

TASA DE CRIMINALIDAD EN USA



MODELOS LÍNEALES TRABAJO FINAL

Ignacio Acosta - Sofía Itté - Mauro Loprete
1er semestre 2021

Índice

1. Introducción	2
2. Análisis Exploratorio de datos	3
2.1. Análisis Univariado	3
2.1.1. Histogramas y Barplots	4
2.1.2. Medidas de resumen	6
2.1.3. Correlación entre variables	6
3. Especificación y selección de modelos inicial	7
3.1. Modelo Inicial	8
4. Diagnóstico	9
4.1. Análisis de multicolinealidad	9
4.2. Método Stepwise	10
4.3. Búsqueda de observaciones influyentes	11
4.4. Heterosedasticidad	11
4.4.1. Test de Breusch-Pagan	11
4.5. Normalidad	12
5. Anexo	13
5.1. Selección de Modelos : StepWise	13
5.2. Script de R	17
6. Bibliografía	17

Índice de figuras

1. Histograma de la Tasa de Criminalidad	2
2. Histogramas (1)	4
3. Histogramas (2)	5
4. Mapa de correlación de variables incluidas	6

1. Introducción

El objetivo de este informe es la construcción de un modelo de regresión lineal múltiple que explique la tasa de criminalidad en USA (número de ofensas reportadas a la policía por habitante).

Para ello se hará uso de una base de datos con un conjunto de variables que en principio se encuentran relacionadas con la variable a explicar.

Haciendo uso de las distintas técnicas estadísticas aprendidas en el curso se buscará descartar variables cuyo aporte no sea suficientemente significativo. Esto busca llegar a un modelo final eficiente (es decir, que explique la tasa de criminalidad de manera acertada haciendo uso de la menor cantidad de variables posibles).

En el transcurso del texto se pondrán a prueba las distintas hipótesis centrales del modelo, tales como la normalidad de los errores y la heteroscedasticidad.

También se trabajará con las observaciones y la existencia de algunas que aporten el mismo nivel de información.

El herramienta gráfico juega un rol fundamental al momento de transmitir la información de manera concisa y entendible. El mismo se encuentra respaldado por tablas que resumen la información de manera más detallada.

Cabe destacar que las observaciones se corresponden con los distintos estados de Estado Unidos.

De manera general se muestra el comportamiento de y de manera resumida:



Figura 1: Histograma de la Tasa de Criminalidad

Es claro que y cuenta con una distribución medianamente asimétrica, tiene un intervalo modal entre los valores de 500 y 800 ofensas (13 estados presentan una tasa de criminalidad comprendida entre esos valores).

2. Análisis Exploratorio de datos

El objetivo de esta sección es presentar las variables a estudiar y como las mismas se relacionan entre sí.

Para ello se hará uso de distintas medidas de resumen univariadas y bivariadas, así como también un herramental gráfico variado que simplificará el entendimiento de las mismas.

Es esta sección fundamental al momento de discutir el modelo final y como a partir de distintas técnicas estadísticas aprendidas en el curso se puede simplificar el *modelo completo* que se presentará en la sección siguiente.

2.1. Análisis Univariado

En esta primer sección se hará especial énfasis en las variables por sí mismas.

Se estudiarán medidas de resumen y a partir de histogramas tendremos un primer acercamiento a la distribución de las mismas y su comportamiento.

Nombre	Descripción	Clasificación
Y	Tasa de criminalidad, número de ofensas reportadas a la policía por habitante	Cuantitativa
M	Número de hombres entre 14 y 24 años cada 1000 habitantes	Cuantitativa
So	Variables indicadora de los estados del sur (0=No, 1=Si)	Cualitativa
Ed	Índice que refleja la escolaridad del estado	Cuantitativa
Po1	Gasto per cápita en policía realizado por el gobierno estatal o local en 1960	Cuantitativa
Po2	Gasto per cápita en policía realizado por el gobierno estatal o local en 1959	Cuantitativa
LF	Tasa de participación en la fuerza laboral civil de sexo masculino entre 14 y 24 años, cada 1000 habitantes	Cuantitativa
M.F	Número de hombres por cada 1000 mujeres	Cuantitativa
Pop	Tamaño de la población del estado cada 100000 habitantes	Cuantitativa
NW	Número de no caucásicos cada 1000 habitantes	Cuantitativa
U1	Tasa de desempleo urbana de hombres entre 14 y 24 años por 1000 habitantes	Cuantitativa
U2	Tasa de desempleo urbana de hombres entre 35 y 39 años por 1000 habitantes	Cuantitativa
GDP	Producto bruto interno per cápita	Cuantitativa
Ineq	Desigualdad del ingreso	Cuantitativa
Prob	Probabilidad de encarcelamiento	Cuantitativa
Time	Tiempo promedio de estadía en cárceles estatales	Cuantitativa

Cuadro 1: Variables a trabajar

2.1.1. Histogramas y Barplots



Figura 2: Histogramas (1)



Figura 3: Histogramas (2)

Como se verá en los histogramas presentados a continuación y haciendo uso de la tabla (más precisamente del **CV**) es claro que las variables, de manera generalizada, presentan una variabilidad baja.

