


## Lesson 2

**A brief review of concepts related to random variables.**

# Lesson 2–Plan

- 1 The concept of random variable and the probability
- 2 Parameters of a random variable
- 3 Parameters in random pairs
- 4 Some discrete models
  - The uniform discrete distribution
  - The binomial distribution
  - The binomial negative distribution
  - The Poisson distribution
- 5 Some continuous models
  - The normal or Gaussian distribution
  - The Central Limit Theorem
  - The uniform continuous distribution
  - The exponential and gamma distributions
  - The beta distribution
- 6 Example of an exercise in the 

# The concept of random variable and the probability

When performing a random experiment, **one (or more) real values** can be associated with each experiment result - we say we have defined a **random variable** or (**a random vector**).

A **random variable** is usually represented by  **$X$** .

A random variable may be:

- **discrete** - for example the number of germinated seeds; registration, at regular intervals, of the number of persons waiting in a queue of a supermarket;
- **continuous** - for example the weight of a subject; the diameter at the height chest of a tree, the length of a sheet.

# Random variable – probability

Associated with each random variable (r.v.) there are:

- a **probability mass function**, if  $X$  **discrete**,

The probability mass function is an **application** that to each value  $x_i \rightarrow p_i = P[X = x_i]$ , satisfying:

$$p_i \geq 0, i = 1, \dots, k \quad \text{e} \quad \sum_{i=1}^k p_i = 1.$$

- or a **density function**, if  $X$  **continuous**.

A function  $f$  is said to be a **density function** if it verifies the conditions:

$$f(x) \geq 0 \quad \forall x \in \mathbb{R}; \quad \int_{-\infty}^{+\infty} f(x) dx = 1$$

# Random variable – probability

Associated with each random variable (r.v.) there is also:

- a real function  $F$ , which is denoted as the **cumulative distribution function** and defined as

$$F(x) = P[X \leq x]$$

If  $X$  is discrete we have  $F(x) = P[X \leq x] = \sum_{x_i \leq x} P[X = x_i]$ , i.e., we have the cumulative probability associated with the variable  $X$  calculated in any  $x \in \mathbb{R}$ .

If  $X$  is continuous we have  $F(x) = P[X \leq x] = \int_{-\infty}^x f(t) dt$   $-\infty < x < \infty$ , where  $f$  is the density function.

Examples of how to calculate a probability, using  $F$ :

- 1  $P(a < X \leq b) = P(X \leq b) - P(X \leq a) = F(b) - F(a)$ ;
- 2  $P(X = a) = F(a) - F(a^-)$  onde  $F(a^-) = \lim_{x \rightarrow a^-} F(x)$
- 3  $P(a < X < b) = P(X < b) - P(X \leq a) = F(b^-) - F(a)$ ;

# Parameters of a random variable

## Expected Value

Given a r.v.  $X$  the **mean value** or **expected value** is denoted as  $E[X]$ ,  $\mu_X$  or simply  $\mu$  and is defined as

$$E[X] = \sum_{i=1}^n x_i p_i \quad X \text{ discrete r.v. with distribution } (x_i, p_i)$$

$$E[X] = \int_{-\infty}^{+\infty} x f(x) dx \quad X \text{ continuous r.v. with density } f(x)$$

## Some properties

- $E[a + bX] = a + b E[X]$ .
- $E[\varphi(X) + \psi(X)] = E[\varphi(X)] + E[\psi(X)]$
- $\inf(X) \leq E[X] \leq \sup(X)$

# Parameters of a random variable

## Variance

The **variance** of a random variable  $X$  is denoted as  $\text{Var}[X]$ ,  $\sigma_X^2$  or  $\sigma^2$  and is defined as

$$\sigma_X^2 = E[(X - \mu)^2]$$

The  $\sigma_X = \sqrt{\text{Var}[X]}$  is **the standard deviation**.

## Some properties

- $\text{Var}[X] = E[X^2] - (E[X])^2$
- $\text{Var}[X] \geq 0$
- $\text{Var}[a + b X] = b^2 \text{Var}[X]$ .

