

Inferencia en el modelo de regresión lineal: intervalos de confianza

Contents

1	Introduccion	1
2	Intervalo de confianza para las β_i	1
3	Intervalo de confianza para σ^2	2
4	Ejemplo	2

1 Introduccion

Un intervalo de confianza para un parámetro es **un rango de valores posibles para dicho parámetro**.

2 Intervalo de confianza para las β_i

Hemos visto que

$$\hat{\beta} \rightarrow N(\beta, Q\sigma^2)$$

donde $Q = (X^T X)^{-1}$. Esto implica que:

$$\hat{\beta}_i \rightarrow N(\beta_i, Q_{ii}\sigma^2), \quad i = 1, 2, \dots, k$$

donde Q_{ij} es el elemento ij de la matriz Q . Aplicando las propiedades de la distribución normal

$$\frac{\hat{\beta}_i - \beta_i}{\sqrt{Q_{ii}\sigma^2}} \rightarrow N(0, 1)$$

Por tanto:

$$\frac{\hat{\beta}_i - \beta_i}{se(\hat{\beta}_i)} \rightarrow t_{n-k-1}$$

donde

$$se(\hat{\beta}_i) = \sqrt{Q_{ii}\hat{\sigma}_R^2}$$

Para deducir la expresión anterior se ha tenido en cuenta que

$$\frac{N(0, 1)}{\sqrt{\frac{\chi_n^2}{n}}} \rightarrow t_n$$

Por tanto, el intervalo de confianza $100(1 - \alpha)\%$ se escribe como

$$\hat{\beta}_i \pm t_{n-k-1; \alpha/2} se(\hat{\beta}_i)$$

3 Intervalo de confianza para σ^2

Partimos de la distribución en el muestreo:

$$\frac{(n-k-1)\hat{s}_R^2}{\sigma^2} \rightarrow \chi_{n-k-1}^2$$

Despejando:

$$\frac{(n-k-1)\hat{s}_R^2}{\chi_{n-k-1; \alpha/2}^2} \leq \sigma^2 \leq \frac{(n-k-1)\hat{s}_R^2}{\chi_{n-k-1; 1-\alpha/2}^2}$$

4 Ejemplo

Vamos a calcular de manera detallada los intervalos de confianza para el modelo $kid_score \sim mom_iq + mom_hs$:

```
d = read.csv("datos/kidiq.csv")
#d$mom_hs = factor(d$mom_hs, labels = c("no", "si"))

m = lm(kid_score ~ mom_iq + mom_hs, data = d)
```

Los parámetros estimados son:

```
coef(m)

## (Intercept)      mom_iq      mom_hs
##   25.731538    0.563906    5.950117

# varianza residual
n = nrow(d)
k = 2 # numero de regresores
(sR2 = sum(resid(m)^2)/(n-k-1))

## [1] 328.9028
```

Vamos a calcular la varianza de los parámetros estimados, es decir $var(\hat{\beta}_i) = Q_{ii}\hat{s}_R^2$:

```
X = cbind(rep(1,n), d$mom_iq, d$mom_hs)
# Q = inv(t(X)*X)
(Q = solve(crossprod(X))) # crossprod es otra manera de calcular t(X) %*% X

##           [,1]      [,2]      [,3]
## [1,]  0.1049491626 -1.025110e-03 -0.0001705848
## [2,] -0.0010251098  1.115594e-05 -0.0001151616
## [3,] -0.0001705848 -1.151616e-04  0.0148740410
```

Por tanto, la matriz de varianzas de los estimadores será

```
(beta_var = sR2 * Q)

##           [,1]      [,2]      [,3]
## [1,] 34.51806922 -0.337161456 -0.05610582
```

```
## [2,] -0.33716146  0.003669219 -0.03787697
## [3,] -0.05610582 -0.037876974  4.89211314
```

Y el standard error de los estimadores, $se(\hat{\beta}_i)$:

```
(beta_se = sqrt(diag(beta_var)))
```

```
## [1] 5.87520802 0.06057408 2.21181218
```

Vamos a calcular ahora el standard error de los estimadores con la matriz de varianzas de los regresores:

```
Xa = cbind(d$mom_iq, d$mom_hs)
(Qa = 1/n*solve(var(Xa)))
```

```
##           [,1]      [,2]
## [1,] 1.113023e-05 -0.0001148963
## [2,] -1.148963e-04  0.0148397690
```

las pequeñas diferencias numéricas se deben a que en las ecuaciones se ha considerado la matriz de covarianzas dividiendo entre n y R calcula dicha matriz dividiendo entre $(n-1)$:

```
(Qa = 1/n*solve(var(Xa)*(n-1)/n))
```

```
##           [,1]      [,2]
## [1,] 1.115594e-05 -0.0001151616
## [2,] -1.151616e-04  0.0148740410
```

El standard error de los estimadores $\hat{\beta}_1$ y $\hat{\beta}_2$ son:

```
sqrt(diag(Qa)*sR2)
```

```
## [1] 0.06057408 2.21181218
```

Para $\hat{\beta}_0$:

```
( xmed = matrix(colMeans(Xa), ncol = 1) )
```

```
##           [,1]
## [1,] 100.0000000
## [2,]  0.7857143
```

```
sqrt( sR2/n*(1 + t(xmed) %*% solve(var(Xa)*(n-1)/n) %*% xmed ) )
```

```
##           [,1]
## [1,] 5.875208
```

Por último, R dispone de una función para calcular la matriz de varianzas de los parámetros estimados, es decir $var(\hat{\beta}) = Q_{ii}\hat{s}_R^2$, mediante:

```
vcov(m)
```

```
##           (Intercept)      mom_iq      mom_hs
## (Intercept) 34.51806922 -0.337161456 -0.05610582
## mom_iq      -0.33716146  0.003669219 -0.03787697
## mom_hs      -0.05610582 -0.037876974  4.89211314
```

Por tanto, el standard error de los estimadores será

```
sqrt(diag(vcov(m)))
```

```
## (Intercept)      mom_iq      mom_hs
## 5.87520802 0.06057408 2.21181218
```

Como vemos, los tres métodos dan el mismo resultado.

El valor de la t con $n-k-1 = 431$ grados de libertad es

```
(t1 = qt(1-0.05/2, df = n-k-1))
```

```
## [1] 1.965483
```

El límite inferior (LI) y el límite superior de los intervalos será:

```
(LI = coef(m) - qt(1-0.05/2, df = n-k-1)*beta_se)
```

```
## (Intercept)      mom_iq      mom_hs  
## 14.1839148    0.4448487    1.6028370
```

```
(LS = coef(m) + qt(1-0.05/2, df = n-k-1)*beta_se)
```

```
## (Intercept)      mom_iq      mom_hs  
## 37.2791615    0.6829634   10.2973969
```

Si lo juntamos todo en una tabla

```
data.frame(estimacion = coef(m), se = beta_se, LI, LS)
```

```
##           estimacion      se      LI      LS  
## (Intercept) 25.731538 5.87520802 14.1839148 37.2791615  
## mom_iq      0.563906 0.06057408 0.4448487 0.6829634  
## mom_hs      5.950117 2.21181218 1.6028370 10.2973969
```

Directamente, mediante la función *confint()* de R se pueden obtener dichos valores:

```
confint(m)
```

```
##           2.5 %      97.5 %  
## (Intercept) 14.1839148 37.2791615  
## mom_iq      0.4448487 0.6829634  
## mom_hs      1.6028370 10.2973969
```

Si queremos otro nivel de confianza, por ejemplo, 90%:

```
confint(m, level = 0.90)
```

```
##           5 %      95 %  
## (Intercept) 16.0468646 35.4162117  
## mom_iq      0.4640559 0.6637562  
## mom_hs      2.3041730 9.5960608
```

- En el caso de la varianza del modelo. Su estimador es:

```
sR2
```

```
## [1] 328.9028
```

Y su intervalo de confianza:

```
c((n-k-1)*sR2/qchisq(1-0.05/2, df = n-k-1), (n-k-1)*sR2/qchisq(0.05/2, df = n-k-1))
```

```
## [1] 289.0557 377.6434
```