



SIMULATION

1. Introduction

- **Basic idea:** To use methods based on random number generation to simulate “complex” processes to get an insight about the random behavior of some important variables;
- The first “real” use of simulation was related to the “Manhattan project” at Los Alamos at the end of World War 2 (nuclear bomb) and was linked to the names of von Neumann, Stanislaw Ulam and Nicholas Metropolis. Since then, simulation techniques are used in many scientific domains and cover a great variety of applications.
- The use of simulation to obtain an approximation to the sampling distribution of a statistic illustrates the process. The process can be summarized in three stages:
 - a. Simulate the production of random samples from the population;
 - b. For each sample generated in the previous stage, get the value of the statistic for which we want to approximate the sampling distribution;
 - c. Now we can use the simulated values of the statistic to get an approximation to its sampling distribution (each replica provides an “observed value” of the statistic).
- Two problems remain: how to generate the random samples and how many replicas should we use.



2. How to generate random samples from a population with a known distribution?

- We follow a two-step procedure:
 - a. First, we generate pseudorandom numbers, i.e. values that can be considered as independent observations of a random variable with a uniform $(0;1)$ distribution;
 - b. Second, using probability theory we transform our pseudorandom numbers to get pseudo random observations from the wanted distribution.

2.1 – Pseudo random numbers generation

- **Random numbers versus pseudorandom numbers.** A random number is a realization of a continuous random variable following a uniform distribution between 0 and 1. All the generated variables should be independent from each other. As it is not possible to generate these numbers (the best we can get is to generate discrete variable but, even that, is quite difficult) we will use pseudorandom numbers.
- A **pseudorandom number generator** is an algorithm for generating a sequence of numbers that look like independent observations of a $U(0;1)$ distribution. Obviously, the sequence is not random as it is completely determined by a relatively small set of initial values.



- One common method to generate pseudorandom numbers is based on the multiplicative linear congruential generator introduced by Lehmer in 1949. The idea is the following

$$n_t = a n_{t+1} \bmod c \quad t = 1, 2, \dots$$

$$u_t = n_t / c$$

where u_t is t -th generated random number ($t = 1, 2, \dots$) and a , c and n_0 are three positive integers with adequate properties. n_0 is called the seed.

- Example

Let $a = 123$, $c = 1000$ and $n_0 = 135$. Generate 25 pseudo random numbers. Verify that $n_{20} = n_0$ and comment (Use Excel or another software).

- As illustrated by the example, the maximum length of the sequence before it begins to repeat can be a major concern. Another important points are to verify if we can assume that the observations **(i)** could have been generated by an independent process and **(ii)** mimic a $U(0;1)$ behavior.

Parameters a , c and n_0 are chosen to answer these points.

- Pseudorandom generation is still an area of active research (one main application of random numbers is cryptography!). The Mersenne-Twister algorithm, introduced in 1997, is replacing linear congruential generator as the “standard”.



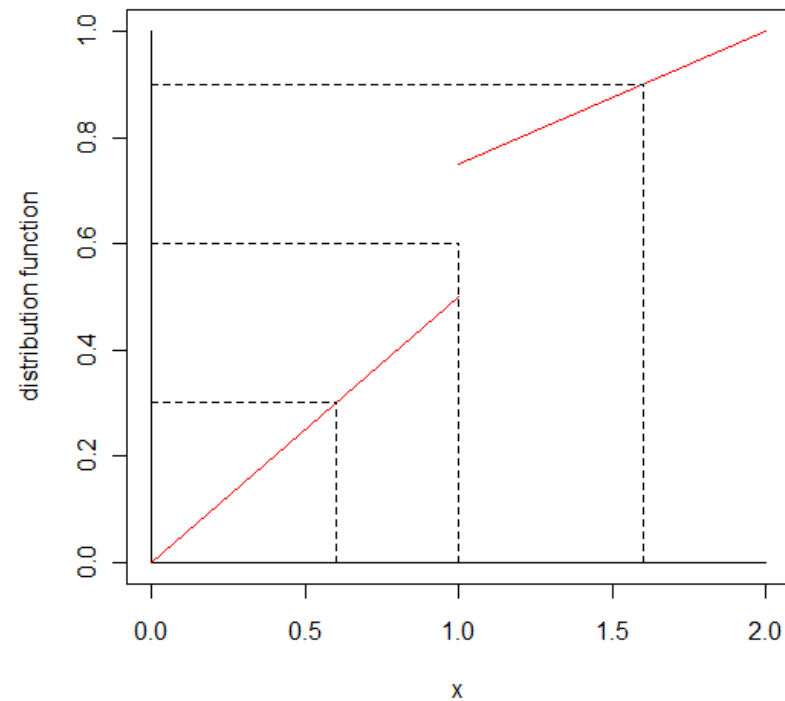
2.2. From a uniform distribution to other distributions