De manera más específica, los histogramas de las variables M, Po1, Nw, Po2, M.F, Pop, U1, U2, Prob y Time cuentan con una distribución asimétrica. La variabilidad entre los valores comprendidos hasta la mediana (aunque baja, como ya se mencionó) es menor que en el resto de las observaciones.

En el caso de la variable GDP y LF, la distribución a diferencia del resto es aproximadamente simétrica. La mediana y la media difieren en un número despreciable.

La variable Ineq también cuenta con una distribución asimétrica pero a diferencia de las demás, cuenta con menor variabilidad entre las observaciones en el tramo central (primer cuartil a tercer cuartil).

2.1.2. Medidas de resumen

Se presenta en forma de tabla el resumen de las variables numéricas. En el mismo se presenta el valor mínimo y máximo de cada variable, medidas de tendencia central tales como lo son el primer y tercer cuartil, junto a la mediana.

A su vez, para estudiar la dispersión se incluye la media aritmética y una medida de variabilidad de la misma, el coeficiente de variación.

Cuadro 2: Medidas descriptivas para variables numéricas

Variable	Min	1er Qu.	Mediana	3er Qu.	Max	Media	CV*100
Número de Hombres 14-24 / 1.000	119.0	130.0	136.0	146.0	177.0	138.6	9.1
Índice Escolaridad	87	98	108	114	122	106	11
Gasto per cápita 1.960	45	62	78	104	166	85	35
Gasto per cápita 1.959	41	58	73	97	157	80	35
Tasa participación masculina 14-24 por 1.000	480.0	530.5	560.0	593.0	641.0	561.2	7.2
Hombres cada 1.000 mujeres	934	964	977	992	1071	983	3
Población cada 100.000	3	10	25	42	168	37	104
Número de no caucásicos cada 1.000 habitantes	2	24	76	132	423	101	102
Tasa desempleo urbana Hombres 14-24 por 1.000	70	80	92	104	142	95	19
Tasa desempleo urbana Hombres 35-39 por 1.000	20	28	34	38	58	34	25
Producto bruto interno per cápita	288	460	537	592	689	525	18
Desigualdad ingreso	126	166	176	228	276	194	21
Probabilidad Encarcelamiento	0.69	3.27	4.21	5.45	11.98	4.71	48.28
Tiempo de estadía en cárceles	12	22	26	30	44	27	27
Tasa de criminalidad	342	658	831	1058	1993	905	43

2.1.3. Correlación entre variables

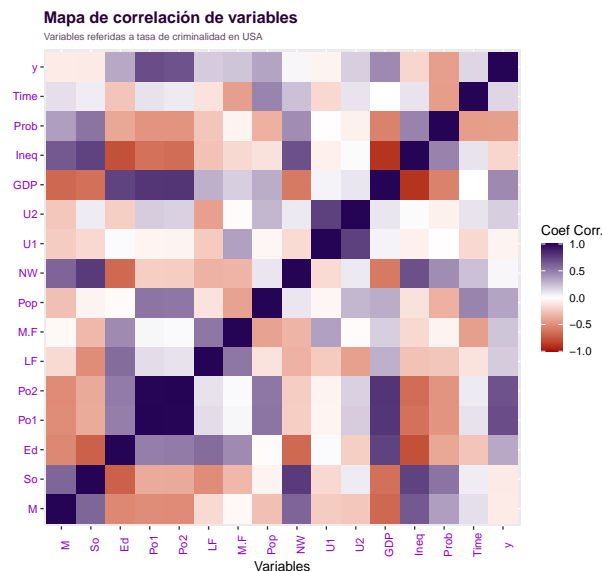


Figura 4: Mapa de correlación de variables incluidas

3. Especificación y selección de modelos inicial

El objetivo de esta sección es la aplicación de las distintas técnicas estadísticas impartidas en el curso para así llegar a un modelo final que no solo sea significativo al momento de estimar a y , sino que también se adecúe a los supuestos y propiedades deseadas (ganando de esta manera fidelidad).

En una primer instancia se planteará un *modelo inicial* constituido por parte de las variables de las cuales se poseen datos.

Claro está, se podría haber planteado en primera instancia un *modelo completo* (es decir, que contenga absolutamente a todas las variables). No es esto errado, pero si desconsiderado con el extenso análisis descriptivo planteado con anterioridad.

Como ya es sabido, variables que presentan una correlación muy alta no son marginalmente significativas al momento de definir la variable de respuesta.

