

Contenidos

- ▶ 5.1: Diagnóstico: Análisis de los residuos
- ▶ 5.2: La descomposición ANOVA (ANalysis Of VAriance)
- ▶ 5.3: Relaciones no lineales y transformaciones para linealización
- ▶ 5.4: El modelo de regresión lineal en forma matricial
- ▶ 5.5: Introducción a la regresión lineal múltiple

Bibliografía

- ▶ Newbold, P. Estadística para los Negocios y la Economía (1997).
 - ▶ Capítulos 12, 13 y 14
- ▶ Peña, D. Regresión y Diseño de Experimentos (2002).
 - ▶ Capítulos 5 y 6

5.1. Diagnóstico en regresión

- ▶ **Supuestos teóricos** del modelo de regresión lineal simple de una var. respuesta y sobre una var. explicativa x :
 - Linealidad: $y_i = \beta_0 + \beta_1 x_i + u_i$, para $i = 1, \dots, n$
 - Homogeneidad: $E[u_i] = 0$, para $i = 1, \dots, n$
 - Homocedasticidad: $\text{Var}[u_i] = \sigma^2$, para $i = 1, \dots, n$
 - Independencia: u_i y u_j son independientes para $i \neq j$
 - Normalidad: $u_i \sim \text{Normal}(0, \sigma^2)$, para $i = 1, \dots, n$
- ▶ Los **métodos de diagnóstico** se utilizan para contrastar si tales supuestos son adecuados para los datos disponibles (x_i, y_i) ; se basan en el **análisis de los residuos** $e_i = y_i - \hat{y}_i$

5.1. Diagnóstico: diagrama de puntos

- ▶ El método más sencillo consiste en la observación visual del **diagrama de puntos** (x_i, y_i)
- ▶ A menudo, este sencillo pero potente método revela pautas que sugieren si el modelo teórico es o no adecuado
- ▶ Ilustraremos su uso con un ejemplo clásico. Consideremos los cuatro conjuntos de datos siguientes

5.1. Diagnóstico: diagrama de puntos

TABLE 3-10
Four Data Sets

DATA SET 1		DATA SET 2	
X	Y	X	Y
10.0	8.04	10.0	9.14
8.0	6.95	8.0	8.14
13.0	7.58	13.0	8.74
9.0	8.81	9.0	8.77
11.0	8.33	11.0	9.26
14.0	9.96	14.0	8.10
6.0	7.24	6.0	6.13
4.0	4.26	4.0	3.10
12.0	10.84	12.0	9.13
7.0	4.82	7.0	7.26
5.0	5.68	5.0	4.74

DATA SET 3		DATA SET 4	
X	Y	X	Y
10.0	7.46	8.0	6.58
8.0	6.77	8.0	5.76
13.0	12.74	8.0	7.71
9.0	7.11	8.0	8.84
11.0	7.81	8.0	8.47
14.0	8.84	8.0	7.04
6.0	6.08	8.0	5.25
4.0	5.39	19.0	12.50
12.0	8.15	8.0	5.56
7.0	6.42	8.0	7.91
5.0	5.73	8.0	6.89

SOURCE: F. J. Anscombe, *op. cit.*

5.1. Diagnóstico: diagrama de puntos

- ▶ Para cada uno de los cuatro conjuntos de datos anteriores, se obtiene el mismo modelo estimado de regresión lineal:
- ▶ $\hat{y}_i = 3.0 + 0.5x_i$
- ▶ $n = 11$, $\bar{x} = 9.0$, $\bar{y} = 7.5$, $r_{x,y} = 0.817$
- ▶ El **error estándar** estimado del estimador $\hat{\beta}_1$,

$$\sqrt{\frac{s_R^2}{(n-1)s_x^2}},$$

toma el valor 0.118. El estadístico T correspondiente toma el valor $T = 0.5/0.118 = 4.237$

- ▶ Sin embargo, los diagramas de puntos correspondientes revelan que los cuatro conjuntos de datos son cualitativamente muy diferentes: ¿Qué conclusiones podemos extraer de estos diagramas?