- Assuming that we got a sequence of independently distributed uniform variables between 0 and 1, the problem is now how to define functions of these variables that follow a given distribution. We can give a general answer to this problem, but in many situations we can design more efficient methods.
- The general answer is called the **inverse transform** and is based on the following theorem:
- **Theorem:** Let $Y \sim U(0;1)$ and $F(x)$ be the distribution function of a continuous random variable (we are assuming that F is strictly increasing on the open interval (a,b) where X has positive density – the interval can be unlimited – with $F(a) = 0$ and $F(b) = 1$). Then the random variable $X = F^{-1}(Y)$ is continuous and $F(x)$ is its distribution function.
- Example: Illustrate the procedure using (i) the exponential distribution and (ii) the normal distribution
- Extensions:
 - Suppose that the distribution function has a jump at $x = c$, i.e. $F(c^-) = a$ and $F(c) = b$ with $b > a$. If the random number u is such that $a \leq u < b$, choose c as the simulated value;
 - Suppose that the distribution function is constant in a given interval, i.e. $F(x) = p$ for $a \leq x \leq b$. If the random number u is equal to p (abnormal case) choose $x = b$;



- **Example 21.2** – Suppose $F_X(x) = \begin{cases} 0.5x & 0 \leq x < 1 \\ 0.5 + 0.25x & 1 \leq x \leq 2 \end{cases}$

Determine the simulated values of x resulting from the uniform numbers 0.3, 0.6 and 0.9.





- **Specific answers** - For many situations one can design a computationally more efficient (and reliable) method to use with specific distributions. The idea is to take advantage of the properties of these distributions.
 - Bernoulli
 - Binomial as the sum of independent Bernoulli variables
 - Normal (Box-Muller transform, Central Limit Theorem)
 - Poisson (Time between events follows an exponential distribution)
- Box-Muller formula: (U_1, U_2) independent uniform variables, then $Z_1 = \sqrt{-2\ln U_1} \cos(2\pi U_2)$ and $Z_2 = \sqrt{-2\ln U_1} \sin(2\pi U_2)$ are independent $n(0;1)$ random variables.
- Example: generate 10 uniforms and then generate (i) 10 Bernoulli variables with parameter 0.2 (ii) 10 normally distributed variables (mean=3, standard deviation=2) using the inverse method and Box-Muller formula.



- **Example 21.1** – generate 10000 pseudo Pareto (with $\alpha = 3$ and $\theta = 1000$) variables and verify that they are indistinguishable from real Pareto observations.

Procedure:

- Generate 10000 (pseudo) uniforms, u_i
- Obtain 10000 (pseudo) Pareto, $x_i = \theta \left((1 - u_i)^{-1/\alpha} - 1 \right)$
- Compare the generated values with the theoretical model:
 - Kolmogorov-Smirnov test
 - χ^2 goodness of fit test (*Loss Models* solution)
 - Other approaches (Anderson-Darling test, graphical techniques, ...)

A possible solution using R

```
> # generate n Pareto (alpha,theta) variables - inverse method
> n=10000; alpha=3; theta=1000
> u=runif(n); x=theta*((1-u)^(-1/alpha)-1)
> # ks test
> Pareto_dist_func=function(x,alpha,theta) {
+   1-(theta/(x+theta))^alpha
+ }
```



```
> ks.test(x, "Pareto_dist_func", alpha=3, theta=1000)
```