Esto se evidencia en los tests de hipótesis en donde se analiza el aporte de cada variable dada las demás variables. Una correlación alta entre variables, podría indicar que parte de la información que aportan una de ellas está también presente en otra y esa cantidad de información se vió cuantificada de manera previa. Lo que tarde o temprano llevaría a descartar alguna de ellas.

En principio, a partir del “arsenal” descriptivo es claro que:

- **Ineq** y **GDP** tienen una correlación negativa altísima (-0.884).
- **P01** y **P02** poseen una correlación negativa casi perfecta (0.994)
- **Ineq** y **Ed** tienen también una correlación negativa bastante alta (-0,794)
- **SO** y **NW** mantienen una correlación positiva y de nivel alto (0,767)

A partir de lo afirmado, se procederá a “descartar”^a alguna de las variables que constituye cada dupla respaldándose en el valor del coeficiente de correlación existente entre las variables explicativas y y .

$$\rho_{y,Ineq} = -0,179$$

$$\rho_{y,P01} = 0,688$$

$$\rho_{y,Ed} = 0,323$$

$$\rho_{y,NW} = 0,03$$

$$\rho_{y,GDP} = 0,441$$

$$\rho_{y,P02} = 0,667$$

$$\rho_{y,SO} = -0,09$$

Se elige aquella variable cuyo coeficiente de correlación con y en valor absoluto sea mayor.

Trás esto, se decide que **Ineq**, **P02** y **NW** no sean incluídos en el modelo inicial ya que se entiende que las mismas no tendrán un aporte significativo en presencia de sus pares.

A manera de resumen podría decirse que este primer acercamiento al modelo sigue fielmente el principio de *parsimonia*¹.

Quedan entonces determinadas las variables a conformar el modelo inicial, que será analizado con detenimiento en la sección siguiente.

¹Frugalidad y moderación en los gastos.

3.1. Modelo Inicial

Como una primera aproximación, se construye un modelo donde se incluyen todas las variables de la tabla de datos, en concreto el siguiente modelo de regresión:

$$\hat{y} = \beta_0 + \beta_1 Time + \beta_2 Prob + \dots \beta_M M$$

Cuadro 3: Test sobre el modelo completo

R^2_{adj}	RSE	F Obs.	P-valor*100	Regresión.gl	Residuos.gl
70.781	209.064	8.429	0	15	31

Recordando que el R^2_a hace referencia al porcentaje de variabilidad de \mathbf{y} que es explicada con el modelo estimado, se considera al mismo como *aceptable*. Por otro lado, haciendo referencia a la significación del modelo, se consideranda el siguiente test de hipótesis y el estadístico F :

$$H_0) \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1) \text{No } H_0$$

$$F_{obs} = \frac{SCE/Regresion.gl}{RSE^2} = \frac{SCE/Regresion.gl}{SCR/Residuos.gl} = \frac{\sum (\hat{y}_i - \bar{y})^2 / Regresion.gl}{\sum (y_i - \hat{y}_i)^2 / Residuos.gl}$$

Siendo SCE la suma de cuadrados explicados por la regresión y RSE^2 el cuadrado del error estándar de los residuos, resulta para este caso particular $SCE = 5.525.982$ y $RSE^2 = 43707,766$, de esta manera se obtiene el F_{obs} que permite rechazar H_0 y así afirmar que el modelo es estadísticamente significativo para explicar a \mathbf{y} .

A continuación se testea la significación de cada variable en forma independiente, los resultados se muestran en el siguiente cuadro:

Cuadro 4: Estimación, error estándar y test individual del modelo completo

Variable	Estimación	Error estándar	Estadístico F	P valor	$(H_0^{\alpha=0.05}) \beta_i = 0$
Intercepto	-5984.288	1628.318	-3.675	0.001	Se rechaza H_0
Número de Hombres 14-24 / 1.000	8.783	4.171	2.106	0.043	Se rechaza H_0
Indicadora Estado Sur	-3.803	148.755	-0.026	0.980	No se rechaza H_0
Índice Escolaridad	18.832	6.209	3.033	0.005	Se rechaza H_0
Gasto per cápita 1.960	19.280	10.611	1.817	0.079	No se rechaza H_0
Gasto per cápita 1.959	-10.942	11.748	-0.931	0.359	No se rechaza H_0
Tasa participación masculina 14-24 por 1.000	-0.664	1.470	-0.452	0.655	No se rechaza H_0
Hombres cada 1.000 mujeres	1.741	2.035	0.855	0.399	No se rechaza H_0
Población cada 100.000	-0.733	1.290	-0.568	0.574	No se rechaza H_0
Número de no caucásicos cada 1.000 habitantes	0.420	0.648	0.649	0.521	No se rechaza H_0
Tasa desempleo urbana Hombres 14-24 por 1.000	-5.827	4.210	-1.384	0.176	No se rechaza H_0
Tasa desempleo urbana Hombres 35-39 por 1.000	16.780	8.234	2.038	0.050	No se rechaza H_0
Producto bruto interno per cápita	0.962	1.037	0.928	0.361	No se rechaza H_0
Desigualdad ingreso	7.067	2.272	3.111	0.004	Se rechaza H_0
Probabilidad Encarcelamiento	-48.553	22.724	-2.137	0.041	Se rechaza H_0
Tiempo de estadía en cárceles	-3.479	7.165	-0.486	0.631	No se rechaza H_0

