

Capítulo 5

Pruebas de Hipótesis para el Modelo Lineal

En este capítulo se estudiará cómo pueden usarse las propiedades de los estimadores de mínimos cuadrados con el fin de realizar inferencias sobre el modelo ajustado.

Todas las pruebas que aquí se mencionarán se basan en las siguientes suposiciones:

- La dependencia entre las variables es lineal. Es decir, el tipo de modelo elegido es correcto, aún cuando las variables incluidas no sean las correctas.
- El vector de errores ε sigue una distribución Normal con media $\mathbf{0}$ y matriz de varianza $\sigma^2 I$. Es decir, los errores son normales, independientes, con media cero y varianza común σ^2 .

5.1 Pruebas sobre los parámetros del modelo.

Nuestro objetivo en esta sección será probar la hipótesis

$$H_0 : \beta_i = 0$$

contra la alternativa

$$H_1 : \beta_i \neq 0$$

Esto equivale a comparar el modelo

$$y_l = \beta_0 + \beta_1 x_{1l} + \dots + \beta_{i-1} x_{i-1,l} + \beta_{i+1} x_{i+1,l} + \dots + \beta_k x_{kl} + \varepsilon_l$$

contra el modelo

$$y_l = \beta_0 + \beta_1 x_{1l} + \dots + \beta_{i-1} x_{i-1,l} + \beta_i x_{i,l} + \beta_{i+1} x_{i+1,l} + \dots + \beta_k x_{kl} + \varepsilon_l$$

es decir, el modelo en el cual no aparece x_i contra el modelo en el cual sí está presente.

Con el fin de construir una prueba para estas hipótesis, necesitaremos el siguiente resultado:

En el capítulo anterior se demostró que $\hat{\beta}_i \sim N(\beta_i, \sigma^2 c_{ii})$ y $\frac{(n-p)S^2}{\sigma^2} \sim \chi_{n-p}^2$. A continuación, probaremos que $\hat{\beta}$ y \mathbf{e} son independientes, lo cual nos garantiza que $\hat{\beta}$ y $S^2 = \frac{1}{n-p} \mathbf{e}'\mathbf{e}$ son independientes.

Recordemos que

$$\begin{aligned} \hat{\beta} &= (X'X)^{-1}X'\mathbf{Y} \\ &= C\mathbf{Y} \end{aligned}$$

y que

$$\begin{aligned} \mathbf{e} &= \mathbf{Y} - X\hat{\beta} \\ &= \mathbf{Y} - (X'X)^{-1}X'\mathbf{Y} \\ &= (I - V)\mathbf{Y}, \end{aligned}$$

donde $V = X(X'X)^{-1}X'$.

Por el Teorema 3.2, $\hat{\beta}$ y \mathbf{e} serán independientes si y sólo si $C(I - V)' = 0$. Veamos, entonces, si esta última igualdad es cierta.

$$\begin{aligned} C(I - V)' &= C(I - V) \\ &= (X'X)^{-1}X'(I - X(X'X)^{-1}X') \\ &= (X'X)^{-1}X' - (X'X)^{-1}X'X(X'X)^{-1}X' \\ &= (X'X)^{-1}X' - (X'X)^{-1}X' \\ &= 0 \end{aligned}$$

Por lo tanto, $\hat{\beta}$ y \mathbf{e} son independientes. Este hecho, como se indicó anteriormente, permite afirmar que $\hat{\beta}$ y S^2 son independientes y que

$$\frac{\frac{\hat{\beta}_i - \beta_i}{\sigma\sqrt{c_{ii}}}}{\sqrt{\frac{(n-p)S^2}{\sigma^2}} \frac{1}{(n-p)}} = \frac{\hat{\beta}_i - \beta_i}{S\sqrt{c_{ii}}} \sim t_{n-p} \quad (5.1)$$

