

# Proyecto 1

Agustín Riquelme y Heriberto Espino

Para este modelo de regresión lineal se usó la base de datos lung cap, la cual recolecta información de un grupo de personas acerca de si son fumadores o no, las edades respectivas, además de la altura, el género y su FEV(Forced Expiratory Volumen), que mide la habilidad de expirar aire de los pulmones. En este caso, la variable que se tomó como predictora es la edad, en embargo, de igual manera se accionará en dos, en fumadores y no fumadores para un mejor análisis de los datos.

Para el caso general, donde se opta por juntar ambos grupos sin diferencia, tenemos los siguientes datos estadísticos, donde podemos observar que las personas que fueron parte de esta muestra van desde los 3 hasta los 19 años de edad con un promedio de 10 años y misma mediana, además, la edad que más se repite es de 9 años y además su variancia su variancia no difiere mucho de la media.

| Edad                |         |
|---------------------|---------|
| Datos estadísticos  | Valores |
| Mínimo              | 3.000   |
| Primer cuartil      | 8.000   |
| Mediana             | 10.000  |
| Media               | 9.931   |
| Tercer cuartil      | 12.000  |
| Máximo              | 19.000  |
| Moda                | 9.000   |
| Varianza            | 8.726   |
| Desviación estándar | 2.954   |

Para el caso del FEV, el mínimo es de 0.8 redondeado con un máximo de 3.6, y la media de este índice es de 2.6 con una variancia mucho más pequeña que la media en términos unitarios y porcentuales. Esta variable de FEV corresponderá a la variable de respuesta, pues se cree que el FEV va a depender de la edad de la persona en cuestión pero igualmente de si es fumador o no, pero esos datos se designarán a continuación.

| FEV                 |         |
|---------------------|---------|
| Datos estadísticos  | Valores |
| Mínimo              | 0.791   |
| Primer cuartil      | 2.169   |
| Mediana             | 2.547   |
| Media               | 2.637   |
| Tercer cuartil      | 3.119   |
| Máximo              | 5.793   |
| Moda                | 1.511   |
| Varianza            | 0.752   |
| Desviación estándar | 0.867   |

Para el caso de los fumadores, estos inician a partir de la edad de 9 años y termina en el mismo rango que los datos en conjunto, siendo 19 con un promedio de 14 años pero con 13 años como la cifra más repetida. Y a pesar de no ser mostrado, el porcentaje de fumadores dentro de la base corresponde a 10%.

| Fumadores           |         |
|---------------------|---------|
| Datos estadísticos  | Valores |
| Mínimo              | 9.000   |
| Primer cuartil      | 12.000  |
| Mediana             | 13.523  |
| Media               | 13.523  |
| Tercer cuartil      | 15.000  |
| Máximo              | 19.000  |
| Moda                | 13.000  |
| Varianza            | 5.472   |
| Desviación estándar | 2.339   |

Para el caso del índice FEV, el mínimo corresponde a 1.7 mientras que el máximo a 4.8, teniendo una media por encima de los datos totales sin embargo una variancia menor.

| No fumadores        |         |
|---------------------|---------|
| Datos estadísticos  | Valores |
| Mínimo              | 1.694   |
| Primer cuartil      | 2.795   |
| Mediana             | 2.465   |
| Media               | 2.566   |
| Tercer cuartil      | 3.751   |
| Máximo              | 4.872   |
| Moda                | 3.297   |
| Varianza            | 0.562   |
| Desviación estándar | 0.750   |

Ahora para el caso de los no fumadores, la edad más pequeña registrada es de 3 años, mientras que la máxima de igual manera corresponde a la máxima conjunta, en este caso, la media es menor, de 9 años al igual que su moda y mediana. Estos datos tienen una mayor variancia que los fumadores pero menor a la conjunta.

| No fumadores        |         |
|---------------------|---------|
| Datos estadísticos  | Valores |
| Mínimo              | 3.000   |
| Primer cuartil      | 8.000   |
| Mediana             | 9.000   |
| Media               | 9.535   |
| Tercer cuartil      | 11.000  |
| Máximo              | 19.000  |
| Moda                | 9.000   |
| Varianza            | 7.511   |
| Desviación estándar | 2.741   |

En el caso del índice FEV, se registró un mínimo de 0.8 y un máximo de 5.8, mayor que el caso de los fumadores. Sin embargo, los demás valores parecen ser más pequeños que los datos de fumadores, por lo que analizarlo de esta manera en conjunto y de dos grupos podrá determinar si el fumar tiene alguna consecuencia en este índice al igual que la edad.

| FEV                 |         |
|---------------------|---------|
| Datos estadísticos  | Valores |
| Mínimo              | 0.791   |
| Primer cuartil      | 1.920   |
| Mediana             | 2.465   |
| Media               | 2.566   |
| Tercer cuartil      | 3.048   |
| Máximo              | 5.793   |
| Moda                | 1.624   |
| Varianza            | 0.723   |
| Desviación estándar | 0.851   |

Entonces, para un mejor análisis, a partir de este momento, los datos serán seleccionados en dos grupos: fumadores y no fumadores. En el primer caso, tendremos los siguientes histogramas y boxplot para el caso de las edades y su índice FEV. Visualmente no observamos ningún outlier en los datos, y no parece haber un comportamiento similar a una distribución normal.

