

UNIVERSIDAD NACIONAL DE ROSARIO
Facultad de Ciencias Económicas y Estadística



Análisis de datos de duración en pacientes con cáncer de mama

Rotterdam tumor bank - 1978-1985

Alumnas: Agustina Mac Kay y Rocio Canteros

Año 2024

Introducción

El cáncer de mama es un tipo de cáncer primario, que se origina en la mama y puede propagarse a otros tejidos u órganos del cuerpo. Es el tipo de cáncer más frecuente y la causa más común de muerte por cáncer en mujeres a nivel mundial.¹

En este estudio se trabajará con información acerca de 583 mujeres que fueron sometidas, entre 1978 y 1985, a una cirugía primaria para extirpar el tumor.

Los datos fueron obtenidos de la base *rotterdam* del paquete *survival* de R. La misma cuenta con el tiempo desde la cirugía hasta la muerte o pérdida de seguimiento de las pacientes, junto a otras covariables basales que se detallan a continuación:

- **Age:** edad al momento de la cirugía (en años).
- **Meno:** estado menopáusico, donde 0 = premenopáusico y 1 = postmenopáusico.
- **Hormon:** variable indicadora de haber recibido un tratamiento hormonal.
- **Chemo:** variable indicadora de haber recibido quimioterapia.
- **Pgr:** receptores de progesterona (en fmol/l).
- **Er:** receptores de estrógeno (en fmol/l).
- **Grade:** grado de diferenciación del tumor, con valores 2 o 3.
- **Size:** tamaño del tumor, con niveles: menos de 20mm, entre 20 y 50mm, 50mm.

De la totalidad de mujeres en estudio, se cuenta con el tiempo exacto hasta la muerte de 377 de ellas y 206 censuras.

Selección del modelo

En primer lugar se compara, para cada variable, su modelo univariado contra un modelo sin covariables. Previo a esto, se definen 2 *dummies* referentes a la variable *Tamaño* y 1 para la variable *Grado*:

size	S_1	S_2	grade	G
< 20	0	0	2	0
20-50	1	0	3	1
> 50	0	1		

Las hipótesis en contraste son:

$$\begin{cases} H_0) & \beta_j = 0 \\ H_1) & \beta_j \neq 0 \end{cases} \quad \text{con } j = \overline{1, 8}$$

Los resultados obtenidos son los siguientes:

Variable	p -value	Decisión
Edad	~ 0	Rechazo H_0
Menopausia	~ 0	Rechazo H_0
Tratamiento hormonal	0.8381	No rechazo H_0
Quimioterapia	0.4217	No rechazo H_0
Receptores de progesterona	0.1316	No rechazo H_0
Receptores de estrógeno	0.0017	Rechazo H_0
Grado de diferenciación	0.0004	Rechazo H_0
Tamaño	~ 0	Rechazo H_0

¹Fuente: [Organización Panamericana de la Salud](#)

Se determina entonces que las variables significativas en esta etapa de la selección de variables son: Edad, Menopausia, Receptores de estrógeno, Grado de diferenciación y Tamaño del tumor.

En segundo lugar, se evaluará si cada una de esas variables siguen siendo significativas en presencia de las demás.

Para cada variable, se compara un modelo aditivo que contenga todas las variables significativas hasta el momento, excepto la variable en cuestión, contra un modelo que contenga todas las variables significativas.

- Modelo aditivo:

$$h_i(t) = h_0(t) \cdot \exp(\beta_1 \cdot \text{age}_i + \beta_2 \cdot \text{meno}_i + \beta_3 \cdot \text{er}_i + \beta_4 \cdot G_i + \beta_5 \cdot S_{1i} + \beta_6 \cdot S_{2i})$$

- Modelos de comparación:

- 1) Sin la variable *Edad*:

$$h_i(t) = h_0(t) \cdot \exp(\beta_2 \cdot \text{meno}_i + \beta_3 \cdot \text{er}_i + \beta_4 \cdot G_i + \beta_5 \cdot S_{1i} + \beta_6 \cdot S_{2i})$$

- 2) Sin la variable *Menopausia*:

$$h_i(t) = h_0(t) \cdot \exp(\beta_1 \cdot \text{age}_i + \beta_3 \cdot \text{er}_i + \beta_4 \cdot G_i + \beta_5 \cdot S_{1i} + \beta_6 \cdot S_{2i})$$

- 3) Sin la variable *Receptores de estrógeno*:

$$h_i(t) = h_0(t) \cdot \exp(\beta_1 \cdot \text{age}_i + \beta_2 \cdot \text{meno}_i + \beta_4 \cdot G_i + \beta_5 \cdot S_{1i} + \beta_6 \cdot S_{2i})$$

- 4) Sin la variable *Grado de diferenciación del tumor*:

$$h_i(t) = h_0(t) \cdot \exp(\beta_1 \cdot \text{age}_i + \beta_2 \cdot \text{meno}_i + \beta_3 \cdot \text{er}_i + \beta_5 \cdot S_{1i} + \beta_6 \cdot S_{2i})$$