Regresemos entonces al problema de probar

$$H_0 : \beta_i = 0$$

$$H_1 : \beta_i \neq 0$$

Parece razonable rechazar H_0 si $\hat{\beta}_i$ está “muy lejos” de cero; este tipo de razonamiento plantea, nuevamente, tenemos el problema de determinar cuándo consideraremos que $\hat{\beta}_i$ está “muy lejos” de cero. Para ello, observemos que si H_0 es cierta, por (5.1)

$$t = \frac{\hat{\beta}_i}{S\sqrt{c_{ii}}} \sim t_{n-p},$$

de modo que podemos usar una t_{n-p} como distribución de referencia, y rechazar H_0 (con un nivel α para la prueba) cuando

$$|t| > t_{n-p, \alpha/2},$$

es decir, cuando t cae en la región indicada en la Figura 5.1.

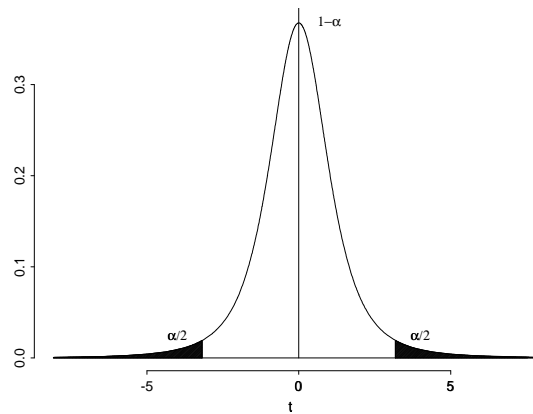


Figura 5.1: Región de rechazo para la prueba t .

Veamos cómo puede usarse esta prueba en el ejemplo 4.4. Comencemos probando la hipótesis de que la distancia recorrida no influye en el tiempo de entrega. Esto equivale a plantear:

$$H_0 : \beta_2 = 0$$

$$H_1 : \beta_2 \neq 0$$

Usando los resultados obtenidos en el capítulo anterior, tenemos

$$t = \frac{\hat{\beta}_2}{\sqrt{S^2 c_{22}}} = \frac{0.456}{\sqrt{9.86 * 0.00218}} = 3.11$$

Si decidimos usar $\alpha = 0.05$, cuando buscamos el valor tabulado de t para un área en la cola de $\alpha/2 = 0.025$ correspondiente a una t con 12 grados de libertad, obtenemos $t_{12,0.025} = 2.18$. Como el valor calculado de t es mayor que el tabulado, rechazamos H_0 . Es decir, podemos afirmar que la distancia recorrida es importante para predecir el tiempo de entrega de la cerveza.

Si queremos hacer la misma prueba para β_1 (es decir, queremos saber si el número de cajas a repartir es importante para determinar el tiempo de entrega), formulamos las hipótesis:

$$H_0 : \beta_1 = 0$$

$$H_1 : \beta_1 \neq 0$$

y obtenemos

$$t = \frac{\hat{\beta}_1}{\sqrt{S^2 c_{11}}} = 5.74$$

Comparando nuevamente contra $t_{12,0.025} = 2.18$, se rechaza H_0 , y podemos afirmar que el modelo $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$ explica mejor la variación de los datos que el modelo $y_i = \beta_0 + \beta_2 x_{2i} + \varepsilon_i$.

Si no es usado con cuidado, este procedimiento puede llevar a conclusiones erróneas, porque los $\hat{\beta}_i$ no son independientes. Este hecho se evidencia en el siguiente ejemplo.

Ejemplo 5.1 *La demanda de un producto es afectada por múltiples factores. En un estudio se tomaron medidas en 9 regiones geográficas sobre urbanización relativa, nivel educativo e ingreso relativo con el fin de determinar su influencia sobre el uso del producto. Los datos recolectados se muestran en la tabla 5.1*

Nivel de Urbanización	Nivel Educativo	Ingreso Relativo	Consumo
X_1	X_2	X_3	Y
42.2	11.2	31.9	167.1
48.6	10.6	13.2	174.4
42.6	10.6	28.7	160.8
39.0	10.4	26.1	162.0
34.7	9.3	30.1	140.8
44.5	10.8	8.5	174.6
39.1	10.7	24.3	163.7
40.1	10.0	18.6	174.5
45.9	12.0	20.4	185.7

Tabla 5.1: Datos sobre el consumo de cierto producto en 9 diferentes regiones (Ejemplo 5.1)

Comencemos ajustando un modelo que contenga las tres variables explicativas, es decir:

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \varepsilon_i$$

Al hacer los cálculos usando un paquete estadístico obtenemos la siguiente tabla:

Coefficiente	Valor	Error Standard	Valor t
	$(\hat{\beta}_i)$	$(S\sqrt{c_{ii}})$	$(\hat{\beta}_i/(S\sqrt{c_{ii}}))$
β_0	60.0143		
β_1	0.2398	1.0121	0.2370
β_2	10.7184	4.5296	0.3663
β_3	-0.7510	0.3950	-1.9014

Comparando los valores de t de la tabla anterior con el valor tabulado $t_{5,0.025} = 2.57$, pareciera no existir evidencia para rechazar $H_0 : \beta_i = 0$ para ningún valor de i .

Antes de tomar la decisión de eliminar todas las variables del modelo, recordemos que la prueba t se basa en los valores de $\hat{\beta}_i$. Tomando en cuenta que $\text{Var}(\hat{\beta}) = \sigma^2(X'X)^{-1}$, sabemos que los $\hat{\beta}_i$, en general, no son independientes, de manera que si eliminamos uno de ellos del modelo, pueden presentarse cambios importantes en los otros que modifiquen la inferencia. En nuestro caso,

$$(X'X)^{-1} = \begin{pmatrix} 33.099 & -0.325 & -1.452 & -0.174 \\ -0.325 & 0.026 & -0.086 & 0.007 \\ -1.452 & -0.086 & 0.519 & -0.020 \\ -0.174 & 0.007 & -0.020 & 0.004 \end{pmatrix}$$

Podemos ver que los elementos que están fuera de la diagonal no son cero, y por tanto existen correlaciones entre los $\hat{\beta}_i$.

Veamos, entonces, qué sucede con el modelo cuando eliminamos el nivel de urbanización X_1 , la variable que corresponde al β con menor valor calculado de t . Es decir, el nuevo modelo que ajustamos es

$$y_i = \beta_0 + \beta_2 x_{2i} + \beta_3 x_{3i} + \varepsilon_i$$

Para este modelo, el paquete estadístico arroja los siguientes resultados:

Coefficiente	Valor	Error Standard	Valor t
β_0	63.021		
β_2	11.517	2.777	4.1469
β_3	-0.8158	0.2614	-3.1206

Al comparar con el valor tabulado $t_{6,0.025} = 2.45$, vemos que existe evidencia suficiente para rechazar $H_0 : \beta_i = 0$ para cualquiera de las dos variables independientes. Es decir, hubiésemos cometido un grave error si hubiésemos eliminado las tres variables del modelo en base a los resultados obtenidos para la prueba t en el primer modelo ajustado.

Este ejemplo nos demuestra que, si no es usada con cuidado, la prueba t puede llevar a resultados erróneos. En general, no es recomendable eliminar más de una variable a la vez cuando aplicamos este procedimiento, pues sólo nos permite comparar modelos que difieren en una variable. En la próxima sección trataremos el problema de comparar modelos que difieren en más de una variable explicativa.

5.2 Comparación de modelos

Consideremos el problema de comparar dos modelos de la forma:

- (1) $y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_l x_{li} + \varepsilon_i$
- (2) $y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_l x_{li} + \beta_{l+1} x_{(l+1)i} + \dots + \beta_k x_{ki} + \varepsilon_i$

Obsérvese que todas las variables explicativas del modelo (1) están contenidas en el modelo (2). En este caso, decimos que el modelo (1) está *anidado* en el modelo (2).

