

ECOMAss2

NathanAung

April 18, 2017

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

```
## Loading required package: car
## Loading required package: lmtest
## Loading required package: zoo
```

```
##
## Attaching package: 'zoo'
## The following objects are masked from 'package:base':
##
##   as.Date, as.Date.numeric
```

```
## Loading required package: sandwich
## Loading required package: survival
```

```
library(stargazer)
```

```
##
## 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
regMCLoop = R6Class(
  public = list(
    nrange=NULL,
    sigmarange=NULL,
    Urange=NULL,
    a1range=NULL,
    a2range=NULL,
    results = NULL,

    initialize = function(nrange,sigmarange,Urange,a2range=0,a1range=1)
    {
      self$nrangle = nrange
      self$sigmarange = sigmarange
      self$Urange = Urange
      self$a2range = a2range
    }
  )
)
```

```

        self$a1range = a1range
        self$results = vector("list",length(nrange)*length(sigmarange)*length(Urange)*length(a2range))#list
    },

    loop = function( a3=1, b1=1, b2=1, b3=0,IV=FALSE)
    {
        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)
                    {
                        for(a1 in self$a1range)
                        {
                            self$results[[i]] = regMClass$new(ns=n,sigma=sigma,dist=dist, a2=a2, a3=a3, b1=b1, b2=b2, b3=b3, IV=IV)
                            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,
    a1 =NULL,
    fittedvalues = NULL,

    initialize = function(ns = 20, n=1000, b1=1, b2=1, b3=0, a1=1, a2=0, a3=1, sigma=1,dist="norm", IV=FALSE)
    {
        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__
  #we dont use a multivariate normal to generate it as this is equivalent.(independent and jointl
  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)
    U=(U-mean(U))/sd(U) #z score rescaled
  }
  #Generate X1 and Y
  X1 = a1*Z1 + a2*X2 + a3*X3 +V
  X1 = X1/sd(X1) #renormalising so all values have sd of 1
  Y = b1*X1 + b2*X2 +b3*X3 + sigma*U
  #Estimates
  if (IV==FALSE){
    hat = lm(Y ~ X1+X2)
  }

  else{
    hat = ivreg(Y~X1+X2|X2+Z1)
  }
  b1list=append(b1list,hat$coefficients["X1"])
  #In confidence interval?
  #b1
  interval = confint(hat,parm = "X1",interval="confidence")
  inInterval = append(inInterval,b1<interval[2]&b1>interval[1])
  intervalrange = append(intervalrange,interval[2]-interval[1])

  #note that this could be probably be optimised greatly by a running mean,sd tally (saves memory

}
self$ns = ns # so we can sort by values used
self$sigma = sigma
self$dist = dist
self$a2 = a2
self$a1 = a1

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. We shall show the code here for this instance but as the code is largely similar, we will not show it for the sake of cluttering in following sections. This probably could have been a method for the regMCLoop class.