For the **standard deviation** we have  $\sigma_{(a+b X)} = |b| \sigma_X$



# Parameters in random pairs

## Brief review of parameter properties in random pairs

If  $(X, Y)$  is a random pair, that can be **discrete** or **continuous**

### Expected value

Given the random pair  $(X, Y)$ , and  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ , we define

$$E[g(X, Y)] = \sum_i \sum_j g(x_i, y_j) p_{ij}, \quad \text{discrete case}$$

$$E[g(X, Y)] = \int \int_{\mathbb{R}^2} g(x, y) f(x, y) dx dy, \quad \text{continuous case.}$$

# Parameters in random pairs

## Properties of the Mean Value

- $E[X \pm Y] = E[X] \pm E[Y]$
- **Desigualdade de Schwarz** If  $E[X^2]$  and  $E[Y^2]$  exist then  $(E[XY])^2 \leq E[X^2]E[Y^2]$ .

**Corollary:**  $(E[X])^2 \leq E[X^2]$

**Remark:** if  $E[X^2]$  exists  $\implies$  then  $E[X]$  also exists.

- If  $X$  and  $Y$  are independent random variables



$$E[XY] = E[X]E[Y]$$

**Remark: The reciprocal is not true**

# The covariance

## The covariance between $X$ e $Y$

Given the random pair  $(X, Y)$  the **covariance** between  $X$  e  $Y$  is

$$\mathbf{Cov}[X, Y] \equiv \sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)]$$

## Exercise

Show that  $\mathbf{Cov}[X, Y] = E[XY] - E[X]E[Y]$

# Variance and covariance properties



$$\text{Var}[X \pm Y] = \text{Var}[X] + \text{Var}[Y] \pm 2\text{Cov}[X, Y]$$

- If  $X$  e  $Y$  are **independent** random variables  $\implies \text{Cov}[X, Y] = 0$ .

**Remark: The reciprocal is not true.**

- If  $X$  e  $Y$  are **independent** random variables

$$\text{Var}[X \pm Y] = \text{Var}[X] + \text{Var}[Y]$$

- $\text{Cov}[a + bX, c + dY] = bd \text{Cov}[X, Y]$ .

- $|\text{Cov}[X, Y]| \leq \sigma_X \sigma_Y$ .

- **Correlation coefficient** is defined as:

$$\rho \equiv \rho_{X,Y} = \frac{\text{Cov}[X, Y]}{\sigma_X \sigma_Y} \quad (\sigma_X > 0; \sigma_Y > 0).$$

# The main discrete and continuous probability models.

# Main discrete models

## The uniform discrete distribution

**Definition** A r.v.  $X$  is said to have a **discrete uniform distribution** if it assumes the values  $x_1, \dots, x_n$  with probabilities  $1/n, \dots, 1/n$ , i.e.  $P(X = x_i) = 1/n, \quad i = 1, \dots, n.$

**Particular case**

$$X = \begin{cases} 1 & 2 & \dots & n \\ 1/n & 1/n & \dots & 1/n \end{cases}$$

### Mean value and variance

$$E[X] = \frac{n+1}{2} \qquad \text{Var}[X] = \frac{n^2-1}{12}$$

Instructions on  to simulate `> sample(v, size, rep=TRUE)`

$v$  vector with the values that the variable can assume

## Function `sample( )`

The function `sample` – allows us to create a random sample from the elements of a vector, **with or without replacement**, with equal probabilities or not.

```
>sample(1:20, 15)
```

15 numbers are randomly selected from 1 to 20 **without replacement**  
the default is “without replacement”.

To select **with replacement** with different probabilities do, for example:

```
>pb<-c(rep(0.1,3), .2, .3, .2);pb  
>sample(1:6,30,rep=T,prob=pb)
```

If the probability is the same it can be omitted.

**Nota:** To generate the same sequence `>set.seed(number)`

# The uniform discrete distribution in R

```
> par(mfrow=c(2,2))
> x1<-sample(1:6,30,rep=T);x1
> dist1<-table(x1);dist1
> plot(dist1)
```

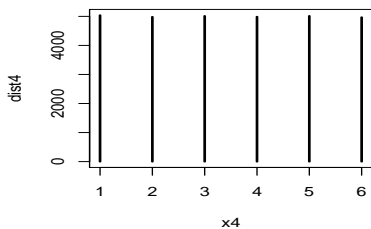
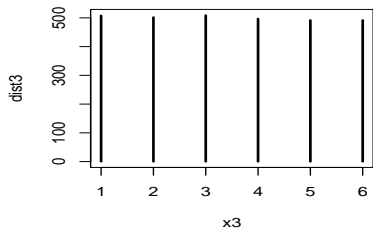
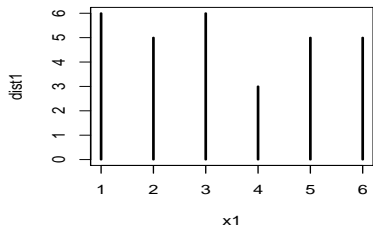
Repeat 300, 3000, and 30000 times (see the graphs of the next *slide* with variables  $x_2$ ,  $x_3$ , and  $x_4$ );

