

Punto 1° Parcial: Ajuste de un modelo de R.L.S

Universidad Nacional de Colombia

Análisis de Regresión 2022-1S

Medellín, Colombia

2022

Daniel Villa 1005087556

Juan Pablo Vanegas 1000640165



UNIVERSIDAD NACIONAL DE COLOMBIA

Contents

1	Objetivos:	4
1.1	Objetivos específicos	4
2	Antecedentes Relevantes	4
3	Variables de respuesta:	4
4	Variable de Control:	4
5	Ajuste del modelo	4
5.1	Intervalos de Confianza	6
5.2	Tabla ANOVA	7
5.3	Coefficiente de determinación	7
6	Análisis de los parámetros del modelo	7
7	Diagnóstico del modelo	7
7.1	Test de normalidad (test de Kolmogorov-Smirnov)	8
8	Transformación del Modelo	9
8.1	Intervalos de Confianza	11
8.2	Tabla ANOVA	12
8.3	Coefficiente de determinación	12
9	Análisis de los parámetros del modelo	12
10	Diagnóstico del modelo	12
10.1	Test de normalidad (test de Kolmogorov-Smirnov)	13
10.2	Autocorrelación (test de Durbin-Watson)	15
11	Predicción	15
11.1	Predicción de nuevas observaciones	15
11.2	Intervalos de confianza para los predictores:	16
12	Modelo N°2 (Peso~Altura)	18
12.1	Ajuste del modelo:	18
12.2	Intervalos de Confianza	20
12.3	Tabla ANOVA	20
12.4	Coefficiente de determinación	20
13	Análisis de los parámetros del modelo	20
14	Diagnóstico del modelo	21
14.1	Test de normalidad (test de Kolmogorov-Smirnov)	22
14.2	Autocorrelación (test de Durbin-Watson)	23
14.3	Valores atípicos:	23
15	Predicción	26
15.1	Predicción de nuevas observaciones	26
15.2	Intervalos de confianza para los predictores	27
16	Apéndice	29
16.1	Lista de figuras	29
16.2	Código:	30

1 Objetivos:

Crear un modelo ajustado de R.L.M. por el cual se pueda predecir la estatura de un individuo (discriminando por genero) sabiendo las estaturas de los padres (madre y padre) utilizando el software estadístico *R*.

1.1 Objetivos específicos

- Plantear el modelo de R.L.M.
- Interpretar los parámetros del modelo.
- Determinar si el efecto de las estaturas de los padres sobre la estatura del sujeto es significativo.
- Interpretar nuestro R^2 .
- Validar los supuestos del modelo.
- Aplicar la prueba de falta de ajuste.

2 Antecedentes Relevantes

La población encuestada, pertenece a estudiantes de la Universidad Nacional de Colombia sede Medellín de diferentes carreras, es decir, la mayoría de los sujetos de la muestra son jóvenes entre los 18 y los 25 años, además decidimos que solamente aquellos que tenían la posibilidad de saber las estaturas de sus padres entraban a nuestra base de datos, ya que el proceso sería más arduo si tomamos datos donde nos faltan llenar valores en las celdas correspondientes.

3 Variables de respuesta:

En nuestro caso será la estatura del sujeto (Hombre o Mujer) para ajustar un modelo para predecir por medio de nuestras variables predictoras la estatura del sujeto.

4 Variable de Control:

En este caso tendremos 3 variables haciendo de este un modelo de R.L.M.

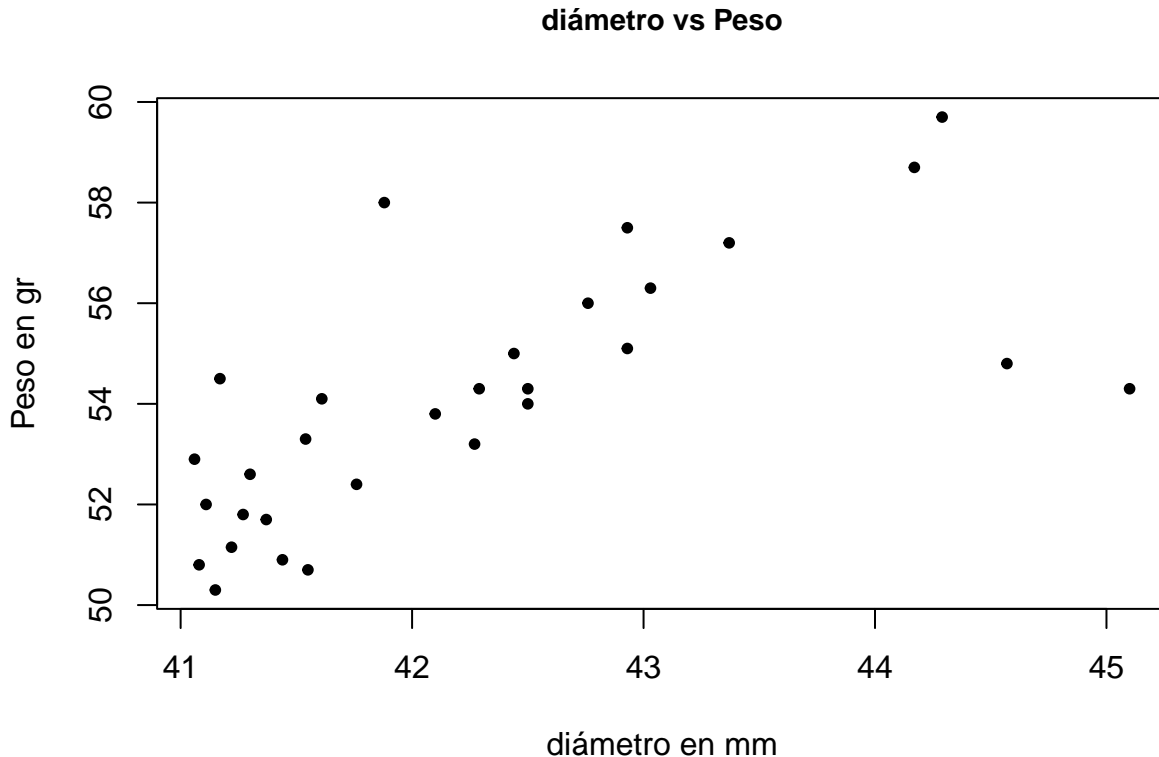
1. Estatura del Padre.
2. Estatura del Madre.
3. Género del sujeto.

5 Ajuste del modelo

Antes de observar o crear un modelo dado unos puntos, primero haremos un test de correlación de los datos para estudiar el grado de variación conjunta entre el diámetro y peso de los huevos:

```
##
## Pearson's product-moment correlation
##
## data:  datos$diametro and datos$peso
## t = 5.1667, df = 28, p-value = 1.758e-05
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.4521558 0.8459674
## sample estimates:
##          cor
## 0.6986211
```

Como podemos ver se rechaza H_0 ya que el $p - valor < 0.05$, es decir nuestros datos peso y el diámetro de los huevos tomados no tienen una correlación no significativa, más bien tienen una correlación positiva entre los datos, que vamos a ver por medio de un grafico de dispersión:

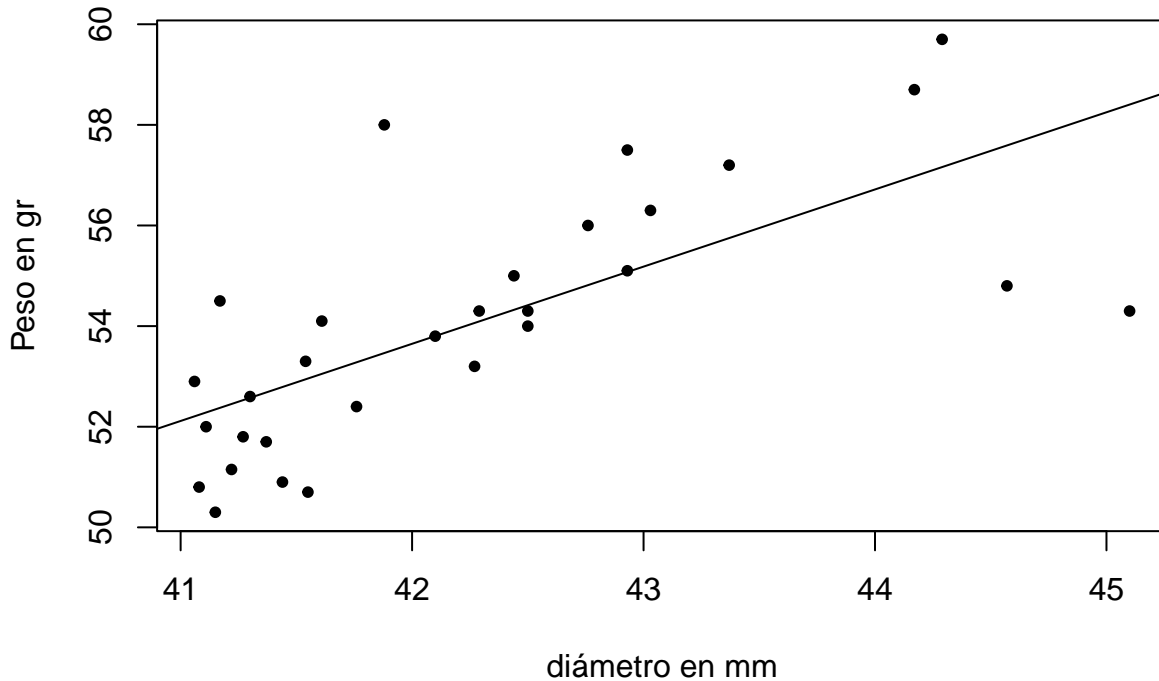


Como podemos ver los datos si tienen en algún grado una correlación positiva (mientras aumentan los valores del diámetro aumenta el peso) apriori.

Una vez visto que existe relación entre las variables pasamos a realizar el ajuste del modelo. Para ello usamos la función `lm()` que toma la forma:

```
##
## Call:
## lm(formula = peso ~ diametro, data = datos)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.1037 -0.9563  0.0118  1.0663  4.5359
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -10.7807    12.5511  -0.859   0.398
## diametro      1.5340     0.2969   5.167 1.76e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.81 on 28 degrees of freedom
## Multiple R-squared:  0.4881, Adjusted R-squared:  0.4698
## F-statistic: 26.7 on 1 and 28 DF,  p-value: 1.758e-05
```

Como podemos ver nuestro intercepto dado el $p - valor > 0.05$ decimos que el intercepto tengan valor a cero, esto es logico ya que seria extraño encontrar huevos con diámetro cero y un valor de peso inicial.



En primer lugar deseamos obtener los estimadores puntuales, errores estándar y p-valores asociados con cada coeficiente

```
##           Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) -10.78066 12.5510919 -0.858942 3.976665e-01
## diametro     1.53402  0.2969033  5.166734 1.758177e-05
```

El resultado del ajuste es:

(12.5510919) (0.2969033)

$$Peso = 0 + 1.53402 * diametro$$

donde los valores entre paréntesis indican los errores estándar de cada coeficiente. Además, puesto que los p-valores son mayores y menores a 0.05, podemos concluir que:

1. En este caso no tiene sentido analizar el valor de la constante para $\text{diámetro} = 0$, ya que pertenecería a un supuesto donde el huevo (imaginariamente) exista, de ahí que el peso del huevo para $\text{diámetro} = 0$ sea de 0, menor que cualquiera de los datos de nuestro conjunto, en conclusión un huevo con $\text{diámetro} = 0$ no existe y más si su peso = 0.
2. Existen evidencias estadísticas suficientes para considerar que hay una relación lineal entre diámetro y peso. Dicha relación es positiva cuando aumenta el diámetro de un huevo dado aumente el peso del mismo. Además vemos que por cada *mm* que aumenta el diámetro de un huevo, aumenta el peso en 1.53 *gramos*.
3. El error estándar residual estimado (*s*) es de 3.1. Este valor es muy importante, es un medidor de la calidad (precisión) del modelo. Además nos vamos a basar en él para calcular los intervalos de confianza para el coeficiente del modelo.

5.1 Intervalos de Confianza

Obtenemos los correspondientes intervalos de confianza para el parámetro m de nuestro modelo = $Y = 0 + m * X$ con nivel de significación al 95%

Parámetro	2.5 %	97.5 %
diametro	0.92584	2.142199

Interpretamos los intervalos: con una probabilidad del 95%, el efecto asociado con diametro se encuentra en el intervalo (0.9258416, 2.142199).

5.2 Tabla ANOVA

Obtenemos la correspondiente tabla ANOVA donde vemos la descomposición de la variabilidad del modelo

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
diametro	1	87.49255	87.492552	26.69514	1.76e-05
Residuals	28	91.76920	3.277471	NA	NA

Observamos que la variabilidad explicada por el modelo, $SSM=87.493$, es inferior a la que queda por explicar (residuos), $SSR=91.769$ y el estadístico $F=26.695$, mayor que 1. Además, volviendo a ver el resumen del modelo

F-statistic: 26.695 on 1 and 28 DF, p-value: 1.758e-05

tenemos que el p – *valor* asociado con el estadístico F es inferior a 0.05.

La conclusión es que hay evidencias suficientes para poder rechazar la hipótesis nula, $H_0 : F = 1$ y por tanto, resulta posible establecer un modelo de regresión lineal para explicar el comportamiento del peso de un huevo en función de su diámetro.

5.3 Coeficiente de determinación

En el `modelo1` el valor de R^2 es **Multiple R-squared: 0.4881**, alrededor del 48.81% de la variabilidad del peso es explicada por la recta ajustada.

6 Análisis de los parámetros del modelo

El test ANOVA significativo nos dice si el modelo tiene, en general, un grado de predicción significativamente bueno para la variable resultado, pero no nos dice nada sobre la contribución individual del modelo. Para encontrar los parámetros del modelo y su significación tenemos que volver a la parte **Coefficients** en el resumen del modelo.

β_1	Estimate	Std. Error	t value	$Pr(> t)$
diametro	1.53402	0.2969033	5.166734	1.7581e-05

Observando la tabla vemos que β_1 es la pendiente de la recta y representa el cambio en la variable dependiente (peso) asociado al cambio de una unidad en la variable predictora. Si nuestra variable predictora incrementa una unidad, nuestro modelo predice que el peso de un huevo se incrementara en 1.534 *gr*, pues en este caso $\beta_1 = 1.534$ Por tanto, la ecuación del modelo queda:

$$Y = 0 + 1.534 * X$$

7 Diagnóstico del modelo

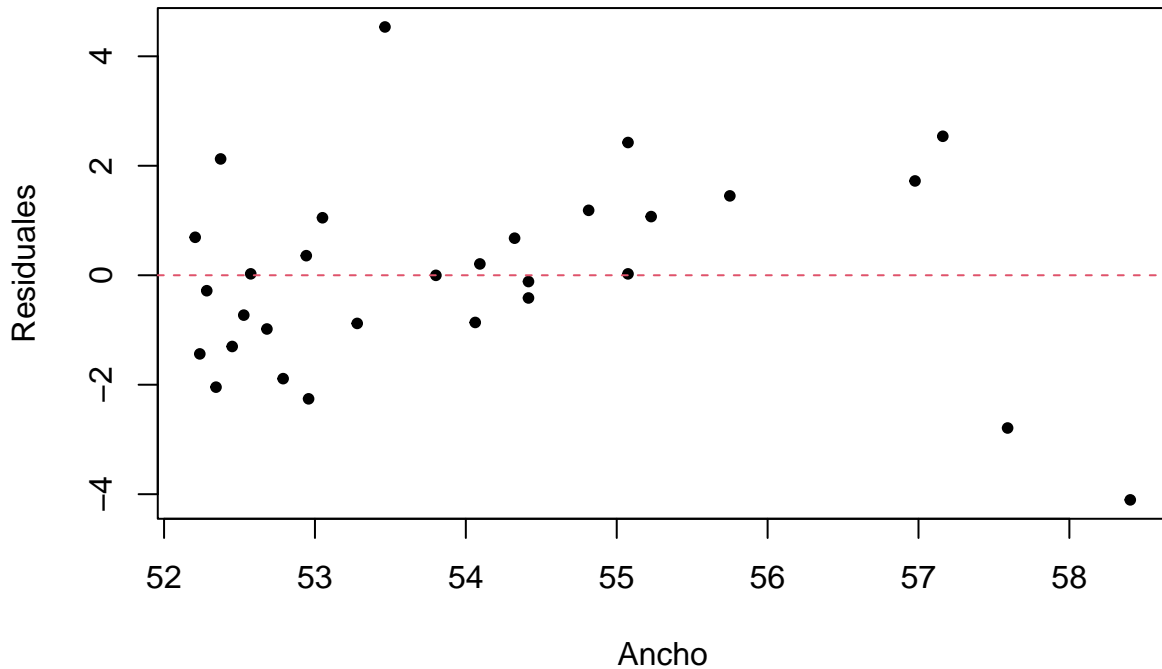
En este apartado hemos hecho uso tanto de J.Faraway (2009) como de Sánchez (2011) para el desarrollo del mismo.

Una vez que tenemos el modelo ajustado procedemos con su diagnóstico, que se realiza a través del **análisis de los residuos ε_i** :

- Las hipótesis de linealidad, homocedasticidad e independencia se contrastan a través de un análisis gráfico que enfrenta los valores de los residuos, ε_i , con los valores ajustados \hat{x}_i .
- Las hipótesis de media cero, varianza constante, incorrelación y normalidad la comprobamos analíticamente

Comenzaremos con el análisis gráfico. Los residuos deberían formar una nube de puntos sin estructura y con, aproximadamente, la misma variabilidad por todas las zonas como se muestra en el gráfico:

Residuales vs. valores ajustados



Continuamos ahora realizando el **diagnóstico analítico**. El primer paso es obtener los residuos, valores ajustados y estadísticos del modelo analizado para poder así estudiar si se cumplen los supuestos del mismo.

Obtención de residuos, valores ajustados y estadísticos necesarios

Para ello, añadimos los correspondientes resultados a nuestros datos a través del siguiente código:

El resultado es la creación de las siguientes variables:

- `fitted.modelo1`: valores ajustados (valores de la variable respuesta) para las observaciones originales de la predictora.
- `residuals.modelo1`: residuos del modelo, esto es, diferencia entre valor observado de la respuesta y valor ajustado por el modelo.
- `rstudent.modelo1`: residuos estudentizados del modelo ajustado.