```
Bias = vector(length=12)
SD = vector(length=12)
CovRate = vector(length=12)
CovLength = vector(length=12)
for(i in 1:12){
  Bias[i] = (dataQ1$results[[i]]$bias)
  SD[i] = (dataQ1$results[[i]]$sd)
  CovRate[i] = (dataQ1$results[[i]]$CovRate)
  CovLength[i]=(dataQ1$results[[i]]$CovLength)
}
ndf=c(rep(20,4),rep(200,4),rep(2000,4))
sigma = c(rep(1,2),rep(2,2),rep(1,2),rep(2,2),rep(1,2),rep(2,2))
Distribution = rep(c("Normal","LogNormal"),6)
Q1=data.frame(ndf,sigma,Distribution,Bias,SD,CovRate,CovLength)
stargazer(Q1,header=FALSE,type="latex",summary=FALSE)
```

Table 1:

	ndf	sigma	Distribution	Bias	SD	CovRate	CovLength
1	20	1	Normal	0.010	0.242	0.949	0.987
2	20	1	LogNormal	-0.003	0.235	0.956	0.997
3	20	2	Normal	0.003	0.477	0.958	1.925
4	20	2	LogNormal	0.018	0.477	0.946	1.996
5	200	1	Normal	-0.005	0.068	0.957	0.280
6	200	1	LogNormal	0.003	0.071	0.954	0.280
7	200	2	Normal	-0.002	0.146	0.948	0.560
8	200	2	LogNormal	0.008	0.148	0.939	0.560
9	2,000	1	Normal	0.001	0.022	0.946	0.088
10	2,000	1	LogNormal	0.001	0.022	0.953	0.088
11	2,000	2	Normal	0.001	0.045	0.944	0.175
12	2,000	2	LogNormal	0.001	0.044	0.950	0.176

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. This is explained by the Gauss Markov Theorem, which guarantees that the OLS is the best unbiased linear estimator. In this case, all the conditions

for the Gauss Markov Theorem are satisfied - the errors have mean 0, are homoskedastic (covariance matrix is identity) and are uncorrelated (generated independently) and so our estimator becomes the BLUE.

b)

For each combination of σ and probabilistic error 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 converges to the true value of the variable asymptotically. Perhaps with a higher number of simulations, this error could be reduced further. Testing across a larger number of sample sizes may have also 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 uniform - little difference is present. While one might expect the log normal case to have significantly less coverage rate, especially because the log normal is assymetric, it may be that one side may compensate for the other. ie. one side may be quite far outside the interval while the other may have a lot more than would be expected.

Similarly, σ appears to have limited effect on coverage rate. Confidence intervals take into account σ to some extent as they rely on an estimate of variance. Hence as standard deviations grow larger, so do confidence intervals. Therefore ~95% of true coefficients are still captured. The number of samples also seems to not have any effect - coverage rate appears to be uniform across n . This is because we replicated over a thousand samples, so from a frequentist perspective, as sample sizes become large or approach infinity, the true β should be within range of the confidence interval 95% of the time.

d)

For all combinations, log normal would have generated much larger 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:

$$\sigma^2$$

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

However because we z-score standardised the data in some cases it may be more accurate than the normal case - z-score scaling skews the data as every lognormal error has mean 0 and standard deviation 1. Obviously, if one draws from a distribution, not every sample will have this present. Initially, the interval length appears to be greater for lognormal distributions however in sample sizes greater than 20, the length is identical with their normal counterpart. This could be an indicator of slower convergence but again, the effect appears to be limited.

For unitary σ , the interval length is far smaller than their $\sigma = 2$ counterparts. This is because confidence intervals take into account the estimated variance which is correlated with sigma - therefore confidence intervals are wider when sigma is larger.

Sample size reduces interval length drastically for all combinations. This is because of the way confidence intervals are calculated - the standard error decreases with n and therefore so does the interval length.

Q2

a)

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

Table 2:

	ndf	sigma	Distribution	Bias	SD	CovRate	CovLength
1	20	1	Normal	0.555	0.323	0.594	1.279
2	20	1	LogNormal	0.583	0.322	0.579	1.281
3	20	2	Normal	0.549	0.514	0.813	2.129
4	20	2	LogNormal	0.576	0.515	0.834	2.150
5	200	1	Normal	0.572	0.099	0	0.361
6	200	1	LogNormal	0.576	0.094	0	0.361
7	200	2	Normal	0.580	0.158	0.035	0.604
8	200	2	LogNormal	0.574	0.153	0.035	0.606
9	2,000	1	Normal	0.577	0.029	0	0.113
10	2,000	1	LogNormal	0.579	0.030	0	0.113
11	2,000	2	Normal	0.578	0.049	0	0.190
12	2,000	2	LogNormal	0.578	0.049	0	0.189

It is clear that omitted variable bias is present here with an apparent consistent bias of ~ 0.577 . This would be ~ 0.33 if X_1 was not rescaled. We can calculate the theoretical bias here:

$$y = X_1\beta_1 + X_2\beta_2 + X_3\beta_3 + U$$

$$X_1 = \alpha_1 Z_1 + \alpha_3 X_3 + V_i$$

OLS estimates:

$$E(\hat{\beta}_1) = E \frac{\text{Cov}(X_1, y)}{\text{Var}(X_1)}$$

We use an approximation here as every X_1 value will be scaled slightly differently due to the method prescribed in rescaling the standard deviation to one.

$$E(\hat{\beta}_1) = E \frac{\text{Cov}(X_1, \frac{1}{\sqrt{3}}(X_1\beta_1 + X_2\beta_2 + X_3\beta_3 + U))}{1}$$

The $1/\sqrt{3}$ is a result of rescaling X_1 to have a standard deviation of one. While it is unlikely that the estimated standard deviation is actually one each time (because it is estimated), this is a good approximation.

$$\begin{aligned} E(\hat{\beta}_1) &= E \frac{\text{Cov}(X_1, X_1\beta_1) + \text{Cov}(X_1, X_2\beta_2) + \text{Cov}(X_1, X_3\beta_3) + \text{Cov}(X_1, U)}{\sqrt{3}} \\ E(\hat{\beta}_1) &= E \frac{\beta_1 \text{Var}(X_1) + \beta_2 0 + \text{Cov}(Z_1 + X_3 + U, X_3\beta_3) + 0}{\sqrt{3}} \\ E(\hat{\beta}_1) &= E(\beta_1) + \frac{\text{Cov}(Z_1 + X_3 + U, X_3\beta_3)}{\sqrt{3}} \\ E(\hat{\beta}_1) &= E(\beta_1) + \frac{\beta_3 \text{Cov}(Z_1, X_3) + \beta_3 \text{Cov}(X_3, X_3) + \text{Cov}(U, X_3)}{\sqrt{3}} \end{aligned}$$

$$E(\hat{\beta}_1) = \beta_1 + \frac{1}{\sqrt{3}}$$

$$\text{Bias} = \frac{1}{\sqrt{3}} \approx 0.5773$$

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 fulfil 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). Because confidence length decreases with sample size and is centered around the mean of $\hat{\beta}$, the true beta is not even nearly within the range of the confidence interval for $\hat{\beta}$. Standard deviation is much larger here where the omitted variable bias is present. This is explained by the additional X_3 factor to y which appears as noise and adds variance within the data. This is more explicitly explained in the next question.

2b)

Bias is expected to disappear as condition 1 is removed. There may be higher standard deviation than question one still as the variance of β can be calculated like so:

$$\text{Var}(\beta) = \text{Var}((X^T X)^{-1} X^T y)$$

$$\text{Var}(\hat{\beta}) = ((X^T X)^{-1} X^T) \sigma I ((X^T X)^{-1} X^T)^T$$

$$\text{Var}(\hat{\beta}) = (X^T X)^{-1} X^T \sigma X (X^T X)^{-1}$$

$$\text{Var}(\hat{\beta}) = \sigma (X^T X)^{-1}$$

And so depends on σ . In this case as X_3 is unobservable, it simply adds to the noise, thereby increasing variance. This in turn increases the length of the confidence interval.

Coverage rate is likely to improve as well due to no omitted variable bias. This means the range of values predicted by the confidence interval matches the true mean of the data.

Q3

- a) Yes, it satisfies the two conditions required for an IV.
 - 1) Z is not correlated with Y other than through X_1
 - 2) Z is correlated relatively strongly with X_1

We know Z is not correlated with Y as it is generated independently.

b)