Remark: Defining a function, for example:

```
> dado<-function(n) sample(1:6,n,replace=T)
> d1<-dado(30);table(d1)
> table(dado(30)) # Do you see any difference?
> dado(300);dado(3000)
```



# Graphics from multiple tosses of a die



# The binomial distribution

When  $n$  Bernoulli **independent** trials are performed, the variable that counts the number of successes that occur is said to have a **binomial distribution** and it is represented by  $X \sim \text{Binom}(n, p)$ , where  $p$  is the **probability of success**. The probability of failure,  $1 - p$ , is usually represented by  $q$ .


$X$  assumes the values  $x = 0, 1, 2, \dots, n$  with probabilities given by

$$P[X = x] = \binom{n}{x} p^x (1 - p)^{n-x}$$

## Mean value and variance

$$E[X] = np$$

$$\text{Var}[X] = np(1 - p) = npq$$

To determine the value of those probabilities, quantiles, or the cumulative distribution function, the  has **pre-defined functions** for many models.

# R functions for existing models

- **d**function ( $x, \dots$ ) - allows to obtain the probability mass function (discrete model) or the density function (continuous model) in  $x$ ;
- **p**function( $q, \dots$ ) - allows to obtain the cumulative distribution function, i.e., returns the probability that the variable is less than or equal to  $q$ ;
- **q**function ( $p, \dots$ ) - allows to calculate the quantile associated to the probability  $p$ ;
- **r**function ( $n, \dots$ ) - allows to generate a sample of  $n$  pseudo-random numbers of the specified model.

Meaning:

**d**ensity, **p**robability, **q**uantile, **r**andom

# Exercises

**Exercise** Let's try to use the functions associated to the binomial model, for example, with  $d, p, q, r$ . Consider a *Binomial* ( $n = 10, p = 0.2$ ).

```
> x<- 0:10
> dbinom(x,size=10,prob=0.2)
> pbinom(3,size=10,prob=0.2,lower.tail = TRUE) # gives P[X<=3]
> qbinom(0.75, size=10, prob=0.2, lower.tail = TRUE)
+ # gives the quantile of probability 0.75
> rbinom(5, size=10, prob=0.2)
> pbinom(3, size=10, prob=0.2, lower.tail = F) #dá P[X>3]
```

The **quantile** is defined as the **smaller value**  $\chi_p$  such that  $F(\chi_p) \geq p$ , being  $F$  the cumulative distribution function.

```
> par(mfrow=c(1,2))
> plot(x,dbinom(x,size=10,prob=0.2),type="h")
> plot(x,dbinom(x,size=10,prob=0.4),type="h")
```

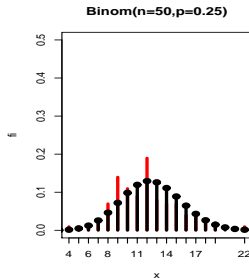
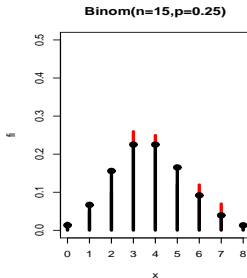
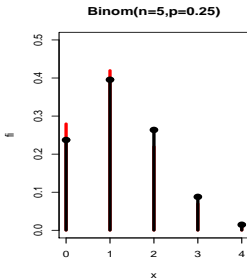
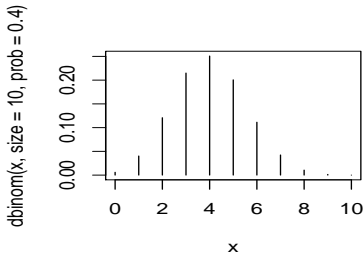
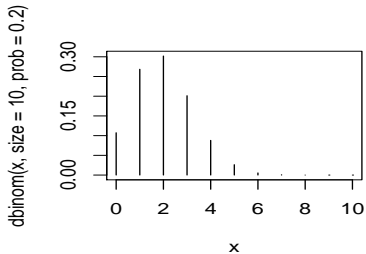
## Exercises (cont.)