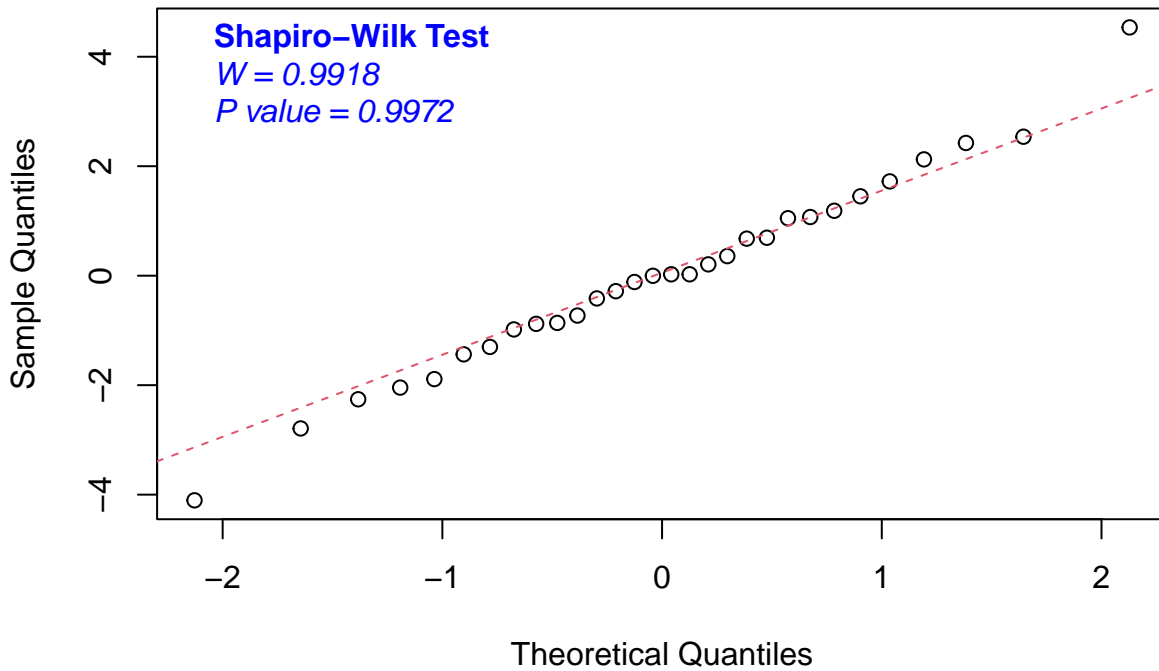
Vamos a utilizar todas estas variables para estudiar si nuestro modelo cumple las hipótesis.

7.1 Test de normalidad (test de Kolmogorov-Smirnov)

Empezamos el análisis con un gráfico `qqplot`, que enfrenta los valores reales a los valores que obtendríamos si la distribución fuera normal. Si los datos reales se distribuyen normalmente, estos tendrán la misma

distribución que los valores esperados y en el gráfico qqplot obtendremos una linea recta en la diagonal.

Normal Q-Q Plot of Residuals



Podemos ver que nuestro $p - \text{valor} > 0.05$ por lo que no se rechaza la hipótesis nula donde los datos se distribuyen normal. además por medio de la grafica nos muestra como nuestros puntos se acomodan bien a la recta.

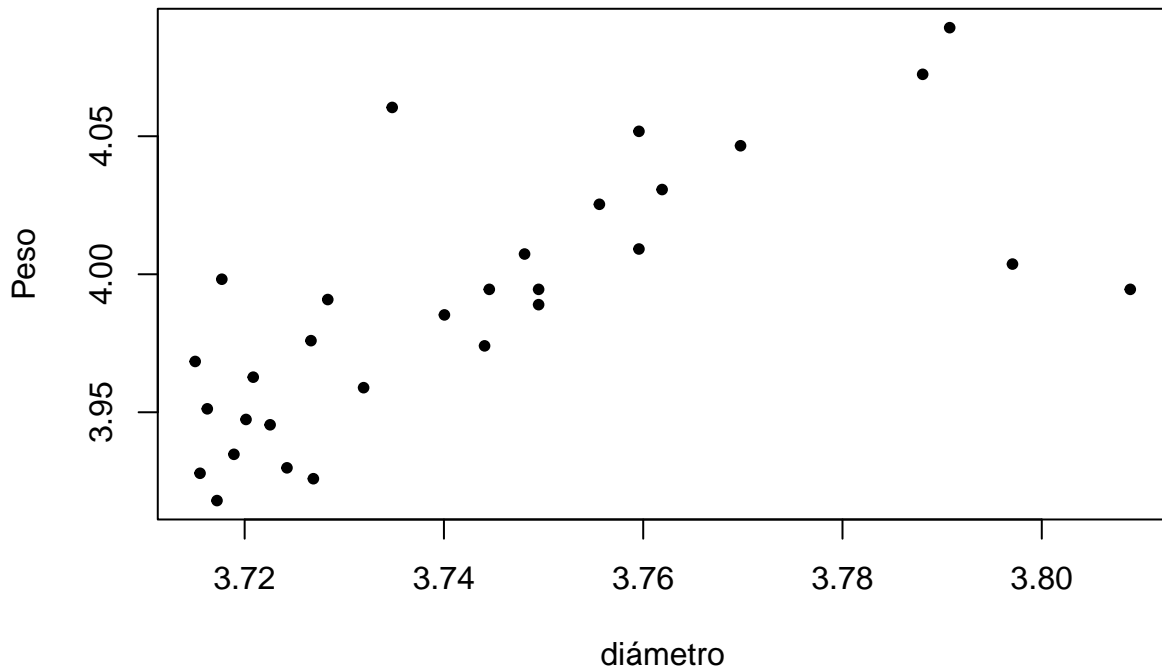
```
##
## studentized Breusch-Pagan test
##
## data:  modelo1
## BP = 4.5466, df = 1, p-value = 0.03298
```

No Existe homogeneidad pues la significación es menor de 0.05, la varianza no es constante a lo largo de la muestra.

8 Transformación del Modelo

Dado que nuestra varianza no es constante, tendremos que hacer una transformación en nuestro modelo, es decir: $\hat{Y}^* = \log(Y)$ $\hat{X}^* = \log(X)$

diámetro vs Peso (Escala Log)

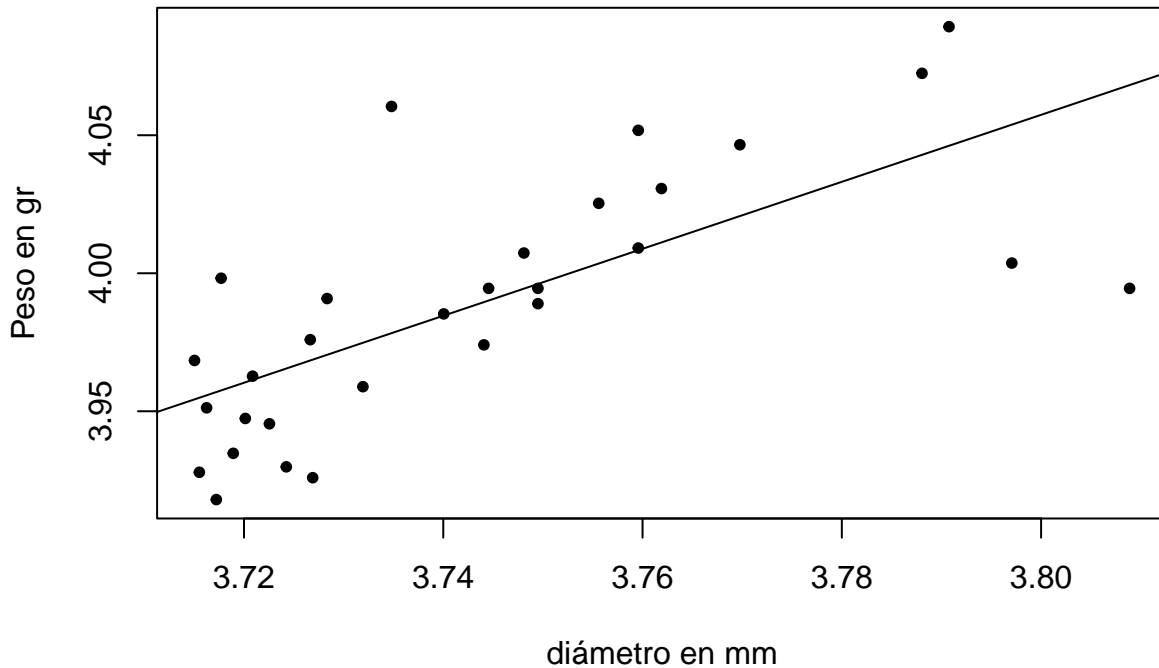


Como podemos ver los datos si tienen en algún grado una correlación positiva (mientras aumentan los valores del diámetro aumenta el peso) apriori.

Una vez visto que existe relación entre las variables pasamos a realizar el ajuste del modelo. Para ello usamos la función `lm(log())` (*log()* <- *logaritmo neperiano*) que toma la forma:

```
##
## Call:
## lm(formula = log(peso) ~ log(diametro), data = datos)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.073651 -0.017460  0.000711  0.020168  0.082143
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -0.5532     0.8645  -0.640   0.527
## log(diametro)  1.2133     0.2309   5.254 1.38e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.03292 on 28 degrees of freedom
## Multiple R-squared:  0.4965, Adjusted R-squared:  0.4785
## F-statistic: 27.61 on 1 and 28 DF, p-value: 1.385e-05
```

Como podemos ver nuestro intercepto dado el $p - \text{valor} > 0.05$ decimos que el intercepto tengan valor a cero, esto es logico ya que seria extraño encontrar huevos con diametro cero y un valor de peso inicial.



En primer lugar deseamos obtener los estimadores puntuales, errores estándar y p-valores asociados con cada coeficiente

```
##           Estimate Std. Error   t value    Pr(>|t|)
## (Intercept) -0.5531989  0.8644818 -0.6399197 5.274268e-01
## log(diametro)  1.2133149  0.2309251  5.2541481 1.384744e-05
```

El resultado del ajuste es:

(0.8644818) (0.2309251)

$$Peso = 0 + 1.2133149 * diametro$$

donde los valores entre paréntesis indican los errores estándar de cada coeficiente. Además, puesto que los p-valores son mayores y menores a 0.05, podemos concluir que:

1. En este caso no tiene sentido analizar el valor de la constante para diámetro = 0, ya que pertenecería a un supuesto donde el huevo (imaginariamente) exista, de ahí que el peso del huevo para diámetro = 0 sea de 0, menor que cualquiera de los datos de nuestro conjunto, en conclusión un huevo con diámetro = 0 no existe y más si su peso = 0.
2. Existen evidencias estadísticas suficientes para considerar que hay una relación lineal entre diámetro y peso. Dicha relación es positiva cuando aumenta el diámetro de un huevo dado aumente el peso del mismo. Además vemos que por cada *mm* que aumenta el diámetro de un huevo, aumenta el peso en 1.21 *gramos*.
3. El error estándar residual estimado (*s*) es de 0.001. Este valor es muy importante, es un medidor de la calidad (precisión) del modelo. Además nos vamos a basar en él para calcular los intervalos de confianza para el coeficiente del modelo.

8.1 Intervalos de Confianza

Obtenemos los correspondientes intervalos de confianza para el parámetro m de nuestro modelo =

$$Y^* = 0 + m * X^* \text{ con nivel significación al 95\%}$$

Parámetro	2.5 %	97.5 %
diametro	0.74028	1.686344

Interpretamos los intervalos: con una probabilidad del 95%, el efecto asociado con diametro se encuentra en el intervalo (0.7402862, 1.686344).

8.2 Tabla ANOVA

Obtenemos la correspondiente tabla ANOVA donde vemos la descomposición de la variabilidad del modelo

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
log(diametro)	1	0.0299145	0.0299145	27.60607	1.38e-05
Residuals	28	0.0303414	0.0010836	NA	NA

Observamos que la variabilidad explicada por el modelo, $SSM=0.029915$, es inferior a la que queda por explicar (residuos), $SSR=0.030341$ y el estadístico $F=27.606$, mayor que 1. Además, volviendo a ver el resumen del modelo

F-statistic: 27.606 on 1 and 28 DF, p-value: 1.385e-05

tenemos que el p – *valor* asociado con el estadístico F es inferior a 0.05.

La conclusión es que hay evidencias suficientes para poder rechazar la hipótesis nula, $H_0 : F = 1$ y por tanto, resulta posible establecer un modelo de regresión lineal para explicar el comportamiento del peso de un huevo en función de su diámetro.

8.3 Coeficiente de determinación

En el `modelo1` el valor de R^2 es **Multiple R-squared: 0.4965**, alrededor del 49.65% de la variabilidad del peso es explicada por la recta ajustada.

9 Análisis de los parámetros del modelo

El test ANOVA significativo nos dice si el modelo tiene, en general, un grado de predicción significativamente bueno para la variable resultado, pero no nos dice nada sobre la contribución individual del modelo. Para encontrar los parámetros del modelo y su significación tenemos que volver a la parte **Coefficients** en el resumen del modelo.

β_1	Estimate	Std. Error	t value	$Pr(> t)$
log(diametro)	1.2133149	0.2309251	5.2541481	1.384744e-05

Observando la tabla vemos que β_1 es la pendiente de la recta y representa el cambio en la variable dependiente (peso) asociado al cambio de una unidad en la variable predictora. Si nuestra variable predictora incrementa una unidad, nuestro modelo predice que el peso de un huevo se incrementara en 1.213 gr, pues en este caso $\beta_1 = 1.213$ Por tanto, la ecuación del modelo queda:

$$Y = 0 + 1.213 * X$$

10 Diagnóstico del modelo

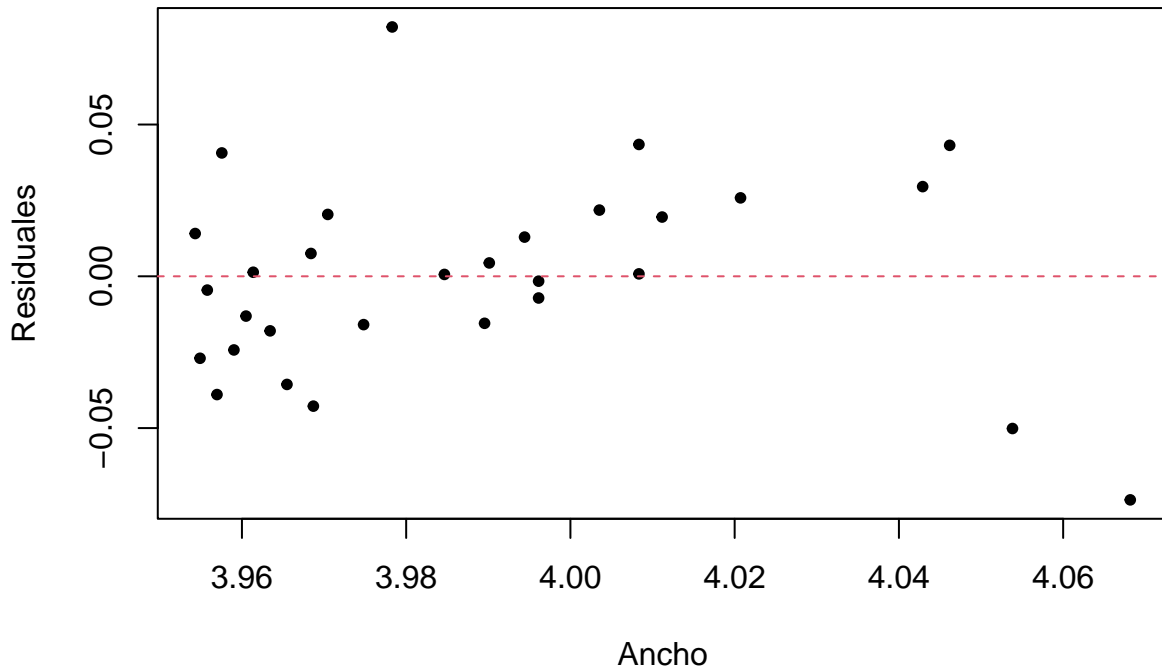
En este apartado hemos hecho uso tanto de *J.Faraway (2009)* como de *Sánchez (2011)* para el desarrollo del mismo.

Una vez que tenemos el modelo ajustado procedemos con su diagnóstico, que se realiza a través del **análisis de los residuos ε_i** :

- Las hipótesis de linealidad, homocedasticidad e independencia se contrastan a través de un análisis gráfico que enfrenta los valores de los residuos, ε_i , con los valores ajustados \hat{x}_i .
- Las hipótesis de media cero, varianza constante, incorrelación y normalidad la comprobamos analíticamente

Comenzaremos con el análisis gráfico. Los residuos deberían formar una nube de puntos sin estructura y con, aproximadamente, la misma variabilidad por todas las zonas como se muestra en el gráfico:

Residuales vs. valores ajustados



Continuamos ahora realizando el **diagnóstico analítico**. El primer paso es obtener los residuos, valores ajustados y estadísticos del modelo analizado para poder así estudiar si se cumplen los supuestos del mismo.

Obtención de residuos, valores ajustados y estadísticos necesarios

Para ello, añadimos los correspondientes resultados a nuestros datos a través del siguiente código:

El resultado es la creación de las siguientes variables:

- `fitted.modelo2`: valores ajustados (valores de la variable respuesta) para las observaciones originales de la predictora.
- `residuals.modelo2`: residuos del modelo, esto es, diferencia entre valor observado de la respuesta y valor ajustado por el modelo.
- `rstudent.modelo2`: residuos estudentizados del modelo ajustado.

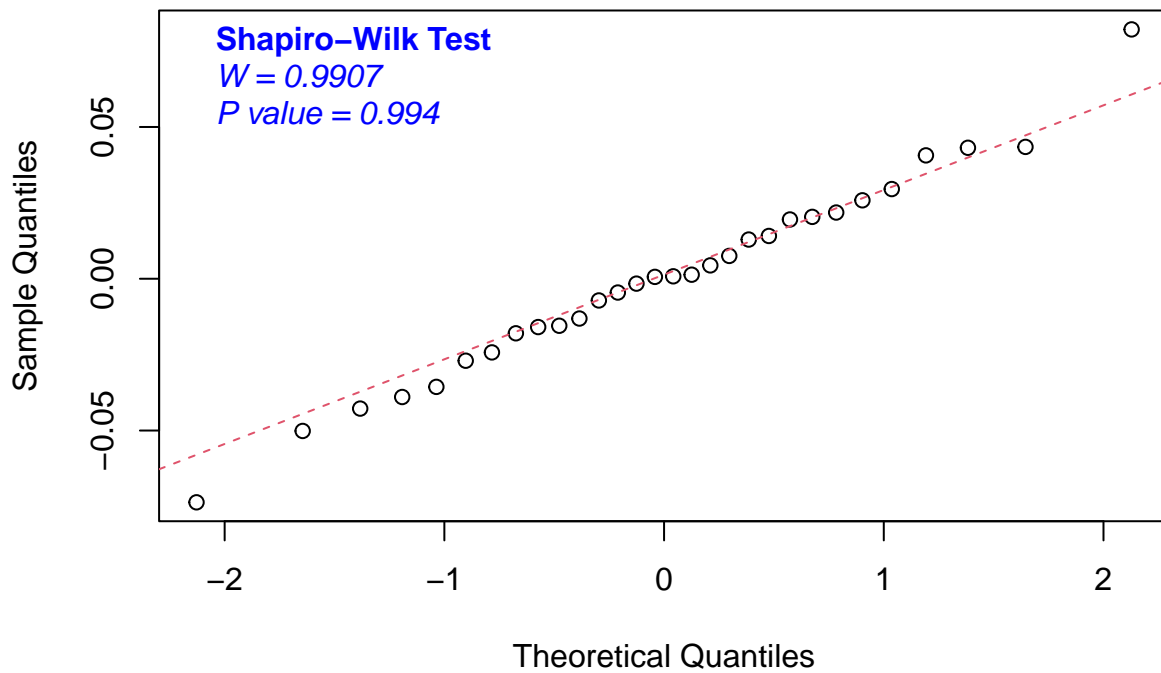
Vamos a utilizar todas estas variables para estudiar si nuestro modelo cumple las hipótesis.

10.1 Test de normalidad (test de Kolmogorov-Smirnov)

Empezamos el análisis con un gráfico `qqplot`, que enfrenta los valores reales a los valores que obtendríamos si la distribución fuera normal. Si los datos reales se distribuyen normalmente, estos tendrán la misma

distribución que los valores esperados y en el gráfico qqplot obtendremos una linea recta en la diagonal.

Normal Q-Q Plot of Residuals

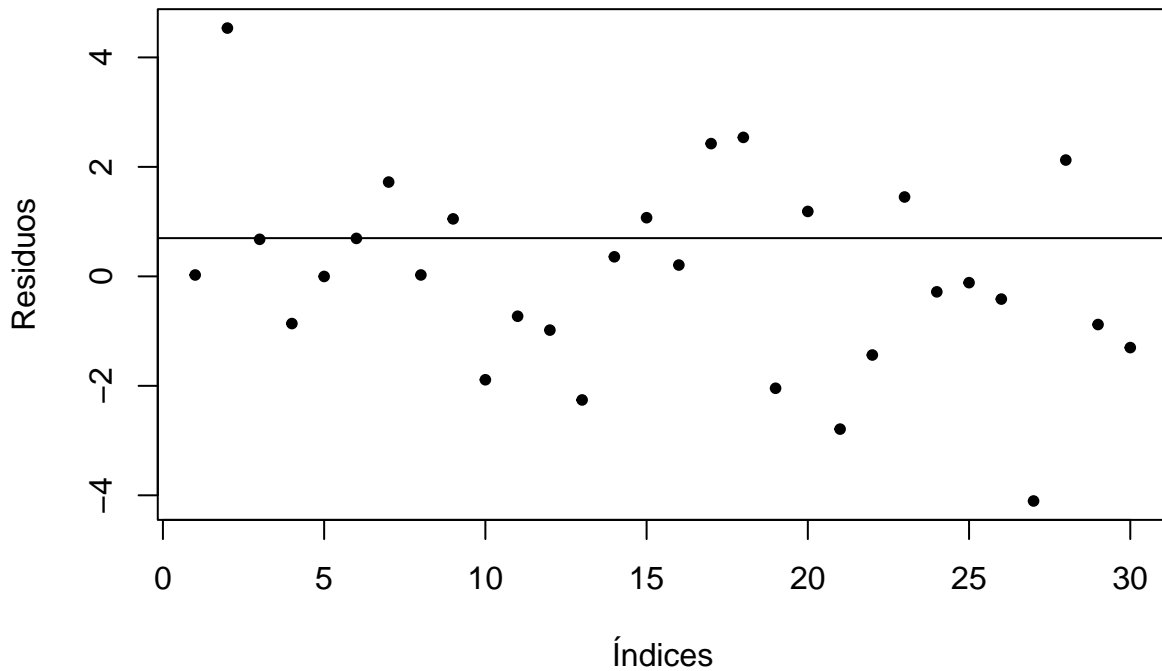


Aqui podemos ver que tambien se cumple el supuesto de normalidad, dejandonos con la siguiente pregunta:
¿ya que nuestros datos están transformados será que nuestra varianza será constante?

```
##  
## studentized Breusch-Pagan test  
##  
## data: modelo2  
## BP = 3.7084, df = 1, p-value = 0.05414
```

Existe homogeneidad pues la significación es mayor de 0.05, la varianza es constante a lo largo de la muestra.

10.2 Autocorrelación (test de Durbin-Watson)



Si hubiera una correlación seria, veríamos picos más largos de residuos por encima y por debajo de la línea de correlación. A menos que estos efectos sean fuertes, puede ser difícil de detectar la autocorrelación, por ello realizamos el contraste de Durbin-Watson.

```
##  
## Durbin-Watson test  
##  
## data: peso ~ diametro  
## DW = 2.0779, p-value = 0.7969  
## alternative hypothesis: true autocorrelation is not 0
```

En el contraste de autocorrelación también aceptamos la hipótesis nula de que no existe correlación entre los residuos con un p – valor superior a 0.05.

11 Predicción

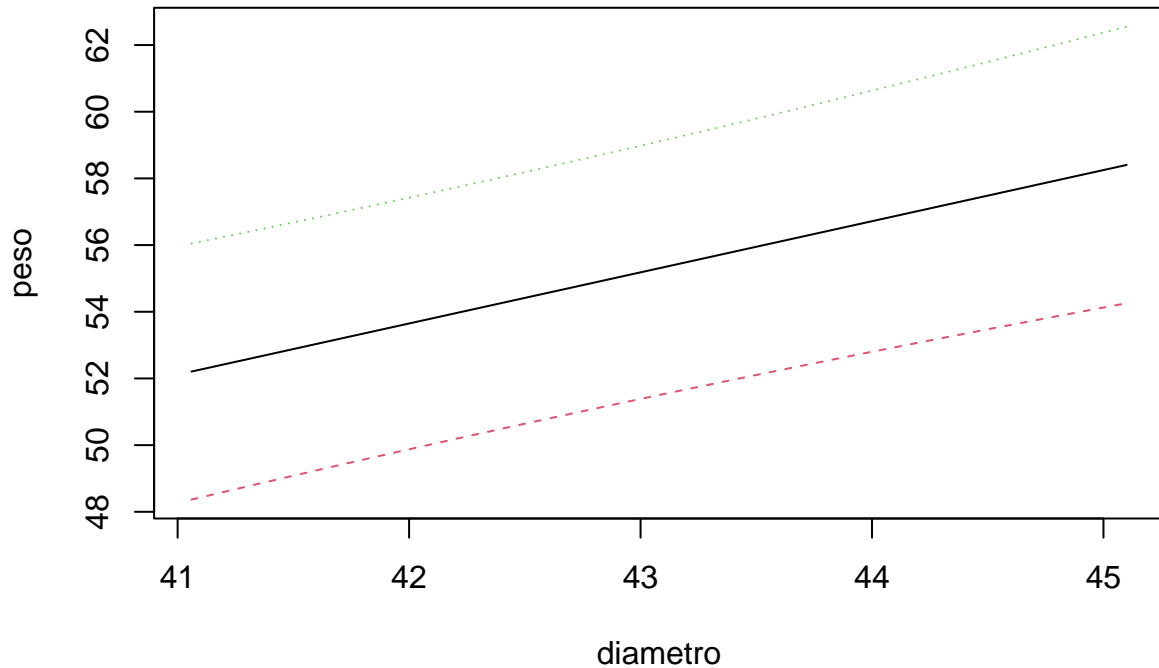
Tenemos un modelo de regresión con la capacidad de relacionar la variable predictora y la variable dependiente. Podemos utilizarlo ahora para predecir eventos futuros de la variable dependiente a través de nuevos valores de la variable predictora.

Para ello debe verificarse alguna de las siguientes condiciones:

- el valor de la predictora está dentro del rango de la variable original.
- si el valor de la predictora está fuera del rango de la original, debemos asegurar que los valores futuros mantendrán el modelo lineal propuesto.

11.1 Predicción de nuevas observaciones

Dibujamos las bandas de predicción, que reflejan la incertidumbre sobre futuras observaciones:



Como nuestros datos estan escaladaos a el `log()` vamos a destransformarlos para obtener nuestras predicciones:

aqui podemos ver en `fit` nuestro diametro ajustado y `lwr` & `upr` nuestros intervalos de prediccion para estos anchos de un huevo dado.

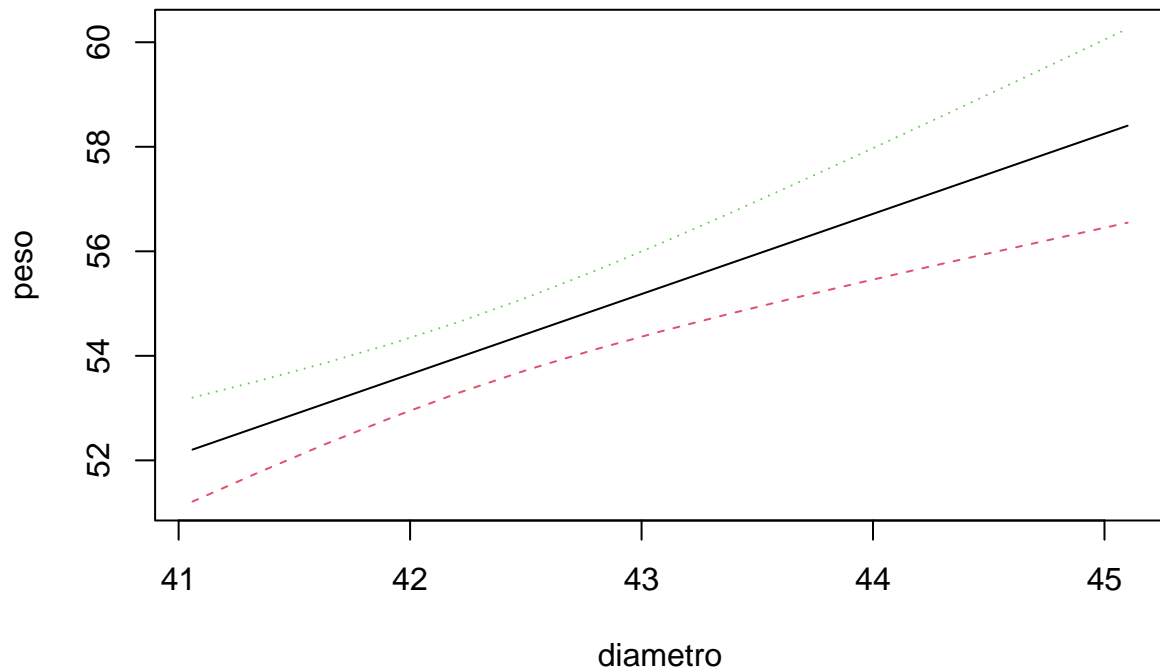
11.2 Intervalos de confianza para los predictores:

Dado un nuevo conjunto de predictores, `x0`, debemos evaluar la incertidumbre en esta prediccion. Para tomar decisiones racionales necesitamos algo más que puntos estimados. Si la prediccion tiene intervalo de confianza ancho entonces entonces los resultados estarán lejos de la estimación puntual.

Nota: Las bandas de confianza reflejan la incertidumbre en la línea de regresión (lo bien que la línea está calculada).

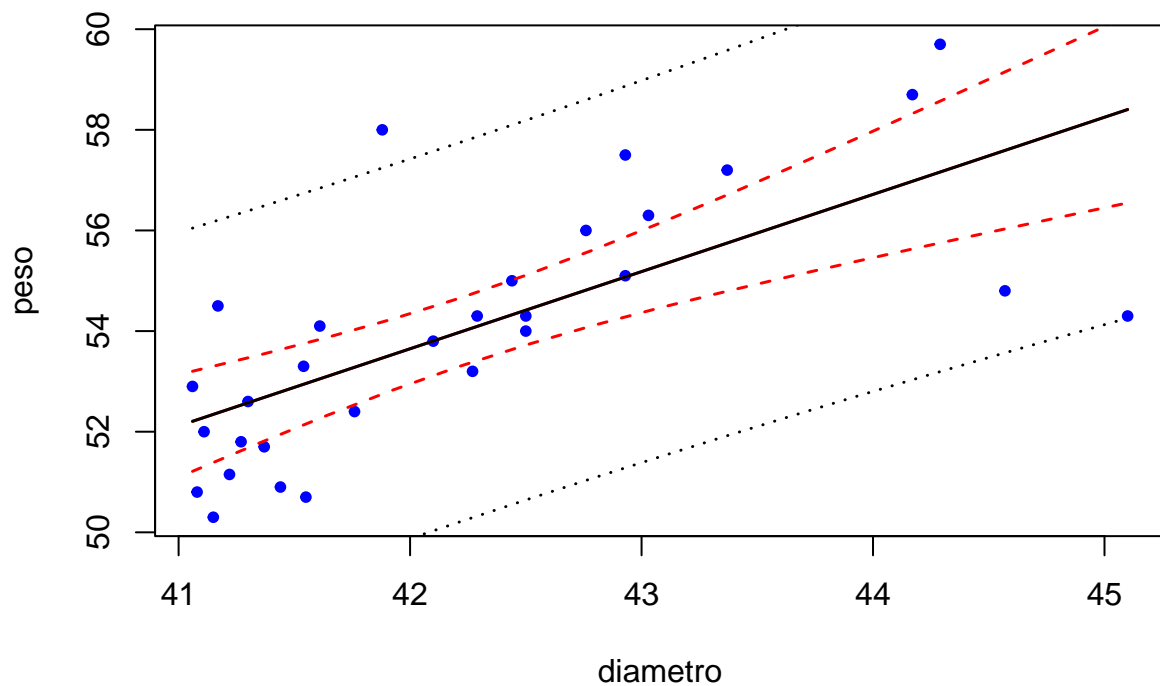
```
##      fit      lwr      upr
## 1 52.20622 51.21131 53.20113
## 2 52.64890 51.77438 53.52341
## 3 53.09157 52.31614 53.86700
## 4 53.53424 52.82755 54.24093
## 5 53.97692 53.29932 54.65451
## 6 54.41959 53.72644 55.11275
```

Dibujamos las bandas de confianza, que además reflejan la incertidumbre sobre futuras observaciones:



Por último podemos hacer un gráfico con la nube de puntos y los dos bandas, la de confianza y la de predicción (*Ferrari & Head, 2010*).

R.L.S. Peso vs Diámetro (IC's & IP's)



Donde:

- Las líneas rojas son Intervalos de Confianza.
- Los puntos azules, son valores predichos.
- y las líneas punteadas negras son el Intervalo de predicción, junto con la línea del modelo quitando la

transformación del `log()`.

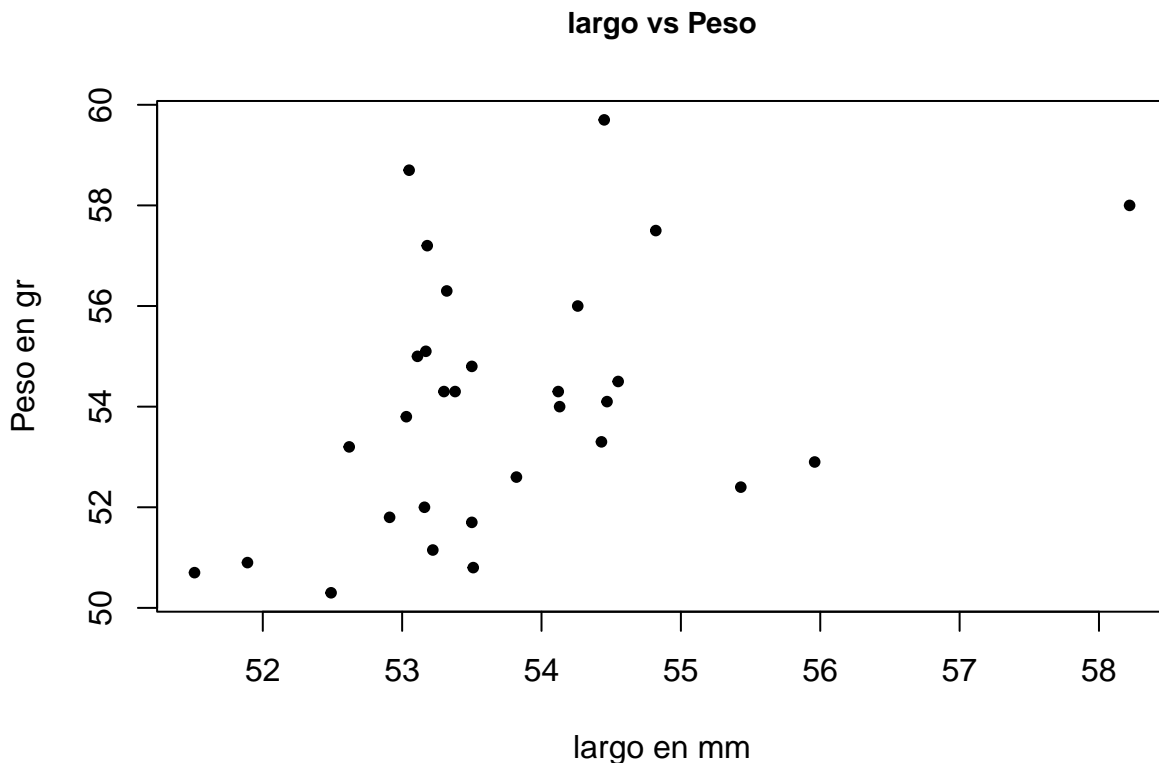
12 Modelo N°2 (Peso~Altura)

12.1 Ajuste del modelo:

Antes de observar o crear un modelo dado unos puntos, primero haremos un test de correlación de los datos para estudiar el grado de variación conjunta entre el diámetro y peso de los huevos:

```
##  
## Pearson's product-moment correlation  
##  
## data:  datos$largo and datos$peso  
## t = 2.4424, df = 28, p-value = 0.02116  
## alternative hypothesis: true correlation is not equal to 0  
## 95 percent confidence interval:  
##  0.06926414 0.67711427  
## sample estimates:  
##      cor  
## 0.4190759
```

Como podemos ver se rechaza H_0 ya que el p -valor < 0.05 , es decir nuestros datos peso y el largo de los huevos tomados tienen una correlación significativa, es decir, una correlación positiva entre los datos, que vamos a ver por medio de un grafico de dispersión:



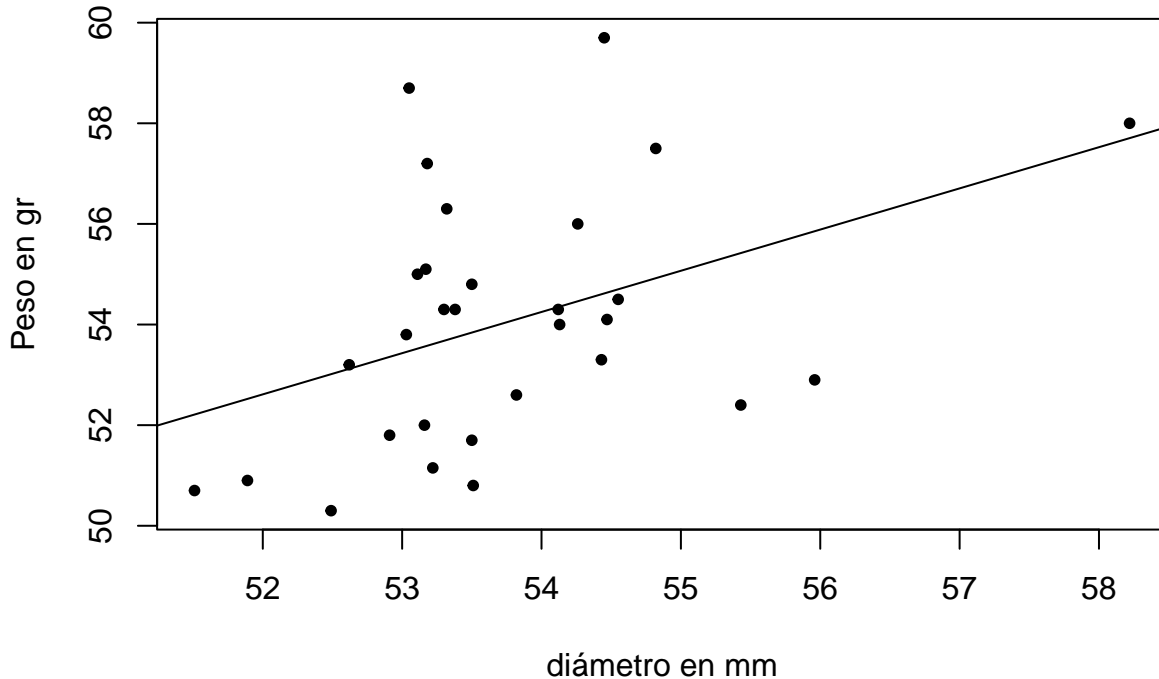
Como podemos ver los datos si tienen en algún grado una correlación positiva (mientras aumentan los valores del diámetro aumenta el peso) apriori.