```
dataQ3 = regMCLoop$new(nrange=c(20,200,2000),sigmarange=c(1,2),Urange=c("norm","lognorm"))
dataQ3$loop(b3=1,IV=TRUE)
```

Asymptotic convergence is slower but bias is reduced eventually. IV estimates are notorious for having finite sample bias and thus for $n = 20$, bias is still very high. As n increases however the bias eventually appears to dissipate.

Table 3:

	ndf	sigma	Distribution	Bias	SD	CovRate	CovLength
1	20	1	Normal	-0.129	1.224	0.945	5.381
2	20	1	LogNormal	-0.341	7.522	0.942	60.594
3	20	2	Normal	-0.173	2.203	0.949	11.399
4	20	2	LogNormal	-0.134	1.446	0.964	5.941
5	200	1	Normal	-0.005	0.172	0.957	0.693
6	200	1	LogNormal	-0.008	0.177	0.949	0.696
7	200	2	Normal	-0.007	0.288	0.953	1.093
8	200	2	LogNormal	-0.012	0.283	0.958	1.097
9	2,000	1	Normal	-0.002	0.056	0.946	0.215
10	2,000	1	LogNormal	-0.001	0.056	0.952	0.215
11	2,000	2	Normal	0.004	0.084	0.961	0.340
12	2,000	2	LogNormal	0.001	0.087	0.960	0.340

Standard deviations are higher due to higher presence of multicollinearity and use of IV estimates. IV estimates produce higher standard errors than ordinary OLS as it relies on more estimations. To demonstrate this:

$$Var(\hat{\beta}_1) = \frac{\sigma^2}{SST_x * R_{x,z}^2}$$

Under OLS though:

$$Var(\hat{\beta}_1) = \frac{\sigma^2}{SST_x}$$

Given $R_{x,z}^2$ is positive and between 0 and 1, the variance will always be smaller for OLS case holding other factors constant.

Interval lengths are also enlarged as a result of increased standard deviations. Coverage rate appears to be as appropriate, as the bias has been reduced. While we may have expected coverage rate to have decreased slightly, the coverage length is quite large due to higher standard deviation and thus are able to capture the bias.

There are some anomalies however in the lognormal values near 1. It seems difficult to explain this, though upon analysing the log normal distribution, we realise that it has a very fat tail. In a sample size of only 20, the presence of say, a positive outlier is less likely to be offset by a negative outlier. Therefore, some $\hat{\beta}$ s are more likely to be extreme and thus estimated standard deviations can become very high under small samples. We suspect further that the IV regression itself does not help due to its higher standard deviation.

Further evidence of this conclusion is that varying the seed often results in one lognormal value becoming very high and one becoming more appropriate. From this, we can attribute this to some random effect, however the likelihood suggests it is due to the heavy tailed distribution.

Q4

```
dataQ4 = regMClloop$new(nrange=200,sigmarange=1,Urange="norm",a1range=c(0.8,seq(0.6,0,by=-0.1)))
dataQ4$loop(b3=1,IV=TRUE)
```

Table 4:

	ndf	sigma	Distribution	Bias	SD	CovRate	CovLength
1	200	1	Normal	-0.002	0.208	0.957	0.824
2	200	1	Normal	-0.023	0.271	0.961	1.059
3	200	1	Normal	-0.053	0.342	0.961	1.313
4	200	1	Normal	-0.018	1.428	0.953	4.428
5	200	1	Normal	-0.481	12.181	0.954	141.050
6	200	1	Normal	-0.967	19.669	0.966	447.843
7	200	1	Normal	0.453	17.775	0.974	754.029
8	200	1	Normal	0.669	25.703	0.990	1,361.399

As the correlation between Z and X decreases, Z becomes takes on more and more of the characteristics of a weak instrument until it becomes an invalid instrument at 0. Bias increases accordingly, and the convergence of the estimator is slower. Because convergence for every given sample size increases with the strength of the IV, standard deviation at every sample size decreases accordingly. This has the follow over effect of increasing confidence interval lengths. Bias is also higher due to finite sample bias in IV regression.

Obviously when α_1 becomes zero it is no longer an instrument and therefore the IVreg method is invalid. However when the instrumental variable is invalid, the bias and standard deviation curiously drops. This can be explained for a number of reasons which were outlined in the previous assignment. We again cite Morgan and Winship. Firstly, IV estimates can always be estimated as sample covariances are never exactly equal to zero (cannot draw an infinite sample). Therefore even for invalid instruments, Iv estimates can be computed. In fact, estimators for the standard errors are constructed under the assumption that the IV is valid and thereby generate “artificially small standard errors” (Morgan and Winship, 2015).

Q5

```
dataQ5 = regMClloop$new(nrange=200,sigmarange=1,Urange="norm",a2range=seq(0,10,by=2))
dataQ5$loop()
```

Table 5:

	ndf	sigma	Distribution	Bias	SD	CovRate	CovLength
1	200	1	Normal	-0.0002	0.073	0.938	0.280
2	200	1	Normal	0.009	0.113	0.941	0.428
3	200	1	Normal	0.009	0.189	0.938	0.705
4	200	1	Normal	-0.005	0.257	0.946	1.013
5	200	1	Normal	0.010	0.339	0.945	1.322

Multicollinearity is present for $\alpha_2 > 0$ here.

We know that X_1 can be expressed as:

$$X_1 = \frac{1}{\sqrt{2 + \alpha_2^2}}(Z_1 + \alpha_2 X_2 + V)$$

This is due to normalising the standard deviation to 1 again.

When multicollinearity is presented in this way, as α_2 increases, X_1 becomes more and more like X_2 (X_2 begins to dominate any other terms within X_1). As such any small variation in X_2 leads to large changes in X_1 . The data is also not well spread and is narrow in the direction of X_2 . Therefore small variations of data may move the span of X around largely. This leads to highly variable estimates and may even begin to affect bias simply due to finite samples. Therefore what we see is with increased values of α_2 we get larger and larger standard deviations (subsequently larger confidence lengths). Bias also increases accordingly simply as a result of having finite samples. Coverage rate is largely unaffected, as there the higher standard deviations are compensated for by the larger confidence lengths. Any small bias that is present is too small to affect coverage rates.