

r uniform → simulations
of a $U(0, 1)$

p uniform → evaluation of the
distribution of

d uniform → evaluation of the density function

help.start()

r uniform(10)

[1]

△ > set.seed(3)
> r uniform(10)
> set.seed(3)
> r uniform(10)

- ? Simulate 10 samples of size 100 from $\mathcal{B}(0.2)$
- a) with the R function \Rightarrow a matrix
 - b) with the runif function
 - \hookrightarrow create a function with several arguments (n : the size of the sample
 R : the number of samples)

a)
 $A = \text{matrix}(\text{data} = \text{rbinom}(1000, 1, 0.2),$
 $nrow = 10)$

100 $\left(\begin{array}{c} \\ \vdots \\ \end{array} \right)^{10}$

$\text{rbinom}(m, n, p) \rightarrow m \times \text{simulation}$
of $\mathcal{B}(n; p)$

$$\mathcal{B}(p) = \mathcal{B}(1; p)$$

$A[1:10,:]$ → extraction of
for rows for column(s) the rows number 1, 2, .., 10

$A[c(1,3,4)]$ → extraction of the
rows number 1, 3 and 4

$c(1,3,4)$ → create a vector of length 3
whose elements are 1, 3, 4

$c(\underbrace{\text{seq}(1, 2, 0.1)}, \underbrace{\text{seq}(f, 5, -0.5)})$

1 vector 1 vector

$A[, 3:4]$ → extraction of the columns number 3 and 4

$A[2:4 | 1:3]$ → extraction of the observations of A which are on line 2, 3, 4 and column 1, 2, 3.

$A[-2, 1] \rightarrow$ we delete the
row number 2

$$E = A[-c(2, 4), 1] \quad a = 2$$

$$J = A[-2, 3] \quad a < -2$$

$$K = J[-4, 1] \quad 2 \rightarrow a$$

$$\Delta K \neq E$$

data frame \Rightarrow quite similar to
a matrix

list \Rightarrow quite similar to
a vector

`l = list(a = seq(1, 10, 0.5), b = f)`

`l[[1]] \Leftrightarrow l$a`

```
{:Sample(c("blue", "red", "yellow"),  
500, replace = TRUE,  
prob = c(0.1, 0.7, 0.2))}
```

P2 = factor(e)

$\ell_2 \rightarrow$ is a factor
associated to qualitative

values : $\overset{\wedge}{\text{red}}$ $\overset{\wedge}{\text{blue}}$ $\overset{\wedge}{\text{yellow}}$

$c(\ell_2, 1:10)$ bad
 $\text{list}(c(\ell_2, 1:10))$

```
K=matrix(data = 0, ncol = 2, nrow = 3)
```

```
K[, 1] = l2[1:3]
```

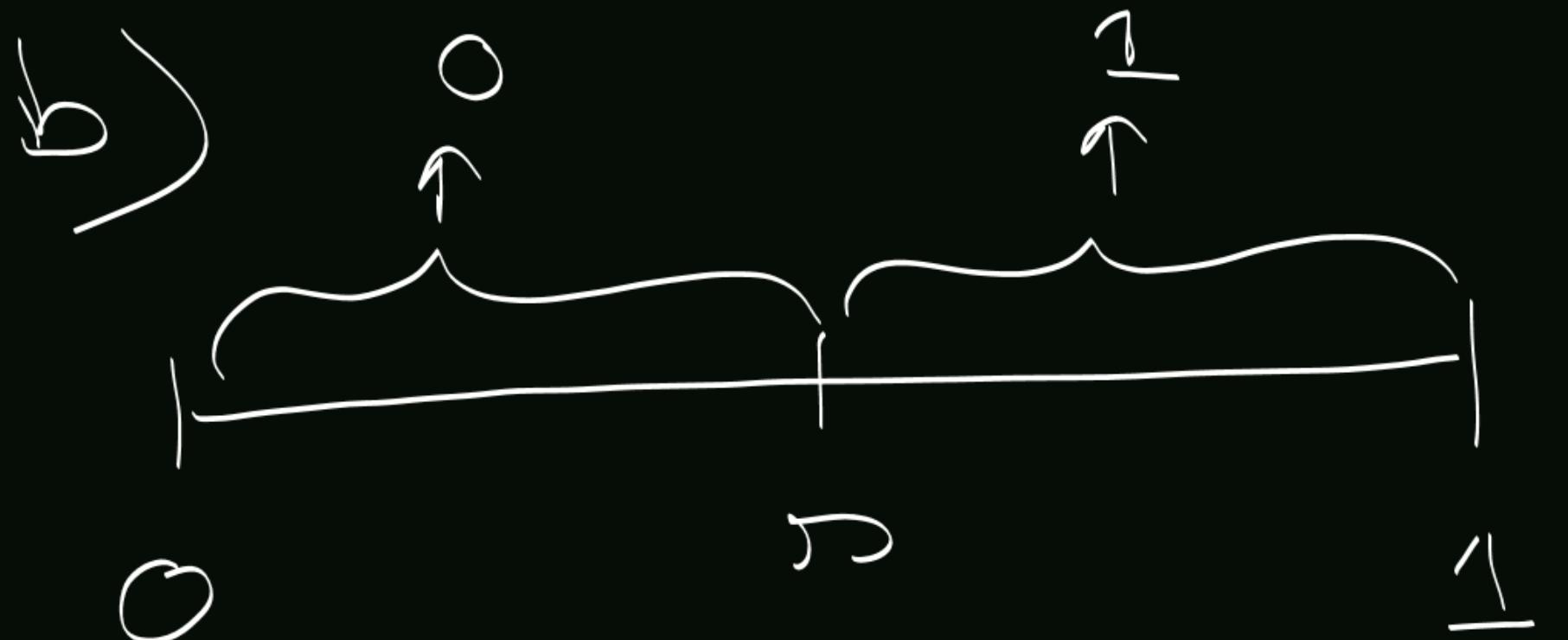
```
K[, 2] = 1:3
```

```
is.matrix(K)
```

```
K
```

R B :

B = matrix (data = 0 , ncol = 10 ,
nrow = 10)



$$P(X=1) = P$$

$$P(X=1) = P(U \geq s) = \begin{cases} 1 & \text{if } U < s \Rightarrow X=0 \\ 1-s & \text{if } U \geq s \Rightarrow X=1 \end{cases}$$