Una vez visto que existe relación entre las variables pasamos a realizar el ajuste del modelo. Para ello usamos la función `lm()` que toma la forma:

```
##
```

```
## Call:
## lm(formula = peso ~ largo, data = datos)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.0482 -1.5604 -0.1238  1.3495  5.2285
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  10.0310    18.0260   0.556  0.5823
## largo         0.8189     0.3353   2.442  0.0212 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.297 on 28 degrees of freedom
## Multiple R-squared:  0.1756, Adjusted R-squared:  0.1462
## F-statistic: 5.965 on 1 and 28 DF,  p-value: 0.02116
```

Como podemos ver nuestro intercepto dado el $p - valor > 0.05$ decimos que el intercepto tengan valor a cero, esto es logico ya que seria extraño encintrar huevos con altura = cero y un valor de peso inicial.



En primer lugar deseamos obtener los estimadores puntuales, errores estándar y p-valores asociados con cada coeficiente

```
##              Estimate Std. Error  t value  Pr(>|t|)
## (Intercept) 10.0309793 18.0260060 0.5564726 0.58230836
## largo        0.8188604  0.3352747 2.4423569 0.02116101
```

El resultado del ajuste es:

(18.0260060) (0.3352747)

$$Peso = 0 + 0.8188604 * largo$$

donde los valores entre paréntesis indican los errores estándar de cada coeficiente. Además, puesto que los

p-valores son mayores y menores a 0.05, podemos concluir que:

1. En este caso no tiene sentido analizar el valor de la constante para largo = 0, ya que pertenecería a un supuesto donde el huevo (imaginariamente) exista, de ahí que el peso del huevo para largo = 0 sea de 0, menor que cualquiera de los datos de nuestro conjunto, en conclusión un huevo con largo = 0 no existe y más si su peso = 0.
2. Existen evidencias estadísticas suficientes para considerar que hay una relación lineal entre largo y peso. Dicha relación es positiva cuando aumenta el largo de un huevo dado aumente el peso del mismo. Además vemos que por cada *mm* que aumenta el diámetro de un huevo, aumenta el peso en 0.82 *gramos*.
3. El error estándar residual estimado (s) es de 4.92. Este valor es muy importante, es un medidor de la calidad (precisión) del modelo. Además nos vamos a basar en él para calcular los intervalos de confianza para el coeficiente del modelo.

12.2 Intervalos de Confianza

Obtenemos los correspondientes intervalos de confianza para el parámetro m de nuestro modelo = $Y = 0 + m * X$ con nivel significación al 95%

Parámetro	2.5 %	97.5 %
diametro	0.13208	1.505639

Interpretamos los intervalos: con una probabilidad del 95%, el efecto asociado con diametro se encuentra en el intervalo (0.1320814, 1.505639).

12.3 Tabla ANOVA

Obtenemos la correspondiente tabla ANOVA donde vemos la descomposición de la variabilidad del modelo

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
largo	1	31.48277	31.482768	5.965107	0.021161
Residuals	28	147.77898	5.277821	NA	NA

Observamos que la variabilidad explicada por el modelo, $SSM=31.4828$, es inferior a la que queda por explicar (residuos), $SSR=5.2778$ y el estadístico $F=5.9651$, mayor que 1. Además, volviendo a ver el resumen del modelo

F-statistic: 5.9651 on 1 and 28 DF, p-value: 0.02116

tenemos que el p -valor asociado con el estadístico F es inferior a 0.05.

La conclusión es que hay evidencias suficientes para poder rechazar la hipótesis nula, $H_0 : F = 1$ y por tanto, resulta posible establecer un modelo de regresión lineal para explicar el comportamiento del peso de un huevo en función de su diámetro.

12.4 Coeficiente de determinación

En el `modelo3` el valor de R^2 es **Multiple R-squared: 0.1756**, alrededor del 17.56% de la variabilidad del peso es explicada por la recta ajustada.

13 Análisis de los parámetros del modelo

El test ANOVA significativo nos dice si el modelo tiene, en general, un grado de predicción significativamente bueno para la variable resultado, pero no nos dice nada sobre la contribución individual

del modelo. Para encontrar los parámetros del modelo y su significación tenemos que volver a la parte **Coefficients** en el resumen del modelo.

β_1	Estimate	Std. Error	t value	$Pr(> t)$
largo	0.8188604	0.3352747	2.4423569	0.02116101

Observando la tabla vemos que β_1 es la pendiente de la recta y representa el cambio en la variable dependiente (peso) asociado al cambio de una unidad en la variable predictora. Si nuestra variable predictora incrementa una unidad, nuestro modelo predice que el peso de un huevo se incrementara en 0.82 gr, pues en este caso $\beta_1 = 0.82$ Por tanto, la ecuación del modelo queda:

$$Y = 0 + 0.82 * X$$

14 Diagnóstico del modelo

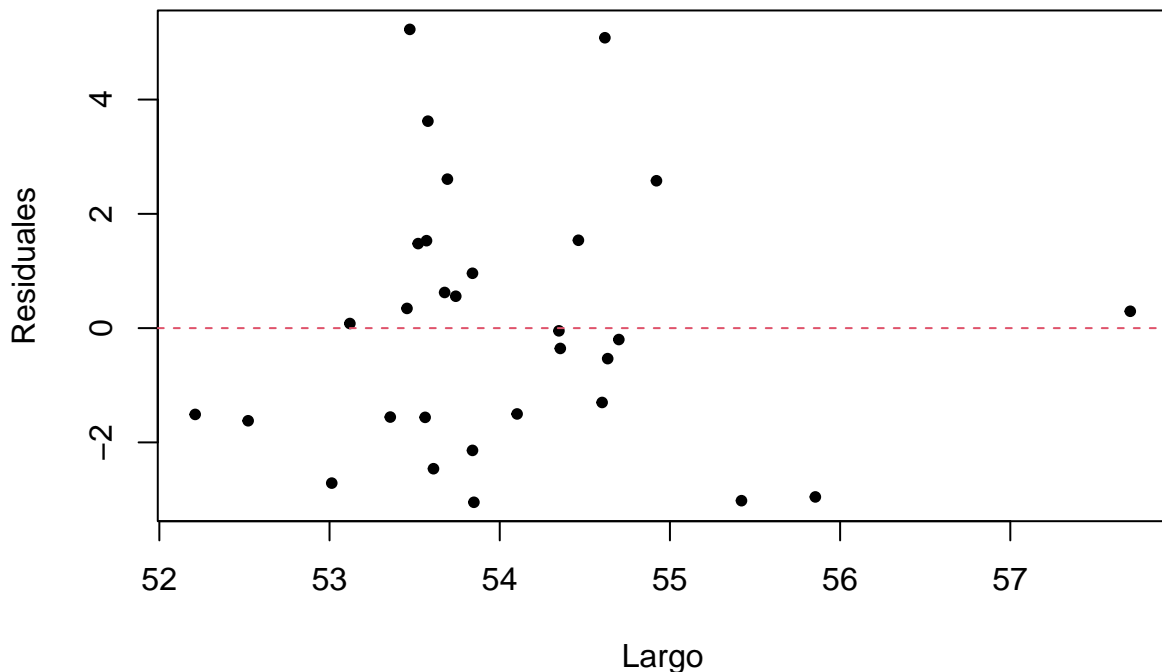
En este apartado hemos hecho uso tanto de J.Faraway (2009) como de Sánchez (2011) para el desarrollo del mismo.

Una vez que tenemos el modelo ajustado procedemos con su diagnóstico, que se realiza a través del **análisis de los residuos ε_i** :

- Las hipótesis de linealidad, homocedasticidad e independencia se contrastan a través de un análisis gráfico que enfrenta los valores de los residuos, ε_i , con los valores ajustados \hat{x}_i .
- Las hipótesis de media cero, varianza constante, incorrelación y normalidad la comprobamos analíticamente

Comenzaremos con el análisis gráfico. Los residuos deberían formar una nube de puntos sin estructura, para esto evaluaremos el siguiente grafico:

Residuales vs. valores ajustados



Continuamos ahora realizando el **diagnóstico analítico**. El primer paso es obtener los residuos, valores ajustados y estadísticos del modelo analizado para poder así estudiar si se cumplen los supuestos del mismo.

Obtención de residuos, valores ajustados y estadísticos necesarios

Para ello, añadimos los correspondientes resultados a nuestros datos a través del siguiente código:

El resultado es la creación de las siguientes variables:

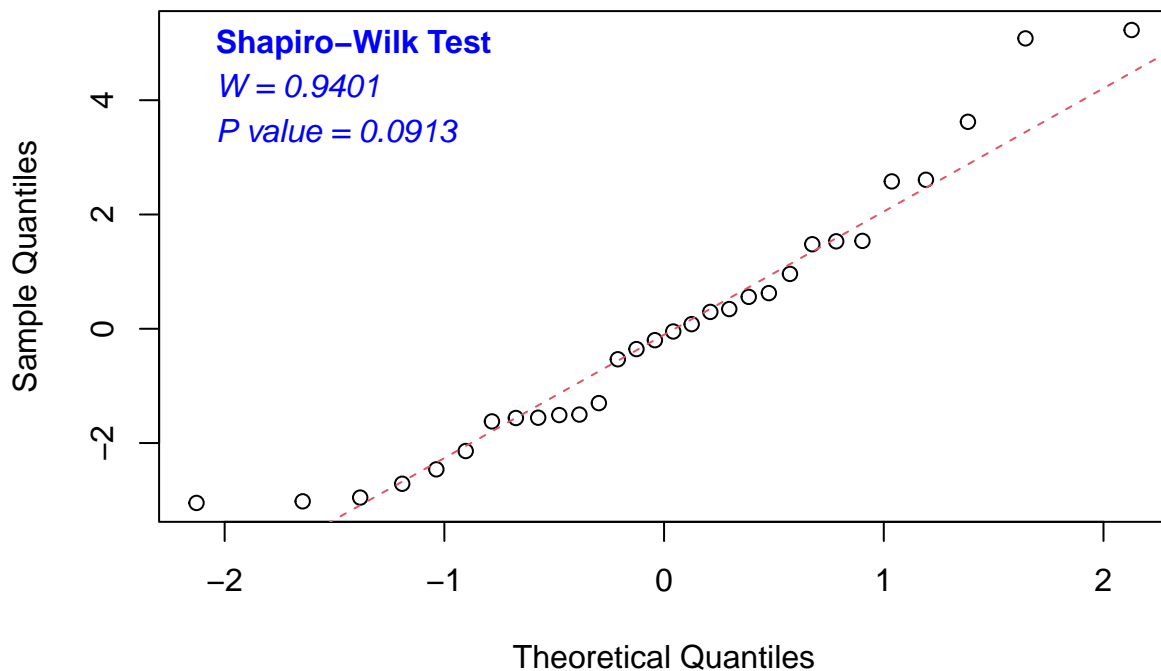
- `fitted.modelo3`: valores ajustados (valores de la variable respuesta) para las observaciones originales de la predictora.
- `residuals.modelo3`: residuos del modelo, esto es, diferencia entre valor observado de la respuesta y valor ajustado por el modelo.
- `rstudent.modelo3`: residuos estudentizados del modelo ajustado.

Vamos a utilizar todas estas variables para estudiar si nuestro modelo cumple las hipótesis.

14.1 Test de normalidad (test de Kolmogorov-Smirnov)

Empezamos el análisis con un gráfico `qqplot`, que enfrenta los valores reales a los valores que obtendríamos si la distribución fuera normal. Si los datos reales se distribuyen normalmente, estos tendrán la misma distribución que los valores esperados y en el gráfico `qqplot` obtendremos una línea recta en la diagonal.

Normal Q-Q Plot of Residuals



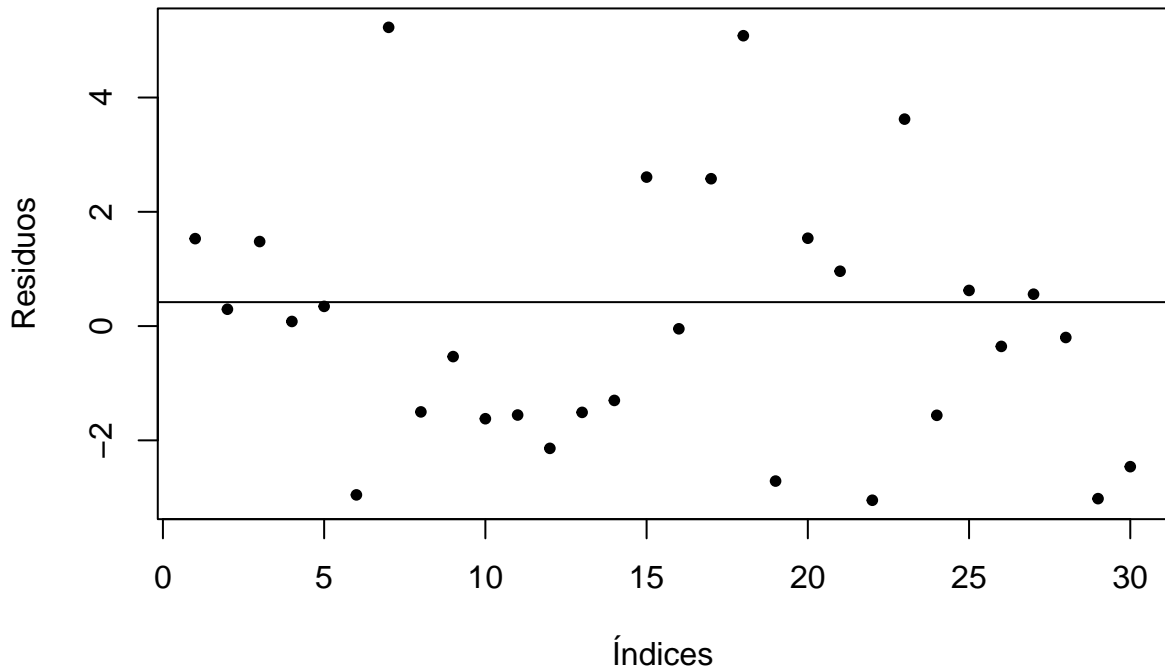
Podemos ver que nuestro $p - \text{valor} > 0.05$ por lo que no se rechaza la hipótesis nula donde los datos se distribuyen normal. además por medio de la grafica nos muestra como nuestros puntos se acomodan bien a la recta.

```
##
## studentized Breusch-Pagan test
##
## data:  modelo3
## BP = 0.00023802, df = 1, p-value = 0.9877
```

Existe homogeneidad pues la significación es mayor de 0.05, la varianza es constante a lo largo de la muestra.

14.2 Autocorrelación (test de Durbin-Watson)

Hemos asumido que los residuos son incorrelados, vamos a comprobarlo.



Si hubiera una correlación seria, veríamos picos más largos de residuos por encima y por debajo de la línea de correlación. A menos que estos efectos sean fuertes, puede ser difícil de detectar la autocorrelación, por ello realizamos el *contraste de Durbin-Watson*.

```
##
## Durbin-Watson test
##
## data: peso ~ largo
## DW = 2.3583, p-value = 0.3051
## alternative hypothesis: true autocorrelation is not 0
```

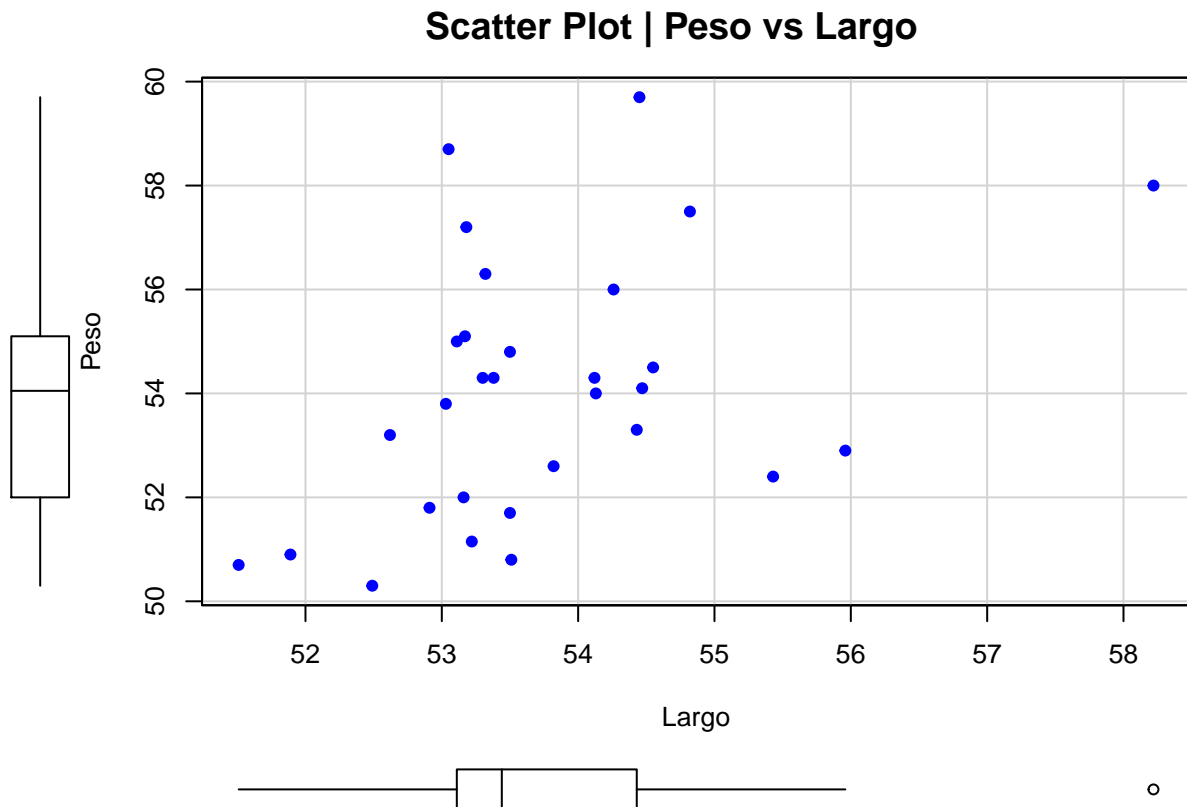
En el contraste de autocorrelación también aceptamos la hipótesis nula de que no existe correlación entre los residuos con un p – *valor* superior a 0.05

14.3 Valores atípicos:

Para observar si hay datos atípicos que hacen que nuestro modelo no sea el más óptimo. Vamos a realizar un test de valores atípicos (Bonferroni).