One-sample Kolmogorov-Smirnov test

```
data: x
```

```
D = 0.0074, p-value = 0.6515
```

```
alternative hypothesis: two-sided
```

```
>
```

```
> # Qui2 test for large samples(n/50 classes, n is multiple of 50)
```

```
> m=n/50; aux1=0:(m-1); lb=theta*((1-(aux1/m))^(-1/alpha)-1);
```

```
> ub=c(lb[2:m], Inf)
```

```
> counts=rep(NA, m)
```

```
> for (j in 1:m) counts[j]=sum((x>=lb[j]) & (x<ub[j]))
```

```
> expected=rep(50, m)
```

```
>
```

```
> chi2=((counts-expected)^2)/expected
```

```
> chi2.test=sum(chi2)
```




```

> result=cbind(lb,ub,counts,expected,chi2)
> result
      lbound      ubound counts expected chi2
[1,] 0.000000  1.672244    45      50 0.50
[2,] 1.672244  3.355730    46      50 0.32
[3,] 3.355730  5.050591    34      50 5.12
[4,] 5.050591  6.756962    46      50 0.32
...
[197,] 2684.031499 3054.801330    51      50 0.02
[198,] 3054.801330 3641.588834    58      50 1.28
[199,] 3641.588834 4848.035476    53      50 0.18
[200,] 4848.035476      Inf    54      50 0.32
> chi2.test
[1] 185.52
> p.value=pchisq(chi2.test,m-1,lower.tail=FALSE)
> p.value
[1] 0.7447192

```



3. How many replicas should be used?

- The answer depends on the problem we want to solve.
- Nowadays we can define a **huge number** of replicas for many situations. This is the usual solution.
- Sometimes, using the Central Limit Theorem we can define an approximate value to the number of replicas to ensure a given precision. Example 21.5 illustrates 3 situations.
- **Example 21.5** – Use simulation to estimate the mean, $F_X(1000)$ and $\pi_{0.9}$, the 90th percentile of the Pareto distribution with $\alpha = 3$ and $\theta = 1000$. In each case, stop the simulation when you are 95% confident that the answer is within $\pm 1\%$ of the true value.
 - As we know the true values, the simulation is useless but we will behave as these values are unknown ($\mu = 500$, $F_X(1000) = 0.875$, $\pi_{0.9} = 1000 \times (0.1^{-1/3} - 1) \approx 1154.435$)
 - We will discuss only the first 2 situations.



- **Mean (μ):** The usual estimator is the statistic $T = \bar{X}$ and $\frac{T - \mu}{\sigma / \sqrt{n}} \overset{\circ}{\sim} n(0;1)$

$$\Pr(|T - \mu| \leq 0.01\mu) = \Pr(|T - \mu| \leq 0.01\mu) = \Pr\left(-\frac{0.01\mu}{\sigma / \sqrt{n}} \leq \frac{T - \mu}{\sigma / \sqrt{n}} \leq \frac{0.01\mu}{\sigma / \sqrt{n}}\right) = 2\Phi\left(\frac{0.01\mu}{\sigma / \sqrt{n}}\right) - 1$$

$$\text{Then } \Pr(|T - \mu| \leq 0.01\mu) \geq 0.95 \Leftrightarrow \Phi\left(\frac{0.01\mu}{\sigma / \sqrt{n}}\right) \geq 0.975 \Leftrightarrow \frac{0.01\mu}{\sigma / \sqrt{n}} \geq 1.96 \Leftrightarrow n \geq \left(\frac{1.96}{0.01}\right)^2 \times \left(\frac{\sigma^2}{\mu^2}\right)$$

As μ and σ^2 are unknown we run a first simulation to estimate them: $\tilde{\mu} = \bar{x}$ and $\tilde{\sigma}^2 = s^2$

Procedure:

- Define m (number of replicas of the first simulation)
- Generate m (pseudo) Pareto observations and compute $\tilde{\mu} = \bar{x}$ and $\tilde{\sigma}^2 = s^2$
- Let n be the smallest integer greater than or equal to $\left(\frac{1.96}{0.01}\right)^2 \times \left(\frac{\tilde{\sigma}^2}{\tilde{\mu}^2}\right)$
- Generate an additional sample with $n - m$ observations and recomputed $\tilde{\mu}$ using all observations (from both samples).