To exemplify the **theoretical binomial distribution** and **the simulated one** (with the generation of pseudo-random numbers)


```
> par(mfrow=c(1,3))
> n<-5;p<-0.25
> x<-rbinom(100,n,p) # 100 random numbers
> ni<-table(x);ni
> fi<-ni/sum(ni);fi
> dbinom(0:n,size=5,prob=0.25)
> plot(fi,type = "h", col = "red",lwd=3,
+      main="Binom(n=5,p=0.25)",ylim=c(0,.5))
> xvals<-0:n;points(xvals,dbinom(xvals,n,p),type="h",lwd=3)
> points(xvals,dbinom(xvals,n,p),type="p",lwd=3)
```

... Repeat with  $n=15$ ,  $n=50$ .

# Examples (cont.)



# More probability models

In the  environment the Negative Binomial model is defined as **the number of failures** that are observed until the  $k$  “**success**” is observed, in a context of independent Bernoulli’s trials.

The variable  $X$ , **number of failures** under the above-mentioned conditions is said to have **Negative Binomial distribution** and it is represented by  $X \sim BN(k, p)$

$p$  is the constant probability of “**success**” from trial to trial

$k$  is the number of “**successes**” that we want to get.

# The binomial negative distribution

Characterizing the r.v.  $X \sim BN(k, p)$ :

Values  $x = 0, 1, 2, \dots$

Probabilities  $P[X = x] = \binom{x+k-1}{x} p^k q^x$

$0 < p < 1$ ,  $q = 1 - p$

Mean value and variance of  $X \sim BN(k, p)$

$$E[X] = \frac{kq}{p}$$

$$\text{Var}[X] = \frac{kq}{p^2}$$

Example in 

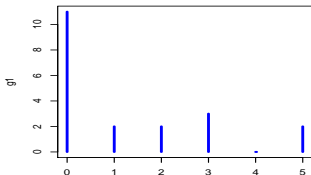
```
> x <- 0:15 #vector of variable values
> dnbinom(x,size=6, prob= 0.4) # probability of 0 to 15 failures
# + until there are 6 successes;
#another parameterization using the above average value
> dnbinom(x, mu = 9, size = 6)
```



# The geometric distribution

If  $k = 1$ , i.e., if we want to determine the **number of failures** to get the **first success**, the variable  $X$  is said to have **geometric distribution**,  $X \sim \text{Geo}(p)$

```
> Ni <- rgeom(20, prob = 1/4)
> g1 <- table(factor(Ni, 0:max(Ni)))
> plot(g1)
```



# The Poisson distribution

## Definition

The r.v.  $X$  that counts the number of successes that occur in a given time interval or domain (independent of the number that occurs in any other disjoint interval or domain) is said to have **Poisson distribution**. It depends only on one parameter  $\lambda \rightarrow$  **average number of successes** that occur in the time interval (or in the specified region).

It is represented by  $X \sim P(\lambda)$  and the law of probability is:

$$P[X = x] = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

# The Poisson distribution

## Mean value and variance

$$E[X] = \lambda \quad \text{Var}[X] = \lambda.$$

Using the 

```
> diff(ppois(c(47, 50), lambda = 50)) # P[47 < X <=50]
> ppois(50,50)-ppois(47,50) # verify that it is the same
```

## The normal or Gaussian distribution

It has a pivotal role in Probability and Statistics because:

- many biometric variables have a form very close to normal;
- sometimes a variable that is not normal can be transformed in a simple way into another with normal distribution;
- the central part of many non-normal models is sometimes reasonably well approximated by a normal distribution.

# Some continuous models

One continuous r.v.  $X$  is said to have a **normal or Gaussian** distribution with parameters  $\mu$  and  $\sigma$  and is represented by  $X \sim \mathcal{N}(\mu, \sigma)$  if the density function is:

$$f(x) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left[ -\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2 \right]$$

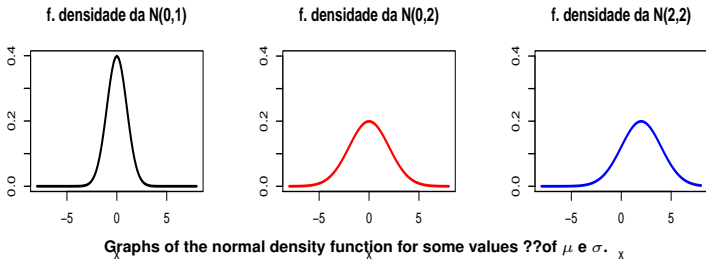
$$-\infty < x < +\infty, \quad -\infty < \mu < +\infty, \quad 0 < \sigma < +\infty$$

# The normal or Gaussian distribution

## Properties of the density curve of the normal distribution

1. It is symmetrical with respect to  $\mu$ .
2. It is an unimodal curve, the mode is  $\mu$ .
3. It has inflection points in  $\mu + \sigma$  e  $\mu - \sigma$ .

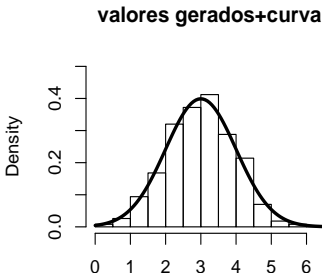
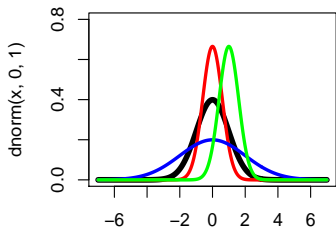
If  $\mu = 0$  and  $\sigma = 1$  the random variable with  $\mathcal{N}(0, 1)$  distribution is called **standard normal** and is usually represented by  $Z$ ,  $Z \sim \mathcal{N}(0, 1)$



```
#calculations and graphs with the normal law
> pnorm(1.96)
> pnorm(-1.96)
> pnorm(3,mean=5,sd=2)
> qnorm(0.75,mean=5,sd=1)
> qnorm(0.75,mean=5,sd=1,lower.tail=T)
> qnorm(0.25,mean=5,sd=1,lower.tail=F)
+           #graficos
> par(mfrow=c(1,2))
> x<-seq(-7,7,.01)
> plot(x,dnorm(x,0,1),type="l",ylim=c(0,.8),lwd=5)
> lines(x,dnorm(x,0,.6),col="red",lwd=3)
> lines(x,dnorm(x,0,2),col="blue",lwd=3)
> lines(x,dnorm(x,1,.6),col="blue",lwd=3)
```

# The normal distribution (graphs)

```
# generating values (cont. exercise)
> y<-rnorm(1000,mean=3,sd=1)
> hist(y,freq=F,ylim=c(0,0.5),
+ main="valores gerados+curva",col=gray(.9))
> curve(dnorm(x,mean=3,sd=1),add=T,lwd=3)
```





# Important results with normal distribution

- Let be  $X \sim \mathcal{N}(\mu, \sigma)$  Then the r.v.  $\frac{X - \mu}{\sigma}$  has a standard normal distribution, i.e.,  $Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1)$ .
- Let  $X_i$   $n$  be r.v. independent, all normal distributed, i.e. having all the same mean value  $\mu$  and the same variance  $\sigma^2$ .

The random variables **sum** and **average**, respectively defined as

$$S_n = \sum_{i=1}^n X_i \quad \text{e} \quad \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

have normal distribution defined as:

$$S_n \sim \mathcal{N}(n\mu, \sigma\sqrt{n}) \quad \text{e} \quad \bar{X}_n \sim \mathcal{N}(\mu, \sigma/\sqrt{n}).$$

# The Central Limit Theorem

We have seen that the sum of independent normal r.v. is still a normal r.v. But **the approximate distribution of the sum of  $n$  random variables** and under certain conditions is also **normal**

## The Central Limit Theorem

Let  $X_1, \dots, X_n$  be independent and identically distributed random variables, with a mean value  $\mu$  and variance  $\sigma^2$  (finite). Se  $n$  'large' the r.v.  $S_n = \sum_{i=1}^n X_i$ , satisfies:

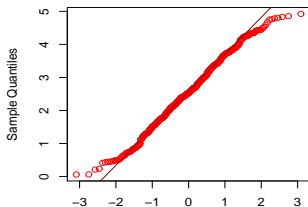
$$\frac{S_n - n\mu}{\sigma\sqrt{n}} \sim \mathcal{N}(0, 1) \quad \text{and we also have} \quad \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1).$$

# The Central Limit Theorem...exercise

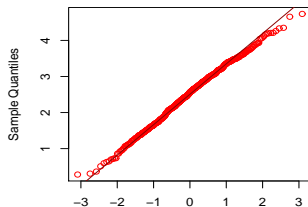
```
> # Uniform distribution(0,5)
> par(mfrow=c(2,2))
> am<-500
> vec.med<-c(rep(0,am))
> n<-c(2,3,10,30)
> for(j in 1:4)
+ {for(i in 1:am)
+ {x<-runif(n[j],0,5)
+ vec.med[i]<-mean(x)}
+ qqnorm(vec.med,main=paste("Q-QPlot Normal, n =",n[j],
+ "n","Médias Pop. U(0,5),"),xlab=" ",
+ col="red")
+ qqline(vec.med,col="darkred")}
```