- 5) Sin la variable *Tamaño del tumor*:

$$h_i(t) = h_0(t) \cdot \exp(\beta_1 \cdot \text{age}_i + \beta_2 \cdot \text{meno}_i + \beta_3 \cdot \text{er}_i + \beta_4 \cdot G_i)$$

Hipótesis: las hipótesis del test se dividen en 2 casos:

- Caso 1) las primeras 4 variables tienen 1 solo coeficiente asociado en el modelo:

$$\begin{cases} H_0) & \beta_j = 0 \\ H_1) & \beta_j \neq 0 \end{cases} \quad \text{con } j = \overline{1, 4}$$

- Caso 2) el tamaño del tumor tiene 2 coeficientes en el modelo, entonces:

$$\begin{cases} H_0) & \beta_5 = \beta_6 = 0 \\ H_1) & \text{Al menos un } \beta_j \neq 0, \quad \text{con } j = \overline{5, 6} \end{cases}$$

Teniendo en cuenta que los resultados de las comparaciones de los modelos fueron los siguientes:

Variable	p -value	Decisión
Edad	0.0229	Rechazo H_0
Menopausia	0.5586	No rechazo H_0
Receptores de estrógeno	0.1276	No rechazo H_0
Grado de diferenciación	0.0056	Rechazo H_0
Tamaño	~ 0	Rechazo H_0

El modelo que se obtendría sería $h_i(t) = h_0(t) \cdot \exp(\beta_1 \cdot \text{age}_i + \beta_2 \cdot G_i + \beta_3 \cdot S_{1i} + \beta_4 \cdot S_{2i})$. Sin embargo, este no es el modelo final dado que falta evaluar si ahora, con otras variables en el modelo, aquellas que no resultaron significativas de manera univariada lo son.

Nuevamente se analizan los resultados brindados por la comparación entre los modelos:

Variable	p -value	Decisión
Tratamiento hormonal	0.5420	No rechazo H_0
Quimioterapia	0.0209	Rechazo H_0
Receptores de progesterona	0.3811	No rechazo H_0

De esta manera, el modelo sería $h_i(t) = h_0(t) \cdot \exp(\beta_1 \cdot \text{age}_i + \beta_2 \cdot \text{chemo}_i + \beta_3 \cdot G_i + \beta_4 \cdot S_{1i} + \beta_5 \cdot S_{2i})$

Una vez definido, se prueban las distintas interacciones dobles entre las variables y resulta que ninguna es significativa. Por lo que no se altera el modelo.

Linealidad

Cuando se tienen variables de tipo continua debe analizarse en qué forma se las incluye en el modelo y hay dos opciones para ello: lineal y no lineal; esta última contempla tanto la idea de categorizar la variable como la de trabajar con una función de ella (logaritmo, al cuadrado, etc). Las hipótesis que se plantean son:

$$\begin{aligned} H_0) & \text{ El efecto de Edad es lineal} \\ H_1) & \text{ El efecto de Edad NO es lineal} \end{aligned}$$

La probabilidad asociada es mayor al 5%, entonces la variable es lineal. Con este resultado, se concluye que el modelo definitivo es el presentado con anterioridad:

$$h_i(t) = h_0(t) \cdot \exp(\beta_1 \cdot \text{age}_i + \beta_2 \cdot \text{chemo}_i + \beta_3 \cdot G_i + \beta_4 \cdot S_{1i} + \beta_5 \cdot S_{2i})$$

Comprobación de supuestos

Como es sabido, todo el trabajo realizado se desarrolló suponiendo que los hazards son proporcionales pero ¿Se cumple esto? Se comprará mediante el uso de los residuos de Schoenfeld obtenidos para cada variable.

Variable	p -value
Edad	~ 0
Quimioterapia	0.0438
Grado	0.2969
Tamaño	0.1809
Global	0.0004

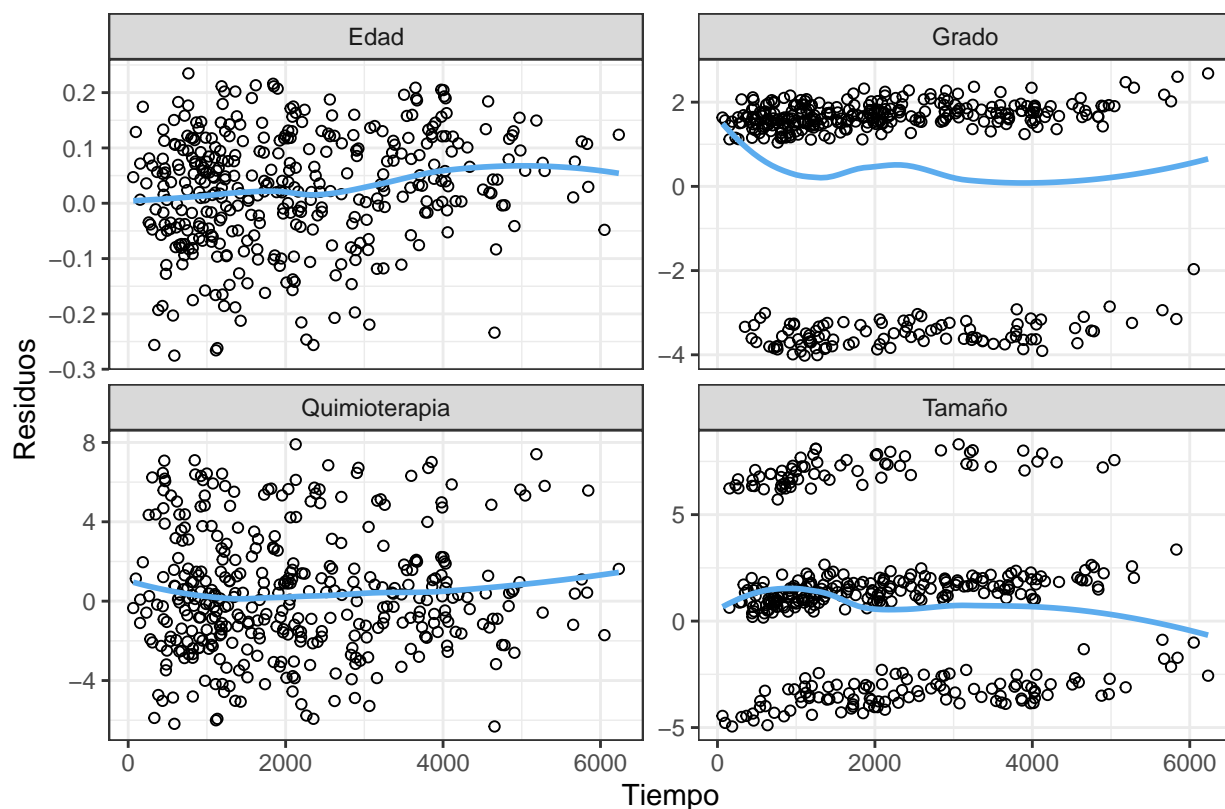


Gráfico 1: Evaluación de la proporcionalidad

La función `cox.zph()`, de manera global, dice que el supuesto no es cierto ya que se rechaza la hipótesis nula. Además, brinda información acerca de qué variables no estarían cumpliendo la proporcionalidad de los hazards: *Edad* y *Quimioterapia*. Esto puede constatarse también de forma gráfica al ver que los residuos de *Quimioterapia* no están divididos en dos grupos (cantidad de categorías que tiene), lo que sí sucede con *Grado* y *Tamaño*.

De la variable *Edad*, por ser continua, se debería observar un patrón aleatorio, pero se observa que los puntos están menos dispersos al avanzar el tiempo.

Para incluir las variables que son significativas pero no cumplen el supuesto de proporcionalidad, se decide lo siguiente:

- Categorizar la variable *Edad* con valores del 1 al 4.
- Plantear un modelo nuevo con la variable *Quimioterapia* estratificada.

Variable	p-value	Variable	p-value
A ₁	0.540	A ₁	0.208
A ₂	0.500	A ₂	0.592
A ₃	0.063	A ₃	0.014
Grado	0.201	Grado	0.259
Tamaño	0.123	Tamaño	0.117
		Quimioterapia	0.065
Global	0.087	Global	0.046

Si bien con un nivel de significación del 5% el solo el modelo estratificado cumple el supuesto de hazards proporcionales (Tabla x), el p-value asociado al test sin estratificar es muy cercano a 0.05 y ese modelo nos da la posibilidad de estimar un parámetro para la Quimioterapia, que es una variable importante para el problema. Por lo tanto, decidimos quedarnos con ese modelo. De esta forma, el modelo final contiene las siguientes variables:

- Edad (categorizada)
- Quimioterapia.
- Grado de diferenciación tumoral.
- Tamaño del tumor.

Por lo que el modelo estimado resulta ser

$$\hat{h}_i(t) = \hat{h}_0(t) \cdot \exp(-0.05 \cdot A_{1i} + 0.5 \cdot A_{2i} + 0.77 \cdot A_{3i} + 0.35 \cdot \text{chemo}_i + 0.31 \cdot G_i + 0.38 \cdot S_{1i} + 0.84 \cdot S_{2i})$$

Interpretación del modelo

Coefficientes

Con solo mirar los coeficientes del modelo se puede concluir:

- El perfil de paciente femenino con mejor pronóstico al realizar la cirugía primaria de extracción del tumor mamario es: mujer entre 40 y 50 años de edad, que no recibió quimioterapia y que posee un tumor de grado 2 con menos de 20cm de diámetro.
- El perfil de paciente femenino con el peor pronóstico es: mujer mayor a 65 años de edad, que recibió quimioterapia y posee un tumor de grado 3 con más de 50cm de diámetro.

Razones de hazards

Otra forma de interpretar el modelo es obteniendo razones de hazards, y sus respectivos intervalos de confianza, para distintos valores de nuestras variables