5.1. Diagnóstico: diagrama de puntos

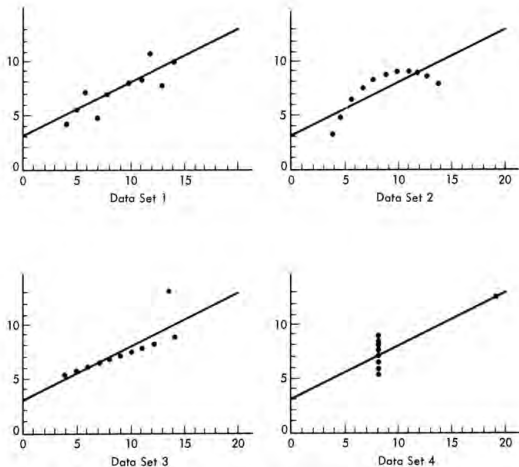


FIGURE 3-29 Scatterplots for the four data sets of Table 3-10
SOURCE: F. J. Anscombe, *op cit*.

5.1: análisis de los residuos

- ▶ Si la observación del diagrama de puntos no basta para descartar el modelo, se utilizan métodos de diagnóstico basados en el **análisis de los residuos** $e_i = y_i - \hat{y}_i$
- ▶ El análisis comienza **tipificando** los residuos (dividiéndolos por la cuasi-desviación típica residual): Las cantidades resultantes se denominan **residuos tipificados**:

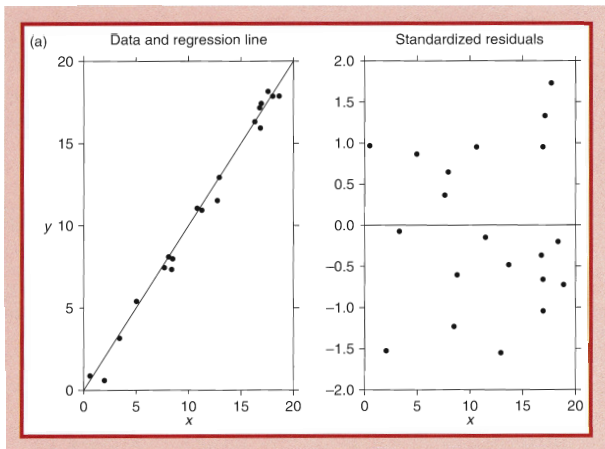
$$\frac{e_i}{s_R}$$

- ▶ Bajo los supuestos del modelo de regresión lineal, los residuos tipificados son aproximadamente variables aleatorias normales estándar independientes
- ▶ Un gráfico de los residuos tipificados no debería mostrar ninguna pauta clara

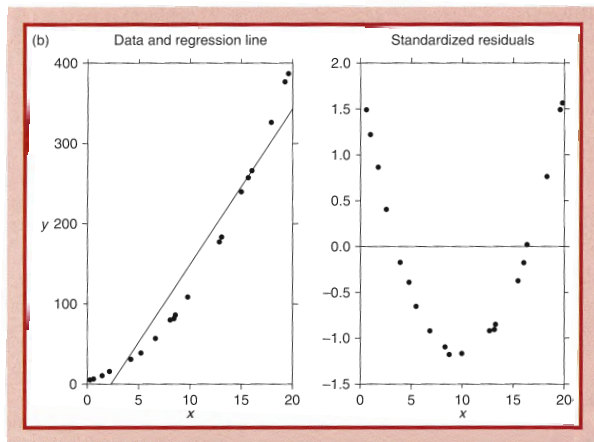
5.1: Diagramas de residuos

- ▶ Hay varios tipos de diagramas de residuos. Los más comunes son:
 - Diagrama de los residuos vs. x
 - Diagrama de los residuos vs. \hat{y}
- ▶ Las desviaciones de los supuestos del modelo dan lugar a pautas, que se pueden identificar visualmente