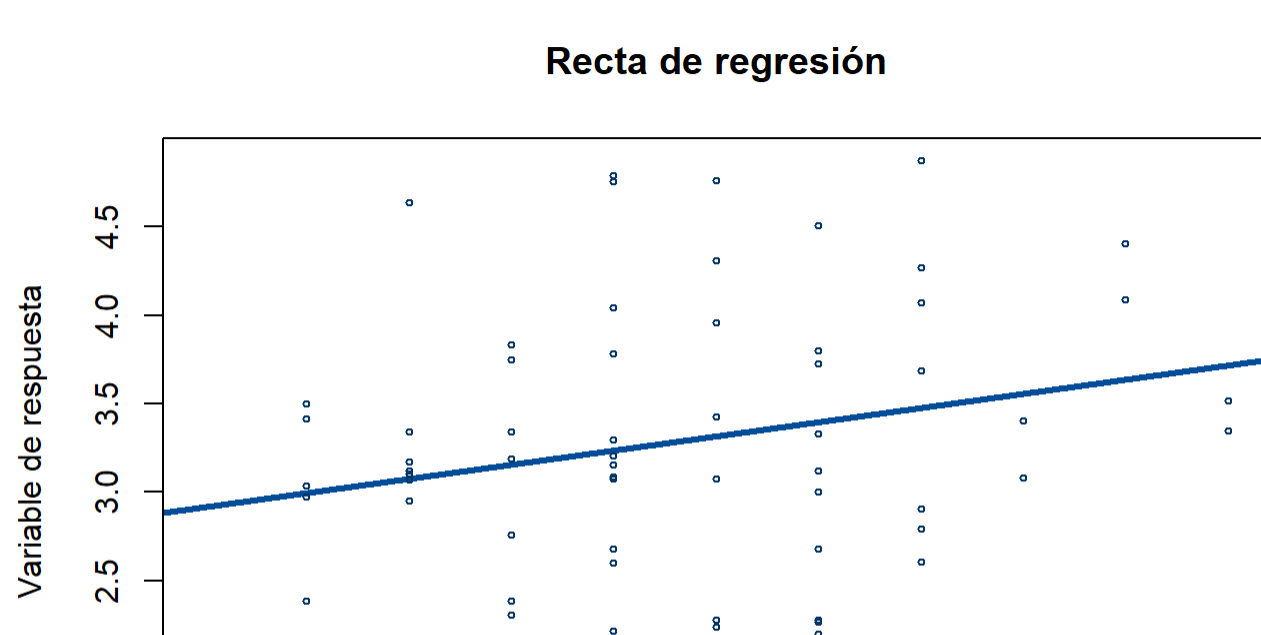
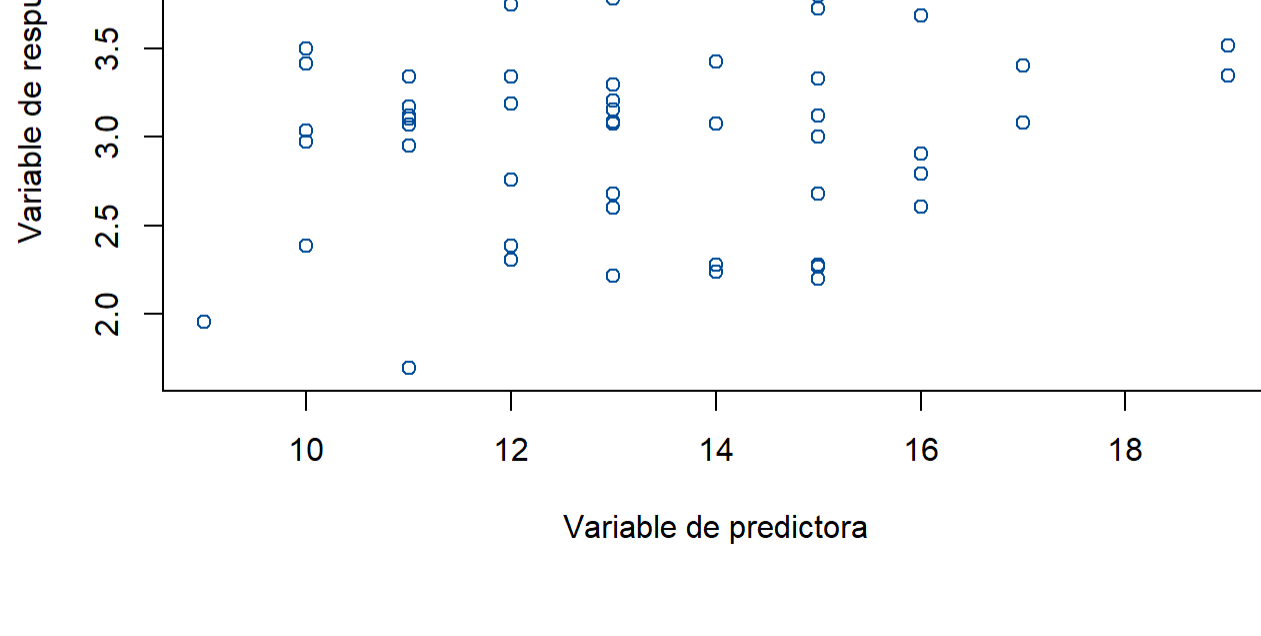
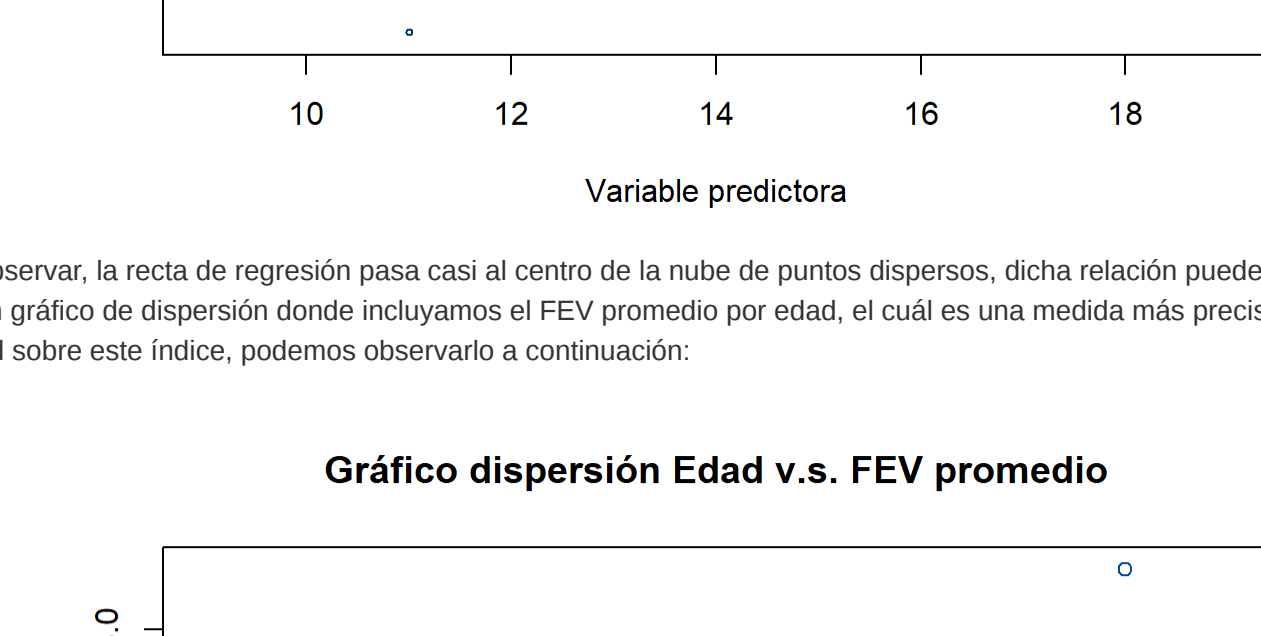
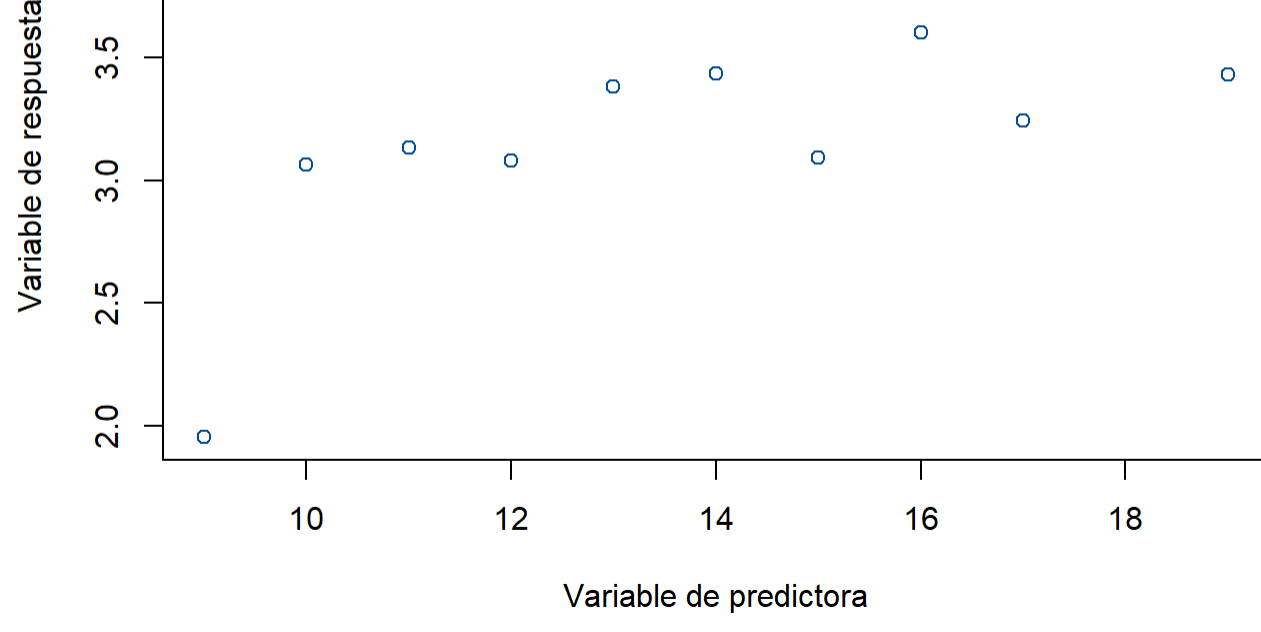