let $U \sim U(0,1)$

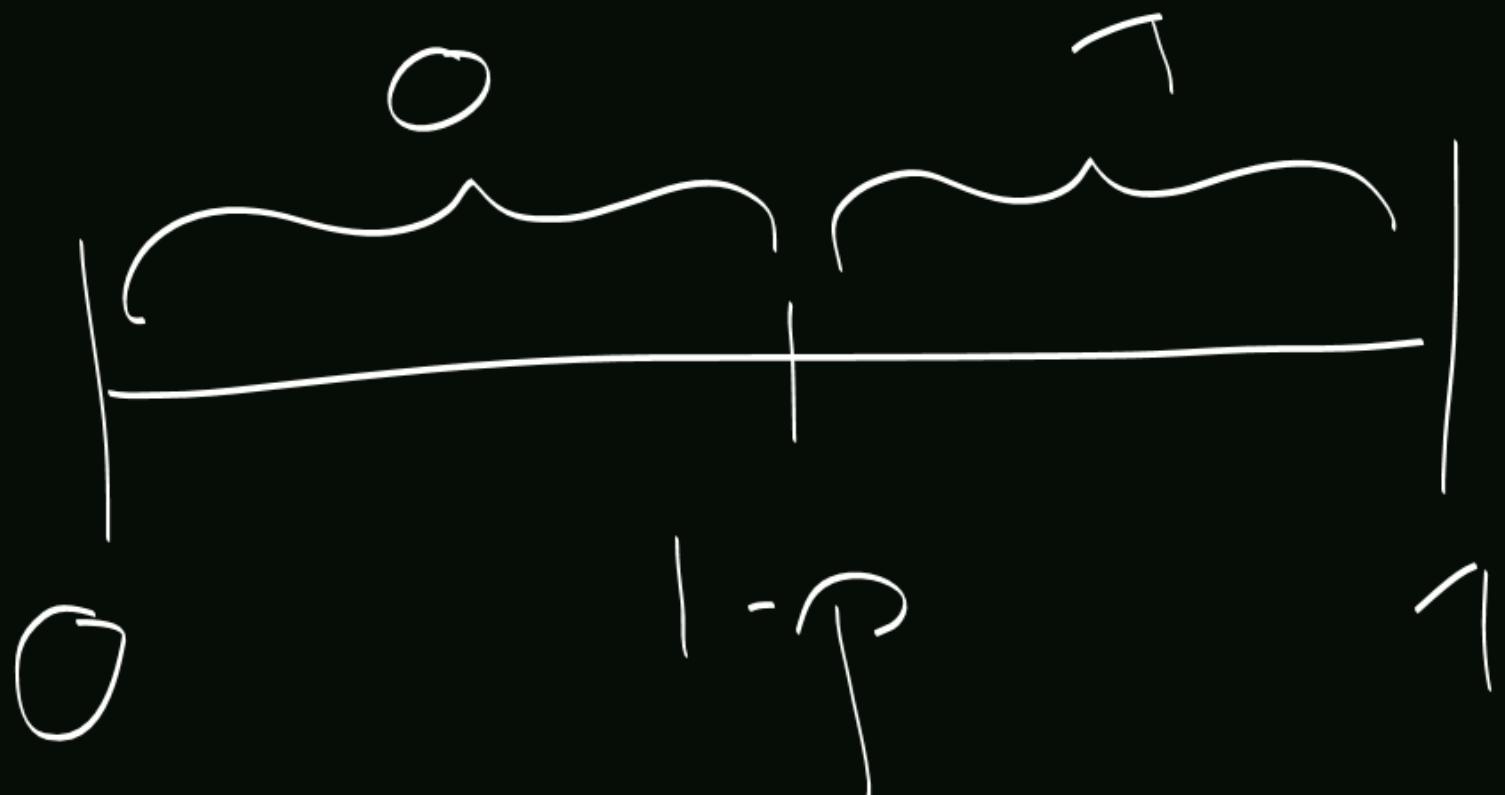
We create X

by:

$\text{if } U < s \Rightarrow X=0$

$\text{if } U \geq s \Rightarrow X=1$

$$\Rightarrow \sigma = 1 - P$$



I simulate u

with respect to

$$a \sim U([0; 1])$$

if $u < 1 - P$, I
create $x = 0$

otherwise I create

$$x = 1$$

```
SimNumber <- function (n, P)
{
  x = c() # empty vector
  for (i in 1:n)
  {
    u = runif(1)
    if (u < (1-P)) {x = c(x, 0)}
    else {x = c(x, 1)}
  }
}
```

simuler = ∞

}

Chapter: Descriptive statistic for variable in dimension 1

I Introduction

We speak about a variable for which we have some observation.

Population : all the "statistical"
individuals that can be
concerned by the study.

Sample : a subset of the population
which is really concerned by
the study

observations \Rightarrow determine the nature
of the variable

qualitative

Quantitative

| categorical | ordinal | discrete | continuous |
|--------------------------------|--|--|---|
| no notion of range | notion of rank in the value | finite set of values or in bijection with \mathbb{N} | The set of values can just be expressed as interval |
| ex: color of eyes, professions | ex: judgment on movie, grade at bachelor (VS, B, TB) | ex: number of children, . | ex: height, weight, |

II Graphics

a) Categorical variable

table of representation

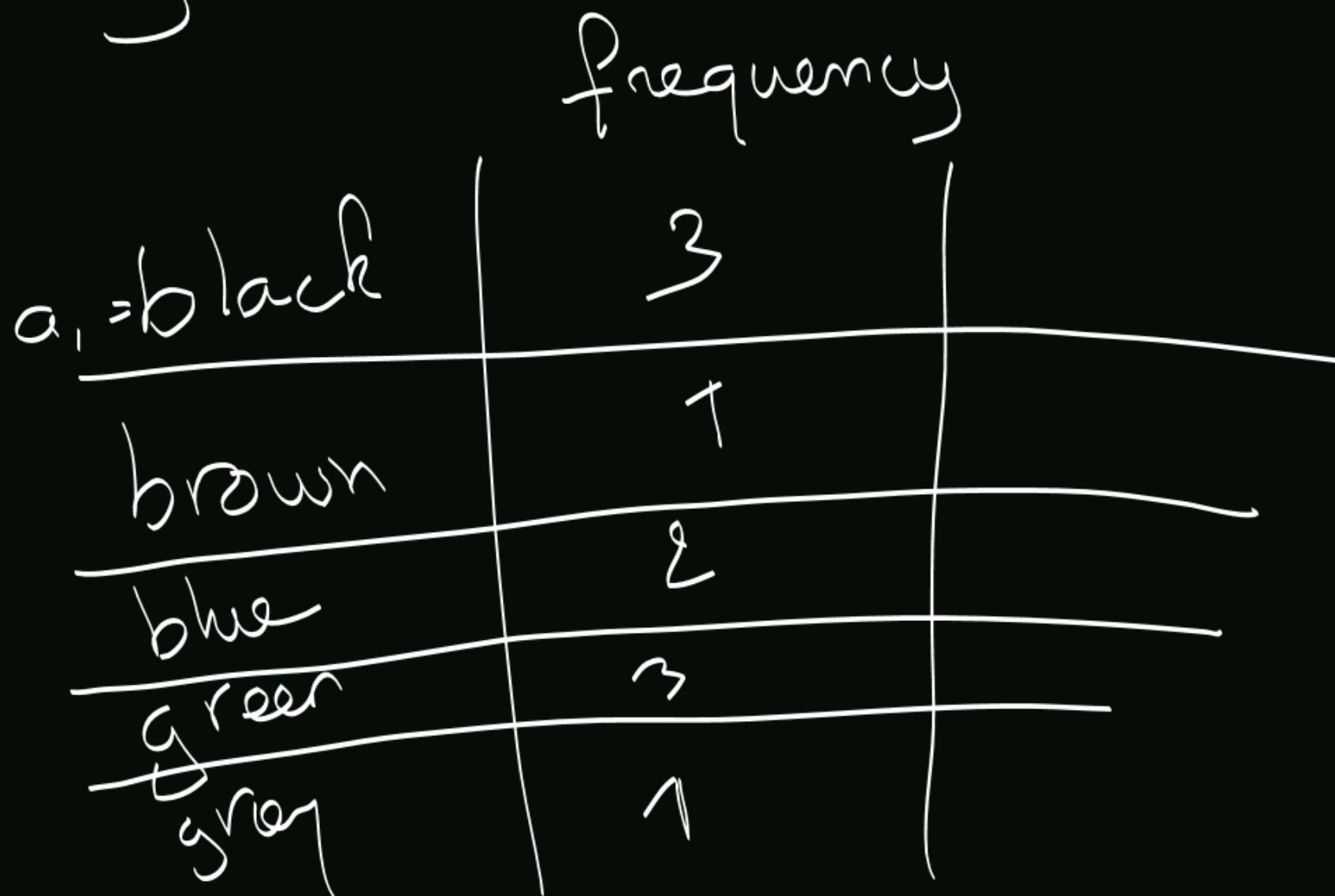
| modalities | Frequency | percent | angle |
|------------|-----------|-----------------|-------------------|
| a_1 | n_1 | f_1 | α_1 |
| a_2 | n_2 | f_2 | α_2 |
| : | | | : |
| a_K | n_K | f_K | α_K |
| | | $\frac{f_K}{n}$ | $\frac{360}{360}$ |