# The Central Limit Theorem...exercise

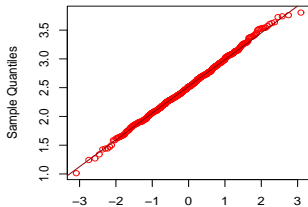
Q-QPlot Normal,  $n = 2$   
Médias Pop.  $U(0,5)$ ,



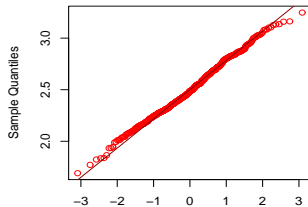
Q-QPlot Normal,  $n = 3$   
Médias Pop.  $U(0,5)$ ,



Q-QPlot Normal,  $n = 10$   
Médias Pop.  $U(0,5)$ ,



Q-QPlot Normal,  $n = 30$   
Médias Pop.  $U(0,5)$ ,



# The Central Limit Theorem – Applications

Let  $X$  be a r.v. with binomial distribution with mean value  $\mu = np$  and variance  $\sigma^2 = npq$ .

$X \sim \mathcal{B}(n, p)$ , i.e., mean value  $\mu = np$  and variance  $\sigma^2 = npq$

$$\frac{X - np}{\sqrt{npq}} \sim \mathcal{N}(0, 1) \quad \text{se} \quad n \rightarrow \infty$$

**Empirical rule** If in the binomial distribution,  $np > 5$  and  $nq > 5 \implies$  the approximation by normal distribution is a good one.

# The Central Limit Theorem – Applications

$$X \sim P(\lambda)$$

If  $\lambda \rightarrow \infty$  then  $\frac{X - \lambda}{\sqrt{\lambda}} \sim \mathcal{N}(0, 1)$ .

## Another convergence

If in the binomial distribution  $n \rightarrow \infty$  and  $p$  is small (let us say  $p < 0.05$  and  $n > 20$ )  $X \sim \mathbf{B}(n, p) \sim \mathbf{P}(np)$

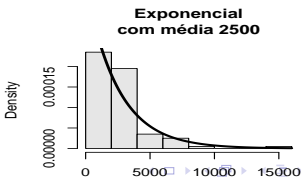
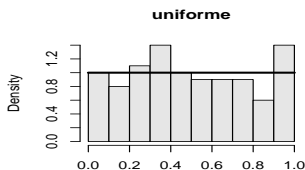
# Other continuous distributions

## Uniform continuous and exponential distribution

```
> u<-runif(100)
> hist(u,freq=F,col=gray(.9),main="uniforme")
> curve(dunif(x),add=T,lwd=3)
```

... and exponential of mean value 2500

```
> x<-rexp(100,1/2500)
> hist(x,probability=TRUE,col=gray(.9),main="Exponencial
+ com média 2500")
> curve(dexp(x,1/2500),add=T)
```



# The gamma distribution

In many areas of sciences there are still many situations in which the Gauss's law does not serve to model the phenomenon.

Let us first briefly refer to the **gamma distribution** who owes his name to **the gamma function**, studied in many areas of mathematics, defined as:

$$\Gamma(\alpha) = \int_0^{+\infty} x^{\alpha-1} e^{-x} dx \quad \text{para } \alpha > 0$$

**Some properties** of the gamma function:

- $\Gamma(\alpha) = (\alpha - 1)\Gamma(\alpha - 1)$  (a recurrence expression)
- When  $\alpha = n$  is a natural number, it is easy to verify that

$$\Gamma(n) = (n - 1)(n - 2)\dots\Gamma(1) = (n - 1)!$$



# The gamma distribution

**Some more properties** of the gamma function:

- $\Gamma(1/2) = \sqrt{\pi}$
- The derivatives of the gamma function are thus defined:

$$\Gamma^{(k)}(\alpha) = \int_0^{\infty} x^{\alpha-1} (\log x)^k e^{-x} dx$$

Some particular values of the derivatives useful in many applications are

$\Gamma'(1) = \gamma = 0.57722\dots$  is the **Euler constant**

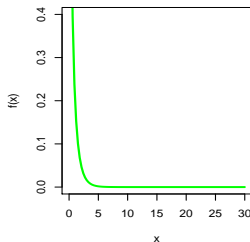
$\Gamma''(1) = \gamma^2 + \pi^2/6 = 1.97811\dots$