Comparar estos modelos equivale a contrastar las hipótesis

$$\begin{aligned} H_0 : & \quad \beta_{l+1} = \beta_{l+2} = \dots = \beta_k = 0 \quad \text{vs.} \\ H_1 : & \quad \beta_j \neq 0, \text{ algún } j = l + 1, \dots, k \end{aligned}$$

Para analizar la geometría de este problema, escribamos ambos modelos en forma matricial como

$$\begin{aligned} (1) \quad \mathbf{Y} &= X_1 \boldsymbol{\beta}_1 + \boldsymbol{\varepsilon} \\ (2) \quad \mathbf{Y} &= X \boldsymbol{\beta} + \boldsymbol{\varepsilon} \end{aligned}$$

y observemos que todas las columnas de X_1 están contenidas en X , de manera que el espacio generado por las columnas de X_1 , $\text{gen}\{X_1\}$, es subespacio del generado por las columnas de X , $\text{gen}\{X\}$.

Sean $\hat{\boldsymbol{\beta}}_1$ y $\hat{\boldsymbol{\beta}}_2$ los estimadores de mínimos cuadrados para $\boldsymbol{\beta}_1$ y $\boldsymbol{\beta}_2$, respectivamente, y definamos

$$\begin{aligned} \mathbf{e}_1 &= \mathbf{Y} - \hat{\mathbf{Y}}_1 = \mathbf{Y} - X_1 \hat{\boldsymbol{\beta}}_1 \\ \mathbf{e}_2 &= \mathbf{Y} - \hat{\mathbf{Y}}_2 = \mathbf{Y} - X \hat{\boldsymbol{\beta}}_2. \end{aligned}$$

Entonces $\mathbf{Y} = X_1 \hat{\boldsymbol{\beta}}_1 + (\mathbf{e}_1 - \mathbf{e}_2) + \mathbf{e}_2$ (ver Figura 5.2), donde estos tres vectores son perpendiculares entre sí. $X_1 \hat{\boldsymbol{\beta}}_1$ está en $\text{gen}\{X_1\}$, cuya dimensión es $l + 1$; $\mathbf{e}_1 - \mathbf{e}_2$ está en el subespacio ortogonal a $\text{gen}\{X_1\}$ en $\text{gen}\{X\}$, que tiene dimensión $k - l$; y finalmente $\mathbf{e}_2 \in \text{gen}\{X\}^\perp$, $\dim(\text{gen}\{X\}) = n - (k + 1)$ (obsérvese que la suma de las dimensiones de los tres subespacios es n).

Supongamos que H_0 es cierta, es decir, que el modelo correcto es el (1). Entonces la diferencia entre \mathbf{e}_1 y \mathbf{e}_2 tiene que haberse producido por causa del azar, así que $\frac{\|\mathbf{e}_1 - \mathbf{e}_2\|^2}{(k-l)}$ debería ser un estimador de la varianza, al igual que $\frac{\|\mathbf{e}_2\|^2}{n-(k+1)} = S^2$, por lo cual ambas cantidades no deberían estar muy lejos. Parece entonces razonable basar una prueba para esta hipótesis en la comparación entre las dos cantidades anteriores, y para ello necesitamos una distribución de referencia.

Esta distribución de referencia puede obtenerse de la siguiente manera: si el modelo correcto es el modelo (1), entonces $E(\mathbf{Y}) = X_1 \boldsymbol{\beta}_1$, y por lo tanto $\frac{1}{\sigma}(\mathbf{Y} - X_1 \boldsymbol{\beta}_1) \sim N(0, I)$.