let x_1, \dots, x_n
be the observations

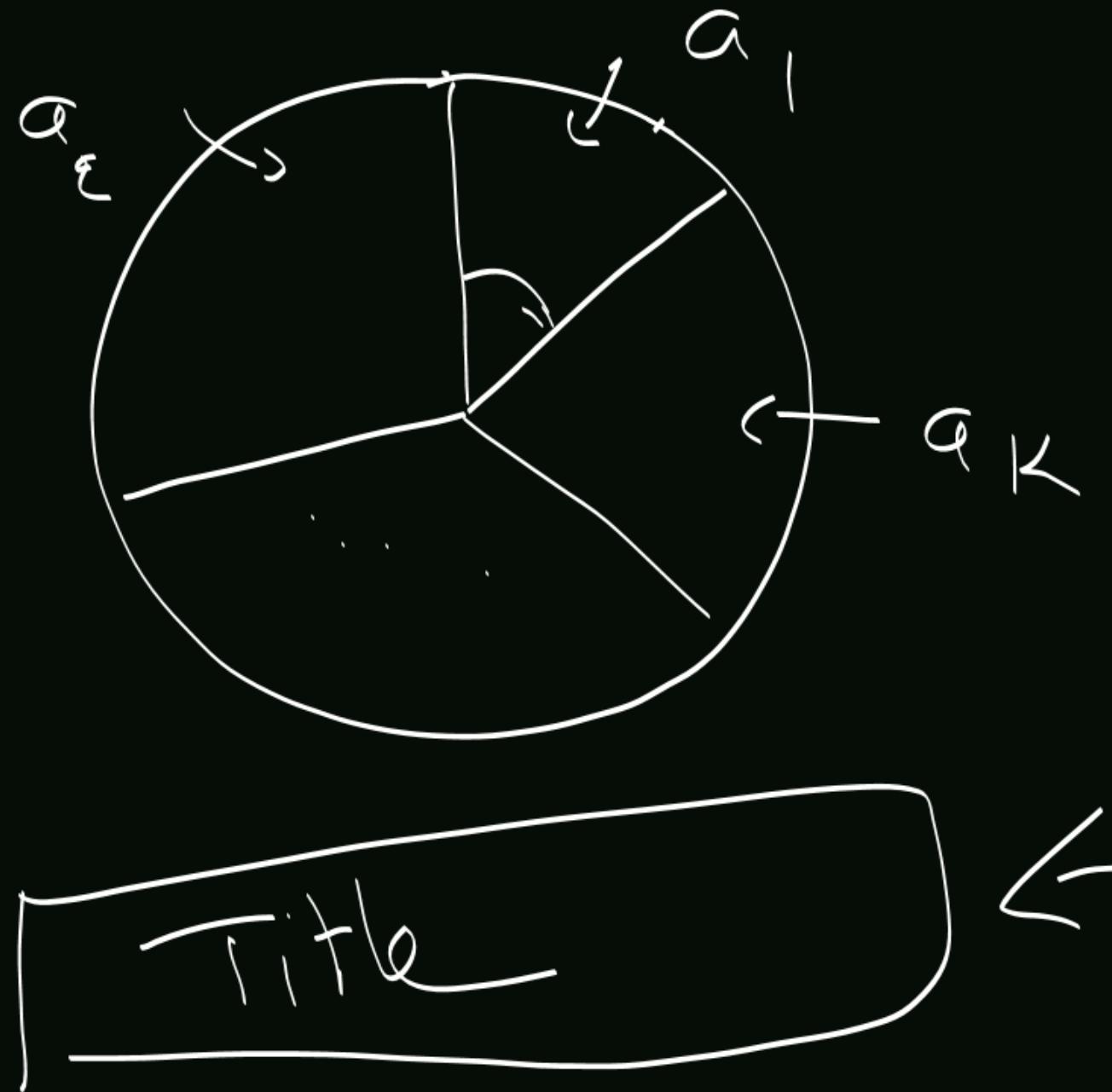
$$\begin{aligned} n_i &= \text{number of} \\ &j \in \{1, \dots, n\} \text{ such} \\ &\text{that } x_j = a_i \\ f_i &= \frac{n_i}{n} \in [0, 1] \\ \alpha_i &= 360 \times f_i \end{aligned}$$

ex:

black black brown blue green
green blue grey black green



The good representation is:



← not forget-

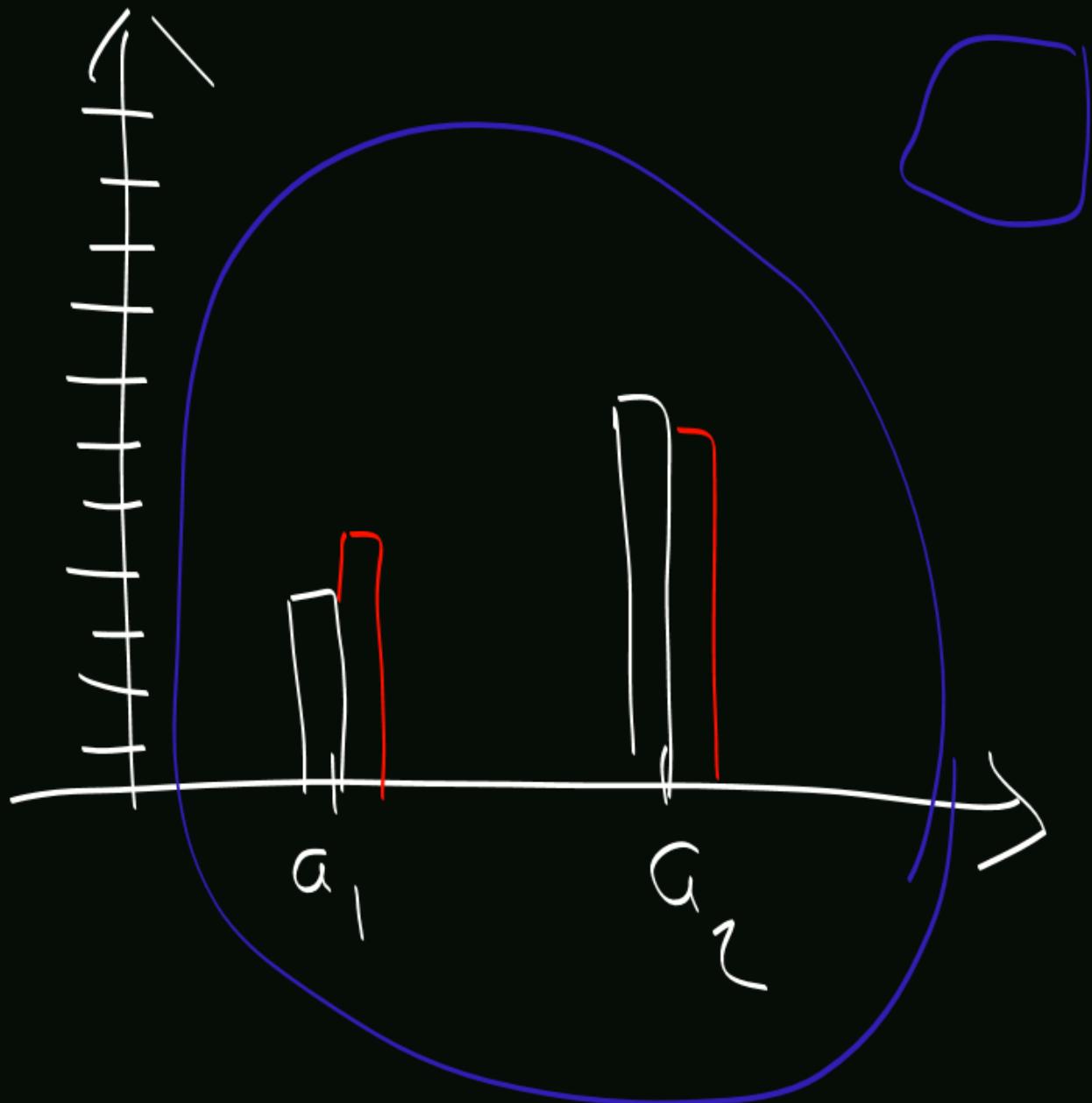
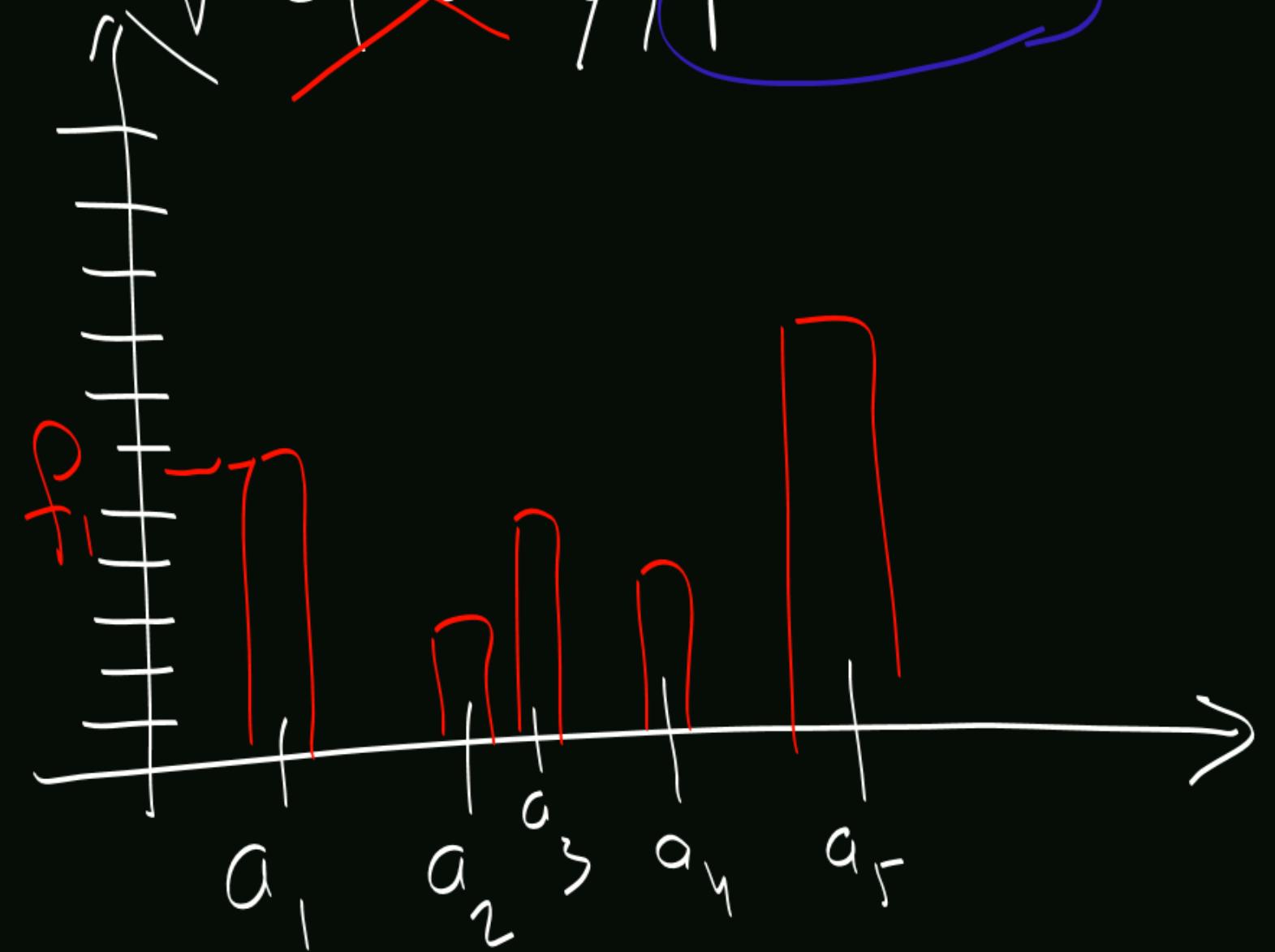
b) ordinal Variable

Table of representation

| modalities | frequency | percent |
|------------|-----------|--------------------|
| a_1 | n_1 | f_1 |
| a_2 | n_2 | f_2 |
| : | : | : |
| a_K | n_K | $\frac{f_K}{\sum}$ |
| | | $\frac{n}{n}$ |