A partir del cuadro presentado y conforme a los tests realizados, se ve claramente que son tan solo 3 las variables que de manera independiente (y **muy importante, en presencia de todas las demás**) logran un aporte significativo al momento de explicar el comportamiento de la tasa de criminalidad.

Ellas son **PO1**: Gasto per cápita en policía en 1960, **U2**: Tasa de desempleo Urbana de hombres entre 35 y 39 años por 1000 habitantes y por último **Prob*1000** Probabilidad de encarcelamiento cada 1000 habitantes.

¿Significa esto que se debe descartar el resto de las variables y plantear un modelo caracterizado por tan solo las 3?, la respuesta es **no**.

Como bien se menciona anteriormente, los tests analizan el aporte dada las demás variables. Una correlación alta entre variables, podría indicar que parte de la información que aportan una de ellas está también presente en otra y esa cantidad de información se vio cuantificada de manera previa.

Con base en esta última afirmación es que se promueve el uso de distintas técnicas que nos permitirán elegir las variables de manera más acertada (y teniendo en cuenta este panorama).

4. Diagnóstico

4.1. Análisis de multicolinealidad

Considerando el problema brevemente mencionado, se analizara la multicolinealidad (aproximada) de las variables independientes del modelo planteado, esto es relevante ya que si existe una relación lineal en la matriz de diseño, esto impactaría directamente en la varianza de los regresores $\beta_k = (X^T X)^{-1} X^T$, haciendo que las estimaciones varíen ante pequeñas variaciones en nuestras observaciones y las predicciones serían menos confiables.

Estamos en frente a un problema de multicolinealidad aproximada cuando es posible afirmar que existe una relación lineal entre las variables explicativas. El término aproximado refiere al hecho que en el caso que se cumpla el fenómeno de forma exacta, la matriz no sería invertible y no existirían estimaciones únicas de los regresores (Teorema de Gauss Markov) y no estaríamos frente a estimadores eficientes (insesgados y de mínima varianza)

Recordando que :

$$\hat{\beta}_k \sim N\left(\beta, \sigma^2 (X^T X)^{-1}\right)$$

Como se menciono anteriormente, ante una posible relación lineal el determinante de la matriz $X^T X$ sería próximo a cero, obteniendo un determinante de la matriz inversa demasiado grande. Es decir para un σ^2 fijo la incertidumbre sería demasiado alta, considerando el hecho que en nuestra primera aproximación a un modelo de regresión es globalmente significativo pero salvo en una cantidad demasiado pequeña se puede afirmar que, de forma independiente existe una relación lineal con la tasa de criminalidad, es por esto que se cuantificara la intensidad de la multicolinealidad con el **Factor de inflación de varianza**.

El **VIF** nos indica en cuantas unidades se incrementa la varianza del estimador ante presencia de colinealidad y se define como :

$$VIF_j = \frac{1}{1 - R_j^2}$$

Donde R_j^2 hace referencia al coeficiente de determinación de una regresión que intenta establecer una relación lineal de X_j con las demás variables explicativas.

Pondremos a prueba las variables explicativas del modelo anteriormente mencionado y diremos que estamos frente a problemas de colinealidad con un $VIF \geq 10$, los resultados se muestran en el cuadro a continuación.

Cuadro 5: Prueba de multicolinealidad : Factor de incremento de Varianza **VIF**

Variable	VIF	Prueba
Número de Hombres 14-24 / 1.000	2.892	No hay problema de colinealidad
Indicadora Estado Sur	5.343	No hay problema de colinealidad
Índice Escolaridad	5.077	No hay problema de colinealidad
Gasto per cápita 1.960	104.659	Problema de colinealidad
Gasto per cápita 1.959	113.559	Problema de colinealidad
Tasa participación masculina 14-24 por 1.000	3.713	No hay problema de colinealidad
Hombres cada 1.000 mujeres	3.786	No hay problema de colinealidad
Población cada 100.000	2.537	No hay problema de colinealidad
Número de no caucásicos cada 1.000 habitantes	4.674	No hay problema de colinealidad
Tasa desempleo urbana Hombres 14-24 por 1.000	6.064	No hay problema de colinealidad
Tasa desempleo urbana Hombres 35-39 por 1.000	5.089	No hay problema de colinealidad
Producto bruto interno per cápita	10.530	Problema de colinealidad
Desigualdad ingreso	8.645	No hay problema de colinealidad
Probabilidad Encarcelamiento	2.809	No hay problema de colinealidad
Tiempo de estadía en cárceles	2.714	No hay problema de colinealidad

En base a esto, podemos afirmar que nuestro modelo presenta problemas con la colinealidad y es por esto que continuaremos con la selección a pasos por el método Stepwise, también que hay que recordar que en el caso de seleccionar variables por la correlación dos a dos entre ellas y quedarse con aquellas que tienen una correlación mas alta con y es un grave error ya que nos estamos olvidando del efecto que puede tener con las demás variables.

4.2. Método Stepwise

El método de Stepwise (basado en el F -Test) comienza seleccionando aquella que tiene una mayor correlación con la variable y , la segunda es aquel modelo que con la variable que se incluyó en el paso anterior maximiza el coeficiente de determinación siguiendo iterando hasta que no se cumpla con el criterio de entrada.

Cuadro 6: Estimación, error estándar y test individual tras aplicar el método Stepwise

Variable	Estimación	Error estándar	Estadístico F	P valor	$(H_0^{\alpha=0.05}) \beta_i = 0$
Intercepto	-5040.505	899.843	-5.602	0.000	Se rechaza H_0
Gasto per capita en policía 1960	11.502	1.375	8.363	0.000	Se rechaza H_0
Desigualdad del ingreso	6.765	1.394	4.855	0.000	Se rechaza H_0
Índice que refleja la escolaridad del estado	19.647	4.475	4.390	0.000	Se rechaza H_0
Número de hombres entre 14 y 24 / 1000	10.502	3.330	3.154	0.003	Se rechaza H_0
Probabilidad de encarcelamiento	-38.018	15.281	-2.488	0.017	Se rechaza H_0
Tasa de desempleo urbana hombres 35-39 años x 1000	8.937	4.091	2.185	0.035	Se rechaza H_0

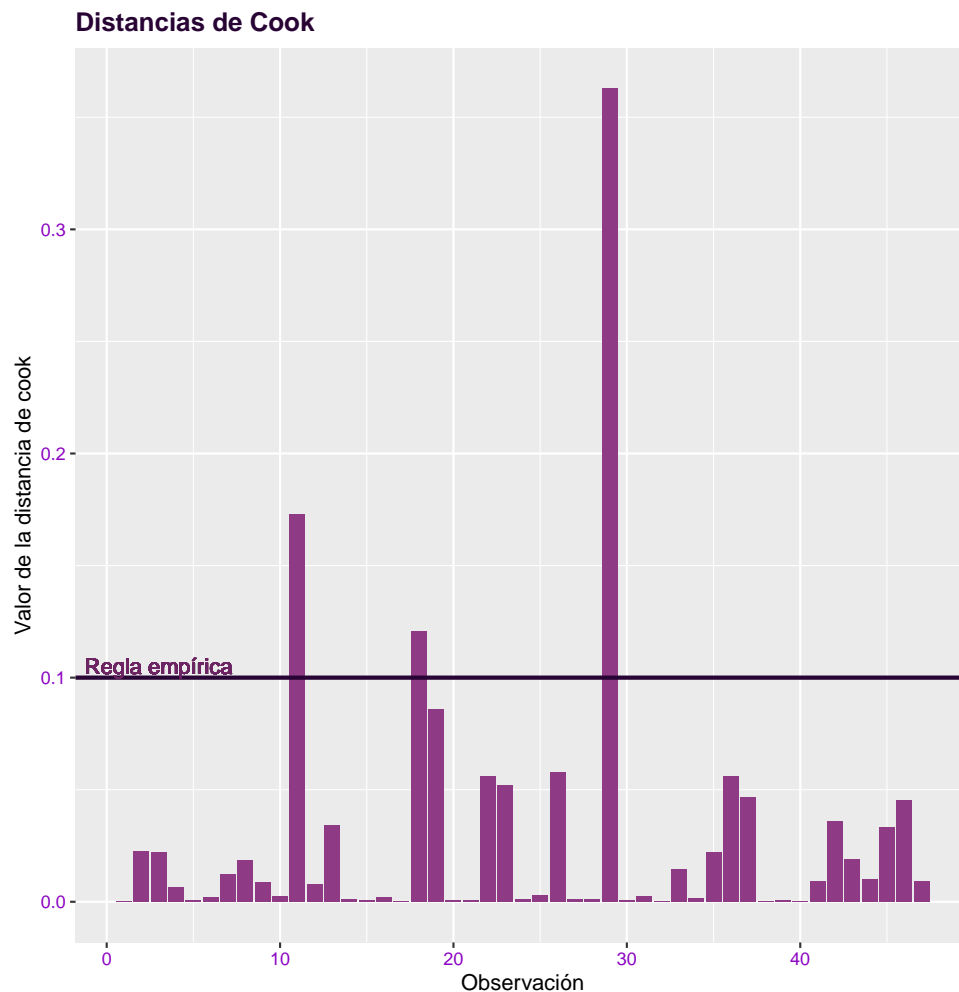
4.3. Búsqueda de observaciones influyentes

En esta sección se buscará estudiar cuales de las observaciones presentan valores influyentes, para ello se hará uso de la Distancia de Cook.

