

ECOMAss2

NathanAung

April 18, 2017

```
knitr::opts_chunk$set(echo = TRUE)
set.seed(42)
library(R6) #using R6 oop system
library(stargazer)
```

```
## Warning: package 'stargazer' was built under R version 3.3.2
```

```
##
```

```
## Please cite as:
```

```
## Hlavac, Marek (2015). stargazer: Well-Formatted Regression and Summary Statistics Tables.
```

```
## R package version 5.2. http://CRAN.R-project.org/package=stargazer
```

Q1

Define two classes that generates results:

```
#Class that generates all combinations in a vector
library(R6) #using R6 oop system
regMCLoop = R6Class(
  public = list(
    nrange=NULL,
    sigmarange=NULL,
    Urange=NULL,
    a2range=NULL,
    results = NULL,

    initialize = function(nrange,sigmarange,Urange,a2range=0)
    {
      self$nrange = nrange
      self$sigmarange = sigmarange
      self$Urange = Urange
      self$a2range = a2range
      self$results = vector("list",length(nrange)*length(sigmarange)*length(Urange)*length(a2range))
    },

    loop = function(a1=1, a3=1, b1=1, b2=1, b3=0)
    {
      i=1#elements of result vector
      for(n in self$nrange)
      {
        for(sigma in self$sigmarange)
        {
          for(dist in self$Urange)
          {
            for(a2 in self$a2range)
            {
```

```

        self$results[[i]] = regMClass$new(ns=n,sigma=sigma,dist=dist, a2=a2, a3=a3, b1=b1, b2=b2,
        i=i+1#next result
    }
}
}
}
}
)
)

```

```

regMClass = R6Class(
  public = list(
    bias = NULL,
    sd = NULL,
    CovRate=NULL,
    CovLength=NULL,
    ns = NULL,
    sigma= NULL,
    dist= NULL,
    a2 =NULL,

initialize = function(ns = 20, n=1000, b1=1, b2=1, b3=0, a1=1, a2=0, a3=1, sigma=1,dist="norm")
{
  intervalrange =vector(mode="numeric",length=0)
  b1list = vector(mode="numeric",length=0)
  inInterval=vector(mode="logical",length=0)#declare vectors so that append method works
  for(j in 1:n)
  {
    #__init__
    Z1 = rnorm(n=ns,mean=0,sd=1)
    X2 = rnorm(n=ns,mean=0,sd=1)
    X3 = rnorm(n=ns,mean=0,sd=1)
    V = rnorm(n=ns,mean=0,sd=1)

    if (dist == "norm")
    {
      U = rnorm(n=ns,mean=0,sd=1)
    }
    else
    {
      U = rlnorm(n=ns,mean=0,sd=1)
    }
    #Generate X1 and Y
    X1 = a1*Z1 + a2*X2 + a3*X3 +V
    Y = b1*X1 + b2*X2 +b3*X3 + sigma*U
    #Estimates
    hat = lm(Y ~ X1+X2)
    b1list=append(b1list,hat$coefficients["X1"])
    #In confidence interval?
    interval = confint(hat,parm = "X1",interval="confidence")
    inInterval = append(inInterval,b1<interval[2]&b1>interval[1])
  }
}

```

```

        intervalrange = append(intervalrange,interval[2]-interval[1])
    }
    self$ns = ns # so we can sort by values used
    self$sigma = sigma
    self$dist = dist
    self$a2 = a2
    self$bias = mean(b1list)-1
    self$sd = sd(b1list)
    self$CovRate = mean(inInterval)
    self$CovLength = mean(intervalrange)
  }
)
)

```

We use a nested loop to generate all combinations needed and store them in a vector of results. These results are then stored within an object. It is possible that a recursive function approach could have been more elegant but for the number of combinations here, nested loop should suffice.

```

dataQ1 = regMCLoop$new(nrange=c(20,200,2000),sigmarange=c(1,2),Urange=c("norm","lognorm"))
dataQ1$loop()

```

Here we create an instance of the class and use the loop methods to generate the results. This does take some time(~20s) as we are looping over quite a large number of iterations.

Table 1:

	ndf	sigma	Distribution	Bias	SD	CovRate	CovLength
1	20	1	Normal	0.001	0.144	0.945	0.590
2	20	1	LogNormal	0.004	0.322	0.951	1.105
3	20	2	Normal	0.004	0.280	0.953	1.170
4	20	2	LogNormal	0.041	0.592	0.945	2.166
5	200	1	Normal	0.001	0.041	0.956	0.162
6	200	1	LogNormal	-0.004	0.097	0.939	0.343
7	200	2	Normal	-0.004	0.083	0.939	0.324
8	200	2	LogNormal	-0.004	0.171	0.951	0.671
9	2,000	1	Normal	-0.0001	0.013	0.964	0.051
10	2,000	1	LogNormal	0.0001	0.028	0.952	0.109
11	2,000	2	Normal	-0.001	0.025	0.954	0.101
12	2,000	2	LogNormal	-0.002	0.055	0.955	0.217

a)