The good representation

~~frequency~~, percent

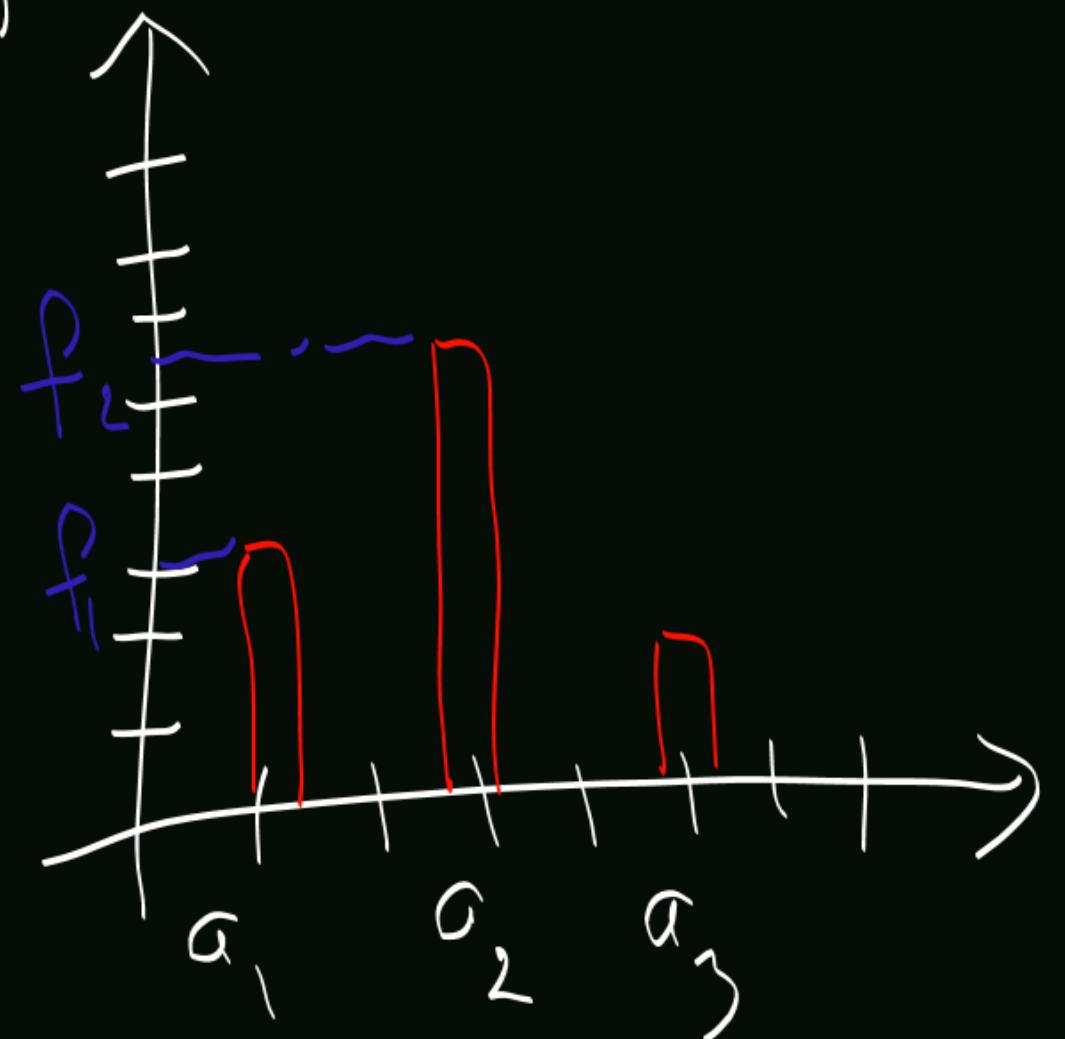


c) Discrete variable

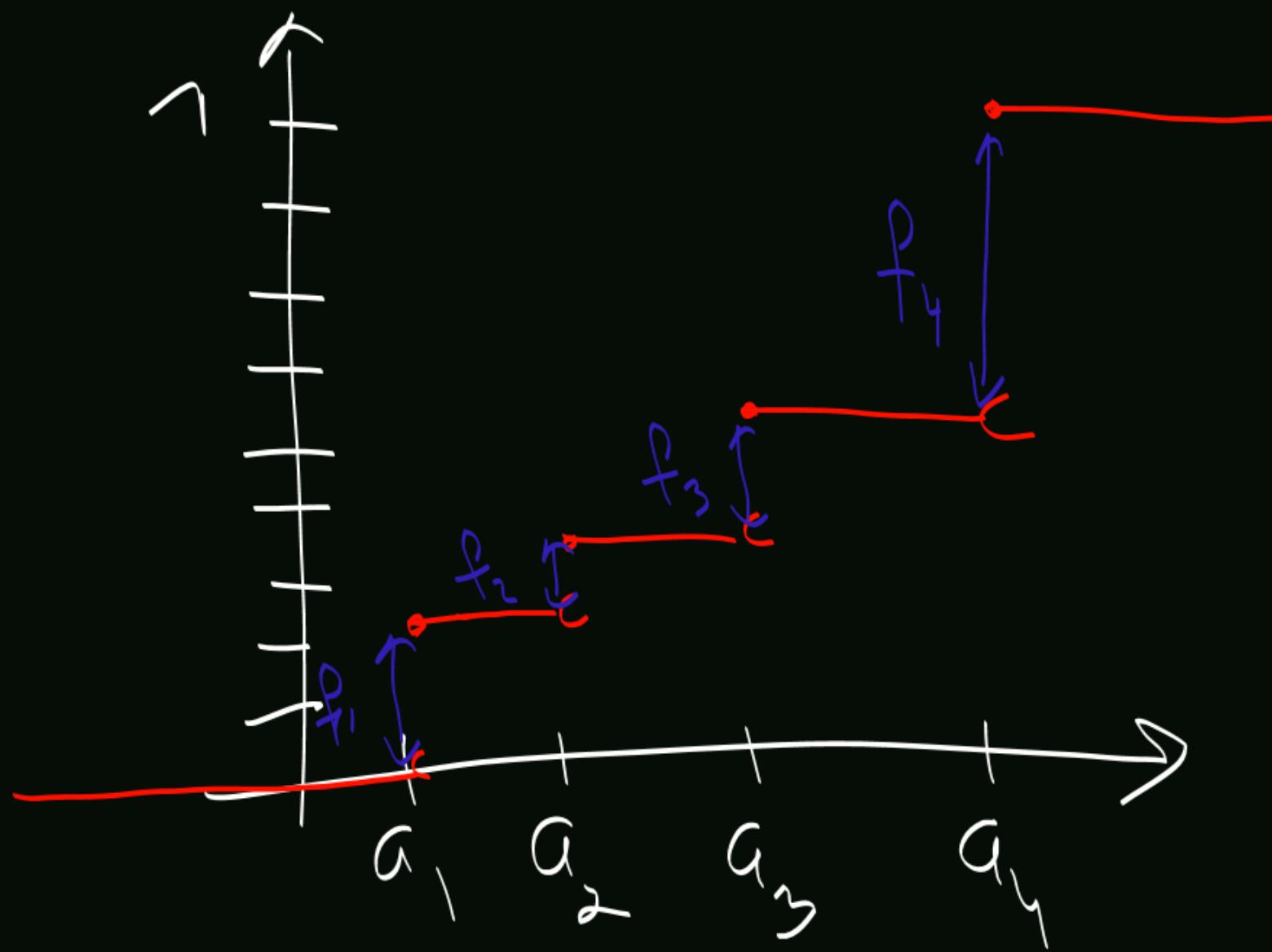
x_1, \dots, x_n : n observations

| Values | Frequency | Percent | cumulative percent |
|--------|-----------|-----------------|--------------------|
| a_1 | n_1 | f_1 | $F_1 = f_1$ |
| | n_2 | f_2 | $F_2 = f_1 + f_2$ |
| : | : | | |
| a_K | n_K | $\frac{f_K}{n}$ | $F_K = 1$ |
| | | | |