# The gamma distribution

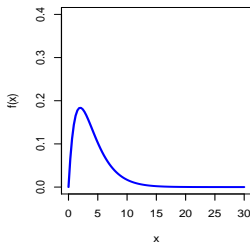
We say that a r.v.  $X$  has a **gamma distribution** with parameters  $\alpha$  e  $\beta$ , ( $\alpha > 0$ ,  $\beta > 0$ ) and we write  $X \sim G(\alpha, \beta)$  with ( $\alpha$  – the shape parameter and  $\beta$  – the scale parameter) if the density function is:

$$f(x) = \begin{cases} \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} & x > 0 \\ 0 & x \leq 0 \end{cases}$$

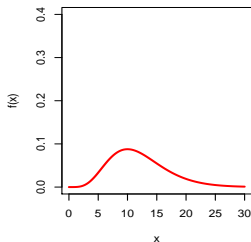
f. dens. da Gamma(0.5,1)



f. dens. da Gamma(2,2)



f. dens. da Gamma(6,2)



Graphs of the density function of a r.v. with distribution  $G(1/2, 1)$ ,  $G(2, 2)$  and  $G(6, 2)$ , from left to right, respectively.

# The gamma distribution

Mean value and variance of  $X \sim G(\alpha, \beta)$

$$E[X] = \alpha \beta \quad \text{Var}[X] = \alpha \beta^2$$

A very important particular case is that one we get by doing  $\alpha = 1$ . The resulting r.v. is said to have **exponential distribution**, is represented by  $X \sim \text{Exp}(\beta)$  and the density function is thus defined,

$$f(x) = \begin{cases} \frac{1}{\beta} e^{-x/\beta} & x > 0 \quad \beta > 0 \\ 0 & x \leq 0 \end{cases}$$

# The exponential distribution

## Mean value and variance

$$E[X] = \beta$$

$$\text{Var}[X] = \beta^2$$

The exponential distribution has been widely used as a model problems related to the duration of life, theory of reliability, waiting times, etc.

### Property

If  $X_i, i = 1, \dots, n$  are independent and identically distributed random variables with  $\text{Exp}(\beta)$ , then

$$\sum_{i=1}^n X_i \sim G(n, \beta).$$

# The exponential distribution

## Remarks:

- There is a very important relationship between the exponential and the Poisson distribution, which often arises in practice. While observing the **occurrence of certain events at time intervals**, we intend to characterize  $T$  the time to the end of which the first occurrence occurs.

## Teorema

Let  $X$  be a Poisson r.v. with parameter  $\lambda$ . Let  $T$  be a r.v. that measures the waiting time for the occurrence of the first event, then  $T$  has an exponential distribution,  $T \sim Exp(\beta)$ , with parameter  $\beta = 1/\lambda$ .

# The beta distribution

One continuous random variable  $X$  is said to have a **beta distribution** with parameters  $(m, n)$  and we write  $X \sim Be(m, n)$  if its density function is of the form

$$f(x) = \begin{cases} \frac{1}{B(m,n)} x^{m-1} (1-x)^{n-1} & 0 < x < 1 \quad m > 0, n > 0 \\ 0 & \text{outros valores de } x \end{cases}$$

where  $B(m, n)$  é a **beta function** so defined

$$B(m, n) = \frac{\Gamma(m)\Gamma(n)}{\Gamma(m+n)} = \int_0^1 x^{m-1} (1-x)^{n-1} dx$$

# The beta distribution

## Properties

1.  $B(m, n) = B(n, m)$

2.  $B(1, 1) = 1$

3.  $B\left(\frac{1}{2}, \frac{1}{2}\right) = \pi$

4.  $B(m, n) = \int_0^{+\infty} \frac{x^{m-1}}{(1+x)^{m+n}} dx$

## Mean value and variance of the beta distribution

$$E[X] = \frac{m}{m+n}$$

$$\text{var}[X] = \frac{mn}{(m+n)^2(m+n+1)}$$

# The beta distribution

The density function of a r.v. with beta distribution presents, as we have said, a great variability of forms.