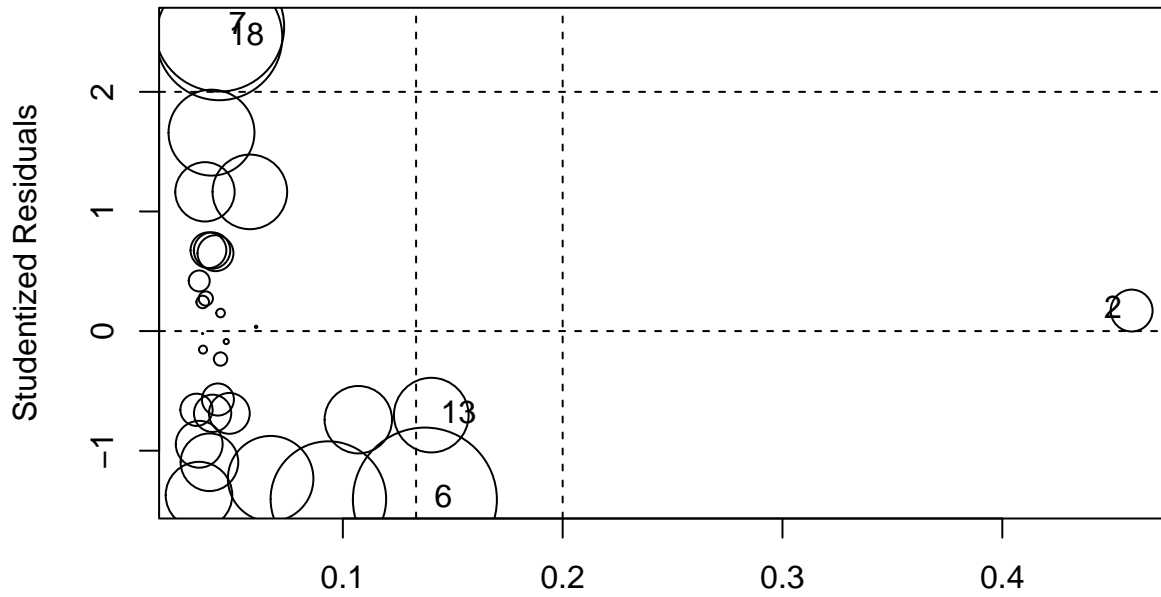
En el caso de observar valores atípicos los pasos a seguir son:

- Descartar que sea un error.
- Analizar si es un caso influyente.
- En caso de ser influyente calcular las rectas de regresión incluyéndolo y excluyéndolo, y elegir la que mejor se adapte al problema y a las observaciones futuras.



En la variable `largo` vemos que la mediana no está centrada en la media y contiene un valor atípico en la muestra, los datos no son uniformes. Con la variable `peso` ocurre casi mismo a excepción de ese dato atípico.

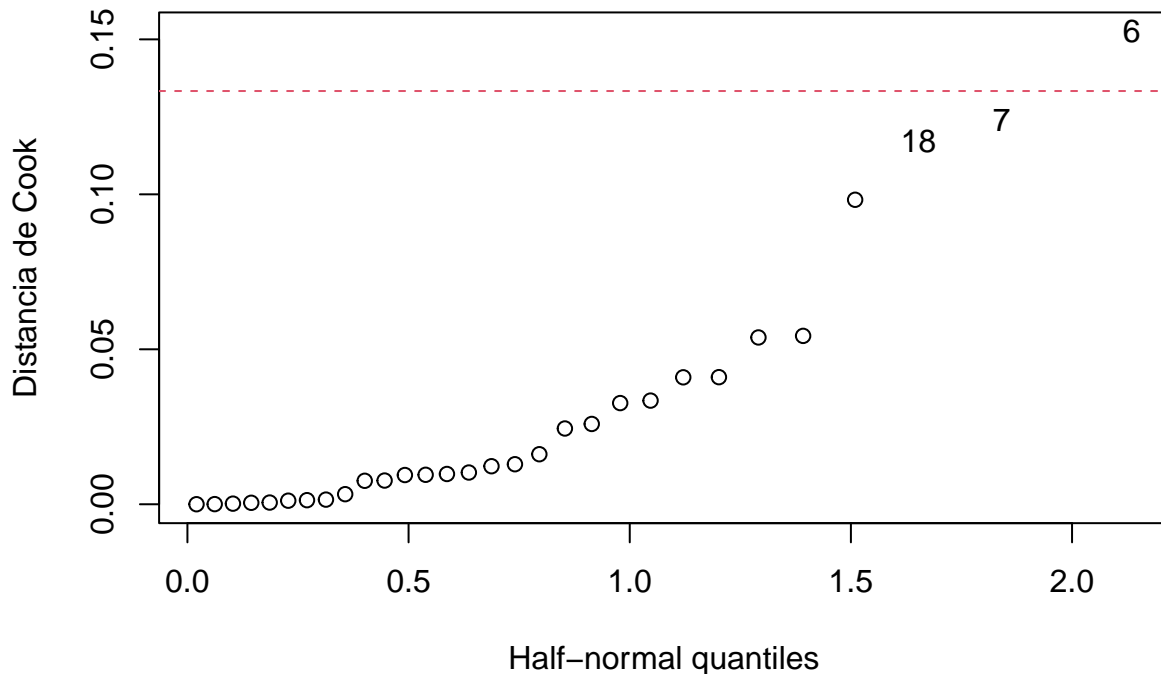
```
## No Studentized residuals with Bonferroni p < 0.05
## Largest |rstudent|:
##   rstudent unadjusted p-value Bonferroni p
## 7 2.544821      0.016962      0.50886
```

Hat-Values

	StudRes	Hat	CookD
2	0.1714820	0.4588309	0.0129136
6	-1.4087161	0.1373254	0.1525856
7	2.5448208	0.0437795	0.1239995
13	-0.7026427	0.1402320	0.0410043
18	2.4573114	0.0437596	0.1170944

El gráfico nos indican que la observación número **6** es un valor influyente dado que $D_i > \frac{4}{n}$ donde D_i es la distancia de cook y n es el tamaño de la muestra. Las observaciones **6,13** y **18** que vemos en el gráfico son medidas influyentes pero solamente se observa en el test que no hay datos atipicos pero por el $|x|$ de los residuales estandarizados se sabe que la observación **7** podria serlo, con un grafico de las distancias de cook miraremos si se cumple lo anterior.



La linea roja representa el corte donde se puede observar si un dato influye o no dentro del modelo, pero con el grafico de influencia nos muestra que la observación #6 no esta alejada de los datos, por lo cual concluimos que no hay datos extraños o atípicos que afecten el modelo.

15 Predicción

Tenemos un modelo de regresión con la capacidad de relacionar la variable predictora y la variable dependiente. Podemos utilizarlo ahora para predecir eventos futuros de la variable dependiente a través de nuevos valores de la variable predictora.

Para ello debe verificarse alguna de las siguientes condiciones:

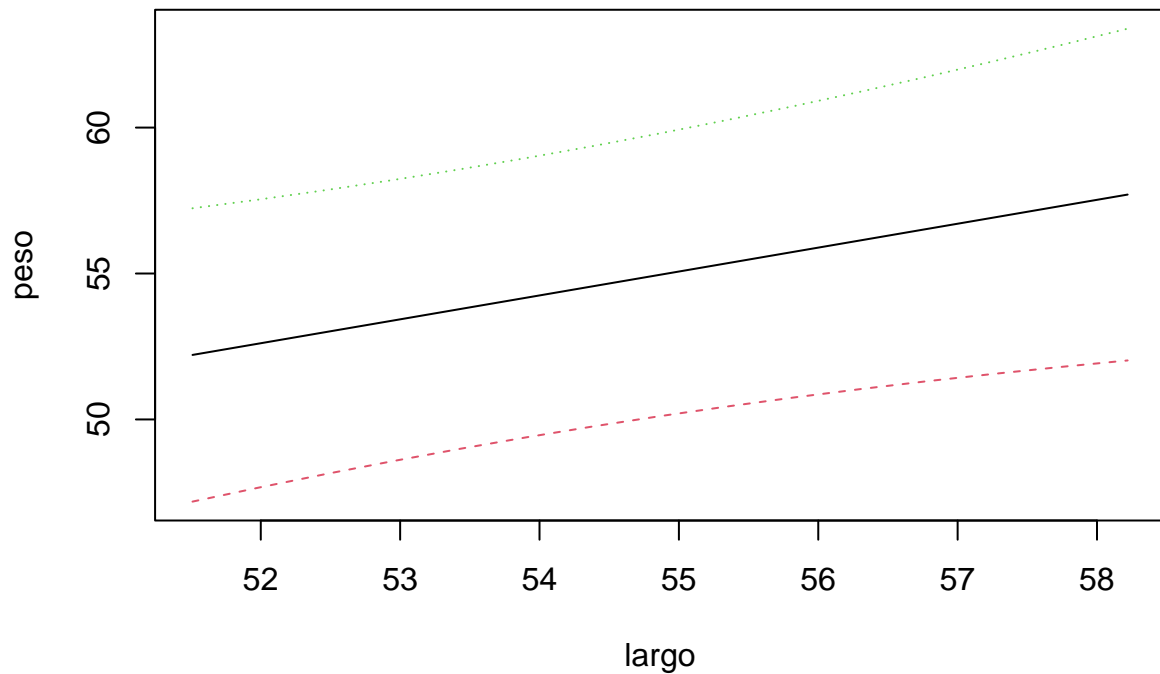
1. el valor de la predictora está dentro del rango de la variable original.
2. si el valor de la predictora está fuera del rango de la original, debemos asegurar que los valores futuros mantendrán el modelo lineal propuesto.

15.1 Predicción de nuevas observaciones

```
##          fit      lwr      upr
## 1 52.21048 47.18543 57.23553
## 2 52.60295 47.66873 57.53717
## 3 52.99542 48.13140 57.85943
## 4 53.38788 48.57254 58.20323
## 5 53.78035 48.99151 58.56920
## 6 54.17282 49.38792 58.95772
```

aquí podemos ver según un largo dado en mm de un huevo caul sera su intervalod e predicción de su peso en las dos columnas lower & upper.

Dibujamos las bandas de predicción, que reflejan la incertidumbre sobre futuras observaciones:

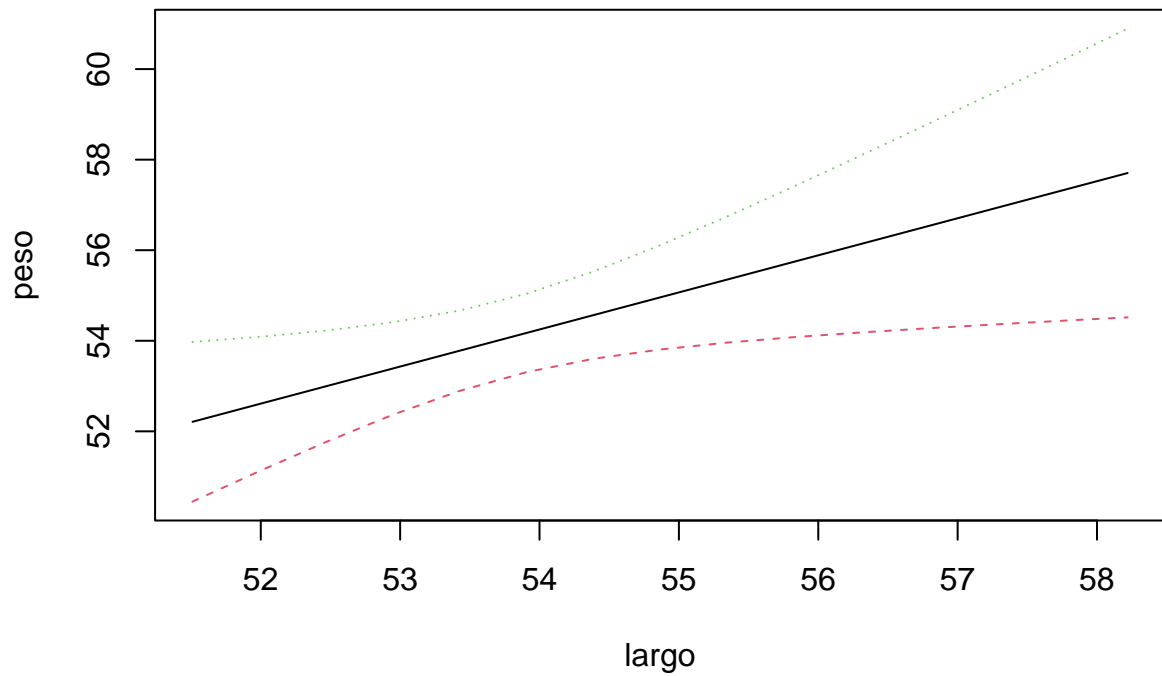


15.2 Intervalos de confianza para los predictores

Además nosotros podemos tener un IC para nuestras x en este caso el largo, es decir, cuanto puede variar el largo de un huevo y según eso dar una estimación de su peso.

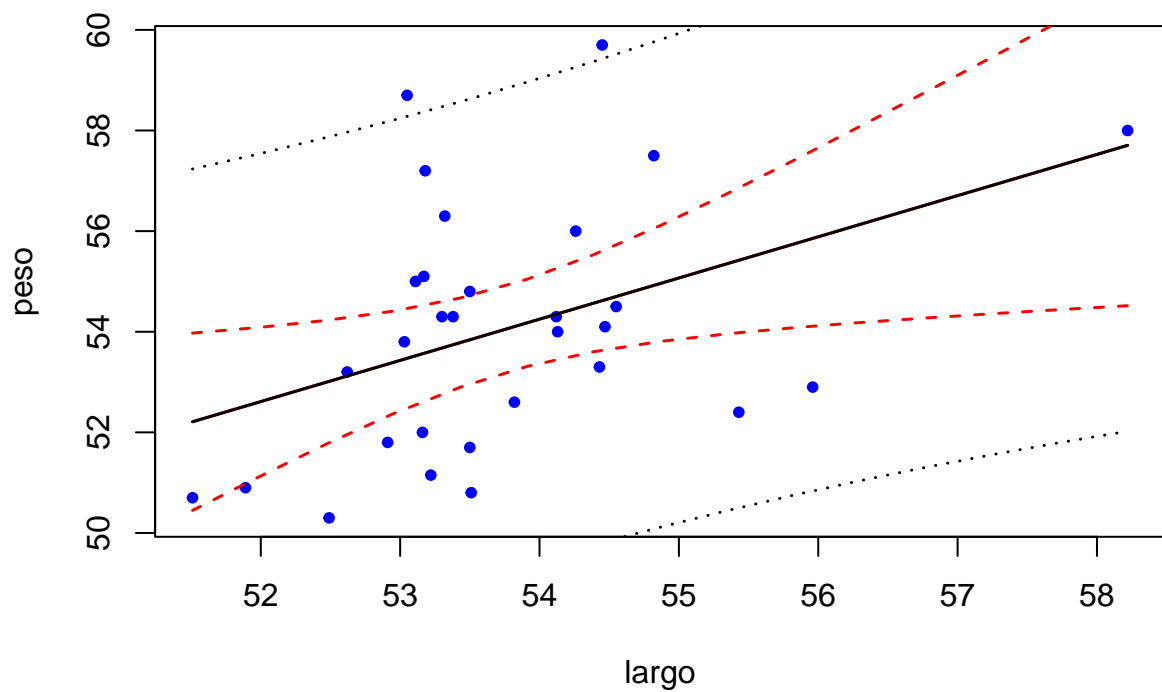
```
##      fit      lwr      upr
## 1 52.21048 50.44823 53.97273
## 2 52.60295 51.11939 54.08651
## 3 52.99542 51.76534 54.22549
## 4 53.38788 52.36714 54.40863
## 5 53.78035 52.89297 54.66774
## 6 54.17282 53.30698 55.03866
```

Dibujamos las bandas de confianza, que además reflejan la incertidumbre sobre futuras observaciones:



Por último podemos hacer un gráfico con la nube de puntos y los dos bandas, la de confianza y la de predicción (Ferrari & Head, 2010).