Podemos escribir

$$\frac{1}{\sigma}(\mathbf{Y} - X_1 \boldsymbol{\beta}_1) = \frac{1}{\sigma}(\hat{\mathbf{Y}}_1 - X_1 \boldsymbol{\beta}_1) + \frac{1}{\sigma}(\mathbf{e}_1 - \mathbf{e}_2) + \frac{1}{\sigma} \mathbf{e}_2$$

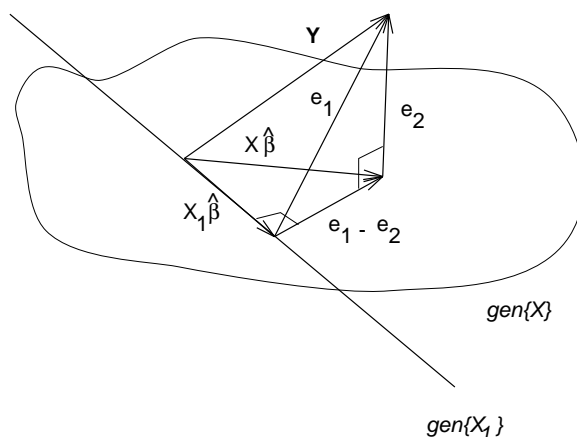


Figura 5.2: Geometría del problema de comparación de modelos lineales.

donde $\frac{1}{\sigma}(\hat{\mathbf{Y}}_1 - X_1\boldsymbol{\beta}_1) \in \text{gen}\{X_1\}$. $\frac{1}{\sigma}(\mathbf{e}_1 - \mathbf{e}_2)$ está en el complemento ortogonal de $\text{gen}\{X_1\}$ en $\text{gen}\{X\}$ y $\frac{1}{\sigma}\mathbf{e}_2$ pertenece al complemento ortogonal de $\text{gen}\{X\}$ en \mathcal{R}^n . Es decir, $\frac{1}{\sigma}(\mathbf{Y} - X_1\boldsymbol{\beta}_1)$ es un vector $N(0, I)$ que puede escribirse como la suma de sus proyecciones en subespacios ortogonales. Esto implica que las normas al cuadrado de dichas proyecciones son independientes, y tienen distribución χ^2 con el número de grados de libertad igual a la dimensión del espacio en el cual se encuentran:

$$\frac{1}{\sigma^2}\|\mathbf{e}_1 - \mathbf{e}_2\|^2 \sim \chi_{k-l}^2$$

$$\frac{1}{\sigma^2}\|\mathbf{e}_2\|^2 \sim \chi_{n-(k+1)}^2$$

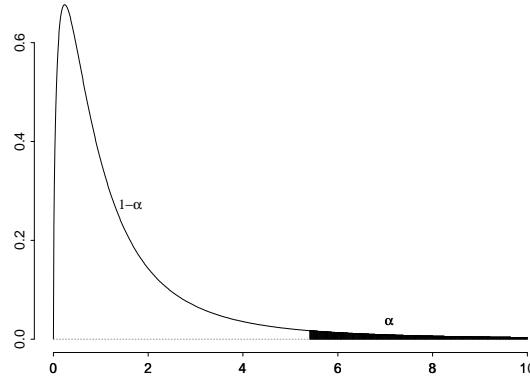
Entonces

$$F = \frac{\frac{\|\mathbf{e}_1 - \mathbf{e}_2\|^2}{k-l}}{\frac{\|\mathbf{e}_2\|^2}{n-(k+1)}} = \frac{\frac{\|\mathbf{e}_1\|^2 - \|\mathbf{e}_2\|^2}{k-l}}{\frac{\|\mathbf{e}_2\|^2}{n-(k+1)}} \sim F_{k-l, n-(k+1)} \quad (5.2)$$

Rechazamos H_0 si $F > F_{k-l, n-(k+1)}^\alpha$, es decir, cuando F cae en la región indicada en la Figura 5.3.

Usando la notación del capítulo anterior, F puede ser escrita como

$$F = \frac{\frac{SSE_1 - SSE_2}{k-l}}{\frac{SSE_2}{n-(k+1)}}$$

Figura 5.3: Región de rechazo para la prueba F .