Using R and choosing $m = 10000$:

```

> alpha=3; theta=1000; m=10000;
> u=runif(m); x1=theta*(((1-u)^(-1/alpha))-1)
> tau2_h=var(x1); miu_h=mean(x1)
> n_min=((1.96/0.01)^2)*(tau2_h/(miu_h^2))
> miu_h; tau2_h; n_min
[1] 502.2063
[1] 878943.7
[1] 133877.9
> m2=123878
> u=runif(m2); x2=theta*(((1-u)^(-1/alpha))-1)
> x=c(x1,x2)
> miu_h=mean(x)
> tau2_h=var(x1); n_min=((1.96/0.01)^2)*(tau2_h/(miu_h^2))
> miu_h; tau2_h; n_min
[1] 498.7444
[1] 737651.6
[1] 113921.9

```



- **Distribution Function** ($F_X(1000)$): Let $p = F_X(1000)$. The usual estimator is the ecdf

$$T = F_n(1000) = \frac{\#\{X_i \leq 1000\}}{n} \text{ and we get } \frac{T - p}{\sqrt{p \times (1-p)} / \sqrt{n}} \overset{\circ}{\sim} n(0;1)$$

$$\Pr(|T - p| \leq 0.01 p) = 2\Phi\left(\frac{0.01 p}{\sqrt{p \times (1-p)} / \sqrt{n}}\right) - 1, \text{ then}$$

$$\Pr(|T - p| \leq 0.01 p) \geq 0.95 \Leftrightarrow n \geq \left(\frac{1.96}{0.01}\right)^2 \times \left(\frac{1-p}{p}\right)$$

As p is unknown we run a first simulation (m replicas) to estimate it: $\tilde{p} = \#\{x_i \leq 1000\} / m$

Procedure:

- Define m (number of replicas of the first simulation)
- Generate m (pseudo) Pareto observations and compute $\tilde{p} = t = \#\{x_i \leq 1000\} / m$
- Let n be the smallest integer greater than or equal to $\left(\frac{1.96}{0.01}\right)^2 \times \left(\frac{1-\tilde{\theta}}{\tilde{\theta}}\right)$
- Generate an additional sample with $n - m$ observations and recomputed \tilde{p} using all observations (from both samples).



Using R and choosing $m = 10000$:

```
> alpha=3; theta=1000; m=10000;
> u=runif(m); x1=theta*(((1-u)^(-1/alpha))-1)
> t=mean(x1<=1000)
> n_min=((1.96/0.01)^2)*(1-t)/t
> t; n_min
[1] 0.881
[1] 5188.994
```

As $10000 > 5189$ we stop.



4. How to use simulation?

- Build a model for the random variable for which we want to approximate the distribution, S . This variable can be function of other random variables.
- Define the number of replicas to be used, NR
- For each replica generate as many pseudo-random variables as we need and compute a value for S using the model from step 1. Let us call this value s_j .
- The cdf of S is then approximated by the ecdf based on s_1, s_2, \dots, s_{NR} . Compute quantities of interest such as the mean, variance, percentiles, probabilities, ... from the ecdf.



3.1 – Approximating the sampling distribution of a statistic

Example A – Consider a normal population with mean 10 and standard deviation 3 from which we observe a sample of size $n=5$. Using simulation, obtain the sampling distribution of \bar{X} and compare with the theoretical result.

Procedure:

- Choose the number of replicas, NR
- For each replica $i, i = 1, 2, \dots, NR$
 - Generate 5 random (pseudo) numbers and obtain 5 r.v. with a $n(10; \sigma = 3)$ distribution
 - Compute the sample average and keep this value as element i of the vector res
- Perform a test to check if the values in vector res can be considered as observations of a normal with mean 10 and standard deviation $3 / \sqrt{5}$.



