

Monte Carlo Simulations

1 Introduction

Doing econometrics means estimating parameters, such as the mean of a population, the coefficients in a linear regression or the autocorrelation of a time series, given a sample of real world data. Besides the point estimate itself, we would like to know how close our estimate is to the true value. In other words we would like to know its “accuracy” or “precision”. An estimator is a (maybe complicated) function of random variables and therefore itself a random variable. The properties of an estimator are fully described by its probability distribution (the so-called sampling distribution). The sampling-distribution can then be used to perform tests against hypothesis. Often we are especially interested in some moments of the sampling distribution, such as the mean and the variance.

In some cases it is possible to calculate the sampling distribution from the econometric model. But sometimes, especially for finite (small) samples, this is either not possible or very difficult. In these cases Monte Carlo experiments are an intuitive way to obtain information about the sampling distribution and hence about the “quality” of the estimator.

2 The Method

The term “Monte Carlo” refers to procedures in which quantities of interest are approximated by generating many random realisations of a stochastic process and averaging them in some way. In statistics, the quantities of interest are the distributions of estimators and test statistics, the size of a test statistic under the null hypothesis, or the power of a test statistic

under some specified alternative hypothesis (see Davidson and MacKinnon 1993, 731). In economic theory, Monte Carlo techniques are used to explore the quantitative properties of models with stochastic elements, for example the correlation between variables in real business cycle models.

How can we use Monte Carlo techniques to find the sampling distribution of an estimator? In the real world, we usually observe just one sample of a certain size N , that will give us just one estimate. The Monte Carlo experiment is a lab situation, where we replicate the real world study many (R) times. Every time, we draw a different sample of size N from the original population. Thus, we can calculate the estimate many times and any estimate will be a bit different. The empirical distribution of these many estimates approximates the true sampling distribution of the estimator.

A Monte Carlo experiment involves the following steps:

- (1) Assume values for the exogenous parts of the model or draw them from their respective distribution function
- (2) Draw a (pseudo) random sample of size N for the error terms in the statistical model from their respective probability distribution functions
- (3) Calculate the endogenous parts of the statistical model
- (4) Calculate the value (e.g. the estimate) you are interested in
- (5) Replicate step 1 to 4 R times
- (6) Examine the empirical distribution of the R values

3 An Example: OLS

Let's explain the above elements in an example: the bivariate ordinary least squares model

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

with $u_i \sim N(0, \sigma^2)$. The stochastic element in the model is u_i , the exogenous part is x_i is either fixed or stochastic. Assuming values for the true parameters b_0 and b_1 and drawing values for the stochastic element, we can simulate the endogenous variable y_t . The values of interest are then the least squares estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ in the simulated data set.

In the core of Monte Carlo experiments is the random number generator. A random number generator produces a sequence of numbers, that are draws from a specific identically and independently distributed random variable. In practice, this is a mathematical algorithm, that produces a sequence of so-called pseudo random numbers. These numbers are in fact not random as the algorithm describes the purely deterministic relationship between the numbers. However, with a good generator, they are indistinguishable from sequences of genuinely random numbers and pass usual statistical tests of independence. Judd (1998) provides a thorough treatment of different pseudo-random number generators.

There is an important limitation of Monte Carlo experiments: We must completely specify the Statistical Model (Data Generating Process DGP). This implies, that we must assume the deterministic parts of the model, the form and the exact parameters of the distribution of the stochastic (error) term and the distribution of exogenous variables. This is a great loss of generality as the results of the experiment apply only to the assumptions made.

4 Implementation in Stata 17

Stata has a built-in random number generator:

```
uniform()
```

returns uniformly distributed pseudo-random numbers on the interval $[0, 1]$. Random numbers for other continuous distributions are calculated using the inverse of the desired distribution, for example

```
generate z = invnormal(uniform())*2+5
```

generates a new variable z with $_N$ (the number of observations in the current dataset) independent draws from a normal distribution with variance $2^2 = 4$ and mean 5. See `help drawnorm` on how to draw a random vector from the *multivariate* normal distribution. You can reset the random number generator with e.g. `set seed 1234` which allows to exactly reproduce your results.

The different steps of a Monte Carlo experiment in Stata are explained by an investigation into the properties of the OLS estimator in a bivariate regression model.

The first task in setting up the Monte Carlo experiment in Stata is to define a program that produces the result of a single experiment, i.e. that performs the steps (1) to (4).

```
capture program drop olssim
program olssim, rclass
    drop _all
    set obs 30
    generate x = uniform() * 10
    generate u = invnormal(uniform()) * 2
    generate y = -2 + 0.5 * x + u
    regress y x
    return scalar b0 = _coef[_cons]
    return scalar b1 = _coef[x]
end
```

The above program clears the data in memory and sets the number of observations in each sample to $N = 30$. Step (1): the realized error terms

u are drawn from the centered normal distribution with variance 2². Step (2): the independent variable x is drawn from a uniform distribution on $[0,10]$. Step (3): the realizations of the dependent variable y are calculated according to the DGP as $y_i = \beta_0 + \beta_1 x_i + u_i$, with true parameter values $\beta_0 = -2$ and $\beta_1 = 0.5$. Step (4): the estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are estimated in the regression of y on x . The last two lines of the program specify the values that are investigated in this Monte Carlo experiment: $\hat{\beta}_0$ and $\hat{\beta}_1$ which will be returned under the names **b0** and **b1**, respectively. The definition of the program can be directly typed into the command window or is part of a do-file. The command **capture program** clears the program from the memory before it is (re-)defined.¹

Step (5) is the replication of the single experiment R times. There is a special Stata command **simulate** that performs this replication and produces a new dataset with the results. For example,

```
simulate b0 = r(b0) b1 = r(b1), reps(1000): olssim
```

performs $R = 1000$ experiments of a bivariate OLS regression on $N = 30$ data points drawn with true parameter values $\beta_0 = -2$ and $\beta_1 = 0.5$. The results are stored in a new dataset with 1000 observations of the two variables $b0$ and $b1$. Each row contains the estimated parameters of a single experiment.

In step (6), we examine the results of the Monte Carlo experiment. This is done by inspecting the new variables $b0$ and $b1$ using the usual command for descriptive statistics. For example

```
sum b1
histogram b1
```

reports mean and standard deviation of the $R = 1000$ estimated slope coefficients $\hat{\beta}_1$ and draws a histogram.

¹More sophisticated implementations of Monte-Carlo experiment would declare the program in a separate so-called ado-file. The program will generally take several arguments that describe the details of the experiments, such as the true parameter values.

The following code runs a Monte Carlo experiment with $R = 1000$ replications for the OLS estimator of the linear model $y_i = \beta_0 + \beta_1 x_i + u_i$ with true parameter values $\beta_0 = -2$ and $\beta_1 = 0.5$ and a sample of size $N = 30$ when the error term u_i follows the Cauchy distribution. It shows the distribution of the t-statistic for the null hypothesis that the true slope coefficient is $\beta_1 = 0.5$. Note that the Cauchy distribution is the distribution of the ratio of two independent normally distributed random variables with mean zero. The Cauchy distribution does not have a finite variance and therefore violates assumption *OLS4* (see the handout on “The Multiple Linear Regression Model”).

```
capture program drop olssim_cauchy
program olssim_cauchy, rclass
    version 10.0
    syntax [, obs(integer 1) beta0(real 0) beta1(real 0) ]
    drop _all
    set obs `obs'
    generate x = uniform()*10
    generate u = invnormal(uniform())/sqrt(invchi2(1,uniform()))
    generate y = `beta0' + `beta1' * x + u
    regress y x
    return scalar t1 = (_coef[x]-`beta1')/_se[x]
end

simulate t1 = r(t1), reps(1000): ///
    olssim_cauchy, obs(30) beta0(-2) beta1(0.5)

sum t1
histogram t1, kdensity ///
    title("Distribution of t-statistic") ///
    xtitle("t-statistic") ///
    legend(order(1 "Histogram" 2 "Kernel density" 3 "N(0,1)")) ///
    plot(function stdnorm = normalden(x,0,1), ra(-4 4) lpattern(dash))
```

5 Implementation in R 4.3.2

R has a built-in random number generator function:

```
x <- runif(N)
```

stores N uniformly distributed pseudo-random numbers on the interval $[0, 1]$ in the vector x . Random numbers for other continuous distributions are calculated using the quantile function of the desired distribution, for example

```
z <- qnorm(x, mean=5, sd=2)
```

generates a vector with z with N (the number of observations in the vector x) independent draws from a normal distribution with variance $2^2 = 4$ and mean 5. You can reset the random number generator by e.g. `set.seed(1234)` which allows to exactly reproduce your results.

