left(\frac{i}{n+1}, F(X_{(i)})\right)$$

lie on the line $y=x$

The same must happen for

$$\left(F^{-1}\left(\frac{i}{n+1}\right), X_{(i)}\right)$$

QQ-plot.

Given a sample (x_1, \dots, x_n)
the plot of the points

$$\left(F^{-1}\left(\frac{i}{n+1}\right), x_{(i)}\right)$$

is called a quantile-quantile plot.
with respect to the distribution F

Points $(F^{-1}\left(\frac{i}{m+1}\right), x_{ci})$

- $\approx y = xc \rightarrow$ sample comes from F
- $\approx y = a + bx \rightarrow$ sample comes from $F_{a,b}$
- different curve \rightarrow sample not from F .

2.3 CORRELATION

Observations x_i can be on multiple variables, for same unit.

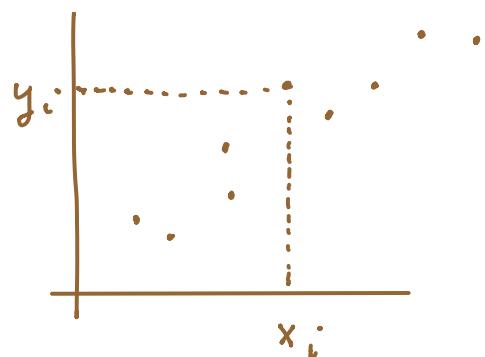
$$\tilde{x}_i = (x_{i1}, x_{i2})$$



$$(x_i, y_i)$$

Ex: $(\text{height of mother}, \text{height of daughter})$

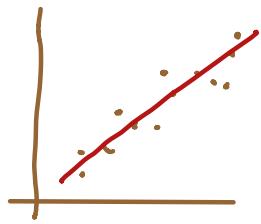
SCATTER - a plot of (x_i, y_i)
 $i = 1, \dots, n$



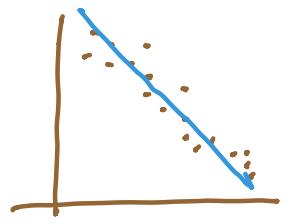
If points lie close to a line



linear correlation
between X, Y



POSITIVE



NEGATIVE

- Association measure:
sample correlation coefficient

$(X_1, Y_1), \dots, (X_n, Y_n)$ i.i.d

$$r_{xy} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(n-1) \sqrt{s_x^2} \sqrt{s_y^2}}$$

- Interpretation: strength of association
 $-1 \leq r_{xy} \leq +1$

$$1) r_{xy} = 1 : (x_i, y_i) \text{ on } y = \bar{y} + \frac{s_y}{s_x} (x - \bar{x})$$

$$2) r_{xy} = -1 : (x_i, y_i) \text{ on } y = \bar{y} - \frac{s_x}{s_y} (x - \bar{x})$$

3) If $(X_1, Y_1), \dots, (X_n, Y_n)$ iid from (X, Y)
and X independent Y

then

$$\rho_{xy} \approx 0$$

4) $|\rho_{xy}| \approx 0 \not\Rightarrow X \perp\!\!\!\perp Y$

See book's example.

• AUTO CORRELATION

In time series $(x_1, \dots, x_h, \dots, x_T)$
observations are not i.i.d
(typically)

this can be verified

1) With a Scatter

SCATTER of $(x_i, x_{i+\tau})$
 $i = 1, 2, \dots, n-1$

should have almost no structure.

2) with an index of autocorrelation

INDEX of autocorrelation of order
 $h = 0, 1, 2, \dots$

$$r_x(h) = \frac{\sum_{i=1}^{n-h} (x_{i+h} - \bar{x})(x_i - \bar{x})}{(n-h) s_x^2}$$

measures correlation
between x_i and x_{i+h} .