Esta es una medida del nivel de influencia de la observación i -ésimas sobre la estimación de $\hat{\beta}$, es decir se busca medir si su presencia o ausencia en el modelo hace que el mismo cambie.

Una distancia de Cook elevada significa que una observación tiene mayor influencia al momento de determinar los $\hat{\beta}$.

$$D_i = \frac{(\hat{\beta} - \hat{\beta}(-i))' X' X (\hat{\beta} - \hat{\beta}(-i))}{(k + 1) \hat{\sigma}^2}$$



Tomando como regla empírica el valor de $\frac{4}{n-k-1} = 4/40$ puede verse que las observaciones **11**, **29** (de manera excesiva) y **17** sobrepasan la regla estipulada.

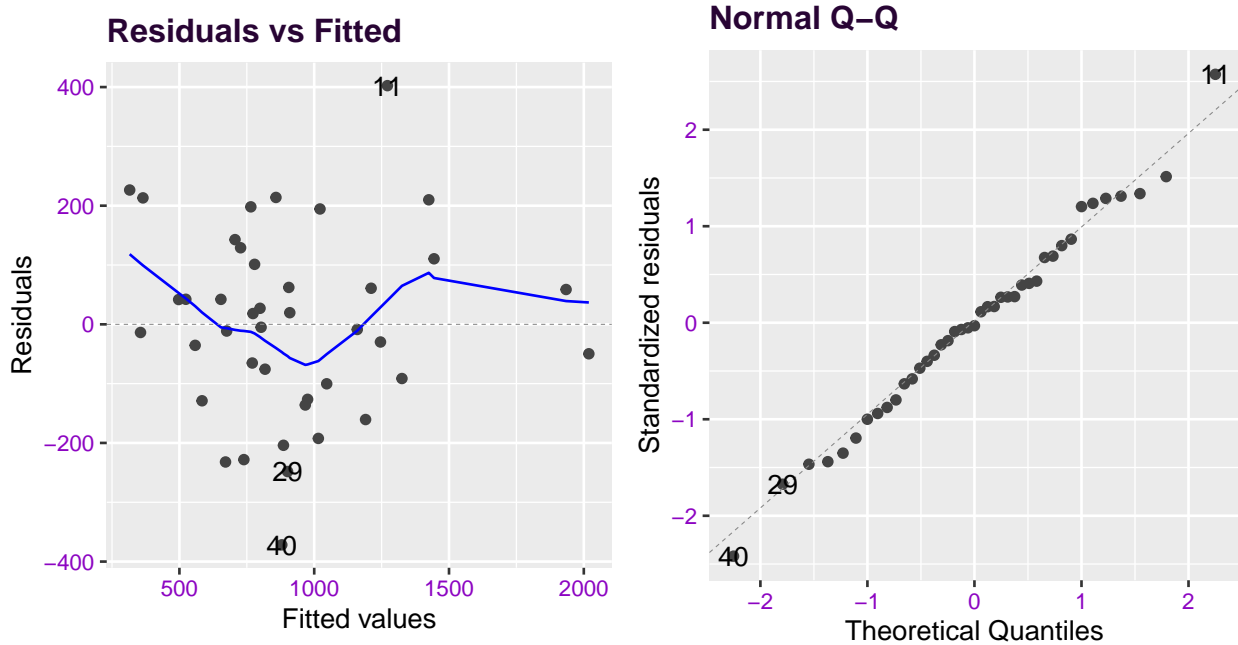
Por ende, se decide retirarlas del modelo reducido ya que las mismas tienen una influencia preponderante en la estimación. Como se vió en clase, observaciones de este tipo pueden llevar a un modelo alejado de la realidad.