Gráfico dispersión Edad v.s. FEV

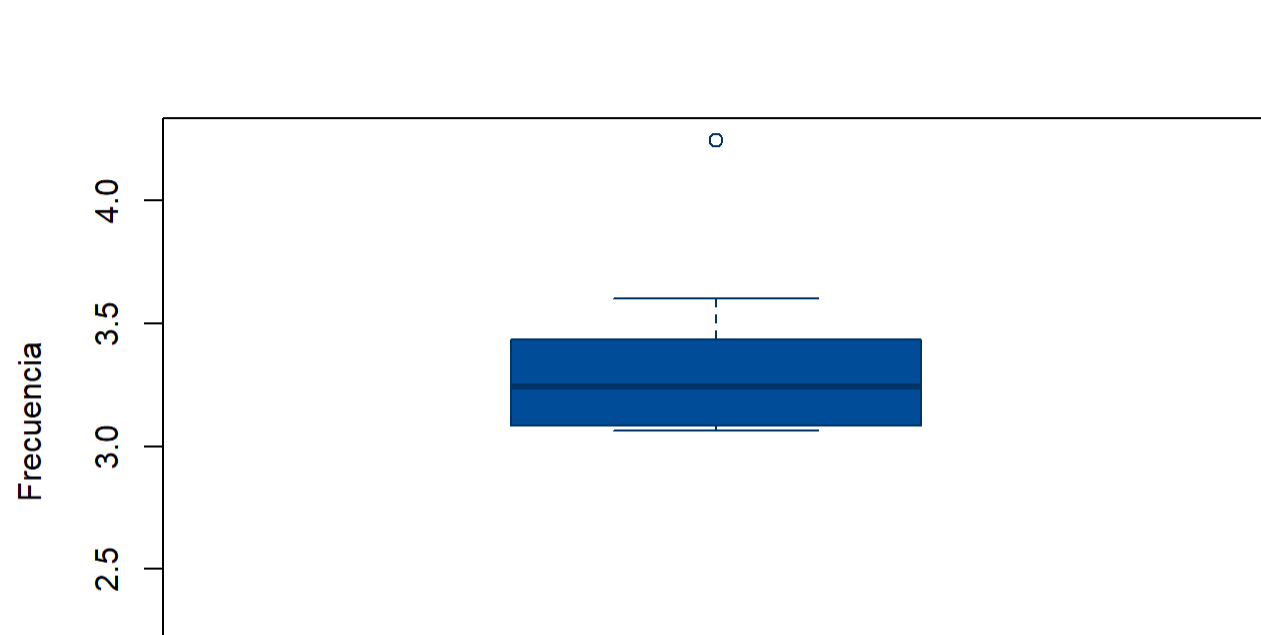


Recta de regresión

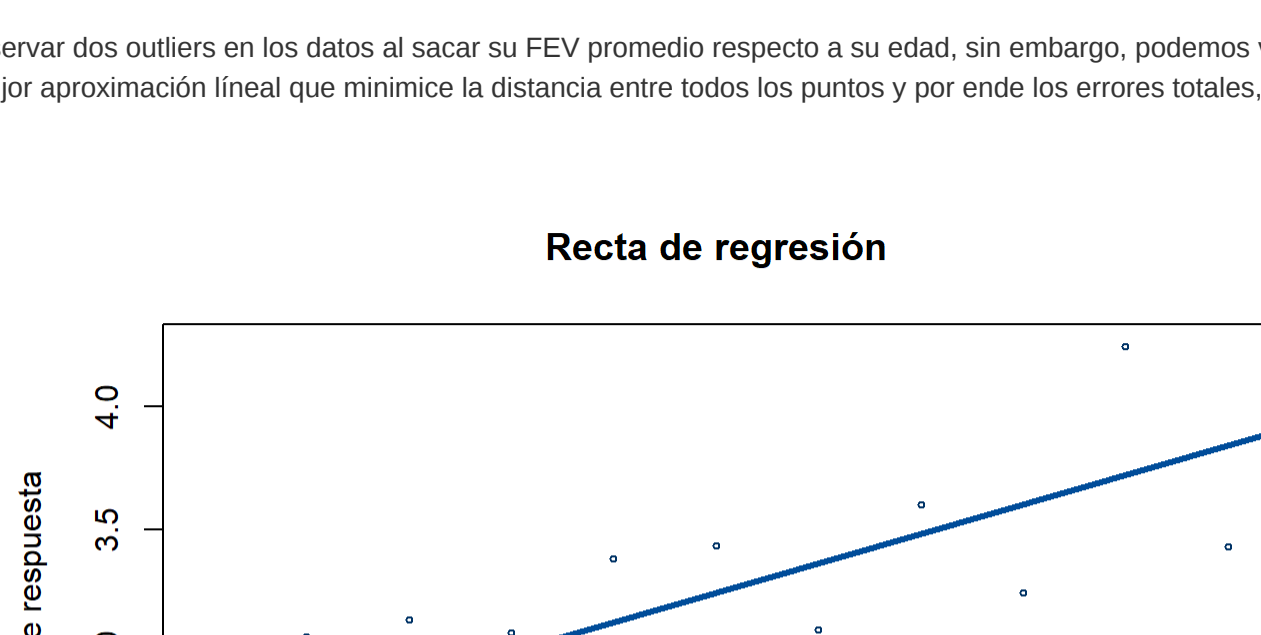


Como podemos observar, la recta de regresión pasa casi al centro de la nube de puntos dispersos, dicha relación puede ser vista de manera más fácil a través de un gráfico de dispersión donde el FEV promedio por edad, el cual es una medida más precisa para poder describir el impacto de la edad sobre este índice, podemos observarlo a continuación:

Gráfico dispersión Edad v.s. FEV promedio

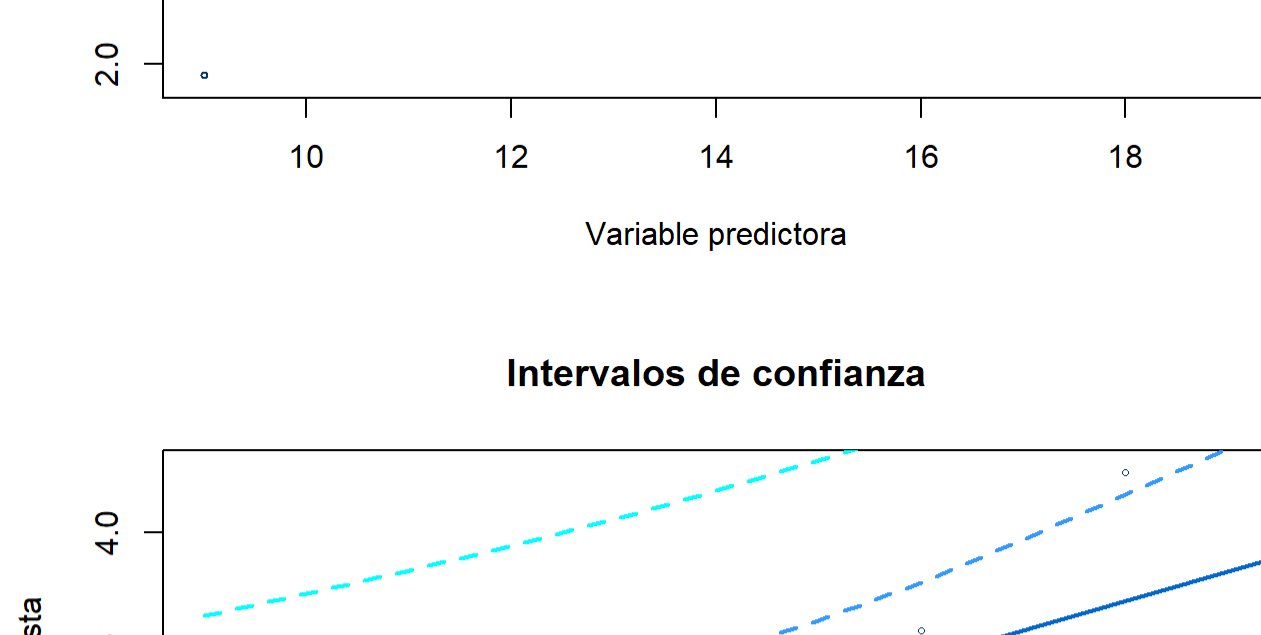


De esta otra manera, con el manejo del promedio de FEV por edad, se puede observar intuitivamente de una mejor manera la forma de la recta de regresión, pudiendo suponer que tiene noción positiva y con una pendiente menor a 1. Además podemos observar un posible outlier a la edad de 9 años, por lo que realizaremos de nuevo un boxplot para su análisis.