Usemos esta prueba en el ejemplo 4.4 para comparar los modelos

$$(1) y_i = \beta_0 + \varepsilon_i$$

$$(2) y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$$

A este tipo de prueba, en la cual se compara un modelo contra el modelo más sencillo posible (constante mas error) se le suele llamar *Prueba de Significancia de la Regresión*.

En terminos de hipótesis, comparar ambos modelos equivale a plantear las hipótesis:

$$H_0 : \beta_1 = \beta_2 = 0 \text{ vs } H_1 : \text{algún } \beta_i \neq 0, i = 1, 2$$

Para el modelo (1), las ecuaciones normales toman la forma

$$X'X = \mathbf{1}'\mathbf{1} = n = X'\mathbf{Y} = \sum_{i=1}^n y_i.$$

donde $\mathbf{1}$ representa un vector de n unos.

Por lo tanto

$$\hat{\beta}_0 = (X'X)^{-1}X'\mathbf{Y} = \frac{\sum_{i=1}^n y_i}{n} = \bar{y}$$

y

$$SSE_1 = \mathbf{Y}'\mathbf{Y} - \hat{\beta}'X'\mathbf{Y} = \sum_{i=1}^n y_i^2 - \bar{y} \sum_{i=1}^n y_i = \sum_{i=1}^n y_i^2 - n\bar{y}^2.$$

Observemos que estas expresiones son válidas para cualquier conjunto de datos.

Para los datos del Ejemplo 4.4, se obtiene $\hat{\beta}_0 = 30.87$ y $SSE_1 = 449.73$.

Por otra parte, sabíamos ya que $SSE_2 = 118.37$. Al estar trabajando con 15 datos, calculamos F como

$$F = \frac{\frac{449.73-118.37}{2-1}}{\frac{118.37}{15-(2+1)}} = \frac{\frac{331.36}{2}}{\frac{118.37}{12}} = 16.8$$

Si fijamos un nivel $\alpha = 0.05$ y comparamos el valor de F que acabamos de calcular con $F_{2,12}^{0.05} = 3.89$, observamos que el valor calculado es mayor que el valor tabulado. Según el procedimiento de prueba desarrollado anteriormente, rechazamos H_0 ; por tanto, decidimos que el modelo (2) es mejor que el modelo (1).

En general, si deseamos hacer una prueba de significancia para un modelo cualquiera, o equivalentemente, comparar los modelos:

$$\begin{aligned} (1) \quad \mathbf{Y} &= \beta_0 \mathbf{1} + \boldsymbol{\varepsilon} \\ (2) \quad \mathbf{Y} &= X\boldsymbol{\beta} + \boldsymbol{\varepsilon} \end{aligned}$$

la prueba F para esta comparación toma la forma

$$\begin{aligned} F &= \frac{\frac{(\mathbf{Y}'\mathbf{Y} - n\bar{y}^2) - (\mathbf{Y}'\mathbf{Y} - \hat{\boldsymbol{\beta}}'X'\mathbf{Y})}{k}}{\frac{\mathbf{Y}'\mathbf{Y} - \hat{\boldsymbol{\beta}}'X'\mathbf{Y}}{n-(k+1)}}} \\ &= \frac{\frac{\hat{\boldsymbol{\beta}}'X'\mathbf{Y} - n\bar{y}^2}{k}}{\frac{\mathbf{Y}'\mathbf{Y} - \hat{\boldsymbol{\beta}}'X'\mathbf{Y}}{n-(k+1)}}} \end{aligned}$$

Si denotamos $SSR = \hat{\boldsymbol{\beta}}'X'\mathbf{Y} - n\bar{y}^2$ (suma de cuadrados de la regresión), podemos escribir

$$F = \frac{\frac{SSR}{k}}{\frac{SSE}{n-(k+1)}} = \frac{MSR}{MSE}$$

Hablando informalmente, SSR representa la variación explicada por las variables independientes y que no explica \bar{y} . Podemos ver, entonces, la prueba de significancia de la regresión como la comparación entre la variación explicada por las variables independientes

(SSR) y la variación debida al ruido (SSE). Si la variación explicada por las variables es sustancialmente mayor que aquella debida al ruido, valdrá la pena complicar el modelo más allá del modelo constante mas error.