4.4. Heterosedasticidad

4.4.1. Test de Breusch-Pagan

```
##
## studentized Breusch-Pagan test
##
## data: .
## BP = 8.9946, df = 5, p-value = 0.1093
```

4.5. Normalidad



Considerando los gráficos anteriormente mencionados podemos ver que la esperanza de los errores se mantiene cercana a cero, a diferencia de algunas observaciones, a su vez, no encontramos un patrón en la dispersión de los errores, en cambio en la gráfica QQ-plot podemos ver que en valores centrales de la distribución se asemeja a una distribución normal a excepción de la observación 40, 11 y 6.

```
## Error: Problem with 'summarise()' column 'Shapiro'.
## i 'Shapiro = shapiro.test(as.numeric())'.
## x 'Shapiro' must be a vector, not a 'htest' object.
## Error in stopifnot(is.numeric(x)): el argumento "x" está ausente, sin valor por omisión
```

En base a los test puedo afirmar que la distribución de los errores es normal.
Me tranque en como presentar una tabla , lo miro luego

5. Anexo



5.1. Selección de Modelos : StepWise

```
## Stepwise regression (forward-backward), alpha-to-enter: 0.15, alpha-to-remove: 0.15
##
## Full model: y ~ M + So + Ed + Po1 + Po2 + LF + M.F + Pop + NW + U1 + U2 +
## GDP + Ineq + Prob + Time
```

```
## <environment: 0x000000001b27a2c0>
##
## --- Step (forward) 1 ---
## Single term additions
##
## Model:
## y ~ 1
##      Df Sum of Sq      RSS      AIC F value      Pr(>F)
## <none>                6426043 492.45
## M      1      41203 6384840 494.19   0.2517   0.618719
## So      1      85596 6340447 493.90   0.5265   0.472414
## Ed      1      647249 5778794 490.10   4.3682   0.043186 *
## Po1     1     3774458 2651585 458.16  55.5154  5.194e-09 ***
## Po2     1     3679620 2746423 459.60  52.2517  1.043e-08 ***
## LF      1      220081 6205962 493.03   1.3830   0.246712
## M.F     1      625471 5800572 490.26   4.2053   0.047059 *
## Pop     1      908112 5517930 488.21   6.4184   0.015426 *
## NW      1         45 6425998 494.45   0.0003   0.986939
## U1      1       3633 6422410 494.43   0.0221   0.882693
## U2      1      219528 6206514 493.03   1.3795   0.247316
## GDP     1     1323717 5102326 485.00  10.1179   0.002878 **
## Ineq    1      238604 6187439 492.90   1.5039   0.227419
## Prob    1     1373492 5052551 484.60  10.6018   0.002341 **
## Time    1      112461 6313582 493.73   0.6947   0.409650
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## --- Step (forward) 2 ---
## Single term additions
##
## Model:
## y ~ Po1
##      Df Sum of Sq      RSS      AIC F value      Pr(>F)
## <none>                2651585 458.16
## M      1      497705 2153880 451.64   8.7808  0.0052297 **
## So      1      173592 2477993 457.38   2.6620  0.1110309
## Ed      1       27826 2623758 459.73   0.4030  0.5293431
## Po2     1       10127 2641458 460.00   0.1457  0.7048187
## LF      1         781 2650804 460.15   0.0112  0.9163072
## M.F     1       50147 2601437 459.38   0.7325  0.3974328
## Pop     1       56487 2595098 459.28   0.8271  0.3688357
## NW      1      198718 2452866 456.97   3.0786  0.0873919 .
## U1      1       1668 2649917 460.14   0.0239  0.8779124
## U2      1       34224 2617361 459.63   0.4969  0.4851746
## GDP     1      360105 2291479 454.18   5.9717  0.0192918 *
## Ineq    1      734344 1917240 446.87  14.5548  0.0004868 ***
## Prob    1     139560 2512025 457.94   2.1111  0.1544380
## Time    1       90065 2561519 458.74   1.3361  0.2549363
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
##
## --- Step (forward) 3 ---
## Single term additions
##
## Model:
## y ~ Po1 + Ineq
##      Df Sum of Sq      RSS      AIC F value    Pr(>F)
## <none>                1917240 446.87
## M      1      116540 1800700 446.29   2.3946 0.130266
## So     1       57785 1859456 447.61   1.1498 0.290531
## Ed     1      356963 1560277 440.42   8.4649 0.006094 **
## Po2    1       14372 1902868 448.56   0.2795 0.600213
## LF     1       66559 1850681 447.42   1.3307 0.256079
## M.F    1      156525 1760715 445.37   3.2893 0.077852 .
## Pop    1        2289 1914952 448.82   0.0442 0.834599
## NW     1       53757 1863483 447.70   1.0674 0.308248
## U1     1        1229 1916012 448.84   0.0237 0.878421
## U2     1        1278 1915963 448.84   0.0247 0.876037
## GDP    1       29366 1887875 448.23   0.5755 0.452875
## Prob   1      334631 1582610 441.00   7.8234 0.008136 **
## Time   1       19419 1897822 448.45   0.3786 0.542128
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## --- Step (forward) 4 ---
## Single term additions
##
## Model:
## y ~ Po1 + Ineq + Ed
##      Df Sum of Sq      RSS      AIC F value    Pr(>F)
## <none>                1560277 440.42
## M      1      173876 1386402 437.57   4.5149 0.040534 *
## So     1         20 1560257 442.42   0.0005 0.982934
## Po2    1       7786 1552491 442.21   0.1806 0.673427
## LF     1      12857 1547420 442.08   0.2991 0.587817
## M.F    1       1162 1559115 442.39   0.0268 0.870784
## Pop    1       1598 1558679 442.38   0.0369 0.848738
## NW     1        107 1560170 442.42   0.0025 0.960585
## U1     1         25 1560252 442.42   0.0006 0.981025
## U2     1      24616 1535661 441.77   0.5771 0.452414
## GDP    1       1951 1558326 442.37   0.0451 0.833070
## Prob   1      275467 1284810 434.45   7.7185 0.008628 **
## Time   1     101432 1458845 439.66   2.5030 0.122374
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## --- Step (forward) 5 ---
## Single term additions
```



```

##
## Model:
## y ~ Po1 + Ineq + Ed + Prob
##      Df Sum of Sq      RSS      AIC F value  Pr(>F)
## <none>                1284810 434.45
## M      1      194411 1090399 429.73   6.2403 0.01734 *
## So     1       42551 1242259 435.07   1.1989 0.28103
## Po2    1       17200 1267610 435.90   0.4749 0.49528
## LF     1       44995 1239815 434.99   1.2702 0.26739
## M.F    1        8010 1276800 436.20   0.2196 0.64227
## Pop    1       13983 1270827 436.01   0.3851 0.53891
## NW     1       21254 1263555 435.77   0.5887 0.44805
## U1     1         421 1284389 436.44   0.0115 0.91535
## U2     1       21109 1263701 435.78   0.5846 0.44962
## GDP    1       11124 1273685 436.10   0.3057 0.58385
## Time   1          0 1284810 436.45   0.0000 0.99777
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## == Step (forward) 6 ==
## Single term additions
##
## Model:
## y ~ Po1 + Ineq + Ed + Prob + M
##      Df Sum of Sq      RSS      AIC F value  Pr(>F)
## <none>                1090399 429.73
## So     1       14101 1076298 431.19   0.4455 0.50900
## Po2    1       16562 1073837 431.10   0.5244 0.47393
## LF     1       58731 1031668 429.46   1.9356 0.17319
## M.F    1         44 1090355 431.73   0.0014 0.97059
## Pop    1       3092 1087307 431.61   0.0967 0.75774
## NW     1        849 1089550 431.70   0.0265 0.87167
## U1     1      15301 1075098 431.15   0.4839 0.49140
## U2     1       92967  997432 428.07   3.1690 0.08399 .
## GDP    1          2 1090397 431.73   0.0000 0.99441
## Time   1       5826 1084572 431.51   0.1827 0.67180
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## == Step (forward) 7 ==
## Single term additions
##
## Model:
## y ~ Po1 + Ineq + Ed + Prob + M + U2
##      Df Sum of Sq      RSS      AIC F value  Pr(>F)
## <none>                997432 428.07
## So     1       9129 988302 429.70   0.3048 0.5846
## Po2    1      19362 978069 429.27   0.6533 0.4247
## LF     1      13984 983448 429.50   0.4692 0.4981