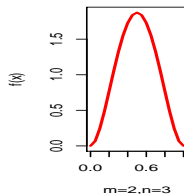
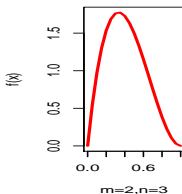
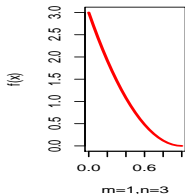
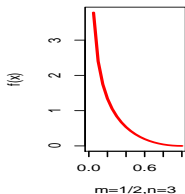
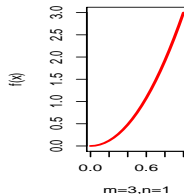
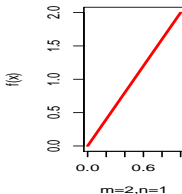
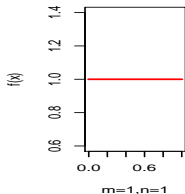
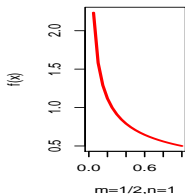
Thus we can characterize the aspect of the density as a function of the parameters.

- se  $m > 1, n > 1 \Rightarrow$  existe uma única moda em  $x = (m - 1)/(m + n - 2)$
- se  $m < 1, n < 1 \Rightarrow$  existe uma antimoda (forma de U)
- se  $(m - 1)(n - 1) \leq 0 \Rightarrow$  forma de J
- se  $m = n \Rightarrow$  symmetry with respect to 0.5.



# The beta distribution

In the following figures, we can see some of these aspects:



# The beta distribution


Let  $X$  and  $Y$  be independent random variables such that  $X \sim G(a_1, b_1)$  e  $Y \sim G(a_2, b_2)$ , then

$$X|(X + Y) \sim Be(a_1, a_2).$$

The beta distribution, just studied, is said to be in the standardized form and is in fact the most widely used form. Its more general form presents four parameters  $(a, b, m, n)$  and the density function is

$$f(x) = \begin{cases} \frac{1}{B(m,n)} \frac{(x-a)^{m-1}(b-x)^{n-1}}{(b-a)^{m+n-1}} & a < x < b \quad m > 0, n > 0 \\ 0 & \text{outros valores de } x \end{cases}$$

# SUMMARY of some distributions in the R

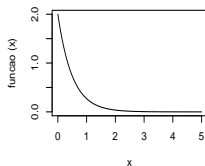
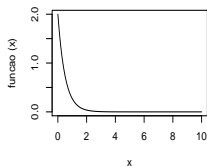
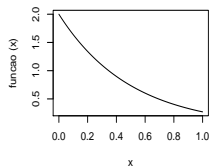
Distribution name in the 	Function	Arguments
Beta	beta	shape1, shape2
Binomial	binom	size, prob
Cauchy	cauchy	location, scale
Chisquare	chisq	df
Exponential	exp	rate
FDist	f	df1, df2
GammaDist	gamma	shape, scale
Geometric	geom	prob
Hypergeometric	hyper	m, n, k
Lognormal	lnorm	meanlog, sdlog
Logistic	logis	location, scale
NegBinomial	nbinom	size, prob
Normal	norm	mean, sd
Poisson	pois	lambda
TDist	t	df
Uniform	unif	min,max
Weibull	weibull	shape, scale

# Example of an exercise in the R

Consider the following function  $f(x) = \begin{cases} 2 e^{-2x} & x > 0 \\ 0 & x \leq 0 \end{cases}$

Let's see that  $f$  is indeed a density function;  
Calculate  $P[X > 1]$  and  $P[0.2 < X < 0.8]$

```
>funcao<-function(x) {  
+   fx<-ifelse(x<0,0,2*exp(-2*x))  
+   return(fx)}  
>par(mfrow=c(1,3))  
>plot(funcao);plot(funcao,0,10);plot(funcao,0,5)
```



# Example of an exercise in theR

```
>integrate(funcao,0,Inf)
>integrate(funcao,1,Inf)
>res<-integrate(funcao,0,1);res;str(res)
>1-res$value

1 with absolute error < 5e-07
0.1353353 with absolute error < 2.1e-05
0.8646647 with absolute error < 9.6e-15
List of 5
 $ value      : num 0.865
 $ abs.error  : num 9.6e-15
 $ subdivisions: int 1
 $ message    : chr "OK"
 $ call       : language integrate(f =funcao,lower = 0,upper = 1)
               attr(*, "class")= chr "integrate"
[1] 0.1353353
```