Solution using R

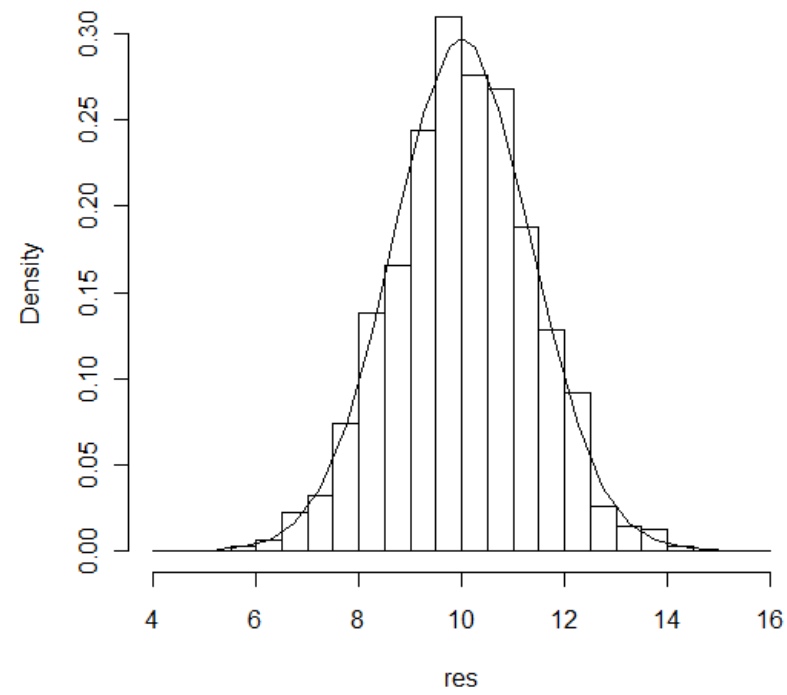
```
> NR=1000; n=5; miu=10; sig=3
> res=rep(NA, NR)
> for(i in 1:NR){
+   x=rnorm(n, miu, sig); res[i]=mean(x)
+ }
> mean(res); sd(res) # can also compute skewness and kurtosis
[1] 10.01766
[1] 1.351131
> ks.test(x, "pnorm", 10, 3/sqrt(5))
```

One-sample Kolmogorov-Smirnov test

```
data: x
D = 0.4427, p-value = 0.2087
alternative hypothesis: two-sided
> breaks=seq(4, 16, 0.5)
> points=c(seq(5+0.5/2, 16-0.5/2, 0.5), 10); points=sort(points)
> dens=dnorm(points, miu, sig/sqrt(n))
> hist(res, breaks, prob=TRUE)
> lines(points, dens, type="l")
```



Histogram of res





Example B – Consider a normal Pareto population ($\alpha = 1.5$, $\theta = 100$) from which we observe a sample of size 10. (i) Explain how to use simulation to get an approximation to the sampling distribution of \bar{X} and compare it with the normal distribution. (ii) Perform the simulation with 1000 replicas. (iii) **Do you think that increasing the sample size will help?**

(i)

Determine NR , the number of replicas to be used.

For each of the NR replicas, $i = 1, 2, \dots, NR$

- Generate 10 pseudo Pareto distributed variables – we generate 10 uniforms(0,1), u_j , $j = 1, 2, \dots, 10$ and using the inverse method we get 10 Paretos, $x_j = \theta \left((1 - u_j)^{-1/\alpha} - 1 \right)$
- Calculate the sample mean, $\bar{x}_i = \sum_{j=1}^{10} x_j$

Now we compare the simulated distribution of \bar{X} using our NR pseudo observations $(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_{NR})$ with a normal distribution and conclude (descriptive statistics, qqplot, ecdf, ...)



(ii)

```
> NR=1000; n=10; alpha=1.5; theta=100
> res=rep(NA, NR)
> for(i in 1:NR){
+   u=runif(n); x=theta*((1-u)^(-1/alpha)-1); res[i]=mean(x)
+ }
> library(moments)
> cbind(mean(res), median(res), sd(res), skewness(res), kurtosis(res))
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] 195.7378 126.9872 339.5334 13.12755 248.533
> qqnorm(res) # result on next slide
```

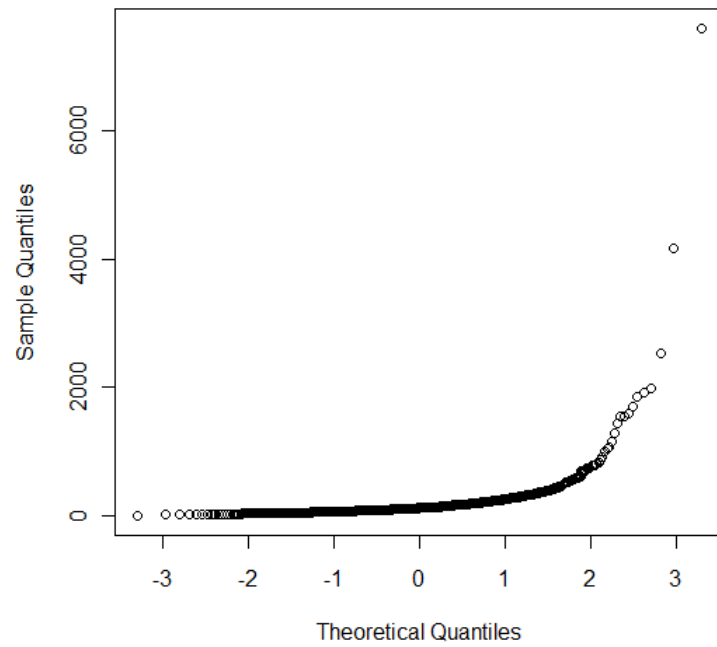
(iii) As $\alpha = 1.5 < 2$ the variance of the Pareto distribution does not exist and consequently the CLT does not apply. The sampling distribution of \bar{X} does not converge to a normal distribution. Using the same R program with $n = 1000$ we get

```
> cbind(mean(res), median(res), sd(res), skewness(res), kurtosis(res))
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] 197.5064 186.5814 115.4582 26.57593 789.149
> > qqnorm(res) # result on next slide
```