```

```
## M.F      1      7210 990222 429.78  0.2403 0.6273
## Pop      1      3226 994206 429.94  0.1071 0.7456
## NW       1       390 997042 430.06  0.0129 0.9103
## U1       1     67397 930034 427.21  2.3914 0.1315
## GDP      1       113 997319 430.07  0.0037 0.9516
## Time     1      6142 991290 429.82  0.2045 0.6541
##
##
## --- Step (forward) 8 ---
## Single term additions
##
## Model:
## y ~ Po1 + Ineq + Ed + Prob + M + U2 + U1
##      Df Sum of Sq    RSS    AIC F value Pr(>F)
## <none>                930034 427.21
## So      1      955.2 929079 429.16  0.0329 0.8572
## Po2     1    12970.8 917064 428.63  0.4526 0.5059
## LF      1    16534.4 913500 428.47  0.5792 0.4522
## M.F     1      293.5 929741 429.19  0.0101 0.9206
## Pop     1     8236.2 921798 428.84  0.2859 0.5965
## NW      1      494.3 929540 429.18  0.0170 0.8970
## GDP     1     1277.4 928757 429.15  0.0440 0.8352
## Time    1    20994.2 909040 428.27  0.7390 0.3964
##
## Call:
## lm(formula = y ~ Po1 + Ineq + Ed + Prob + M + U2 + U1, data = Datos)
##
## Coefficients:
## (Intercept)      Po1      Ineq      Ed      Prob      M
## -4614.686    12.366     6.266    19.596   -39.840     9.284
##      U2      U1
##    16.785   -4.745
```

5.2. Script de R

6. Bibliografía