Graphics:
barplot
percent



plot of cumulative percent



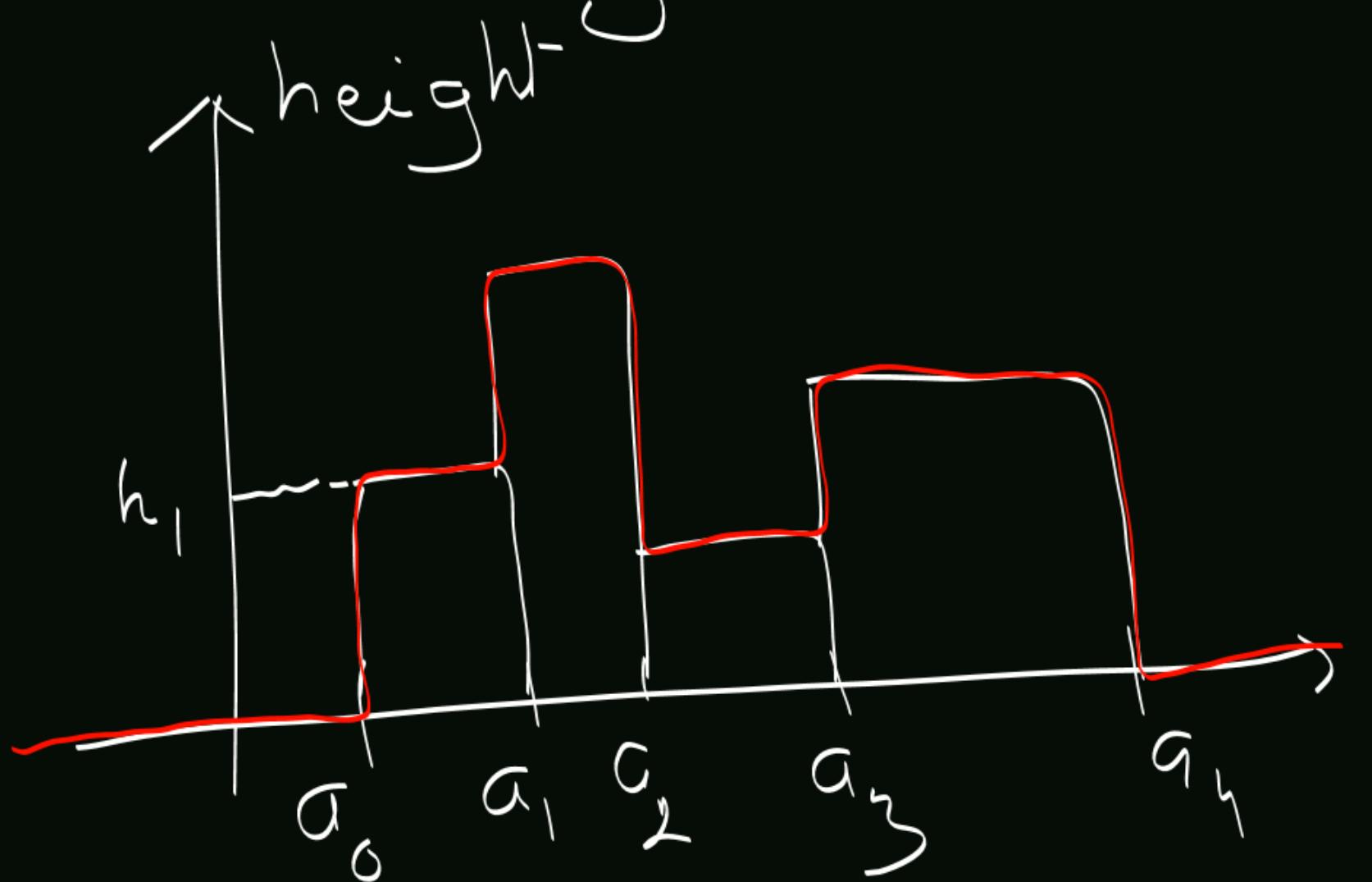
d) Continuous

| classes | frequency | percent | cumulative percent | height |
|------------------|----------------|----------|--------------------|----------|
| $[a_0; a_1[$ | n_1 | f_1 | $F_1 = f_1$ | h_1 |
| $[a_1; a_2[$ | n_2 | f_2 | $F_2 = F_1 + f_2$ | h_2 |
| \vdots | \vdots | \vdots | \vdots | \vdots |
| $[a_{K-1}; a_K[$ | n_K | f_K | $F_K = 1$ | h_K |
| | \overline{n} | | | |

Graphics

histogram

height



area of the histogram

$$= \sum_{i=1}^K (a_i - a_{i-1}) \times h_i$$

$$= \sum_{i=1}^K (a_i - a_{i-1}) \times \frac{f_i}{(a_i - a_{i-1})}$$

$$= 1$$

2 questions:

- How many classes to create?
- How to determine the limits of the classes?

Sturge's method :

i) if n denotes the number of observations, the number K of classes is defined by:

$$K \approx 1 + 3.22 \times \log_{10}(n)$$

$$\text{range} = \frac{a_K - a_0}{K}$$

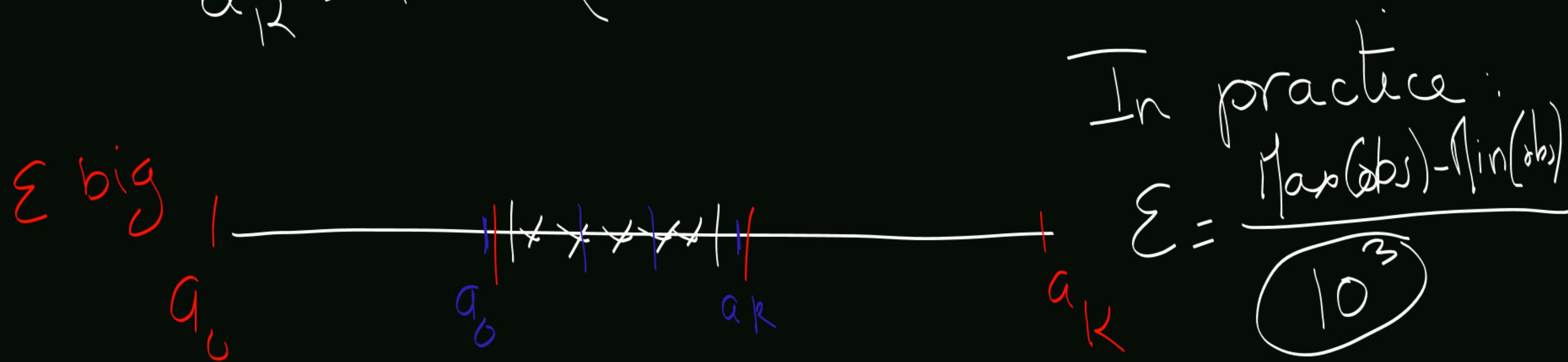
$$\forall i \in \{1, \dots, K-1\}, a_i = a_0 + i \times \text{range}$$



2) We choose ε , with ε small
with respect to the data.

$$a_0 = \min(\text{observations}) - \varepsilon$$

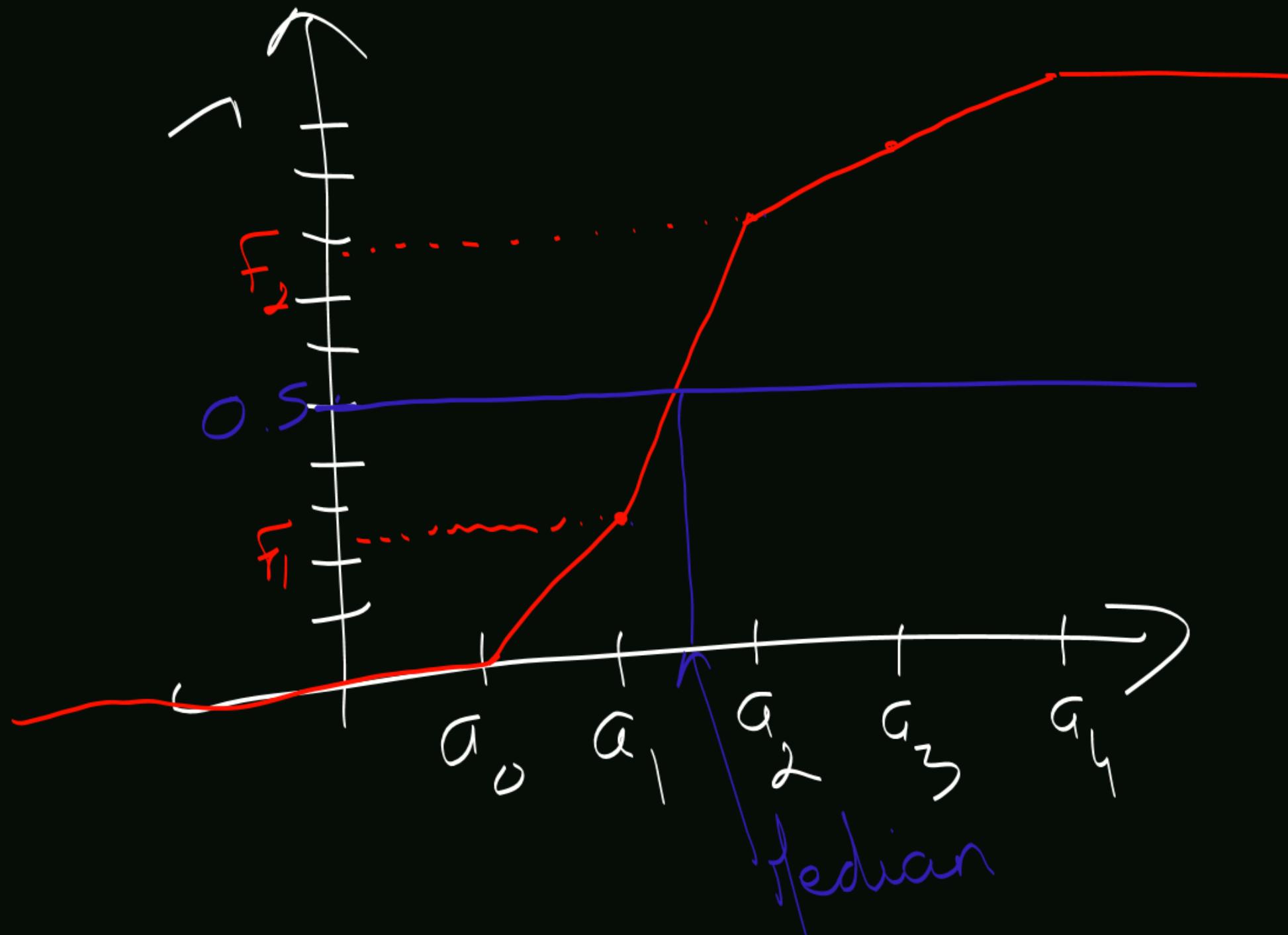
$$a_K = \max(\text{observations}) + \varepsilon$$



Rk:

The density f^{cl^-} is an extrapolation
of the histogram.

plot of cumulative percent



\Rightarrow assumption:
The observations
are uniform
on each class!