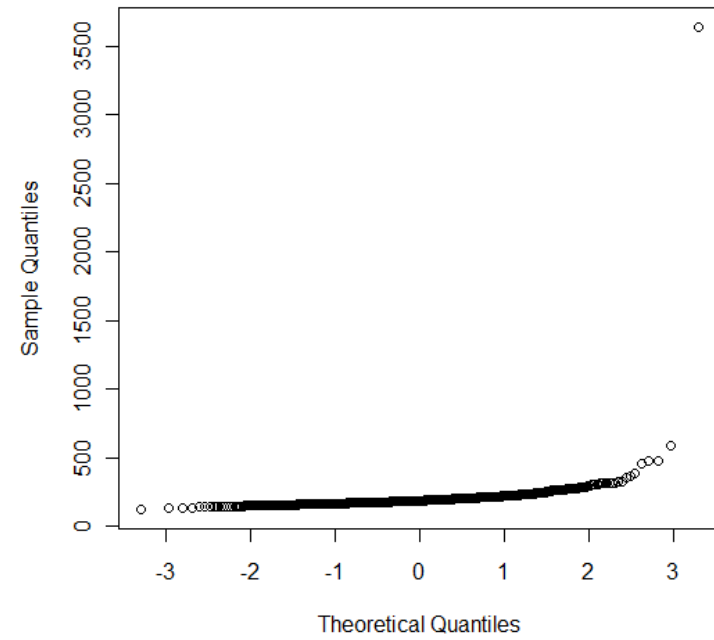
Sample size=10

Normal Q-Q Plot



Sample size=1000

Normal Q-Q Plot





Example C – Using simulation determine the p-value of the Kolmogorov-Smirnov test to test if the sample (11.79,11.25,6.83,10.47,13.60,10.60,17.40,10.99,16.45,12.47,8.19,13.46,13.82,11.93,7.47,10.09,14.26, 12.04,12.13,9.66) came from a normal distribution with mean 10 and standard deviation 3.

Procedure

- Compute the value of the test statistic using the observed sample, D
- Determine NR , the number of replicas to be used
- For each of the NR replicas, $i = 1, 2, \dots, NR$
 - Generate 20 pseudo normal (mean 10 and standard deviation 3) variables – generate 20 uniforms(0,1), $u_j, j = 1, 2, \dots, 20$, then, using the Box-Muller method, get 20 $n(0,1)$, z_j , and finally get $x_j = 10 + 3 \times z_j, j = 1, 2, \dots, 20$.
 - Using the generated sample, compute the value of the K-S test statistic, KS_i
- The estimated p-value is given by the proportion of values of KS greater than D



Using R and 10000 replicas we get

```
> x=c(11.79,11.25,6.83,10.47,13.60,10.60,17.40,10.99,16.45,12.47,8.19,
+     13.46,13.82,11.93,7.47,10.09,14.26,12.04,12.13,9.66)
>
> # test if x follows a normal distribution with mean 10 and stdev=3
> a=ks.test(x,"pnorm",10,3)
> D=a$statistic
> a
      One-sample Kolmogorov-Smirnov test
data:  x
D = 0.3122, p-value = 0.03132
alternative hypothesis: two-sided
>
> NR=10000; n=length(x)
> res=rep(NA, NR)
> for(i in 1:NR){
+   y=rnorm(20,10,3); a=ks.test(y,"pnorm",10,3)
+   res[i]=a$statistic
+ }
> p.value=mean(res>=D); p.value
[1] 0.0326
```