The different steps of a Monte Carlo experiment in R are explained by an investigation into the properties of the OLS estimator in a bivariate regression model.

The first task in setting up the Monte Carlo experiment in R is to define a function that produces the result of a single experiment.

```
olssim <- function(N, b1, b0){
  u <- qnorm(runif(N)) * 2
  x <- runif(N, min=0, max=10)
  y <- b0 + b1 * x + u
  ols <- lm(y ~ x, data=data.frame(y,x))
  b <- ols$coefficients
  return(b)
}
```

The above function draws performs the steps (1) to (4) for a bivariate regression model with true parameters β_0 and β_1 for N observations. Step (1): the realized error terms u are drawn from the centered normal distribution with variance $2^2 = 4$. Step (2): the independent variable x is drawn from a uniform distribution on $[0, 10]$. Step (3): the realizations of the dependent variable y are calculated according to the DGP as $y_i =$

$\beta_0 + \beta_1 x_i + u_i$, with true parameter values β_0 and β_1 . Step (4): the estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are estimated in the regression of y on x . The last line of the function specifies that $\hat{\beta}_0$ and $\hat{\beta}_1$ will be returned.

Step (5) is the replication of the single experiment R times. The function `replicate()` performs this replication. For example,

```
b_mc <- replicate(1000, olssim(N=30, b0=-2, b1=0.5))
```

performs $R = 1000$ experiments of a bivariate OLS regression on $N = 30$ data points drawn with true parameter values $\beta_0 = -2$ and $\beta_1 = 0.5$. The results are stored in a new matrix called `b_mc`. Each column contains the estimated parameters of a single experiment.

In step (6), we examine the results of the Monte Carlo experiment. This is done by inspecting the R estimates of both variables using the usual functions for descriptive statistics. For example,

```
b1_mc <- b_mc[,2]  
mean(b1_mc)  
sd(b1_mc)  
hist(b1_mc)
```

reports mean and standard deviation of the $R = 1000$ estimated slope coefficients $\hat{\beta}_1$ and draws a histogram.

The following code runs a Monte Carlo experiment with $R = 1000$ replications for the OLS estimator of the linear model $y_i = \beta_0 + \beta_1 x_i + u_i$ with true parameter values $\beta_0 = -2$ and $\beta_1 = 0.5$ and a sample of size $N = 30$ when the error term u_i follows the Cauchy distribution. It shows the distribution of the t-statistic for the null hypothesis that the true slope coefficient is $\beta_1 = 0.5$. Note that the Cauchy distribution is the distribution of the ratio of two independent normally distributed random variables with mean zero. The Cauchy distribution does not have a finite variance and therefore violates assumption *OLS4* (see the handout on “The Multiple Linear Regression Model”).

```
olssim_cauchy <- function(N, b0, b1){
  u <- qnorm(runif(N))/qnorm(runif(N))
  x <- runif(N, min=0, max=10)
  y <- b0 + b1 * x + u
  ols <- lm(y ~ x, data=data.frame(y,x))
  output <- summary(ols)
  t1 <- (output$coefficients[2,1]-b1)/output$coefficients[2,2]
  return(t1)
}

N <- 30
R <- 1000
t1_mc <- replicate(R, olssim_cauchy(N, b0=-2, b1=0.5))

hist(t1_mc, breaks=40, prob=TRUE, xlim = c(-4,4), ylim = c(0,0.4),
     xlab="t-statistic", main="Distribution of t-value")
curve(dnorm(x), lty=2, lwd=2, add=TRUE)
lines(density(t1_mc), lwd=2)
legend(-4, 0.4, legend=c("Kernel density", "N(0,1)"), lty=1:2, lwd=2)
```

References

- Cameron, A. Colin and Pravin K. Trivedi (2005), *Microeconometrics: Methods and Applications*, Cambridge University Press. Section 7.7.
- Davidson, Russell and James G. MacKinnon (1993), *Estimation and Inference in Econometrics*, Oxford University Press, chapter 21.
- Judd, Kenneth L. (1998), *Numerical Methods in Economics*, MIT Press, chapter 8.