? Take the datasets data1, data2,
data3, data4
and do the good graphics

? Take the file named extra and
try to illustrate the fact that
the density is the limit of an histogram

To load the data in the

R console, use the
read.table function

> setwd() → to change the work directory
> getwd()

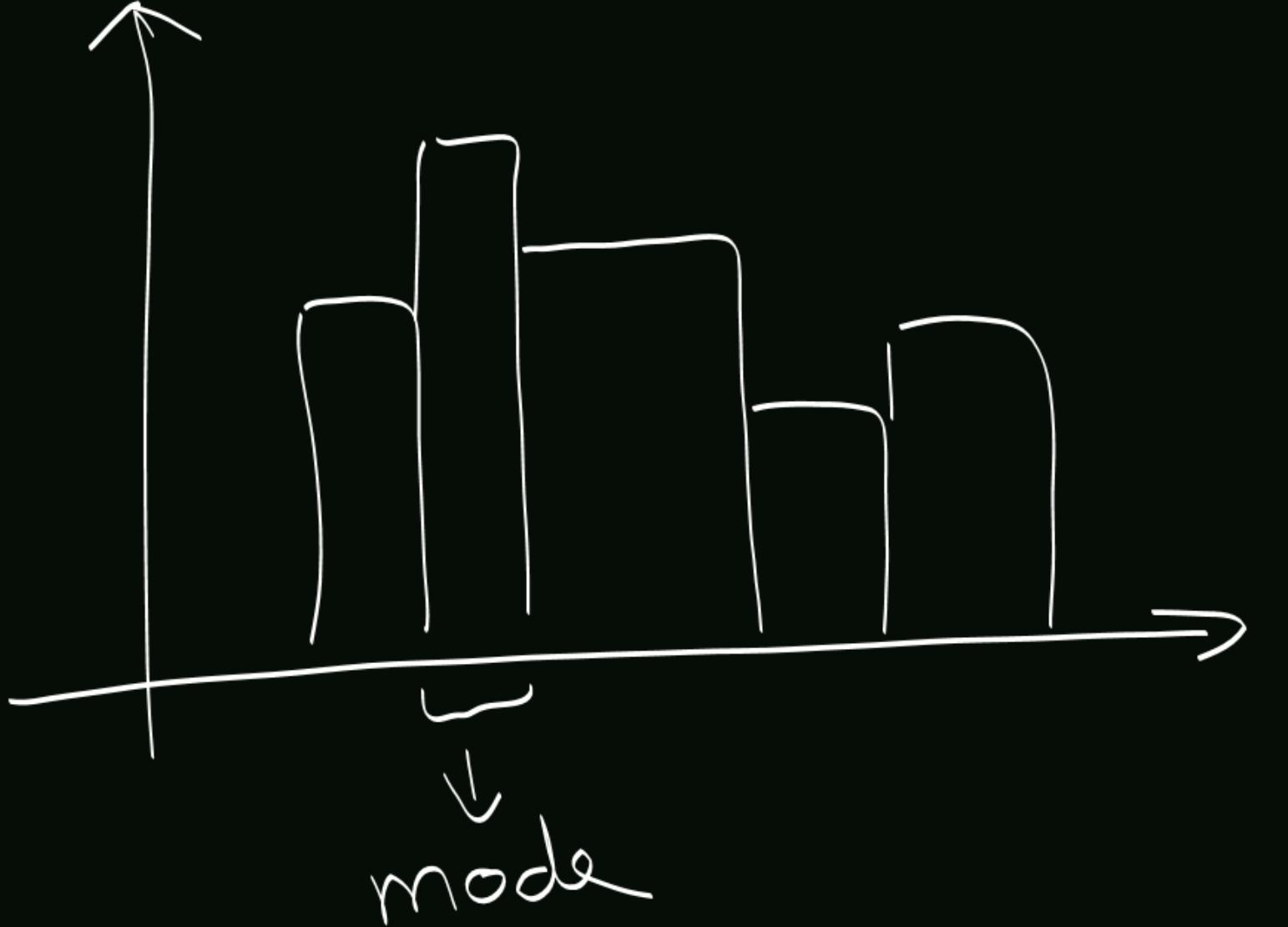
IV Statistical parameters

- parameters of position (quantile)
- dispersion (variance/sd, range, coefficient of variation)
- centrality (mean, median, mode)