3.1 – More complex analysis

- Among the examples presented in *Loss Models* we will analyze Example 21.17. Skip section 21.2.4 (unless you are familiar with copulas) and return to section 21.2.6 after bootstrap has been presented.
- **Example 21.17** – An insurance company offers the following product to individuals age 40. A single premium of 10000 is paid (an administrative fee has already been deducted). In return, there are two possible benefits. The 10000 is invested in a mutual fund. If the policyholder dies during the next four years, the fund value is paid to the beneficiary. If not, the fund value is returned to the policyholder at the end of the four years. The policyholder may purchase a guarantee. If the fund has earned less than 5% per year at the time of payment, the payment will be based on a 5% per year accumulation rather than the actual fund value. Determine the 90th percentile of the cost of providing this guarantee. Assume that the force of mortality is constant at 0.02 and that 50000 policies will be sold. Also assume that the annual fund increase has a lognormal distribution with $\mu = 0.06$ and $\sigma = 0.02$.
We also assume that payments to the beneficiaries (policyholders who died) are made at the end of each year.



Constant force of mortality at 0.02 $\rightarrow q = 1 - e^{-0.02} = 0.0198$

n_0 - number of policies at the beginning of year 1. $n_0 = 50000$

C_0 - Capital value for each policy at the beginning of year 1. $C_0 = 10000$

Let us describe the simulation for replica k

- Years $i = 1, 2, 3$
 - Simulate the number of death m_i from $M_i \sim b(n_{i-1}, 0.0198)$ and calculate the number of survivors $n_i = n_{i-1} - m_i$
 - Simulate x_i , the fund increase during year i from $X_i \sim \text{lnormal}(0.06; 0.02)$ and calculate the value of the policy at the end of year i , $C_i = X_i \times C_{i-1}$
 - If the value of the fund is smaller than the guaranteed capital the insurance fund has to pay the difference to the beneficiaries of the policyholders who died during the year, i.e.
$$\text{if } (C_i < C_0 \times 1.05^i) \text{ then } P_i = (1.05^i C_0 - C_i) m_i \text{ else } P_i = 0$$
- Years $i = 4$ (similar to the previous years except that we have to reimburse all policyholders)
 - Simulate x_i , the fund increase during year i from $X_i \sim \text{lnormal}(0.06; 0.02)$ and calculate the value of the policy at the end of year i , $C_i = X_i \times C_{i-1}$



- If the value of the fund is smaller than the guaranteed capital the insurance fund has to pay the difference to all policyholders (including the beneficiaries of the) i.e.

$$\text{if } (C_i < C_0 \times 1.05^i) \text{ then } P_i = (1.05^i C_0 - C_i) n_{i-1} \text{ else } P_i = 0$$

- Compute the present value of the payments due to the guarantee. Let us assume that the discount rate is given by the increase of the fund, i.e. $Paym_k = \sum_{i=1}^k P_i / v_i$ where $v_i = \prod_{j=1}^i x_j$.

Repeat the procedure for $k = 1, 2, \dots, NR$ where, for instance, we can define $NR = 10000$.

Sort array $Paym_k$ and estimate the 90th percentile using the usual method. We can calculate the corresponding value per policy (just divide by 50000) to price the insurance option. We can also estimate the mean value or others statistics of interest.

See file example21.17.xlsx to see a realization of the simulation.

Challenging question: Can you develop the simulation using R?



```

q=1-exp(-0.02); C0=10000; N0=50000;
G=rep(1.05,4); G=cumprod(G)*C0

NR=1000; pp=rep(NA,NR)
for(j in 1:NR){
  pay=rep(0,4)
  x=rlnorm(4,0.06,0.02); discount=1/cumprod(x); C=C0*cumprod(x)
  NReimb=rep(NA,4)
  n=N0
  for(i in 1:3){
    NReimb[i]=rbinom(1,n,q); n=n-NReimb[i]
  }
  NReimb[4]=n
  pay=NReimb*(G-C)*(G>C)
  pp[j]=sum(pay*discount)
  # G;C;NReimb;pay;
}

pp_p=pp/N0
mean(pp_p); sd(pp_p); quantile(pp_p,0.9,type=6)
hist(pp_p)

```