La información necesaria para la prueba de significancia de la regresión se suele resumir en una tabla a la cual se le da el nombre de *Tabla de Análisis de Varianza* o *Tabla ANOVA* (ANalysis Of VAriance), de la siguiente forma:

Fuente	Grados de Libertad	SS	MS	F
Modelo	k	$SSR = \hat{\beta}' X' Y - n\bar{y}^2$	$MSR = \frac{SSR}{k}$	$\frac{MSR}{MSE}$
Error	$n - (k + 1)$	$SSE = Y' Y - \hat{\beta}' X' Y$	$MSE = \frac{SSE}{n - (k + 1)}$	
TOTAL	$n - 1$	$SST = Y' Y - n\bar{y}^2$		

En el caso del ejemplo 4.4 (para el cual ya hicimos la prueba de significancia) la tabla ANOVA toma la siguiente forma

Fuente	Grados de Libertad	SS	MS	F
Modelo	2	331.36	165.68	16.8
Error	12	118.37	9.86	
TOTAL	14	449.37		

Si deseamos comparar dos modelos anidados, la información necesaria se encuentra en las tablas ANOVA de cada modelo. Para respaldar esta afirmación, supongamos que en el Ejemplo 5.1 deseamos comparar el modelo según el cual el consumo (Y) depende sólo del nivel educativo (X_2) contra el modelo que explica el consumo en términos del nivel de urbanización (X_1), el nivel educativo (X_2) y el ingreso relativo (X_3).

$$(1) y_i = \beta_0 + \beta_2 x_{2i} + \varepsilon_i$$

$$(2) y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \varepsilon_i$$

Usando un paquete estadístico, obtenemos las tablas ANOVA para ambos modelos. Para el modelo (1):

Fuente	Grados de Libertad	SS	MS	F
Modelo	1	754.41	754.41	10.06
Error	7	524.8	74.97	
TOTAL	8	1279.21		

Para el modelo (2):

Fuente	Grados de Libertad	SS	MS	F
Modelo	3	1081.35	360.45	9.11
Error	5	197.85	39.57	
TOTAL	8	1279.21		

Comparar los modelos anteriores corresponde a probar las hipótesis: $H_0 : \beta_1 = \beta_3 = 0$ vs $H_1 : \beta_1 \neq 0$ ó $\beta_3 \neq 0$. Usando una prueba F para esta comparación, obtenemos:

$$F = \frac{\frac{SSE_1 - SSE_2}{k-l}}{\frac{SSE_2}{n-(k+1)}} = \frac{524.8 - 197.85}{7-5} = 4.13$$

Fijando un nivel $\alpha = 0.05$ para la prueba, cuando comparamos el valor anterior con $F_{2,5}^{0.05} = 5.79$, concluimos que los datos no proporcionan suficiente evidencia para rechazar H_0 , es decir, el nivel de urbanización y el ingreso relativo no parecen ser importantes para predecir el consumo cuando comparamos con el modelo que sólo contiene el nivel educativo.

Observemos que este resultado no coincide con el obtenido usando la prueba t; en ese caso, después de eliminar el nivel de urbanización del modelo, concluíamos que tanto el nivel educativo como el ingreso relativo son significativos. Este tipo de contradicciones es muy común cuando aplicamos procedimientos estadísticos, y en general nos indican la necesidad de usar otros elementos de decisión. En última instancia, es necesario destacar el importantísimo papel que juega el sentido común a la hora de realizar un análisis estadístico.