R.L.S. Peso vs Largo (IC's & IP's)



16 Apéndice

16.1 Lista de figuras

Figura:

1. Grafico de Dispersión diametro vs peso
2. R.L.S de diametro vs peso
3. Residuales vs valores ajustados (diametro vs peso)
4. QQplot de los residuales
5. Grafico de Dispersión diametro vs peso (Escala Log)
6. R.L.S de diametro vs peso Transformada ($\log()$)
7. Residuales vs valores ajustados (diametro vs peso) Transformada ($\log()$)
8. QQplot de los residuales Transformada ($\log()$)
9. Residuales de cada observación
10. Grafico de Intervalos de predicción de la R.L.S $\log()$
11. Grafico de Intervalos de confianza de la R.L.S $\log()$
12. R.L.S. Peso vs Largo (IC's & IP's)
13. Grafico de Dispersión largo vs peso
14. R.L.S de largo vs peso
15. Residuales vs valores ajustados (largo vs peso)
16. QQplot de los residuales
17. Residuales de cada observación
18. Scatter plot Peso vs largo
19. Influence plot (Residuos estandarizados)
20. Grafico de Intervalos de predicción de la R.L.S
21. Grafico de Intervalos de confianza de la R.L.S
22. R.L.S. Peso vs Largo (IC's & IP's)

16.2 Código:

```
## ----message=FALSE, warning=FALSE, include=FALSE-----
library(readr)
library(tidyverse)
library(kableExtra)
library(magrittr)
library(ggExtra)
library(GGally)
library(janitor)
library(tidystats)
library(car)
library(faraway)
library(lmtest)
library(graphics)
datos <- read_delim("eggs.csv", delim = ";") %>% clean_names()
names(datos)[3] <- "diametro"

myQQnorm <- function(modelo, student = F, ...){
  if(student){
    res <- rstandard(modelo)
    lab.plot <- "Normal Q-Q Plot of Studentized Residuals"
  } else {
    res <- residuals(modelo)
    lab.plot <- "Normal Q-Q Plot of Residuals"
  }
  shapiro <- shapiro.test(res)
  shapvalue <- ifelse(shapiro$p.value < 0.001, "P value < 0.001",
    paste("P value = ", round(shapiro$p.value, 4), sep = ""))
  shapstat <- paste("W = ", round(shapiro$statistic, 4), sep = "")
  q <- qqnorm(res, plot.it = FALSE)
  qqnorm(res, main = lab.plot, ...)
  qqline(res, lty = 2, col = 2)
  text(min(q$x, na.rm = TRUE), max(q$y, na.rm = TRUE)*0.95,
    pos = 4, 'Shapiro-Wilk Test', col = "blue", font = 2)
  text(min(q$x, na.rm = TRUE), max(q$y, na.rm = TRUE)*0.80,
    pos = 4, shapstat, col = "blue", font = 3)
  text(min(q$x, na.rm = TRUE), max(q$y, na.rm = TRUE)*0.65,
    pos = 4, shapvalue, col = "blue", font = 3)
}

## ----echo=FALSE, message=FALSE, warning=FALSE-----
cor.test(datos$diametro, datos$peso)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(datos$diametro, datos$peso, xlab = "diámetro en mm",
  ylab = "Peso en gr", main = "diámetro vs Peso",
  cex.main = 0.95, pch=20)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
modelo1 <- lm(peso~diametro, data=datos)
```

```

summary(modelo1)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(datos$diámetro, datos$peso, xlab = "diámetro en mm",
      ylab = "Peso en gr", pch=20)
abline(modelo1)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
summary(modelo1)$coefficients

## ----message=FALSE, warning=FALSE, include=FALSE-----
MSR. ancho <- mean(summary(modelo1)$residuals^2)

## ----message=FALSE, warning=FALSE, include=FALSE-----
confint(modelo1, level = 0.95)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
anova(modelo1)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(fitted(modelo1), residuals(modelo1), xlab = "Ancho",
      ylab = "Residuales", main = "Residuales vs. valores ajustados", pch=20)
abline(h = 0, lty = 2, col = 2)

## ----message=FALSE, warning=FALSE, include=FALSE-----
datos$fitted.modelo1 <- fitted(modelo1)
datos$residuals.modelo1 <- residuals(modelo1)
datos$rstudent.modelo1 <- rstudent(modelo1)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
modelo1 %>% myQQnorm()

## ----echo=FALSE, message=FALSE, warning=FALSE-----
bptest(modelo1)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(log(datos$diámetro), log(datos$peso), xlab = "diámetro",
      ylab = "Peso", main = "diámetro vs Peso (Escala Log)",
      cex.main = 0.95, pch=20)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
modelo2 <- lm(log(peso)~log(diámetro), data=datos)
summary(modelo2)

```

```

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(log(datos$diametro), log(datos$peso), xlab = "diámetro en mm",
      ylab = "Peso en gr", pch=20)
abline(modelo2)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
summary(modelo2)$coefficients

## ----message=FALSE, warning=FALSE, include=FALSE-----
MSR.log.ancho <- mean(summary(modelo2)$residuals^2)

## ----message=FALSE, warning=FALSE, include=FALSE-----
confint(modelo2, level = 0.95)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
anova(modelo2)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(fitted(modelo2), residuals(modelo2), xlab = "Ancho",
      ylab = "Residuales", main = "Residuales vs. valores ajustados", pch=20)
abline(h = 0, lty = 2, col = 2)

## ----message=FALSE, warning=FALSE, include=FALSE-----
datos$fitted.modelo2 <- fitted(modelo2)
datos$residuals.modelo2 <- residuals(modelo2)
datos$rstudent.modelo2 <- rstudent(modelo2)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
modelo2 %>% myQQnorm()

## ----echo=FALSE, message=FALSE, warning=FALSE-----
bptest(modelo2)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(datos$residuals.modelo1, pch = 20,
      ylab = "Residuos", xlab = "Índices")
abline(h = cor(datos$peso, datos$diametro))

## -----
dwtest(peso-diametro, alternative = "two.sided", data = datos)

```



```

## ----message=FALSE, warning=FALSE, include=FALSE-----
x0 <- seq(min(datos$diametro), max(datos$diametro), length = 15)
dfp <- data.frame(diametro = x0)
pred.ip <- predict(modelo1, dfp, interval = "prediction", se.fit = TRUE, data = datos)
head(pred.ip$fit)

## -----
matplot(x0, pred.ip$fit, type = "l", xlab = "diametro", ylab = "peso")

## ----message=FALSE, warning=FALSE, include=FALSE-----
x0 <- seq(min(datos$diametro), max(datos$diametro), length = 15)
dfp <- data.frame(diametro = x0)
pred.ip <- predict(modelo2, dfp, interval = "prediction",
se.fit = TRUE, data = datos)
pred.ip1 <- predict(modelo1, dfp, interval = "prediction",
se.fit = TRUE, data = datos)
head(pred.ip$fit)
newpred <- exp(pred.ip$fit)
head(newpred)

## -----
pred.ic <- predict(modelo1, dfp, interval = "confidence",
se.fit = TRUE, data = datos)
head(pred.ic$fit)

## -----
matplot(x0, pred.ic$fit, type = "l", xlab = "diametro", ylab = "peso")

## -----
plot(datos$diametro, datos$peso, pch = 20, xlab = "diametro",
ylab = "peso", col="blue")

# Añadimos las bandas
matlines(dfp$diametro, pred.ic$fit, lty = c(1, 2, 2),
        lwd = 1.5, col = "red")

matlines(dfp$diametro, pred.ip1$fit, lty = c(1, 3, 3),
        lwd = 1.5, col = "black")
title(main= "R.L.S. Peso vs Diámetro (IC's & IP's)")

## ----echo=FALSE, message=FALSE, warning=FALSE-----
cor.test(datos$largo, datos$peso)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(datos$largo, datos$peso, xlab = "largo en mm",
      ylab = "Peso en gr", main = "largo vs Peso",
      cex.main = 0.95, pch=20)

```

```

## ----echo=FALSE, message=FALSE, warning=FALSE-----
modelo3 <- lm(peso~largo, data=datos)
summary(modelo3)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(datos$largo, datos$peso, xlab = "diámetro en mm",
      ylab = "Peso en gr", pch=20)
abline(modelo3)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
summary(modelo3)$coefficients

## ----message=FALSE, warning=FALSE, include=FALSE-----
MSR.largo <- mean(summary(modelo3)$residuals^2)

## ----message=FALSE, warning=FALSE, include=FALSE-----
confint(modelo3, level = 0.95)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
anova(modelo3)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(fitted(modelo3), residuals(modelo3), xlab = "Largo",
     ylab = "Residuales", main = "Residuales vs. valores ajustados", pch=20)
abline(h = 0, lty = 2, col = 2)

## ----message=FALSE, warning=FALSE, include=FALSE-----
datos$fitted.modelo3 <- fitted(modelo3)
datos$residuals.modelo3 <- residuals(modelo3)
datos$rstudent.modelo3 <- rstudent(modelo3)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
modelo3 %>% myQQnorm()

## ----echo=FALSE, message=FALSE, warning=FALSE-----
bptest(modelo3)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
plot(datos$residuals.modelo3, pch = 20, ylab = "Residuos", xlab = "Índices")
abline(h = cor(datos$peso, datos$largo))

```

```

## ----echo=FALSE, message=FALSE, warning=FALSE-----
dwtest(peso~largo, alternative = "two.sided", data = datos)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
scatterplot(peso~largo,data = datos,smooth = F, pch=19,
            regLine = F, xlab = "Largo", ylab = "Peso")
title(main = "Scatter Plot | Peso vs Largo")

## ----echo=FALSE, message=FALSE, warning=FALSE-----
outlierTest(modelo3, cutoff = 0.05, n.max = 10, order = TRUE)
influencePlot(modelo3, id.n = 2)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
cook3 <- cooks.distance(modelo3)
labels3 <- rownames(datos)
halfnorm(cook, 3, labs = labels, ylab = "Distancia de Cook")
abline(h=4/30, lty = 2, col = 2)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
x0.l <- seq(min(datos$largo), max(datos$largo), length = 15)
dfp.l <- data.frame(largo = x0.l)
pred.ip.l <- predict(modelo3, dfp.l, interval = "prediction",
                    se.fit = TRUE, data = datos)
head(pred.ip.l$fit)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
matplot(x0.l, pred.ip.l$fit, type = "l", xlab = "largo", ylab = "peso")

## ----echo=FALSE, message=FALSE, warning=FALSE-----
pred.ic.l <- predict(modelo3, dfp.l, interval = "confidence",
                    se.fit = TRUE, data = datos)
head(pred.ic.l$fit)

## -----
matplot(x0.l, pred.ic.l$fit, type = "l", xlab = "largo", ylab = "peso")

## -----
plot(datos$largo, datos$peso, pch = 20, xlab = "largo",
     ylab = "peso", col="blue")

# Añadimos las bandas
matlines(dfp.l$largo, pred.ic.l$fit, lty = c(1, 2, 2),
        lwd = 1.5, col = "red")

matlines(dfp.l$largo, pred.ip.l$fit, lty = c(1, 3, 3),
        lwd = 1.5, col = "black")

title(main= "R.L.S. Peso vs Largo (IC's & IP's)")

```

17 Referencias

- [1] Coleman, D. E., & Montgomery, D. C. (1993). A Systematic Approach to Planning for a Designed Industrial Experiment. *Technometrics*, 35(1), 1–12. <https://doi.org/10.2307/1269280>
- [2] The jamovi project (2021). jamovi. (Version 2.2) [Computer Software]. Retrieved from <https://www.jamovi.org>.
- [3] R Core Team (2021). R: A Language and environment for statistical computing. (Version 4.0) [Computer software]. Retrieved from <https://cran.r-project.org>. (R packages retrieved from MRAN snapshot 2021-04-01).
- [4] Fox, J., & Weisberg, S. (2020). car: Companion to Applied Regression. [R package]. Retrieved from <https://cran.r-project.org/package=car>.
- [5] Ali S. Hadi, S. C. &. (2006). Linear models with r (4th edition.). John Wiley & Sons. Retrieved from http://samples.sainsburysebooks.co.uk/9780470055458_sample_381725.pdf
- [6] Ferrari, D., & Head, T. (2010). Regression in r. part i: Simple linear regression. UCLA Department of Statistics Statistical Consulting Center. Retrieved October 13, 2014, from http://scc.stat.ucla.edu/page_attachments/0000/0139/reg_1.pdf
- [7] Field, A., Miles, J., & Field, Z. (2012). Discovering statistics using r (1st edition.). Sage Publications Ltd.
- [8] J.Faraway, J. (2009). Linear models with r (1st edition.). Taylor & Francis e-Library. Retrieved from <http://home.ufam.edu.br/jcardoso/PPGMAT537/Linear%20Models%20with%20R.pdf>
- [9] Kabacoff, R. (2014). Creating a figure arrangement with fine control. Retrieved October 13, 2014, from <http://www.statmethods.net/advgraphs/layout.html>
- [10] Pérez, J. L. (2014). La estadística: Una orqueta hecha instrumento. Retrieved October 13, 2014, from <http://estadisticaorquestainstrumento.wordpress.com/>
- [11] Sánchez, J. G. P. (2011). Regresión lineal simple. Universidad Politécnica de Madrid. Retrieved October 13, 2014, from <http://ocw.upm.es/estadistica-e-investigacion-operativa/introduccion-a-la-estadistica-basica-el-diseno-de-experimentos-y-la-regresion-lineal/contenidos/Material-de-clase/Regresion.pdf>
- [12] (SCG), S. S. C. G. (2013). Multiple linear regression (r). San Diego State University. Retrieved October 13, 2014, from <http://scg.sdsu.edu/mlr-r/>
- [13] SPSS. (2007). Análisis de regresión lineal: El procedimiento regresión lineal. IBM SPSS Statistics. Retrieved October 13, 2014, from http://pendientedemigracion.ucm.es/info/socivmyt/paginas/D_departamento/materiales/analisis_datosyMultivariable/18reglin_SPSS.pdf

Punto 2° Parcial: Ajuste de un modelo de R.L.M

Universidad Nacional de Colombia
Análisis de Regresión 2022-1S
Medellín, Colombia
2022

Daniel Villa 1005087556

Juan Pablo Vanegas 1000640165



UNIVERSIDAD NACIONAL DE COLOMBIA

Contents

1	Objetivos:	3
1.1	Objetivos específicos	3
2	Antecedentes Relevantes	3
3	Variables de respuesta:	3
4	Variable de Control:	3
5	Creación del Modelo	3
6	Análisis de correlación	4
7	Ajuste del modelo	5
8	Comparación de modelos	6
9	Selección del “mejor” modelo	9
9.1	Resumen de las variables de interes	9
9.2	Correlación	9
10	Elección de los predictores	11
11	Condiciones para la regresión múltiple lineal	12
11.1	Relación lineal entre los predictores numéricos y la variable dependiente:	12
11.2	Distribución normal de los residuos:	13
11.3	Variabilidad constante de los residuos:	13
11.4	Autocorrelación:	14
11.5	Identificación de posibles valores atípicos o influyentes	14
12	Apendice:	15
12.1	listas de Figuras:	15
12.2	Codigo:	16
13	Referencias:	19

1 Objetivos:

Crear un modelo ajustado de R.L.M. por el cual se pueda predecir la estatura de un individuo (discriminando por genero) sabiendo las estaturas de los padres (madre y padre) utilizando el software estadístico *R*.

1.1 Objetivos específicos

- Plantear el modelo de R.L.M.
- Interpretar los parámetros del modelo.
- Determinar si el efecto de las estaturas de los padres sobre la estatura del sujeto es significativo.
- Interpretar nuestro R^2 .
- Validar los supuestos del modelo.
- Aplicar la prueba de falta de ajuste.

2 Antecedentes Relevantes

La población encuestada, pertenece a estudiantes de la Universidad Nacional de Colombia sede Medellín de diferentes carreras, es decir, la mayoría de los sujetos de la muestra son jóvenes entre los 18 y los 25 años, además decidimos que solamente aquellos que tenían la posibilidad de saber las estaturas de sus padres entraban a nuestra base de datos, ya que el proceso sería más arduo si tomamos datos donde nos faltan llenar valores en las celdas correspondientes.

3 Variables de respuesta:

En nuestro caso será la estatura del sujeto (Hombre o Mujer) para ajustar un modelo para predecir por medio de nuestras variables predictoras la estatura del sujeto.

4 Variable de Control:

En este caso tendremos 3 variables haciendo de este un modelo de R.L.M.

1. Estatura del Padre.
2. Estatura del Madre.
3. Genero del sujeto.

5 Creación del Modelo

Nuestro modelo tendrá la forma de:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \varepsilon_i, \quad i = 1, \dots, 30 \quad \text{donde } \varepsilon_i \sim N(0, \sigma^2)$$

Para llegar a lo anterior primero miraremos nuestra base de datos y como se comportan estas variables (*algunos datos...*).

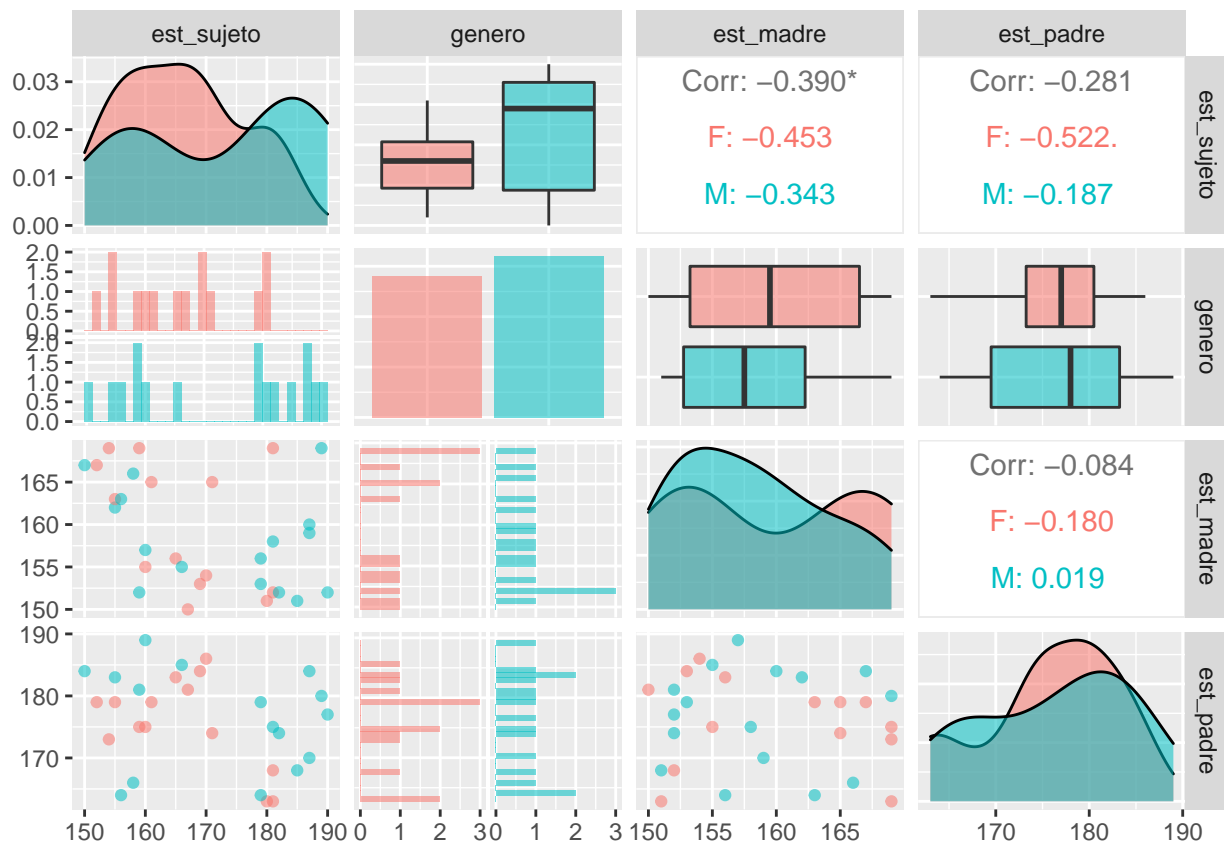
est_sujeto	genero	exp_sujeto	est_madre	est_padre	ced_sujeto	ced_madre	ced_padre	exp_padre	exp
180	F	1/11/2017	151	163	2567	1262	8323	18/10/1965	1/1
181	M	22/6/2016	158	175	6044	5350	0316	21/6/1984	15/
187	M	26/9/2020	159	170	6946	4240	0659	31/12/1977	6/
155	M	8/10/2020	162	183	7023	2121	4144	21/1/1970	4/
159	F	21/6/2016	169	175	1791	7832	3727	11/6/1987	19/
.
.
.
150	M	17/11/2016	167	184	8126	9626	9357	30/6/1981	13/
154	F	5/2/2020	169	173	9538	2408	4497	12/9/1982	30/
167	F	7/9/2021	150	181	2894	7612	4657	22/3/1980	31/
156	M	18/6/2020	163	164	8849	4763	8827	9/6/1989	23/
170	F	4/8/2017	154	186	9227	4016	1264	29/7/1968	10/

Definimos nuestras variables:

- **est_sujeto**: Estatura en cm del sujeto encuestado.
- **genero**: Genero del sujeto encuestado.
- **exp_sujeto**: Fecha de expedición de la cedula del sujeto.
- **est_madre**: Estatura de la madre del sujeto en cm.
- **est_padre**: Estatura del padre del sujeto en cm.
- **ced_sujeto**: Ultimos cuatro dígitos de la cedula del sujeto encuestado.
- **ced_madre**: Ultimos cuatro dígitos de la cedula de la madre.
- **ced_padre**: Ultimos cuatro dígitos de la cedula del padre.
- **exp_madre**: Fecha de expedición de la cedula de la madre.
- **exp_padre**: Fecha de expedición de la cedula del padre.

6 Análisis de correlación

Comenzamos representando los datos en una nube de puntos múltiple, donde vemos la relación entre cada par de variables.



Según la siguiente tabla veremos que tan debil o fuerte son los datos respecto a otros:



Figure 1: corr_guide

como podemos apreciar nuestras 3 correlaciones entre cada variable predictora y nuestra Y es una *Correlación Negativa débil* ya que ninguno supera ± 0.5 .

7 Ajuste del modelo

```
##
## Call:
## lm(formula = est ~ est_pa + est_ma + genero, data = stature)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -19.854  -6.954   0.101   4.956  26.325
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 381.5306    73.4427   5.195 2.01e-05 ***
## est_pa      -0.5417     0.2766  -1.959  0.0609 .
## est_ma      -0.7519     0.3204  -2.346  0.0269 *
## generoM       5.7221     4.1315   1.385  0.1778
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.2 on 26 degrees of freedom
## Multiple R-squared:  0.3025, Adjusted R-squared:  0.222
## F-statistic: 3.758 on 3 and 26 DF,  p-value: 0.02294
```

El error típico residual es 11.2, la $R^2 = 0.3025$, aunque para el modelo múltiple es mejor fijarnos en su valor ajustado $R_a^2 = 0.222$. Esto que significa que la recta de regresión explica el 22.22% de la variabilidad del modelo. Además, $F = 3.758$ con una significación $p < 0.05$, lo que nos dice que nuestro modelo de regresión resulta un poco mejor que el modelo básico.

8 Comparación de modelos

Pretendemos seleccionar el “mejor” subconjunto de predictores por varias razones:

1. Explicar los datos de la manera más simple. Debemos eliminar predictores redundantes.
2. Predictores innecesarios añade ruido a las estimaciones.
3. La causa de la multicolinealidad es tener demasiadas variables tratando de hacer el mismo trabajo. Eliminar el exceso de predictores ayuda a la interpretación del modelo.
4. Si vamos a utilizar el modelo para la predicción, podemos ahorrar tiempo y/o dinero al no medir predictores redundantes.

Puesto que tenemos dos variables explicativas disponemos de 6 modelos posibles:

modelo 1 : $est \sim est_{ma} + est_{pa} + genero$

modelo 2 : $est \sim est_{ma} + genero$

modelo 3 : $est \sim est_{pa} + genero$

modelo 4 : $est \sim est_{ma}$

modelo 5 : $est \sim est_{pa}$

modelo 6 : $est \sim genero$

Vamos a ajustar cada uno de los modelos

```
##
## Call:
## lm(formula = est ~ est_ma + genero, data = stature)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -18.075  -9.035  -1.057   8.030  23.859
```

```

##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 278.2976    53.7713   5.176 1.9e-05 ***
## est_ma      -0.7020     0.3358  -2.091  0.0461 *
## generoM      5.4878     4.3413   1.264  0.2170
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.77 on 27 degrees of freedom
## Multiple R-squared:  0.1995, Adjusted R-squared:  0.1402
## F-statistic: 3.365 on 2 and 27 DF,  p-value: 0.04957
##
## Call:
## lm(formula = est ~ est_pa + genero, data = stature)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -22.784  -9.021   2.940   8.162  18.059
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 252.2777    52.4691   4.808 5.1e-05 ***
## est_pa      -0.4902     0.2978  -1.646  0.111
## generoM      6.9006     4.4298   1.558  0.131
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.1 on 27 degrees of freedom
## Multiple R-squared:  0.1548, Adjusted R-squared:  0.09216
## F-statistic: 2.472 on 2 and 27 DF,  p-value: 0.1033
##
## Call:
## lm(formula = est ~ est_ma, data = stature)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.884 -10.284  -3.723   7.131  26.948
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 289.6144    53.5839   5.405 9.18e-06 ***
## est_ma      -0.7548     0.3367  -2.242  0.0331 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.9 on 28 degrees of freedom
## Multiple R-squared:  0.1522, Adjusted R-squared:  0.1219
## F-statistic: 5.025 on 1 and 28 DF,  p-value: 0.0331
##
## Call:
## lm(formula = est ~ est_pa, data = stature)

```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -19.3440 -10.9008   0.4737   9.9393  21.2097
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 252.7696    53.7882   4.699 6.31e-05 ***
## est_pa      -0.4721     0.3051  -1.548   0.133
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.4 on 28 degrees of freedom
## Multiple R-squared:  0.0788, Adjusted R-squared:  0.0459
## F-statistic: 2.395 on 1 and 28 DF,  p-value: 0.1329
##
## Call:
## lm(formula = est ~ genero, data = stature)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -22.688 -11.821   1.929  11.562  17.312
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  166.071     3.330  49.874 <2e-16 ***
## generoM       6.616     4.560   1.451   0.158
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 12.46 on 28 degrees of freedom
## Multiple R-squared:  0.06994, Adjusted R-squared:  0.03672
## F-statistic: 2.106 on 1 and 28 DF,  p-value: 0.1579
```

Para evitar la elección subjetiva del mejor modelo, podemos comparar todos los modelos mediante una tabla ANOVA conjunta para cada par de modelos.

Nota: Se escoje el de menor error estandar (RSS)

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
27	3740.781	NA	NA	NA	NA
27	3949.951	0	-209.1702	NA	NA

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
27	3740.781	NA	NA	NA	NA
28	3962.171	-1	-221.3902	1.597938	0.2170003

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
27	3740.781	NA	NA	NA	NA
28	4304.956	-1	-564.1753	4.072073	0.0536355

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
27	3740.781	NA	NA	NA	NA
28	4346.366	-1	-605.5851	4.370958	0.0461022

```
## [1] "Mejor modelo del 2 al 6 con menor RSS"
```

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
27	3740.781	NA	NA	NA	NA
26	3259.715	1	481.066	3.837058	0.0609385

Comparando ambas tablas anova deducimos que el modelo que mejor se ajusta a los datos es el `modelo2` pues reduce el error estándar.

9 Selección del “mejor” modelo

Existen distintos métodos a la hora de construir un modelo complejo de regresión con varios predictores

- El **método jerárquico** en el que se seleccionan los predictores basándose en un trabajo anterior y el investigador decide en qué orden introducir las variables predictoras al modelo.
- El **método de entrada forzada** en el que todas las variables entran a la fuerza en el modelo simultáneamente.
- Los **métodos paso a paso** que se basan en un criterio matemático para decidir el orden en que los predictores entran en el modelo.

Nosotros vamos a utilizar en R los métodos paso a paso

9.1 Resumen de las variables de interes

```
##      est      genero      est_ma      est_pa
## Min.   :150.0   F:14   Min.   :150.0   Min.   :163.0
## 1st Qu.:159.0   M:16   1st Qu.:153.0   1st Qu.:170.8
## Median :168.0           Median :157.5   Median :178.0
## Mean   :169.6           Mean   :159.0   Mean   :176.2
## 3rd Qu.:181.0           3rd Qu.:165.0   3rd Qu.:182.5
## Max.   :190.0           Max.   :169.0   Max.   :189.0
```

En todas las variables explicativas los valores de la media y la mediana son muy cercanos, lo cual es muy bueno.

9.2 Correlación

```
##
## Pearson's product-moment correlation
##
## data:  df$est and df$est_ma
## t = -2.2416, df = 28, p-value = 0.0331
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.65788264 -0.03466696
## sample estimates:
##      cor
## -0.3900645
##
## Pearson's product-moment correlation
##
## data:  df$est and df$est_pa
## t = -1.5476, df = 28, p-value = 0.1329
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.58211099  0.08850857
```

```
## sample estimates:
##      cor
## -0.2807117

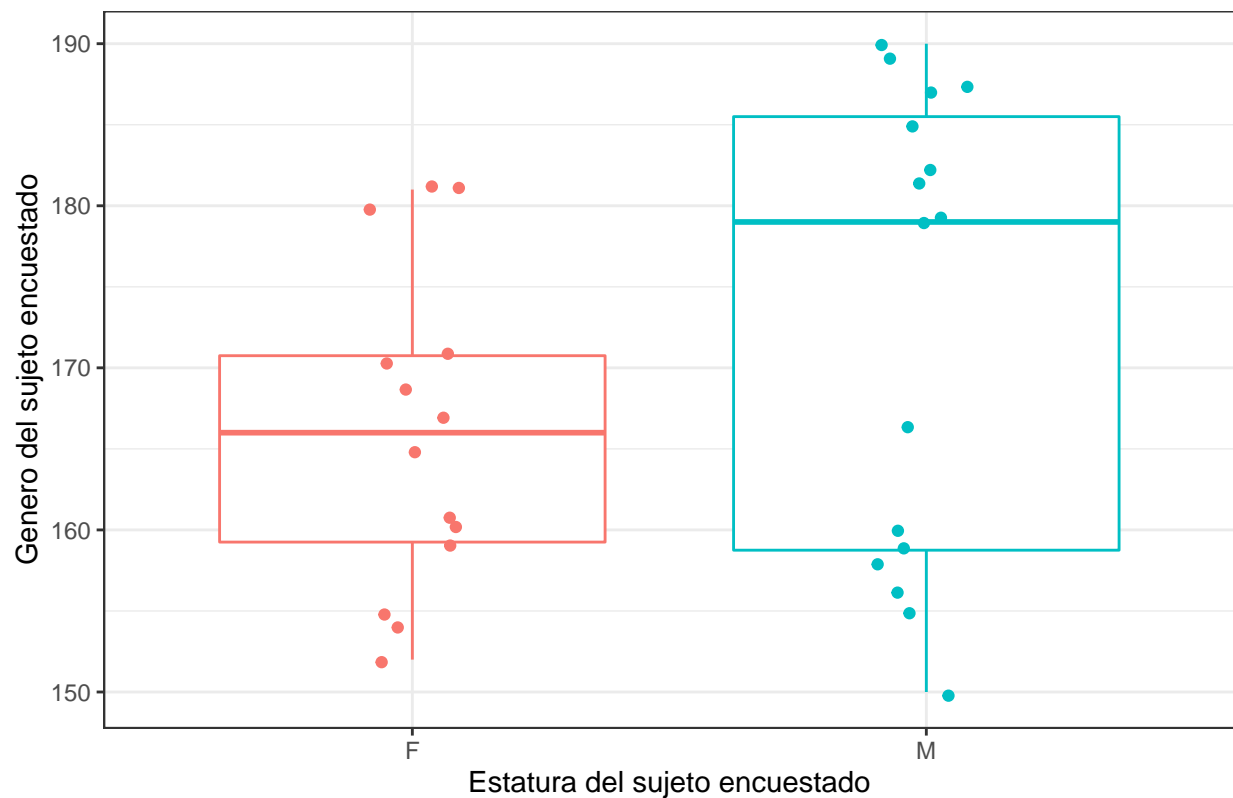
##
## Pearson's product-moment correlation
##
## data:  df$est_pa and df$est_ma
## t = -0.44376, df = 28, p-value = 0.6606
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.4308663  0.2852892
## sample estimates:
##      cor
## -0.08356935
```

Como podemos ver nuestras variables están correlacionadas negativamente, además tenemos que no superan el umbral de ± 0.5 lo que conlleva a que tenga una correlación débil.

Si miramos sus p -value podemos encontrar que solamente para el caso de las madres se rechaza la hipótesis H_0 donde se asume que la correlación entre este tipo de variables puede ser igual a cero, es decir, que de alguna forma los datos están correlacionados negativamente entre la estatura del sujeto y de la madre.

Vamos a aplicar los tres métodos a nuestros modelos para cómo funciona cada uno de ellos

Boxplot Genero vs Estatura



El análisis gráfico y de correlación muestran una relación lineal significativa entre la variable `est` y `est_ma`. La variable `genero` parece influir de forma significativa en la estatura. Ambas variables pueden ser buenos predictores en un modelo lineal múltiple para la variable dependiente `est`.

Esto confirma que nos quedaremos con el `modelo2` por ende a este trabajaremos los datos.

Volvemos a refrescar un poco la memoria con el **modelo2**:

```
##           Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) 278.2976107 53.7712777  5.175581 1.900584e-05
## est_ma      -0.7020405  0.3357947 -2.090684 4.610222e-02
## generoM      5.4877921  4.3412792  1.264096 2.170003e-01
```

se nota puede notar viendo el valor $Pr(|t|)$ que la variable genero es la unica que se rechaza, dentro de la hipotesis que el $\beta_i = 0$, es decir, nuestro modelo solamente se complementaria de los datos de la estatura de la madre.

pero observando el $R^2 = 0.1402$ podemos ver que solamente el kodelo explica el 14.02% de la variabilidad observada en la estatura de los sujetos; nuestro modelo junto con el valor $p - value = 0.04957 \approx 0.05$ nos demuestra que nuestros datos no son significativos.

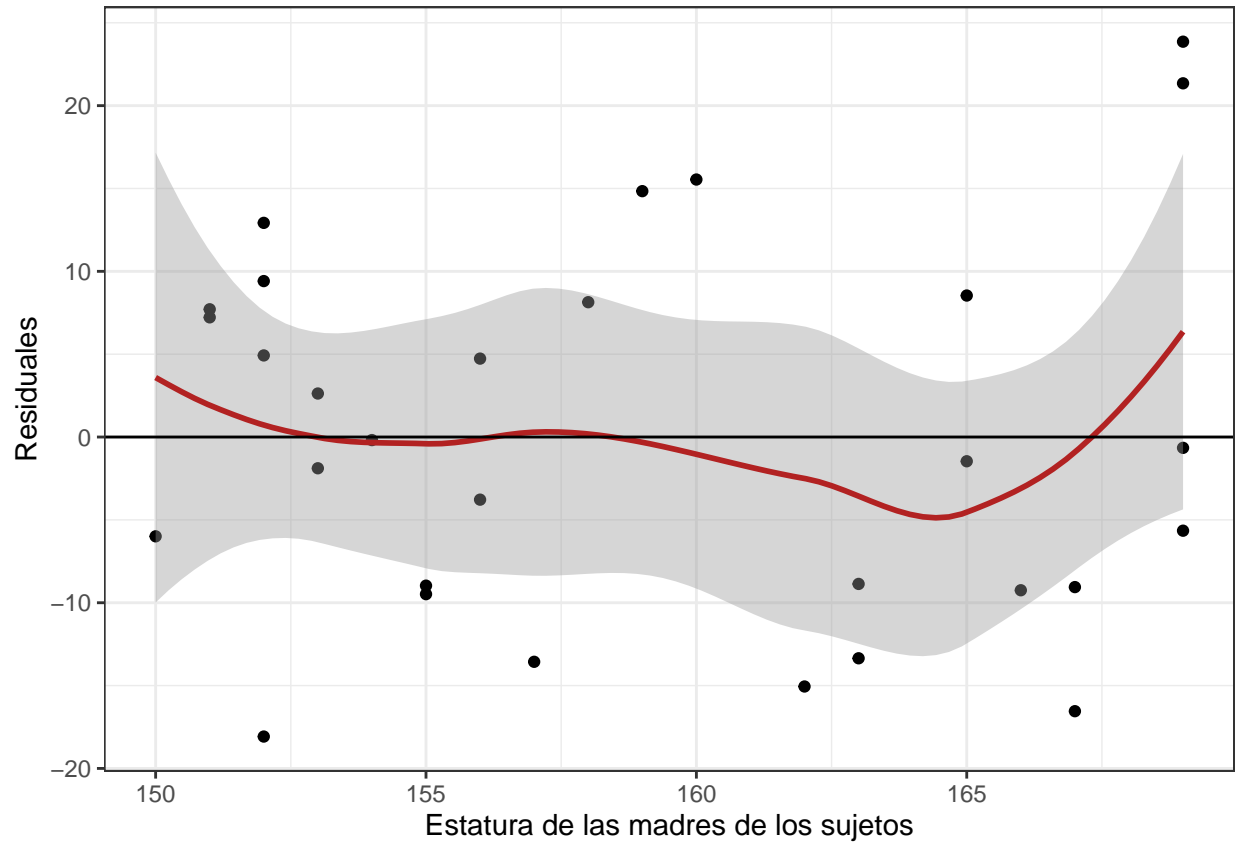
10 Elección de los predictores

En este caso, al solo haber dos predictores, a partir del summary del modelo se identifica que solo la variable `est_ma` es importante.

```
##
## Call:
## lm(formula = est ~ est_ma + genero, data = stature)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -18.075  -9.035  -1.057   8.030  23.859
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 278.2976    53.7713   5.176 1.9e-05 ***
## est_ma      -0.7020     0.3358  -2.091  0.0461 *
## generoM      5.4878     4.3413   1.264  0.2170
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11.77 on 27 degrees of freedom
## Multiple R-squared:  0.1995, Adjusted R-squared:  0.1402
## F-statistic: 3.365 on 2 and 27 DF,  p-value: 0.04957
```

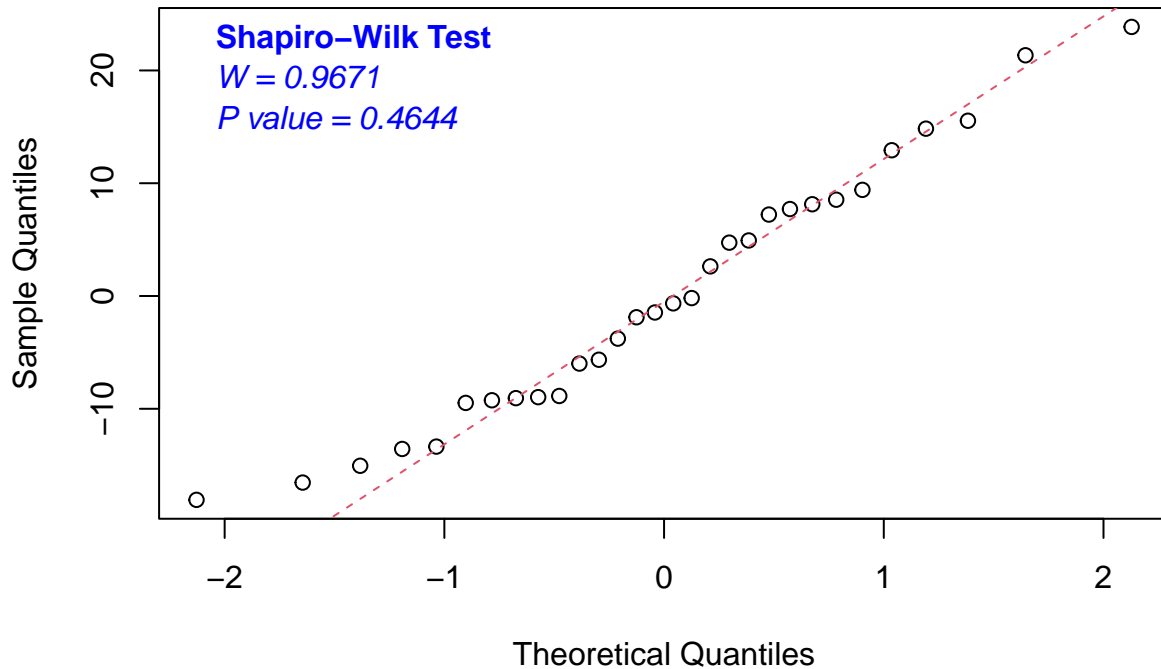
11 Condiciones para la regresión múltiple lineal

11.1 Relación lineal entre los predictores numéricos y la variable dependiente:



Se satisface la condición de linealidad. Se aprecian posibles datos atípicos.

11.2 Distribución normal de los residuos: Normal Q-Q Plot of Residuals



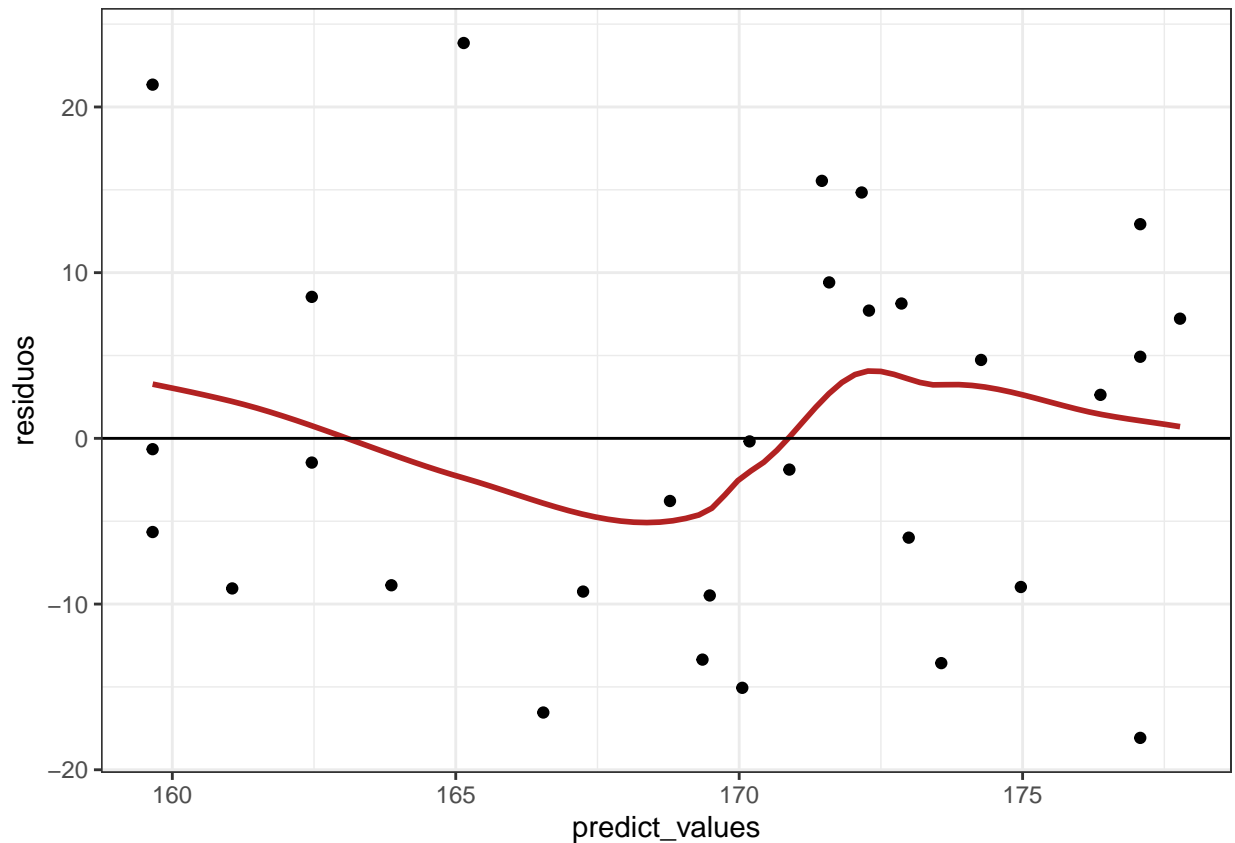
apesar de los datos atípicos se observa que cumple con la normalidad de los residuos.

Nos surge una duda, ¿estos datos atípicos estarán influenciando nuestros datos de alguna manera?

11.3 Variabilidad constante de los residuos:

```
ggplot(data = data.frame(predict_values = predict(modelo2),  
                          residuos = residuals(modelo2)),  
       aes(x = predict_values, y = residuos))+  
  geom_point()+  
  geom_smooth(color = "firebrick", se = FALSE)+  
  geom_hline(yintercept = 0)+  
  theme_bw()
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```



```
bptest(modelo2)
```

```
##
## studentized Breusch-Pagan test
##
## data: modelo2
## BP = 8.9341, df = 2, p-value = 0.01148
```

Como podemos ver nuestros ε_i no tienen un varianza constante, apra ello miraremos estos datos atípicos y si hay necesidad de transformar los datos o aplicar un modelo donde las varianzas no sean constantes.

Nota: No multicolinealidad: Dado que solo hay un predictor cuantitativo no se puede dar colinealidad.

11.4 Autocorrelación:

```
dwt(modelo2, alternative = "two.sided")
```

```
## lag Autocorrelation D-W Statistic p-value
## 1 -0.03978963 2.063677 0.844
## Alternative hypothesis: rho != 0
```

No hay evidencia de autocorrelación

11.5 Identificación de posibles valores atípicos o influyentes

```
outlierTest(modelo2)
```

```
## No Studentized residuals with Bonferroni p < 0.05
```

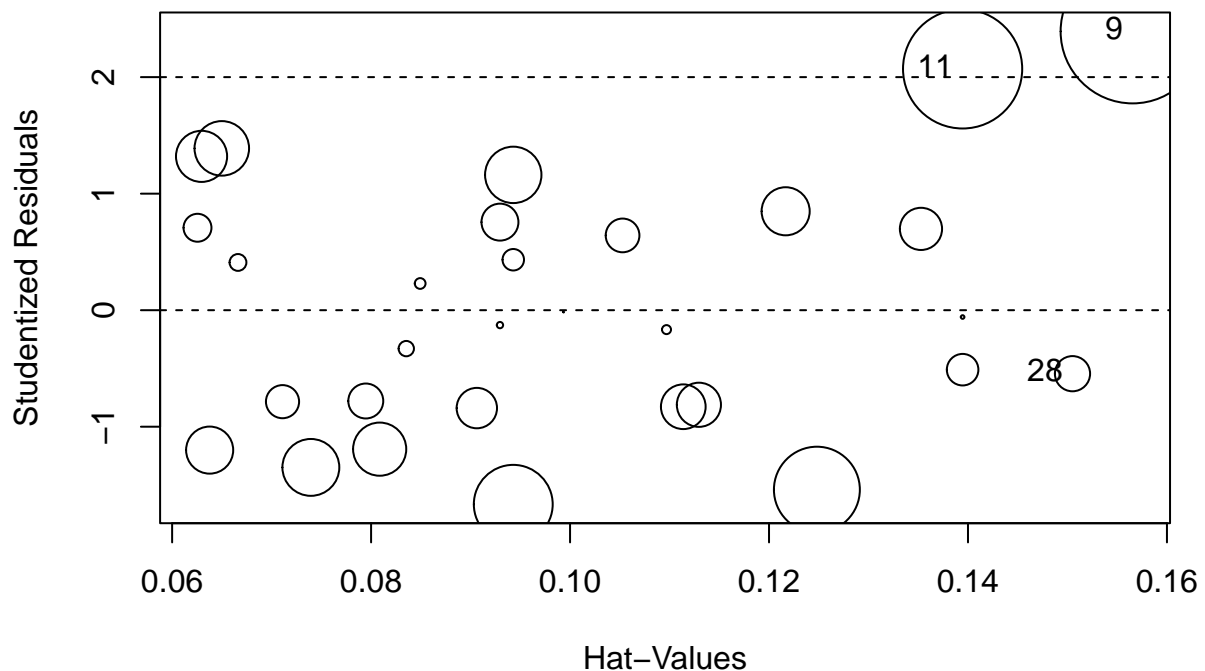
```
## Largest |rstudent|:
##   rstudent unadjusted p-value Bonferroni p
## 9 2.392447          0.024254      0.72762
```

Tal como se apreció en el estudio de normalidad de los residuos, la observación 13 tiene un residuo estandarizado >2 (más de 2 veces la desviación estándar de los residuos) por lo que se considera un dato atípico. El siguiente paso es determinar si es influyente.

```
summary(influence.measures(modelo2))
```

```
## Potentially influential observations of
##   lm(formula = est ~ est_ma + genero, data = stature) :
##
##   dfb.1_ dfb.est_ dfb.gnrM dffit   cov.r cook.d hat
## 9 -0.80   0.80    0.54    1.03_* 0.73 0.30 0.16
```

```
influencePlot(modelo2)
```



	StudRes	Hat	CookD
9	2.3924473	0.1565516	0.3013984
11	2.0706534	0.1394605	0.2064776
28	-0.5450432	0.1505058	0.0180131

El análisis muestran varias observaciones influyentes aunque ninguna excede la distancia de $Cook > 1$, pero como lo vimos antes, el dato en nuestra modelo tiene influencia si $Cook > \frac{4}{n}$, donde $n = 30$, es decir, $Cook > 0.1333$; nuestros datos influenciadores y que superan, pero al no superar las distancias del HAT, en fin. no hay datos atípicos en los datos, porque se llega a la conclusión que este modelo no sirve para predecir o hacer inferencia sobre los datos de la madre y su hijo.

12 Apendice:

12.1 listas de Figuras:

1. Correlation plot.

2. Residuales vs Est_Madre.
3. QQnorm norm Residuales.
4. Predict vs Residuos
5. Infuence plot.

12.2 Código:

```
## ----message=FALSE, warning=FALSE, include=FALSE-----
library(readr)
library(tidyverse)
library(kableExtra)
library(magrittr)
library(ggExtra)
library(GGally)
library(janitor)
library(tidystats)
library(car)
library(faraway)
library(lmtest)
library(graphics)
# lectura de la base de datos:
stature <- read_csv("estaturas.csv")

## ----echo=FALSE, message=FALSE, warning=FALSE-----
kable(rbind(head(stature, n = 5),rep(".", ncol(stature)),
              rep(".", ncol(stature)),rep(".", ncol(stature)),
              tail(stature, n = 5)),digits = 30, align = "c")

## ----echo=FALSE, message=FALSE, warning=FALSE-----
stature$genero %<>% as.factor()
stature$exp_sujeto %<>% as.Date(format="%d/%m/%Y")
stature$exp_padre %<>% as.Date(format="%d/%m/%Y")
stature$exp_madre %<>% as.Date(format="%d/%m/%Y")
stature$ced_madre %<>% as.character()
stature %>% ggpairs(.,columns=c(1,2,4,5), aes(color=genero,alpha=0.5))

## ----message=FALSE, warning=FALSE, include=FALSE-----
colnames(stature) <- c("est","genero","exp","est_ma","est_pa","ced",
                      "ced_ma","ced_pa","exp_pa","exp_ma")

## ----echo=FALSE, message=FALSE, warning=FALSE-----
modelo.m <- lm(est ~ est_pa + est_ma + genero, data = stature)
summary(modelo.m)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
modelo2 <- lm(est~est_ma+genero, data=stature)
modelo3 <- lm(est~est_pa+genero, data=stature)
modelo4 <- lm(est~est_ma, data=stature)
```

```

modelo5 <- lm(est~est_pa, data=stature)
modelo6 <- lm(est~genero, data=stature)

summary(modelo2)
summary(modelo3)
summary(modelo4)
summary(modelo5)
summary(modelo6)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
anova(modelo2,modelo3)
anova(modelo2,modelo4)
anova(modelo2,modelo5)
anova(modelo2,modelo6)
print("Mejor modelo del 2 al 6 con menor RSS")
anova(modelo2,modelo.m)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
# Eliminamos la variables que no nos interesan en el modelo:
df <- stature[, c(1,2,4,5)]
summary(df)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
cor.test(df$est, df$est_ma)
cor.test(df$est, df$est_pa)
cor.test(df$est_pa, df$est_ma)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
df %>% ggplot(., aes(x=genero, y=est, color=genero))+
  geom_boxplot()+
  geom_jitter(width = 0.1)+
  theme_bw()+theme(legend.position = "none")+
  labs(x="Estatura del sujeto encuestado",
       y="Genero del sujeto encuestado",
       title = "Boxplot Genero vs Estatura")

## ----echo=FALSE, message=FALSE, warning=FALSE-----
summary(modelo2)$coefficients

## ----echo=FALSE, message=FALSE, warning=FALSE-----
summary(modelo2)

## ----echo=FALSE, message=FALSE, warning=FALSE-----
df %>% ggplot(.,aes(x=est_ma, y=modelo2$residuals))+
  geom_point()+
  geom_smooth(color="firebrick")+
  geom_hline(yintercept = 0)+

```

```

theme_bw()+
labs(x="Estatura de las madres de los sujetos",
     y="Residuales")

## ----echo=FALSE, message=FALSE, warning=FALSE-----
myQQnorm <- function(modelo, student = F, ...){
  if(student){
    res <- rstandard(modelo)
    lab.plot <- "Normal Q-Q Plot of Studentized Residuals"
  } else {
    res <- residuals(modelo)
    lab.plot <- "Normal Q-Q Plot of Residuals"
  }
  shapiro <- shapiro.test(res)
  shapvalue <- ifelse(shapiro$p.value < 0.001, "P value < 0.001", paste("P value = ", round(shapiro$p.v
shapstat <- paste("W = ", round(shapiro$statistic, 4), sep = "")
  q <- qqnorm(res, plot.it = FALSE)
  qqnorm(res, main = lab.plot, ...)
  qqline(res, lty = 2, col = 2)
  text(min(q$x, na.rm = TRUE), max(q$y, na.rm = TRUE)*0.95, pos = 4, 'Shapiro-Wilk Test', col = "blue",
  text(min(q$x, na.rm = TRUE), max(q$y, na.rm = TRUE)*0.80, pos = 4, shapstat, col = "blue", font = 3)
  text(min(q$x, na.rm = TRUE), max(q$y, na.rm = TRUE)*0.65, pos = 4, shapvalue, col = "blue", font = 3)
}

modelo2 %>% myQQnorm()

## -----
ggplot(data = data.frame(predict_values = predict(modelo2),
                          residuos = residuals(modelo2)),
       aes(x = predict_values, y = residuos))+
  geom_point()+
  geom_smooth(color = "firebrick", se = FALSE)+
  geom_hline(yintercept = 0)+
  theme_bw()

bptest(modelo2)

## -----
dwt(modelo2, alternative = "two.sided")

## -----
outlierTest(modelo2)

## -----
summary(influence.measures(modelo2))
influencePlot(modelo2)

```

13 Referencias:

- [1] Coleman, D. E., & Montgomery, D. C. (1993). A Systematic Approach to Planning for a Designed Industrial Experiment. *Technometrics*, 35(1), 1–12. <https://doi.org/10.2307/1269280>
- [2] The jamovi project (2021). jamovi. (Version 2.2) [Computer Software]. Retrieved from <https://www.jamovi.org>.
- [3] R Core Team (2021). R: A Language and environment for statistical computing. (Version 4.0) [Computer software]. Retrieved from <https://cran.r-project.org>. (R packages retrieved from MRAN snapshot 2021-04-01).
- [4] Fox, J., & Weisberg, S. (2020). car: Companion to Applied Regression. [R package]. Retrieved from <https://cran.r-project.org/package=car>.
- [5] Ali S. Hadi, S. C. &. (2006). Linear models with r (4th edition.). John Wiley & Sons. Retrieved from http://samples.sainsburysebooks.co.uk/9780470055458_sample_381725.pdf
- [6] Ferrari, D., & Head, T. (2010). Regression in r. part i: Simple linear regression. UCLA Department of Statistics Statistical Consulting Center. Retrieved October 13, 2014, from http://scc.stat.ucla.edu/page_attachments/0000/0139/reg_1.pdf
- [7] Field, A., Miles, J., & Field, Z. (2012). Discovering statistics using r (1st edition.). Sage Publications Ltd.
- [8] J.Faraway, J. (2009). Linear models with r (1st edition.). Taylor & Francis e-Library. Retrieved from <http://home.ufam.edu.br/jcardoso/PPGMAT537/Linear%20Models%20with%20R.pdf>
- [9] Kabacoff, R. (2014). Creating a figure arrangement with fine control. Retrieved October 13, 2014, from <http://www.statmethods.net/advgraphs/layout.html>
- [10] Pérez, J. L. (2014). La estadística: Una orqueta hecha instrumento. Retrieved October 13, 2014, from <http://estadisticaorquestainstrumento.wordpress.com/>
- [11] Sánchez, J. G. P. (2011). Regresión lineal simple. Universidad Politécnica de Madrid. Retrieved October 13, 2014, from <http://ocw.upm.es/estadistica-e-investigacion-operativa/introduccion-a-la-estadistica-basica-el-diseno-de-experimentos-y-la-regresion-lineal/contenidos/Material-de-clase/Regresion.pdf>
- [12] (SCG), S. S. C. G. (2013). Multiple linear regression (r). San Diego State University. Retrieved October 13, 2014, from <http://scg.sdsu.edu/mlr-r/>
- [13] SPSS. (2007). Análisis de regresión lineal: El procedimiento regresión lineal. IBM SPSS Statistics. Retrieved October 13, 2014, from http://pendientedemigracion.ucm.es/info/socivmyt/paginas/D_departamento/materiales/analisis_datosyMultivariable/18reglin_SPSS.pdf