FEV

Aquí podemos observar dos outliers en los datos al sacar su FEV promedio respecto a su edad, sin embargo, podemos ver menos dispersos los datos para una mejor aproximación lineal que minimice la distancia entre todos los puntos y por ende los errores totales, esto se verá a continuación:



Intervalos de confianza

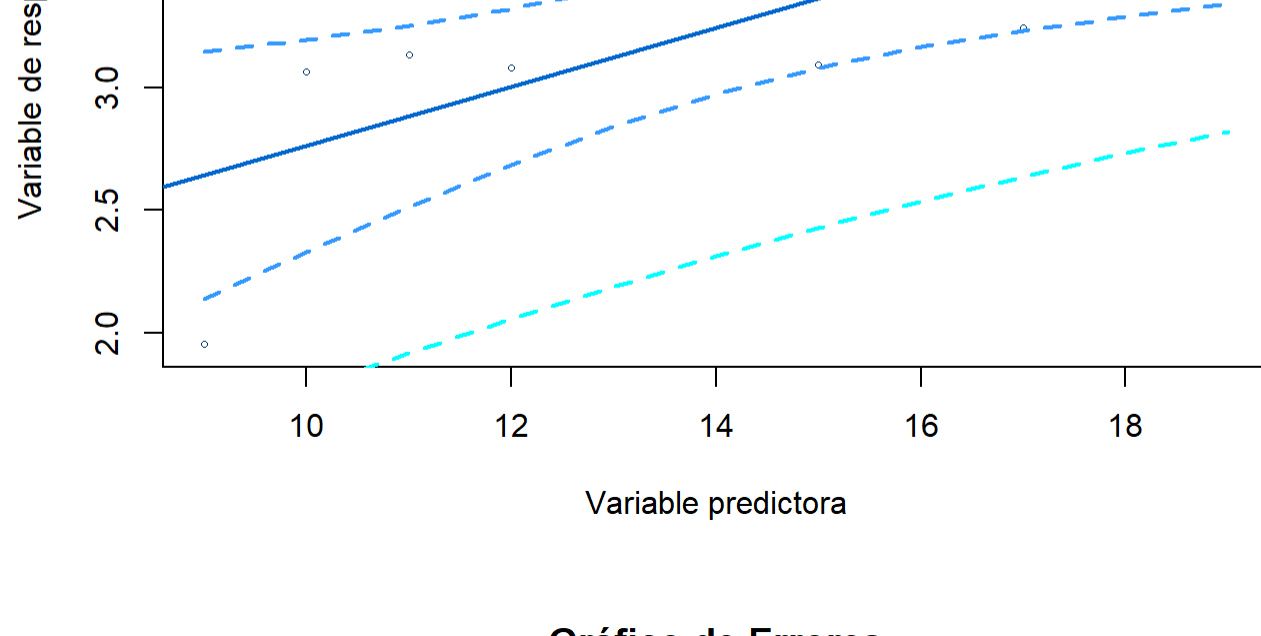
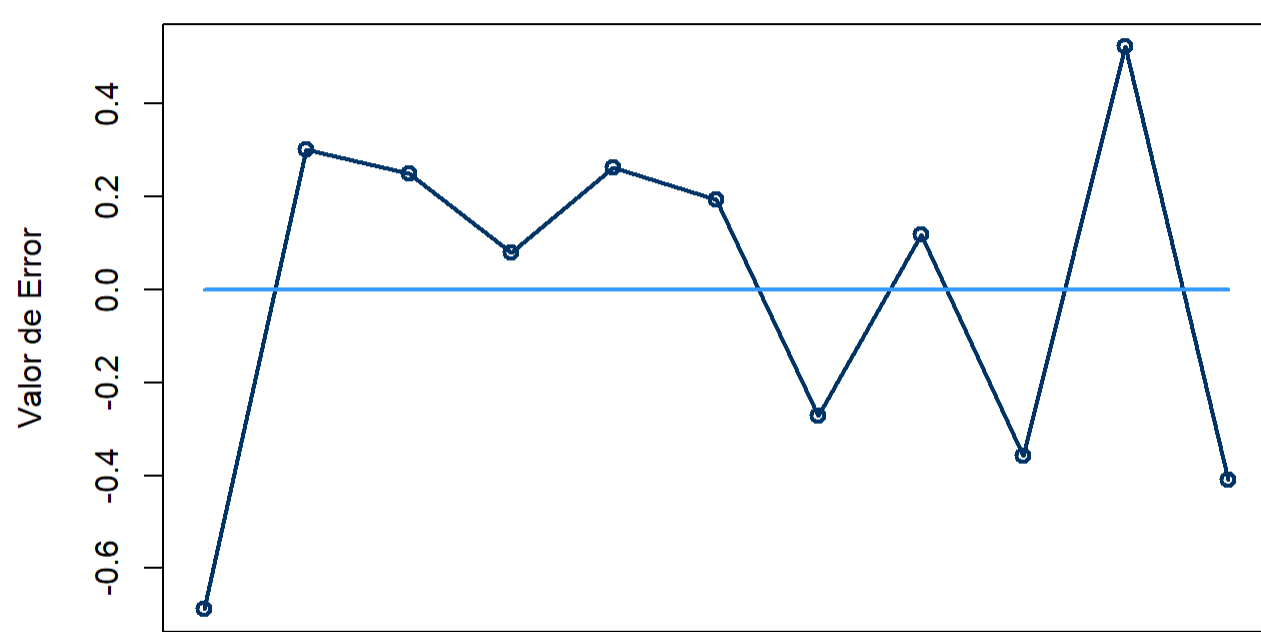