In all combinations, $|\text{bias}|$ is unanimously within ± 0.005 of 0. Considering that even under all combinations in which the number of samples is 20 that this is present, these results strongly indicate that the estimator is or very close to unbiased. It is impossible to say definitively due to the stochastic nature of the simulation, however the weak law of large numbers and central limit theorem point to this conclusion. This seems appropriate considering the theoretical aspects of the estimator. (more explanation needed here?)

b)

For each combination, of σ and distribution, we vary the number of samples and observe the results. In nearly all cases, bias decreases monotonically throughout. In the few cases where it does not, bias is very close to zero and is still within the same order of magnitude of the previous estimate. This is to be expected even with a consistent estimator as the estimator becomes closer and closer to the true value of the variable. Perhaps with a higher number of simulations, this error would be reduced. Further, testing across a larger number of sample sizes may have graphically indicated the monotonicity of the bias as sample size increased. With only three data points, there could be any number of small changes in between.

c)

Coverage rate is consistently within 0.15 of the expected 95% coverage. Coverage rates between normal and log normal distributions appear to be inconclusive. Similarly, σ appears to have limited effect on coverage rate. The number of samples also inconclusive. One could of course perform an anova test to show these and reveal possible interaction effects.

d)

For all combinations, log normal generates much large confidence intervals than its normal counterpart.

This can be seen by examining the variance of normal and log normal distributions. With standard deviation σ , the variance for normal and lognormal respectively is:

$$\sigma^2$$
$$(e^{\sigma^2} - 1)(e^{2\mu + \sigma^2})$$

One can verify that for all values the variance will be greater for the log normal case . Therefore the estimated standard deviation will be much larger for the distribution and thus parameter estimates will have a larger confidence interval for the same coverage rate.

This can be verified again as for unitary σ , the coverage rate is far smaller than their $\sigma = 2$ counterparts.

Sample size reduces interval length drastically for all combinations. 2a)

```
dataQ2a = regMClloop$new(nrange=c(20,200,2000),sigmarange=c(1,2),Urange=c("norm","lognorm"))
dataQ2a$loop(b3=1,a3=1)
```

It is clear that omitted variable bias is present here with an apparent consistent bias of ~0.33. ‘work out theoretical bias here’. Standard deviation is also uniformly higher compared to a) and as such confidence intervals are larger. Coverage rate is especially poor here indicating a severe deviation from the requirements and assumptions of linear regression.

Whilst X_3 determines X_1 in both 1) and 2a), it does not determine Y in 1). As such, it does not fulfill the requirements for an omitted variable bias as endogeneity does not occur:

- 1) The omitted variable must be a determinant of the dependant variable(ie. regression coefficient is non-zero)
- 2) Omitted variable is correlated with the explanatory variable.

Whilst condition 2) is present through both questions, condition 1 only holds for question 2a).

2b)

Table 2:

	ndf	sigma	Distribution	Bias	SD	CovRate	CovLength
1	20	1	Normal	0.342	0.185	0.562	0.766
2	20	1	LogNormal	0.356	0.327	0.714	1.237
3	20	2	Normal	0.334	0.310	0.824	1.275
4	20	2	LogNormal	0.311	0.622	0.902	2.271
5	200	1	Normal	0.331	0.053	0	0.210
6	200	1	LogNormal	0.333	0.100	0.082	0.371
7	200	2	Normal	0.331	0.088	0.038	0.350
8	200	2	LogNormal	0.339	0.176	0.488	0.695
9	2,000	1	Normal	0.334	0.017	0	0.065
10	2,000	1	LogNormal	0.334	0.030	0	0.116
11	2,000	2	Normal	0.333	0.028	0	0.109
12	2,000	2	LogNormal	0.334	0.058	0.002	0.222

```
dataQ2b = regMCloop$new(nrange=c(20,200,2000),sigmarange=c(1,2),Urange=c("norm","lognorm"))
dataQ2b$loop(b3=1,a3=0)
```

Table 3:

	ndf	sigma	Distribution	Bias	SD	CovRate	CovLength
1	20	1	Normal	0.002	0.252	0.944	1.024
2	20	1	LogNormal	0.006	0.390	0.957	1.546
3	20	2	Normal	0.003	0.403	0.952	1.606
4	20	2	LogNormal	-0.005	0.827	0.950	2.816
5	200	1	Normal	-0.001	0.070	0.945	0.281
6	200	1	LogNormal	-0.008	0.112	0.963	0.458
7	200	2	Normal	-0.004	0.110	0.952	0.444
8	200	2	LogNormal	0.005	0.230	0.951	0.858
9	2,000	1	Normal	-0.0004	0.023	0.944	0.088
10	2,000	1	LogNormal	0.001	0.037	0.954	0.148
11	2,000	2	Normal	0.0003	0.036	0.949	0.139
12	2,000	2	LogNormal	0.002	0.069	0.952	0.274

As expected bias disappears as condition 1 is removed. The another anomalies can be explained as follows: