

## Una alternativa a *Stata*: usando *R* para estimación de Modelos de Regresión

Bolívar Morales-Oñate<sup>1</sup>, Carlos Jiménez-Mosquera<sup>2</sup>, Paúl Méndez<sup>3</sup>

<sup>1</sup> Facultad de Ciencias, Escuela Superior Politécnica de Chimborazo, Riobamba

<sup>2</sup> Colegio de Ciencias e Ingenierías, Universidad San Francisco de Quito, Quito

<sup>3</sup> Departamento de Ingeniería Matemática, Universidad de Concepción, Concepción

\*Autor para correspondencia / Corresponding author [bolivar.morales@esepoch.edu.ec](mailto:bolivar.morales@esepoch.edu.ec), [cjimenez@usfq.edu.ec](mailto:cjimenez@usfq.edu.ec), [pmendez@ing-mat.udec.cl](mailto:pmendez@ing-mat.udec.cl)

## An alternative to *Stata* Using *R* for estimation of Regression Models

### Resumen

*RegUtils* es un paquete implementado para el entorno y lenguaje de programación R, que contiene un conjunto de funciones destinadas a facilitar el uso de este software como una alternativa al software comercial *Stata* para la estimación de Modelos de Regresión. Las funciones aquí implementadas, están pensadas para facilitar la migración de modelos estimados en *Stata* al software R, utilizando una sintaxis en muchos sentidos equivalente a la de los comandos de *Stata*, pero que al mismo tiempo sea compatible con paquetes relevantes de R y los paradigmas de programación en R.

**Palabras clave:** Paquete *RegUtils*, *Stata*, R, migración de modelos.

Código JEL: C13,C87.

### Abstract

*RegUtils* is a package implemented for the environment and programming language R, which contains a set of functions to facilitate the use of this software as an alternative to commercial software *Stata* for estimating regression models. Functions within this package are designed to help with the migration of models developed in *Stata* to R. These are constructed keeping a syntax equivalent to *Stata* counterparts, while maintaining compatibility with R relevant packages and programming paradigms.

**Keywords:** *RegUtils* package, *Stata*, R, migration models.

## INTRODUCCIÓN

El análisis estadístico en general, y en particular, la estimación de regresiones lineales y no lineales, como muchas otras áreas, se apoya en la actualidad en el uso de herramientas computacionales para este análisis. Dos de las herramientas más usadas son: *R* y *Stata*. Estas facilitan, entre otras cosas, el uso de procesos de estimación de parámetros que de otra forma serían inalcanzables, por la complejidad e intensidad del cálculo involucrado. Trabajos como [1], ilustran este punto en el contexto del manejo de grandes bases de datos multinivel longitudinales.

En muchos problemas existe una relación inherente entre dos o más variables, y resulta necesario explorar la naturaleza de esta relación. El análisis de regresión es una técnica estadística para el modelado y la investigación de la relación entre dos o más variables. Unas aplicaciones generales serían: líneas de tendencia en series de tiempo, relaciones causales en medicina, modelos económicos, entre otros [2].

*R* es un lenguaje y entorno de programación para análisis estadístico y gráfico. Se trata de un proyecto de software libre, resultado de la implementación GNU del premiado lenguaje *S*. *R* y *S-Plus* –versión comercial de *S*– son, probablemente, los dos lenguajes más utilizados en investigación por la comunidad estadística, siendo además muy populares en el campo de la investigación biomédica, la bioinformática y las matemáticas financieras. A esto contribuye la posibilidad de cargar diferentes bibliotecas o paquetes con finalidades específicas de cálculo o procedimientos gráficos [3].

*Stata* es un paquete de software estadístico creado en 1985 por StataCorp. Es utilizado principalmente por instituciones académicas y empresariales dedicadas a la investigación, especialmente en economía, sociología, ciencias políticas, biomedicina y epidemiología. *Stata* permite, entre otras funcionalidades, la gestión de datos, el análisis estadístico, el trazado de gráficos y las simulaciones [4].

El objetivo principal de este trabajo es describir la implementación y sintaxis de un conjunto de comandos paralelos a la oferta de *Stata*, a fin de facilitar el uso de *R* para quienes están familiarizados o vienen de un entorno en el que se usaba principalmente *Stata*. Estos comandos se escogieron luego de analizar la oferta de modelos de regresión tanto en *Stata* como en *R*. La sección 2 de este trabajo ilustra el modelo teórico de cada procedimiento contenido en *RegUtils*. Luego, en la sección 3 se describe brevemente los pasos para la instalación y uso de *RegUtils*. La sección 4 presenta los resultados de la implementación de *RegUtils* donde se evidencia la equivalencia entre *R* y *Stata*. Finalmente, el trabajo concluye presentando las conclusiones y posibles desarrollos posteriores.

## MARCO TEÓRICO

Se han revisado 35 comandos *Stata* (ver anexo 1), 5 de ellos (14%) no tuvieron una equivalencia directa en *R*. Estos son: **areg**, comando que aborda los modelos de regresiones lineales de grandes conjuntos en variables binarias. **boxcox**, es un comando que maneja modelos de regresión con sus transformaciones. **eivreg** se utiliza para modelos de regresión de errores en las variables. **etregress**, maneja modelos de regresión lineal con efectos de tratamiento endógeno. **ivtobit** se usa para la estimación de modelos de regresión *tobit* con variables endógenas.



Las siguientes subsecciones están organizadas de tal manera que el lector pueda familiarizarse con la descripción y el modelo asociado a cada comando implementado en **RegUtil**s.

## Regresión Lineal con conjuntos grandes de variables binarias (**areg**)

### Descripción

Este modelo se utiliza cuando se desea ajustar un modelo de regresión lineal que tenga factores dentro del conjunto de variables explicativas de modo que el ajuste implique el uso de un gran número de variables binarias al momento de generar la matriz de diseño.

### Modelo

Suponga que desea estimar una regresión con un gran número de variables binarias a partir de  $n$  individuos:

$$y = \mathbf{X}\beta + \mathbf{d}_1\gamma_1 + \mathbf{d}_2\gamma_2 + \dots + \mathbf{d}_k\gamma_k + \epsilon \quad (1)$$

donde  $y$  es de dimensión  $n \times 1$ ,  $X$  (de dimensión  $n \times p$  donde  $p$  es el número de columnas de la matriz de diseño) es la matriz de variables explicativas (sin incluir las variables binarias),  $d_i$  son las variables binarias,  $\gamma_i$  son los coeficientes de las variables binarias y  $\epsilon$  es el error [5].

## Modelos de Regresión con transformaciones de Box Cox (**BOXCOX**)

### Descripción

Este modelo se utiliza para estimar –usando el método de máxima verosimilitud– los parámetros de regresión de variables a las cuales se les ha aplicado la transformación de Box-Cox.

### Modelo

En su trabajo seminal, [6] proponen una transformación de los datos que es útil para reducir el sesgo, estabilizar la varianza y concebir que los datos tengan una distribución más parecida a una normal, entre otras [7]. De entre los modelos *lhsonly*, *rhsonly*, *lambda* y *theta* que ofrece Box-Cox, para el ejemplo se utilizó el modelo *theta*:

$$y_j^{(\theta)} = \beta_0 + \beta_1 x_{1j}^{(\lambda)} + \beta_2 x_{2j}^{(\lambda)} + \dots + \beta_k x_{kj}^{(\lambda)} + \gamma_1 z_{1j} + \gamma_2 z_{2j} + \dots + \gamma_l z_{lj} + \epsilon_j \quad (2)$$

donde  $\epsilon \sim N(0, \sigma^2)$ , y está sujeta a una transformación Box-Cox con parámetro  $\theta$  y cada  $x_1, x_2, \dots, x_k$  son transformadas por Box-Cox con parámetro  $\lambda$ . Las variables  $z_1, z_2, \dots, z_l$  son covariables no transformadas [4].

## Modelos de regresión con errores en las variables (**eivreg**)

### Descripción

Este modelo implica el ajuste de regresiones en las cuales una o más variables independientes son medidas con ruido aditivo. El sesgo introducido por este ruido trata de ser compensando mediante el uso de un coeficiente de confiabilidad.

**Modelo**

Si una covariable del modelo tiene un error de medición, una regresión tradicional no estimaría adecuadamente su efecto. Además, los coeficientes de las demás covariables del modelo podrían estar sesgados debido a la presencia en el modelo de esa variable. Se puede ajustar el sesgo si se conoce la confiabilidad (*reliability*):

$$reliability = 1 - \frac{noise\ variance}{total\ variance} \tag{3}$$

Esto es, dado el modelo

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{u} \tag{4}$$

para alguna variable  $x_i$  en  $\mathbf{X}$ ,  $x_i$  es observada con error,  $x_i = x_i^* + e$  y la *varianza de ruido* es la varianza de  $e$ . La varianza total es la varianza de  $x_i$  [4]

**Modelos de regresión lineal con efectos de tratamiento endógeno (etregress)**

**Descripción**

Este modelo estima los parámetros de una regresión aumentada con una variable binaria endógena, donde el principal objetivo generalmente es parámetro que captura el efecto promedio en el tratamiento (*average treatment effect - ATE*). Esta variable endógena debe estar correlacionada con el tratamiento pero no con el error ni las variables explicativas del modelo principal.

**Modelo**

El modelo de regresión de efectos de tratamiento endógeno está compuesto de una ecuación para el resultado  $y_j$  y una ecuación para el tratamiento endógeno  $t_j$ ,

$$y_j = \mathbf{x}_j\beta + \delta t_j + \epsilon_j \tag{5}$$

$$t_j = \begin{cases} 1 & \text{si } \mathbf{w}_j\gamma + u_j > 0 \\ 0 & \text{en otro caso} \end{cases} \tag{6}$$

Donde  $\mathbf{x}_j$  son las covariables del modelo principal,  $\mathbf{w}_j$  son covariables usadas para modelar el tratamiento, y los términos de error  $\epsilon_j$  y  $u_j$  tienen una distribución normal bivariada con media cero y matriz de covarianza:

$$\begin{bmatrix} \sigma^2 & \rho\sigma = \lambda \\ \rho\sigma = \lambda & 1 \end{bmatrix} \tag{7}$$

Las covariables  $\mathbf{x}_j$  y  $\mathbf{w}_j$  son exógenas, esto es, no están correlacionadas con el término de error [4].



## Modelo de regresión *tobit* con variables endógenas (*ivtobit*)

### Descripción

Corresponde a modelos *tobit* donde una o más de las variables regresoras son determinadas endógenamente.

**Modelo** 
$$y_{1i}^* = \mathbf{y}_{2i}\boldsymbol{\beta} + \mathbf{x}_{1i}\boldsymbol{\lambda} + u_i \quad (8)$$

El modelo es:

$$\mathbf{y}_{2i} = \mathbf{x}_{1i}\boldsymbol{\Gamma}_1 + \mathbf{x}_{2i}\boldsymbol{\Gamma}_2 + \mathbf{v}_i \quad (9)$$

donde  $i = 1, \dots, N$ ;  $\mathbf{y}_{2i}$  es un vector  $1 \times p$  de variables endógenas;  $\mathbf{x}_{1i}$  es un vector  $1 \times k_1$  de variables exógenas;  $\mathbf{x}_{2i}$  es un vector  $1 \times k_2$  de instrumentos adicionales; y la ecuación  $y_{2i}$  está escrita en forma reducida. El modelo supone que  $(u_i, \mathbf{v}_i) \sim N(\mathbf{0})$ .  $\boldsymbol{\beta}$  y  $\boldsymbol{\lambda}$  son vectores de parámetros instrumentales,  $\boldsymbol{\Gamma}_1$  y  $\boldsymbol{\Gamma}_2$  son matrices de parámetros de forma reducida. y  $y_{1i}^*$  es observada en forma censurada [7].

## INSTALACIÓN

Como se ha mencionado, el anexo 1 contiene la información de todos los comandos investigados y sus paralelos entre Stata y R. La tabla 7 muestra los comandos que han sido desarrollados exclusivamente en RegUtils.

**Tabla 1.** Comandos implementados en RegUtils. Fuente: [3] y [4]. Elaboración propia

Comando en Stata	Comando en R	Descripción
<i>areg</i>	alm	Ajusta regresiones con variables dummy
<i>boxcox</i>	boxcox.r	Modelos de regresión Box–Cox
<i>eivreg</i>	eivlm	Regresión con errores en las variables
<i>etregress</i>	etreg	Regresión lineal con efectos endógenos
<i>ivtobit</i>	ivtobit	Regresión tobit con variables endógenas

Actualmente el paquete se encuentra en el repositorio Github: <https://github.com/bolimorales/RegUtils>.

Para utilizarlo se deben realizar los siguientes pasos:

1. Dependencias: **RegUtils** depende de los paquetes **MaxLik** [8], **Formula** [9], **car** [10], **sandwich** [11] y **censReg** [12]. Por lo tanto, es necesario instalar estos paquetes previamente.
2. Instalar el paquete **devtools** [13]. Entre otras cosas, este paquete permite instalar paquetes disponibles en el repositorio *GitHub*.
3. Ejecutar los siguientes comandos:
  - (a) `library(devtools)`
  - (b) `install_github("bolimorales/RegUtils")`
4. El comando `library(RegUtils)` le permite usar el paquete

Como es usual en los paquetes de R, se puede acceder a la ayuda de cada una de las funciones de la tabla 7. Por ejemplo, para acceder a la ayuda de **etreg** se ejecuta: **help(etreg)**.

## ESTIMACIÓN DE MODELOS

A continuación se presentan los resultados obtenidos de cada una de las funciones del paquete **RegUtils**. Note que se mantiene consistencia con el marco teórico. Es decir, cada uno de los modelos presentados en la sección 2 es ilustrado mediante ejemplos y su respectivo contraste entre R y Stata.

### Regresión Lineal con conjuntos grandes de variables binarias (*areg*)

#### Ejemplo

Se dispone de un conjunto de datos que describen las características de 74 autos. Las variables de interés se muestran en la tabla 2:

**Tabla 2.** Datos para la regresión con variables binarias. Fuente: [4]. Elaboración propia.

Variable	Descripción	Tipo
<i>mpg</i>	millas por galón	Numérica
<i>Weight</i>	peso en libras	Numérica
<i>gear ratio</i>	proporción entre plato y piñón	Numérica
<i>rep78</i>	medida (de 1 al 5) del registro de reparaciones del auto donde 1 es la peor y 5 es la mejor	Factor

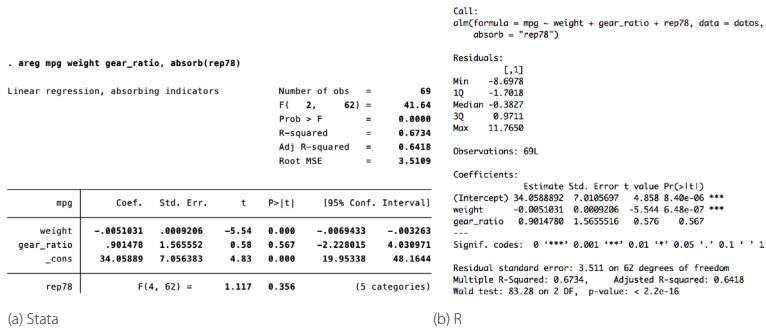
Se desea explicar las millas por galón en función de las demás variables presentadas en la tabla anterior. Una forma de resolver el problema es ajustar un modelo de regresión tradicional, lo cual precisaría de 4 niveles para el factor *rep78*. Sin embargo, **areg** permite, en lugar de realizar una prueba *t* para cada coeficiente, hacer una prueba *F* conjunta para todos los niveles de un determinado factor. Este proceso permite trabajar de forma eficiente sobre conjuntos de datos donde la cantidad de variables binarias generadas es muy grande, pero también nos permite centrar el análisis en un conjunto de variables independientes al margen de los grupos creados por las variables binarias [5].

#### Aplicación

Aquí se ajusta el modelo descrito en la parte anterior tanto en Stata (figura 1a) como en R (figura 1b).



Figura 1. Resultados de la regresión con variables binarias en Stata y R. Fuente y elaboración: autores



Note que el comando especifica **absorb(rep78)**, esto indica que se aplica una prueba *F* sobre el conjunto de los niveles de la variable *rep78*. Su valor *p* indica que, en conjunto, los niveles de *rep78* no son significativos a un  $\alpha = 0.05$ . Es decir que los niveles de *rep78* no influyen en las millas por galón. Además, los coeficientes de las demás variables se ajustan en forma estándar.

El resultado puede escribirse de la siguiente manera:

$$\text{mpg} = 34.058 - 0.0051\text{weight} + 0.901\text{gear\_ratio}. \quad (10)$$

Se aprecia que los niveles de la variable *rep78* no se imprimen debido al método utilizado para obtener más eficiencia, es decir, la prueba *F*. Note que los coeficientes para *weight*, *gear ratio* y la constante son exactamente los mismos al igual que el error estándar, valores *t* y, consecuentemente, el valor *p* (ver tablas de coeficientes). En la salida estándar de R no se reportan los intervalos de confianza para los coeficientes pero estos se pueden obtener a través del comando **confint** (usando el modelo guardado en *fit1*) con el que se obtiene los mismos valores (con una pequeña diferencia en el intercepto).

El primer estadístico *F* (igual a 41.64) mostrado en la salida de Stata corresponde al Chi-cuadrado del test de Wald (83.28) (ubicado al final de la salida de R en la figura 1b), dividido para los grados de libertad, es decir  $41.64 = 83.28/2$ . En ambos casos el Chi-cuadrado por lo que la conclusión de la prueba es la misma.

También se observa que coinciden el  $R^2$ , el  $R^2$  ajustado y los errores estándar residuales ( $MSE = 3.511$ ), los cuales suelen utilizarse para evaluar la bondad de ajuste de la regresión. La prueba *F* para *rep78*, no se presenta de forma estándar en R, pero se puede calcular mediante el comando **anova** con el que se obtiene el mismo valor, esto es, 1.117.

## Regresión con transformaciones de Box Cox (BOX- COX)

### Ejemplo

A continuación se utiliza un subconjunto de datos de la Segunda Encuesta Nacional de Salud y Nutrición (NHANES II), para crear un modelo de estimación del nivel individual de presión arterial. Las variables a utilizar en el modelo se muestran en la tabla 3.

**Tabla 3.** Datos para la regresión con transformaciones de BoxCox. Fuente: [4]. Elaboración propia

Variable	Descripción	Tipo
<i>bpdiast</i>	Presión arterial	Numérica
<i>bmi</i>	Índice de masa corporal	Numérica
<i>tcresult</i>	Colesterol sérico	Numérica
<i>age</i>	Edad	Numérica
<i>sex</i>	Sexo	Factor

Con el objeto de corregir la no linealidad en la relación con las variables se calcula la transformación de Box-Cox para las variables *bmi*, y *tcresult*, así como la transformación de la variable dependiente dentro del modelo de regresión.

### Aplicación

Para estimar este modelo se utiliza en R el comando **boxcox**, parte del paquete RegUtils; en Stata se utiliza el comando **boxcox**. Estos comandos permiten estimar varios tipos de modelos que incluyen el modelo *lambda*, donde se usa el mismo parámetro de transformación sobre todas las variables, modelos *theta*, que utilizan un parámetro diferente para la variable independiente, y modelos que mezclan variables transformadas y no transformadas. Los dos comandos permiten especificar un subconjunto de variables independientes a las cuales no se aplicará la transformación. En la figura 2b se muestran los resultados del ejemplo.





**Figura 2.** Resultados de la regresión con transformaciones de BoxCox en Stata y R.  
Fuente y elaboración: autores

Estimates of scale-variant parameters				
	Coef.	chi2(df)	P>chi2(df)	df of chi2
<b>Notrans</b>				
age	.003011	319.060	0.000	1
sex	-.1054007	243.284	0.000	1
_cons	5.835555			
<b>Trans</b>				
bmi	.0072041	1369.235	0.000	1
tresult	.004734	81.177	0.000	1
<b>/sigma</b>				
	.3348267			

Test	Restricted	log likelihood	chi2	Prob > chi2
H0:				
theta=lambda = -1		-40362.898	775.82	0.000
theta=lambda = 0		-39709.945	31.92	0.000
theta=lambda = 1		-39928.686	307.24	0.000

Coefficients:	Estimate	Chi sq	df	Chi sq	Pr(> t )
(Intercept)	5.8355e+00	NA	NA	NA	NA
bmi	8.7203e-02	1.3692e+03	1	< 2.2e-16	***
tresult	4.7339e-03	8.1177e+01	1	< 2.2e-16	***
age	3.8111e-03	3.1906e+02	1	< 2.2e-16	***
sex	-1.0540e-01	2.4328e+02	1	< 2.2e-16	***

(a) Stata

(b) R

Se escribe los resultados de la siguiente manera:

$$y_j^{(0.198)} = 5.83 + 0.087bmi^{(0.638)} + 0.0047tresult^{(0.638)} + 0.003age - 0.105sex \quad (11)$$

Se aprecia que todos los coeficientes son significativos a un  $\alpha = 0.05$  como se verifica en la figura 2b. Los coeficientes para lambda y theta coinciden, así como sus pruebas de hipótesis. Se presenta además el mismo número de observaciones, la misma estimación para sigma y de la verosimilitud,  $-39775$ . Para realizar la prueba del ratio de verosimilitudes en R se puede usar el comando `logLik.ratio.test` del paquete RegUtils.

## Regresión con errores en las variables (eivreg)

### Ejemplo

Se utiliza datos de una industria automotriz. Se asume que el peso de los autos es medido con ruido aditivo que puede ser aproximado por una confiabilidad de 0.85. Bajo este supuesto se realiza un modelo de regresión lineal simple para el precio utilizando las variables que se detallan en la tabla 4:

**Tabla 4.** Datos para la regresión con errores en las variables. Fuente: [4].Elaboración propia.

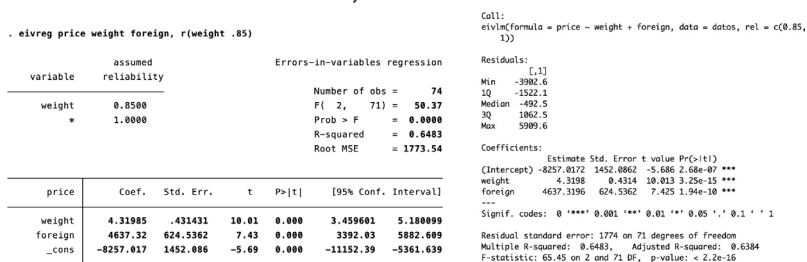
Variable	Descripción	Tipo
<i>price</i>	Precio	Numérica
<i>weight</i>	peso en libras	Numérica
<i>foreign</i>	Extranjero	Binaria

## Aplicación

En Stata el comando `eivreg` ajusta modelos de regresión con errores en las variables. En R se implementó el comando `eivlm` como parte del paquete RegUtils. En la figura 3 se muestran los resultados en Stata y R.

**Figura 3.** Resultados de la regresión con errores en las variables en Stata y R.

Fuente y elaboración: autores



(a) Stata

(b) R

Se puede escribir los resultados de la siguiente manera:

$$\text{price} = -8257.01 + 4.319\text{weight} + 4637.3\text{foreign}. \quad (12)$$

Los coeficientes habrían sido menores de no haber tomado en cuenta la medición con error de la variable *weight* (Ver [14] para más detalles de este tipo de modelos). El coeficiente de confiabilidad afecta el coeficiente y el error estándar obtenidos para la variable peso. Además se observa que los coeficientes de los parámetros y las pruebas de hipótesis coinciden.

## Modelos de regresión lineal con efectos de tratamiento endógeno (etregress)

### Ejemplo

Se utiliza un subconjunto de la base de datos sobre ingresos de mujeres de Estados Unidos en 1972 con edades entre 18 y 30 años, para modelar los efectos promedios en el tratamiento para la variable *union* (pertenece o no a un sindicato) sobre el ingreso. Las variables a ser incluidas en el modelo se detallan en la tabla 5.

**Tabla 5.** Datos para la regresión lineal con efectos de tratamiento endógenos. Fuente: [4].

Variable	Descripción	Tipo
<i>wage</i>	Ingreso	Númérica
<i>grade</i>	Años de instrucción	Númérica
<i>smsa</i>	Variable indicativa de pertenencia a un distrito estadístico	Binaria
<i>black</i>	Variable indicativa para Afro-Americanos	Binaria
<i>tenure</i>	Permanencia en el trabajo actual	Númérica
<i>south</i>	Variable indicativa para residencia en el sur	Binaria

Elaboración propia

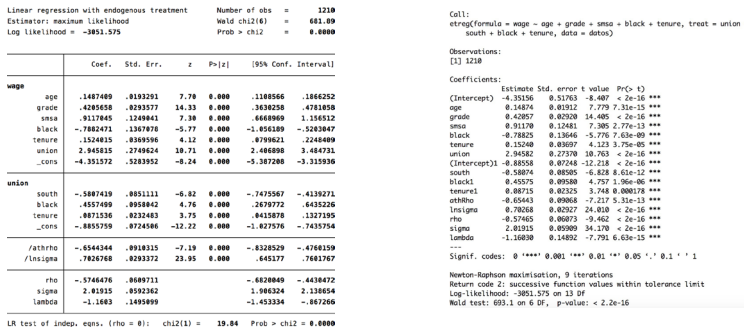


De estas variables se utiliza *south*, *black* y *tenure* para modelar la variable endógena unión.

### Aplicación

Para estimar el modelo se emplea el comando **etregress** en Stata y el comando **etreg** del paquete RegUtils en R. En la figura 5 se muestran los resultados obtenidos.

**Figure 4.** Resultados de la regresión lineal con efectos de tratamiento endógenos en Stata y R.



Fuente y elaboración: autores

El modelo ajustado del ejemplo sería:

$$\text{wage} = -4.35 + 0.148\text{age} + 0.420\text{grade} + 0.911\text{smsa} - 0.788\text{black} + 0.152\text{tenure} + 2.94\text{union} \quad (13)$$

$$\text{union} = \begin{cases} 1 & \text{si } -0.885 - 0.580\text{south} + 0.455\text{black} + 0.087\text{tenure} > 0 \\ 0 & \text{en otro caso} \end{cases} \quad (14)$$

$$\begin{bmatrix} 2.01^2 & -0.574 * 2.01 = -1.160 \\ -0.574 * 2.01 = -1.160 & 1 \end{bmatrix} \quad (15)$$

Se observa en la figura 8 que los coeficientes son significativos y que el efecto de tratamiento endógeno (coeficiente de union) es 2.94

## Modelo de regresión *tobit* con variables endógenas (*ivtobit*)

### Ejemplo

Se utiliza a continuación datos de ingresos de un grupo de mujeres asumiendo que todas las mujeres que deciden no trabajar reciben \$ 10,000 en pagos de asistencia social y la manutención de los hijos. En el modelo se incluirán las variables que se muestran en la tabla 6:

**Tabla 6.** Datos para la regresión tobit con variables endógenas. Fuente: [4]. Elaboración propia.

Variable	Descripción	Tipo
<i>fem_inc</i>	Ingreso	Númerica
<i>fem_educ</i>	Años de instrucción formal	Númerica
<i>kids</i>	Número de hijos	Númerica
<i>other_inc</i>	Otros ingresos del hogar	Númerica
<i>male_educ</i>	Años de instrucción formal del esposo	Númerica

En el modelo se considera la variable *other\_inc* como endógena, motivo por el cual se utiliza la variable *male\_educ* como variable instrumental.

### Aplicación

Para estimar el modelo se utiliza el comando `ivtobit` de Stata. En R se usa el comando del mismo nombre implementado como parte del paquete `RegUtils`. La figura 5 muestra los resultados obtenidos.

**Figura 5.** Resultados de la regresión tobit con variables endógenas en Stata y R. Fuente y elaboración: autores

```
Tobit model with endogenous regressors      Number of obs =      500
Log likelihood = -3226.0845                 Wald chi2(3) =      117.42
                                           Prob > chi2 =      0.0000
```

	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
other_inc	-.9045399	.1329762	-6.80	0.000	-1.165168 - .6439114
fem_educ	3.272391	.3968708	8.25	0.000	2.494538 4.050243
kids	-3.322357	.7218628	-4.59	0.000	-4.737332 -2.097322
_cons	19.24735	7.372391	2.61	0.009	4.797725 33.69697
/alpha	.2987654	.1379965	2.11	0.035	.0202972 .5612336
/lns	2.874031	.0586672	56.72	0.000	2.774725 2.973337
/lnv	2.813383	.0316228	88.97	0.000	2.751404 2.875363
s	17.78026	.097228			16.03422 19.55707
v	16.66621	.5270318			15.66461 17.73186

```
Instrumented: other_inc
Instruments: fem_educ kids male_educ

Wald test of exogeneity (/alpha = 0): chi2(1) = 4.44 Prob > chi2 = 0.0351

Obs. summary:      272 left-censored observations at fem_inc=10
                  228 uncensored observations
                   0 right-censored observations
```

(a) Stata

```
Call:
ivtobit(formula = fem_inc ~ other_inc + fem_educ + kids | male_educ +
fem_educ + kids, data = datos, left = 10)
```

```
Observations:
Total      Left-censored      Uncensored Right-censored
500        272                228                0
```

```
Coefficients:      Estimate Std. error t value Pr(> |t|)
(Intercept) 19.2528 6.9474 2.771 0.00558 **
other_inc -0.9047 0.1380 -6.913 4.75e-12 ***
fem_educ 3.2727 0.4170 7.849 4.20e-15 ***
kids -3.3128 0.6778 -4.888 1.02e-06 ***
sigma_u 18.3606 1.1331 16.204 < 2e-16 ***
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

BHHH maximisation, 6 iterations
Return code 2: successive function values within tolerance limit
Log-likelihood: -3226.085 on 7 DF

Wald test: 103.8 on 3 DF, p-value: < 2.2e-16

alpha: 0.2989144
Sigma:
[1,] 1.1 1.2
[1,] 337.11146 80.80733
[2,] 80.80733 277.77013
```

(b) R

El modelo estimado es:

$$fem\_inc = 19.24 - 0.9045other\_inc + 3.27fem\_educ - 3.31kids \quad (16)$$

En ambos comandos el estimador se calcula utilizando el método de máxima verosimilitud. Todas las variables del modelo son estadísticamente significativas a un  $\alpha = 0.05$  y se aprecia que son los mismos valores.



## CONCLUSIONES

El uso del paquete RegUtils, sin duda es una herramienta muy útil para usuarios que desean migrar o usar simultáneamente R y Stata dado que permite obtener modelos similares a aquellos ajustados con los comandos de Stata, con una sintaxis que resulta muy fácil de equiparar con la de Stata.

El paquete facilita que estos modelos se integren con otras funcionalidades y comandos de R. Sin embargo, se recomienda que una vez que se alcance el nivel adecuado de familiaridad con R, se explore otros paquetes y rutinas propias del lenguaje que, sin tener las mismas salidas de Stata pueden conducirnos al ajuste de modelos de utilidad.

Utilizar los comandos implementados en R para el manejo de estos modelos conlleva algunas ventajas sobre el software propietario, se podría mencionar entre estas el hecho de que no tendría costo para el usuario. Además, R es versátil para el desarrollo de rutinas, acceso al código para aprender de él y modificarlo en función de las necesidades del investigador y facilidades para el trabajo colaborativo al poder compartir análisis sin preocupación de que el receptor disponga de las licencias.

Como todos los proyectos realizados en R, se recomienda que todo usuario interesado se motive en generar retroalimentación para su perfeccionamiento. Por ejemplo podría colaborar con paquetes o correcciones de paquetes existentes en el Comprehensive R Archive Network – CRAN. También se puede contribuir a mejorar las rutinas del paquete RegUtils, ya que el código se encuentra disponible en el sitio github.

## REFERENCES

- [1] D. F. McCaffrey, J. R. Lockwood, K. Mihaly, T. R. Sass, et al., *A review of stata commands for fixed-effects estimation in non-linear models*, *Stata Journal*, 12 (2012), p.406.
- [2] B. Baltagi, *Econometric analysis of panel data*, John Wiley & Sons, 2002.
- [3] R Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, Vienna, Austria, 2017.
- [4] StataCorp, *Stata 15 Base Reference Manual*, Stata Press, College Station, TX, 2017.
- [5] M. C. Lovell, *A simple proof of the fwl (frisch-waugh-lovell) theorem*. 2006.
- [6] G. E. Box and D. R. Cox, *An analysis of transformations*, *Journal of the Royal Statistical Society. Series B (Methodological)*, (1964), pp.211– 252.
- [7] J. M. Wooldridge, *Introducción a la econometría: un enfoque moderno*, Editorial Paraninfo, 2006.
- [8] A. Henningsen and O. Toomet, *maxlik: A package for maximum likelihood estimation in R*, *Computational Statistics*, 26 (2011), pp.443– 458.
- [9] A. Zeileis and Y. Croissant, *Extended model formulas in R: Multiple parts and multiple responses*, *Journal of Statistical Software*, 34 (2010), pp. 1–13.
- [10] J. Fox and S. Weisberg, *An R Companion to Applied Regression*, Sage, Thousand Oaks CA, second ed., 2011.
- [11] A. Zeileis, *Econometric computing with hc and hac covariance matrix estimators*, *Journal of Statistical Software*, 11 (2004).
- [12] A. Henningsen, *censReg: Censored Regression (Tobit) Models*, 2017. R package version 0.5-26.
- [13] H. Wickham and W. Chang, *devtools: Tools to Make Developing R Packages Easier*, 2017. R package version 1.13.4.
- [14] N. R. Draper and H. Smith, *Applied regression analysis*, vol. 326, John Wiley & Sons, 2014.
- [15] B. Morales-Oñate, C. Jiménez-Mosquera, P. Méndez, and V. Morales-Oñate, *RegUtils: Tools for STATA users in Regression Models Estimation*, 2015. R package version 0.1.
- [16] A. Ghalanos, *rugarch: Univariate GARCH models.*, 2018. R package version 1.4-0.
- [17] D. Wuertz, T. Setz, and Y. Chalabi, *fArma: Rmetrics - Modelling ARMA Time Series Processes*, 2017. R package version 3042.81.
- [18] T. Coelli and A. Henningsen, *frontier: Stochastic Frontier Analysis*, 2017. R package version 1.1-2.
- [19] P. Chaussé, *Computing generalized method of moments and generalized empirical likelihood with R*, *Journal of Statistical Software*, 34 (2010), pp. 1–35.
- [20] O. Toomet and A. Henningsen, *Sample selection models in R: Package sampleSelection*, *Journal of Statistical Software*, 27 (2008).
- [21] O. Toomet, *intReg: Interval Regression*, 2015. R package version 0.2-8.
- [22] C. Kleiber and A. Zeileis, *Applied Econometrics with R*, Springer-Verlag, New York, 2008. ISBN 978-0-387-77316-2.
- [23] A. Henningsen and J. D. Hamann, *systemfit: A package for estimating systems of simultaneous equations in r*, *Journal of Statistical Software*, 23 (2007), pp. 1–40.
- [24] R. Koenker, *quantreg: Quantile Regression*, 2017. R package version 5.34.



- [25] W. N. Venables and B. D. Ripley, *Modern Applied Statistics with S*, Springer, New York, fourth ed., 2002. ISBN 0-387-95457-0.
- [26] Y. Rosseel, *lavaan: An R package for structural equation modeling*, *Journal of Statistical Software*, 48 (2012), pp. 1–36.
- [27] Y. Croissant and A. Zeileis, *truncreg: Truncated Gaussian Regression Models*, 2016. R package version 0.2-4.
- [28] G. Millo, *Robust standard error estimators for panel models: A uniting approach*, *Journal of Statistical Software*, 82 (2017), pp. 1–27.
- [29] K. Kashin, *panelAR: Estimation of Linear AR(1) Panel Data Models with Cross-Sectional Heteroskedasticity and/or Correlation*, 2014. R package version 0.1.
- [30] J. Pinheiro, D. Bates, S. DebRoy, D. Sarkar, and R Core Team, *nlme: Linear and Nonlinear Mixed Effects Models*, 2017. R package version 3.1-131.

## ANEXO

Tabla 7. Cuadro comparativo de comandos entre Stata y R. Fuente: [20] y [22]. Elaboración propia

Comando en Stata	Librería	Comando en R	Descripción
areg	RegUtils[18]	alm	Ajusta regresiones con variables dummy
arch	rugarch[8]	ugarchFit	Modelos de regresión con errores ARCH
arima	fArma[28]	armaFit	Modelos ARIMA
boxcox	RegUtils	boxcox.r	Modelos de regresión Box?Cox
cnsreg	stats[20]	lm	Regresión lineal restringida
eivreg	RegUtils	eivlm	Regresión con errores en las variables
etregress	RegUtils	etreg	Regresión lineal con efectos endógenos
frontier	frontier[4]	sfa	Modelos estocásticos frontier
gmm	gmm[3]	gmm	método generalizado de momentos de estimación
heckman	sampleSelection[24]	selection	Modelos de selección Heckman
intreg	intReg[23]	intReg	Regresión intervalo
ivregress	AER[13]	ivreg	Regresión de variables instrumentales simples
ivtobit	RegUtils	ivtobit	Regresión tobit con variables endógenas
newey	sandwich[29]	NeweyWest	Regresión de Newey ?West con errores estándar
nl	stats	nls	Estimación no lineal de mínimos cuadrados
nlsur	systemfit[10]	nlsystemfit	Estimación de sistemas de ecuaciones no lineales
qreg	quantreg[14]	rq	Regresión cuantil (incluye mediana)
reg3	systemfit	systemfit	Regresión de mínimos cuadrados en tres estados (3SLS)
rreg	MASS[25]	rlm	Un tipo de regresión robusta
gsem	lavaan[21]	cfa	modelos de ecuaciones estructurales generalizadas
sem	lavaan	cfa	modelos de ecuaciones estructurales lineales
sureg	systemfit	systemfit	regresión aparentemente no relacionada
tobit	censReg[9]	censReg	Regresión tobit
truncreg	truncreg[5]	truncreg	Regresión truncada
xtabond	plm[17]	pgmm	Estimación de datos de panel dinámicos de Arellano?Bond
xtdpd	plm	pgmm	Estimación de datos de panel dinámicos lineales
xtfrontier	frontier	sfa	Modelos frontier de datos de panel dinámicos
xtgls	panelAR[12]	panelAR	Modelos GLS de datos de panel
xhtaylor	plm	pht	Modelos de estimación de errores de Hausman?Taylor
xtintreg	plm	survreg	Modelos de regresión de datos de panel en intervalos
xtivreg	plm	plm	Regresión de datos de panel con variables instrumentales (2SLS)
xtpcse	panelAR	panelAR	Regresión lineal con errores estándar
xtreg	plm	plm	Modelos lineales con efectos aleatorios
xtregar	nlme[19]	lme	Modelos lineales con efectos aleatorios AR(1)
xttobit	censReg	censReg	Modelos tobit de datos de panel