Gráfico de Errores



Histograma de Errores



Boxplot de errores

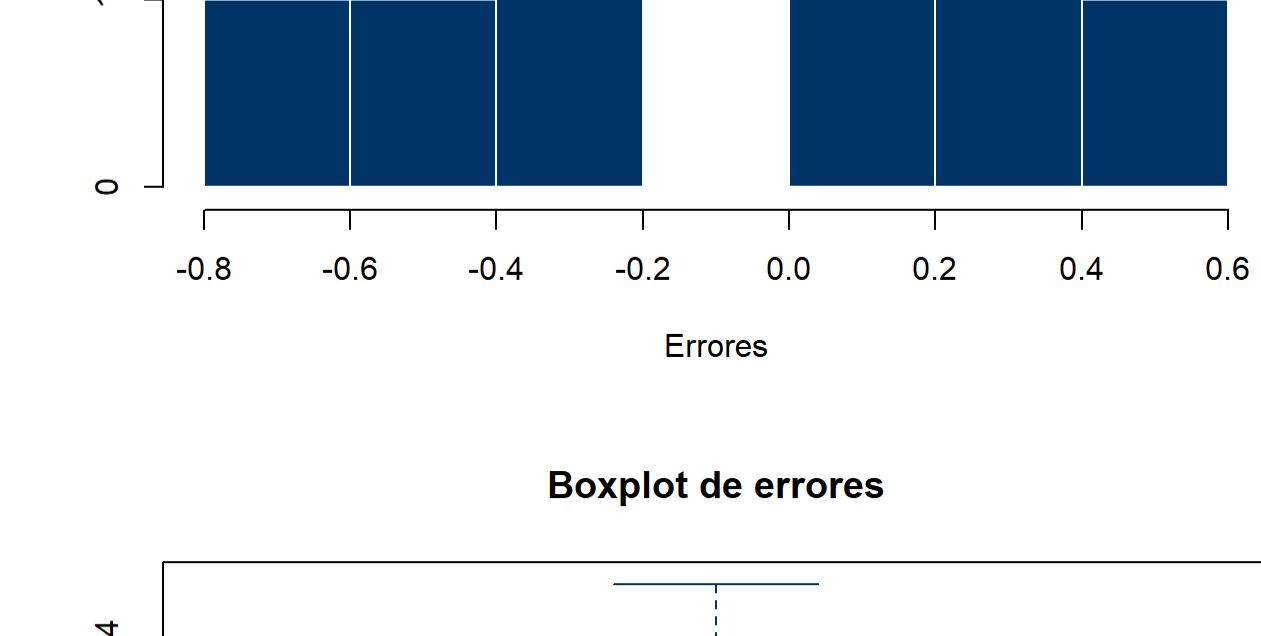
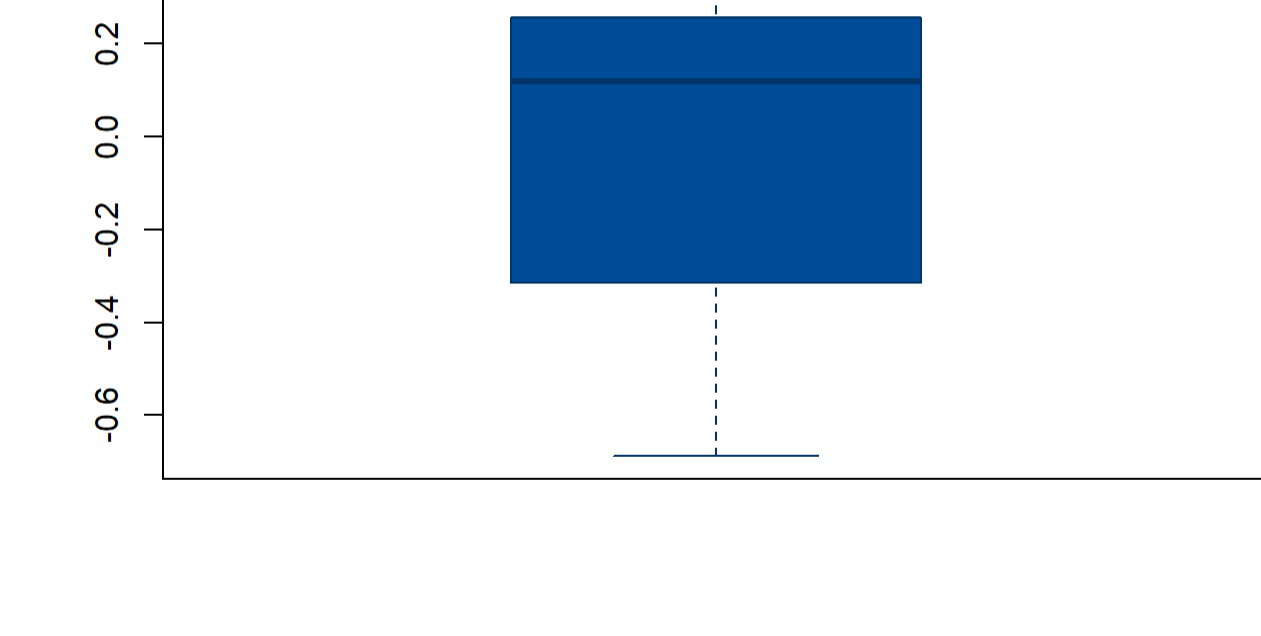


Gráfico QQ de Errores



Información del modelo

| Datos            | Valores             |
|------------------|---------------------|
| a                | 0.05                |
| b0               | 1.560207005217      |
| b1               | 0.12009024100224    |
| S                | 0.393903024150109   |
| S^2              | 0.155159592434601   |
| CV               | 12.1513270369503    |
| df               | 9                   |
| qti              | -2.2021571527982    |
| qtd              | 2.2621571657982     |
| tc               | 3.19777003936184    |
| P valor          | 0.010870134695964   |
| Min. b1          | 0.03513898964712    |
| Máx. b1          | 0.205059407253977   |
| Prom. err.       | -3.532527056255e-17 |
| Var. err.        | 0.13964363191141    |
| Shap. P valor    | 0.366242866344028   |
| Anderson P valor | 0.218479236573656   |
| Rho Pearson      | 0.7202926566454     |
| Rho Spearman     | 0.772727272727273   |
| dW               | 2.35620794814095    |

Inferencia para b1 H0: b1 = 0 vs H1: b1 ≠ 0

Prueba de hipótesis Conclusión

tc Rechazamos H0, b1 ≠ 0

p valor Rechazamos H0, b1 ≠ 0

Coefficiente de correlación

Rho Conclusión

Pearson Hay una correlación fuerte

Spearman Hay una correlación fuerte

Test de normalidad H0: Los datos son normales vs H1: Los datos no son normales

Test de normalidad Valor

Normalidad Shapiro p valor > α, no rechazamos H0

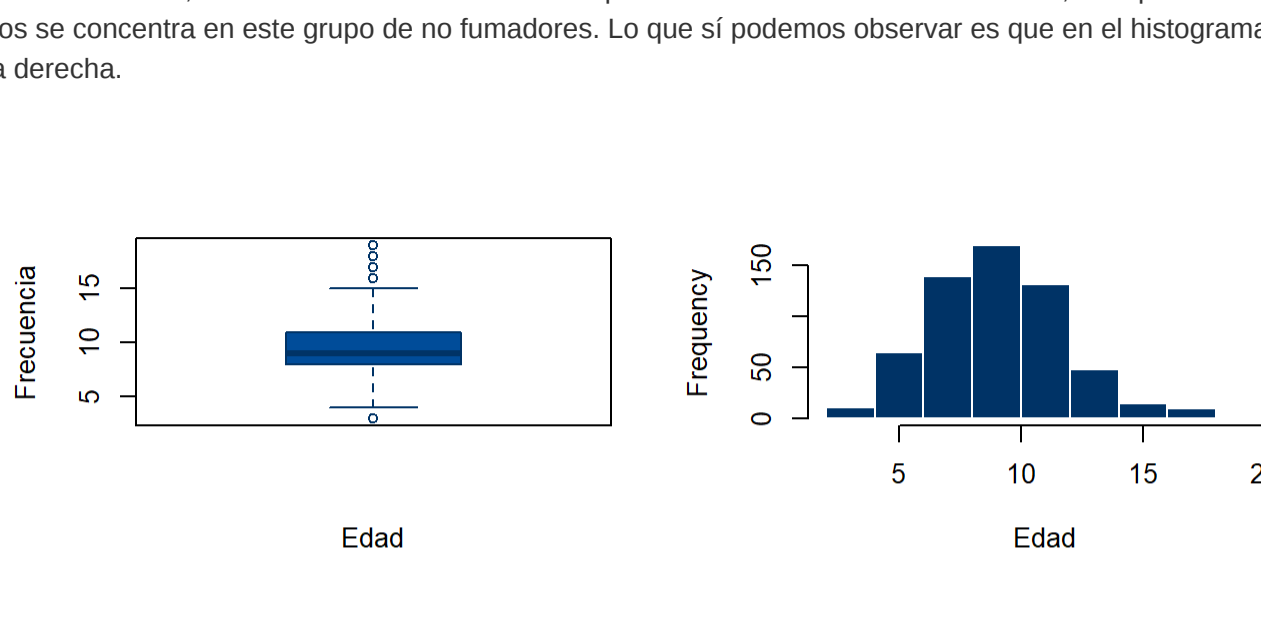
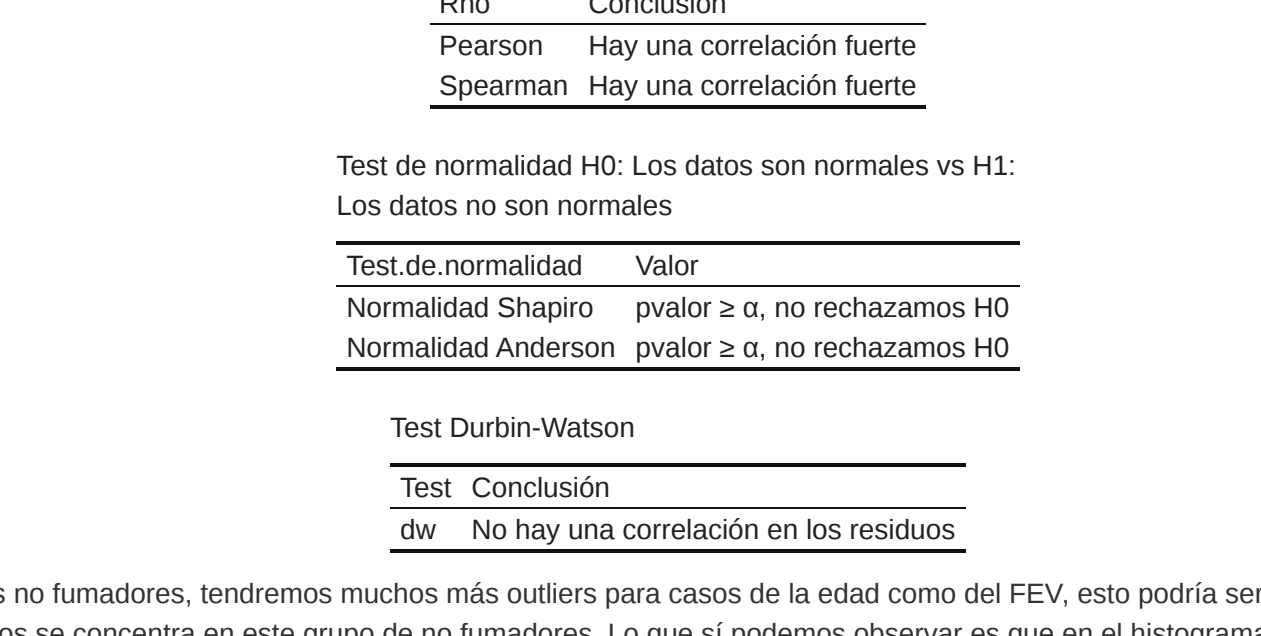
Normalidad Anderson p valor > α, no rechazamos H0

Test Durbin-Watson

Test Conclusión

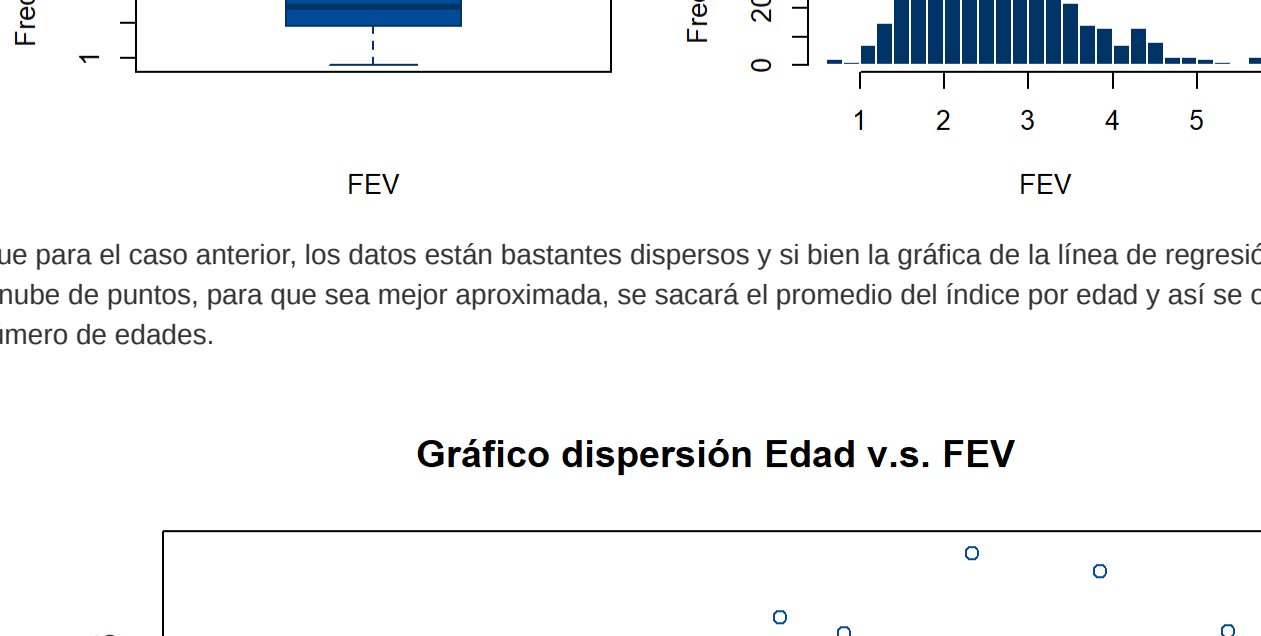
dW No hay una correlación en los residuos

Para el caso de los no fumadores, tendremos los mismos datos más outliers para el caso de la edad como del FEV, esto podría ser esperado al saber que la mayoría de los datos se concentra en este grupo de no fumadores. Lo que si podemos observar es que en el histograma ambos datos parecen tener un sesgo a la derecha.

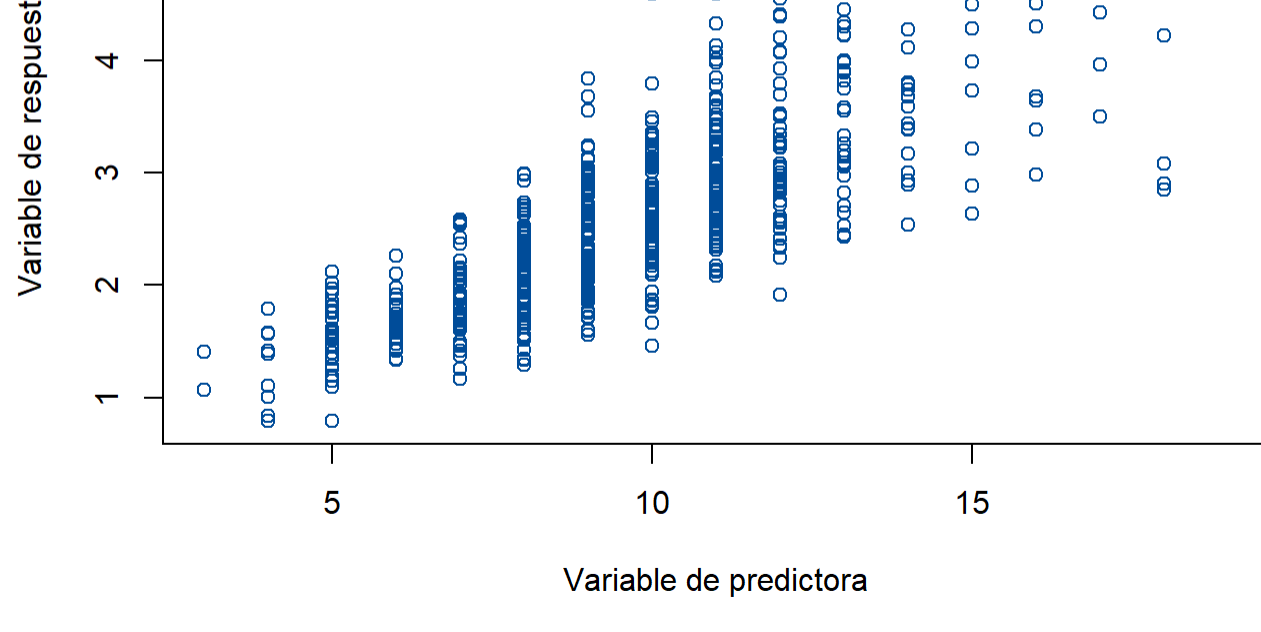


De igual manera que para el caso anterior, los datos están bastante dispersos y si bien la gráfica de la línea de regresión se encuentra en un punto medio de la nube de puntos, lo que se puede indicar que los errores de este modelo podrían ser menores que el anterior y sugiriendo que el fumar pueda alterar los resultados normales a los que se debería de comportar el índice FEV según la edad.

Gráfico dispersión Edad v.s. FEV



Recta de regresión



Intervalos de confianza

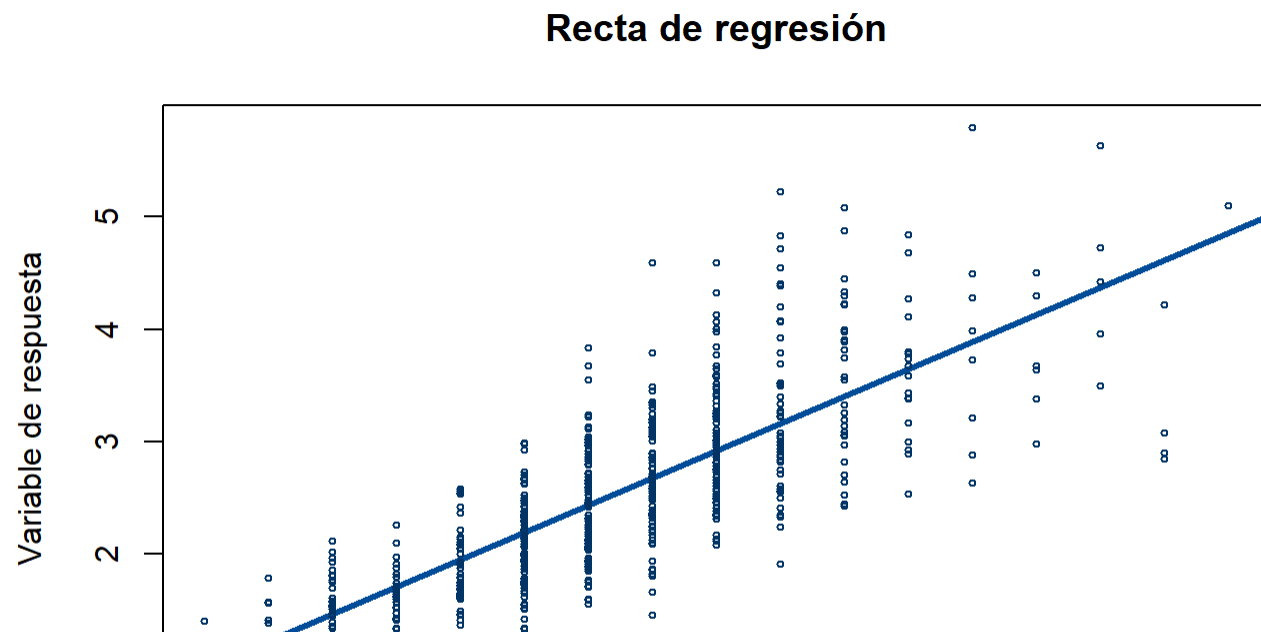
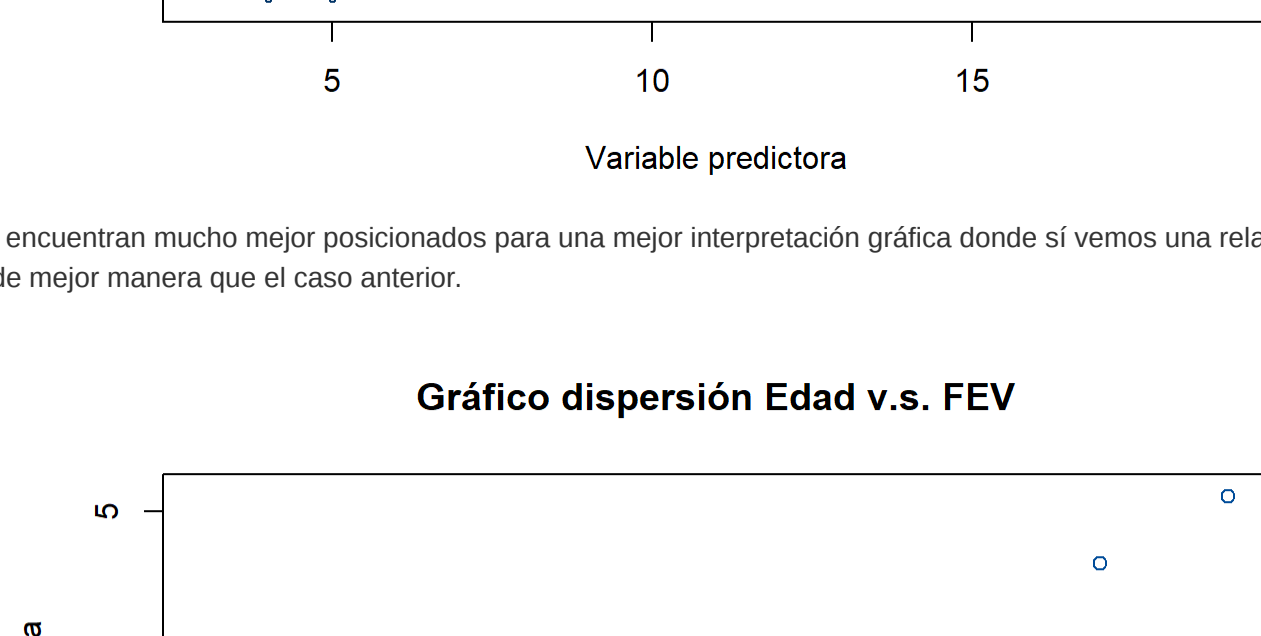
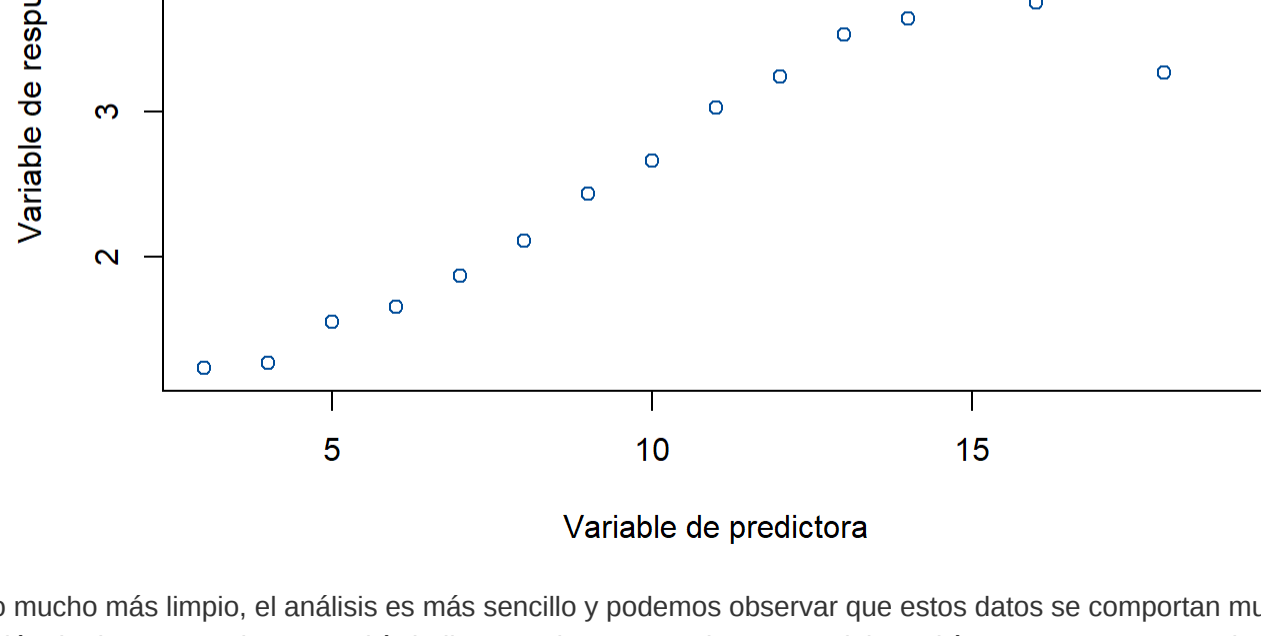


Gráfico de Errores



Histograma de Errores



Boxplot de errores

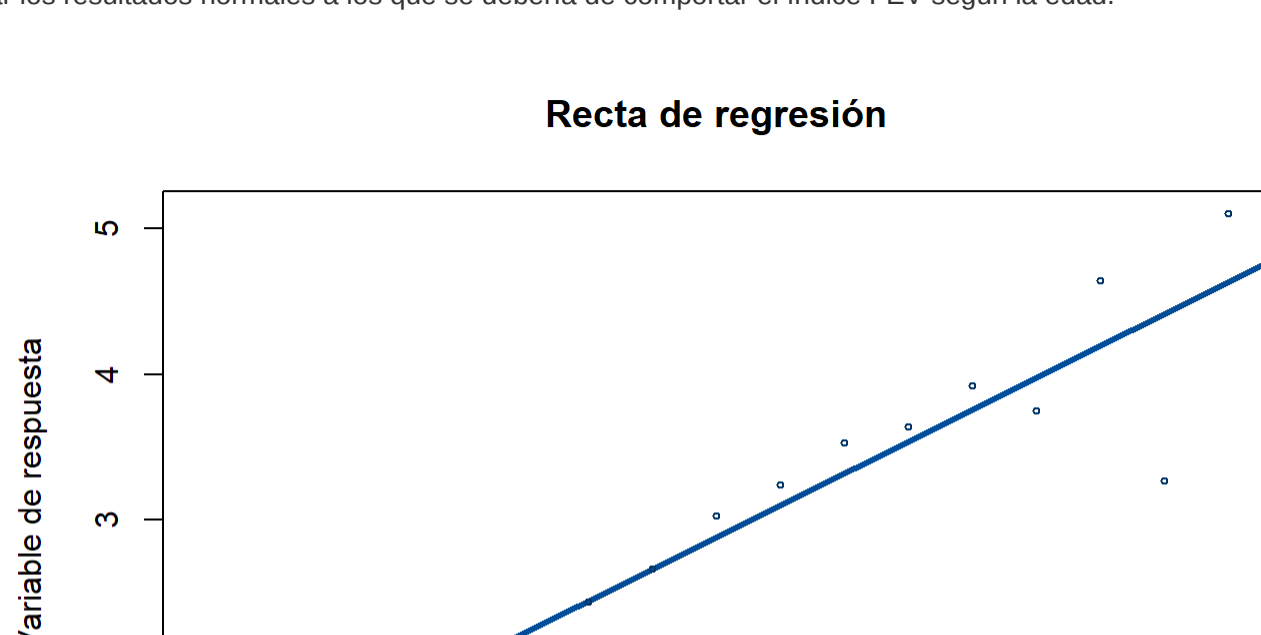
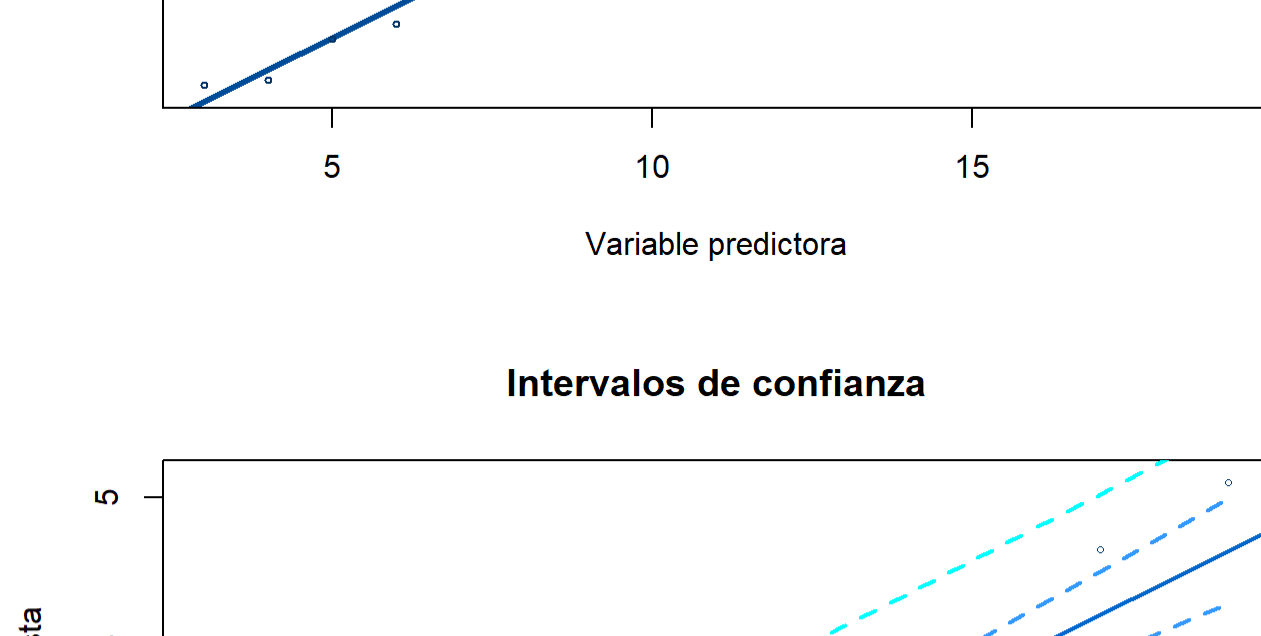


Gráfico QQ de Errores



Información del modelo

| Datos            | Valores                          |
|------------------|----------------------------------|
| a                | 0.05                             |
| b0               | 0.468253267996954                |
| b1               | 0.2191226501407799               |
| S                | 0.303813904046725                |
| S^2              | 0.132214472285774                |
| CV               | 12.6315879123559                 |
| df               | 15                               |
| qti              | -2.1314495455978                 |
| qtd              | 2.1314495455978                  |
| tc               | 12.17446865173                   |
| P valor          | 0.00000003559972489347           |
| Min. b1          | 0.180753297392903                |
| Máx. b1          | 0.257462020388895                |
| Prom. err.       | -0.00000000000000239848695484323 |
| Var. err.        | 0.123951067767913                |
| Shap. P valor    | 0.00124133691555556              |
| Anderson P valor | 0.003356766036633                |
| Rho Pearson      | 0.952927310022899                |
| Rho Spearman     | 0.95078431373549                 |
| dW               | 2.97416176068524                 |

Inferencia para b1 H0: b1 = 0 vs H1: b1 ≠ 0

Prueba de hipótesis Conclusión

tc Rechazamos H0, b1 ≠ 0

p valor Rechazamos H0, b1 ≠ 0

Coefficiente de correlación

Rho Conclusión

Pearson Hay una correlación fuerte

Spearman Hay una correlación fuerte

Test de normalidad H0: Los datos son normales vs H1: Los datos no son normales

Test de normalidad Valor

Normalidad Shapiro p valor < α, rechazamos H0, No son normales

Normalidad Anderson p valor < α, rechazamos H0, No son normales

Test Durbin-Watson

Test Conclusión

dW No hay una correlación en los residuos