a) Parameters of locality

→ all is done for a quantitative variable

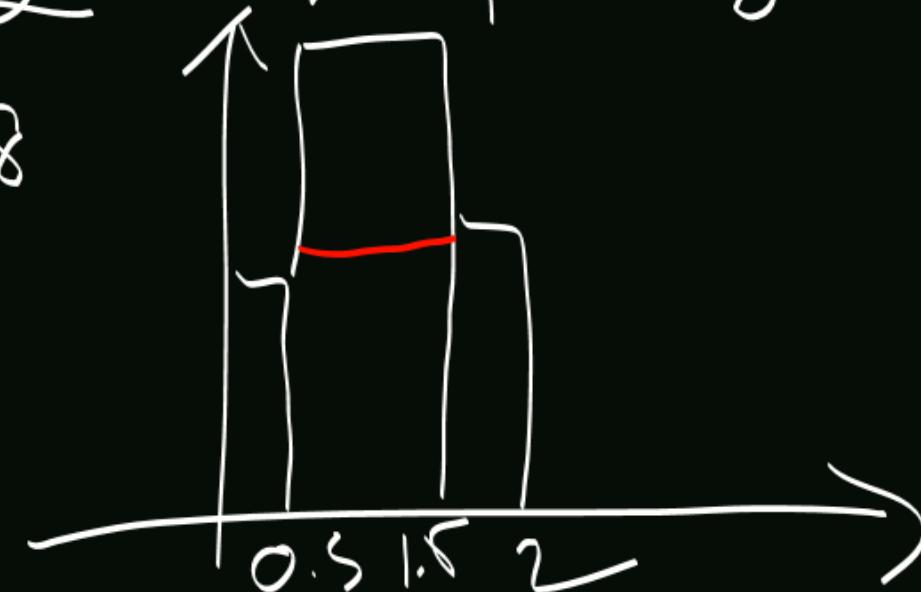
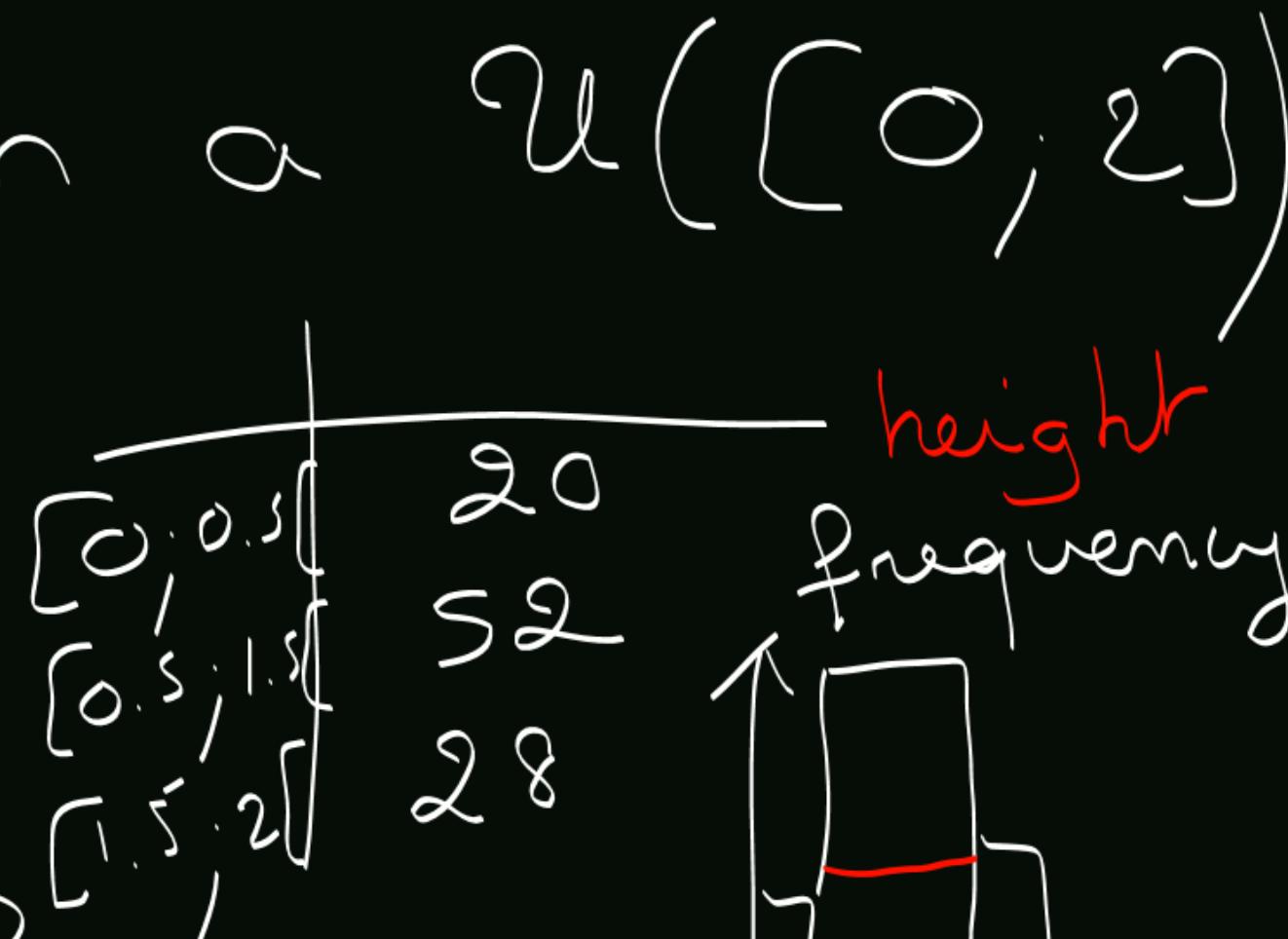
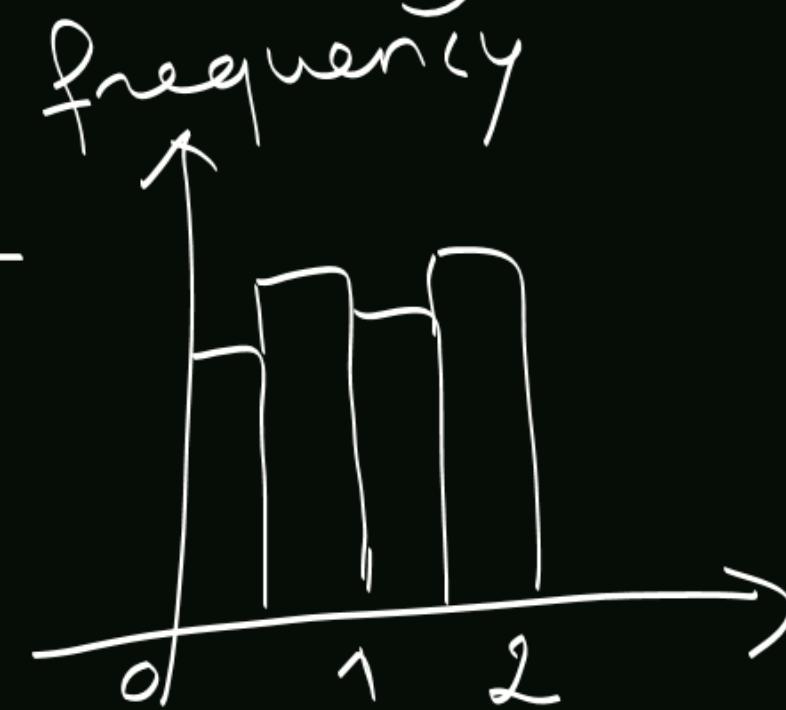
mode = $\begin{cases} \text{value associated to the biggest frequency} & \text{discrete} \\ \text{class associated to the biggest height} & \text{continuous} \end{cases}$



Rk:

Suppose you have observations
that come from a $U([0, 2])$

| | |
|------------|----|
| $[0; 0.5[$ | 20 |
| $[0.5; 1[$ | 28 |
| $[1; 1.5[$ | 24 |
| $[1.5; 2]$ | 28 |



- mean

2 formula:

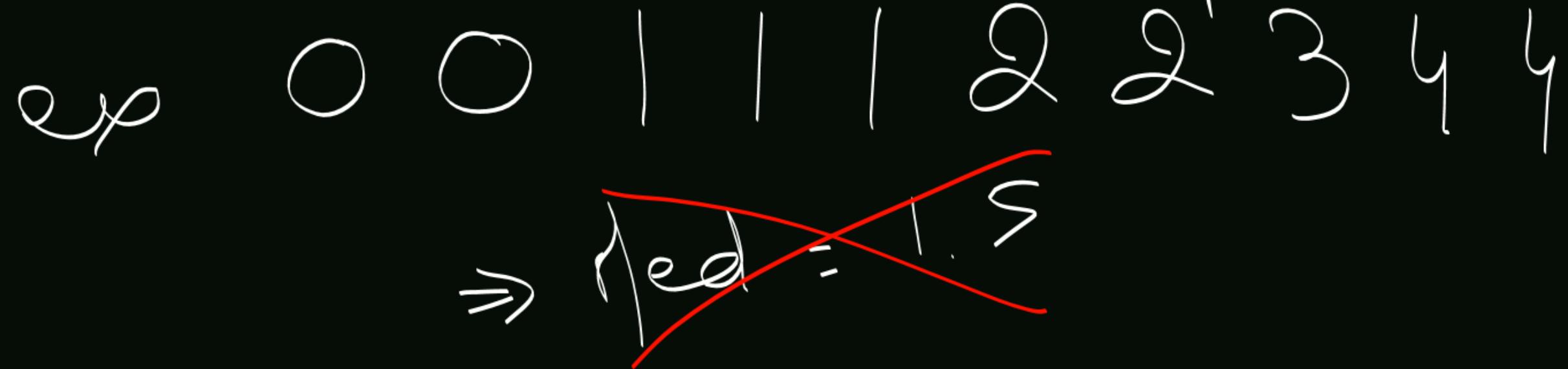
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{i=1}^K n_i c_i$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^K n_i c_i$$

where c_i is the middle
of the class $[a_{i-1}, a_i]$

- median

Property : The median should
be a possible value.



- discrete

Median is the value associated
to cumulative percent which is
immediately bigger or equal to 0.5.

| values | cumulative values |
|--------|----------------------|
| a_1 | 0.1 |
| a_2 | 0.4 |
| a_3 | 0.6 |
| a_4 | 0.7 |
| a_5 | 1 |

median = a_3

- continuous

- * to find the find the class where the median is
This class is the one associated to the cumulative percent which is immediately bigger or equal to 0.5

let F be the function associated to the cumulative percent.

Let $[a, b]$ the class where the

median is

$$F(a) < 0.5 \leq F(\text{Median})$$

- to evaluate the median,
we perform some interpolation
by a linear function on $[a, b[$

$$\begin{cases} SF(a) = P \times a + q & P \text{ and } q \\ F(b) = P \times b + q \end{cases}$$

P and q
are
unknown!

We solve :

$$O.S = P \times \text{Median} + q$$

We find the Median

$$\Rightarrow \text{Median} = a + (b - a) \times \frac{O.S - F(a)}{F(b) - F(a)}$$