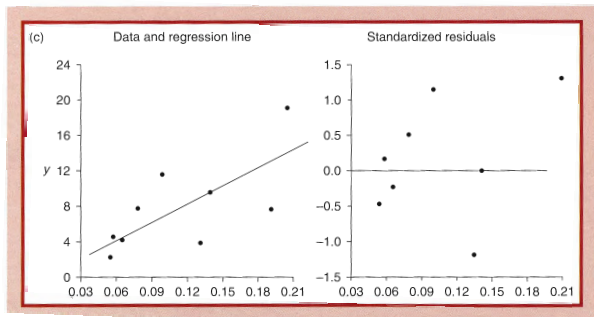
5.1: Ej: consistencia con el modelo teórico



5.1: Ej: No linealidad



5.1: Ej: Heterocedasticidad



5.1: Datos atípicos

- ▶ A partir del gráfico de la recta de regresión podemos observar **datos atípicos**, que presentan desviaciones sustanciales de la recta de regresión
- ▶ Los estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$ de los parámetros de la recta de regresión son muy sensibles a tales datos atípicos
- ▶ Por ello, es importante identificar tales datos y comprobar si son válidos
- ▶ Veremos que Statgraphics permite mostrar los datos que producen “*Unusual Residuals*”, así como “*Influential Points*”

5.1: Normalidad de los errores

- ▶ Recordemos que uno de los supuestos teóricos del modelo de regresión lineal es que los errores tienen una distribución normal
- ▶ Podemos comprobar este supuesto visualmente a partir de la observación y análisis de los residuos e_i , empleando varios métodos:
 - ▶ Observación del histograma de frecuencias de los residuos
 - ▶ Observación de un “*Normal Probability Plot*” para los residuos (desviaciones importantes de los datos de la línea recta en este gráfico indican desviaciones sustanciales del supuesto de normalidad)

5.2: La descomposición ANOVA

- ▶ ANOVA: *ANalysis Of VAriance*
- ▶ Al ajustar un modelo de regresión lineal $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ a un conjunto de datos (x_i, y_i) , para $i = 1, \dots, n$, podemos distinguir tres **fuentes de variación** en las respuestas:
 - **variación debida al modelo:** $SCM = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$, donde las siglas “SC” se refieren a “suma de cuadrados”, y la “M” se refiere al “Modelo”
 - **variación residual:** $SCR = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2$
 - **variación total:** $SCT = \sum_{i=1}^n (y_i - \bar{y})^2$
- ▶ La descomposición ANOVA indica que $SCT = SCM + SCR$

5.2: El coeficiente de determinación R^2

- ▶ La descomposición ANOVA indica que $SCT = SCM + SCR$
- ▶ Notemos que: $y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$
- ▶ $SCM = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ mide la variación de las respuestas debida a la regresión (explicada por los valores predichos \hat{y})
- ▶ Por lo tanto, el cociente SCR/SCT es la proporción de variación de la respuesta no explicada por la regresión
- ▶ El cociente $R^2 = SCM/SCT = 1 - SCR/SCT$ es la proporción de variación de las respuestas explicada por la regresión; se conoce como **coeficiente de determinación**
- ▶ Resultado: $R^2 = r_{xy}^2$ (coef. de correlación al cuadrado)
- ▶ Ej: si $R^2 = 0.85$, la variable x explica un 85% de la variación de la variable y

5.2: Tabla ANOVA

Fuente de variación	SC	G.L.	Media	Cociente F
Modelo	SCM	1	SCM/1	SCM/ s_R^2
Residuos/Errores	SCR	$n - 2$	SCR/ $(n - 2) = s_R^2$	
Total	SCT	$n - 1$		

5.2: Contraste de hipótesis ANOVA

- ▶ Contraste de hipótesis $H_0: \beta_1 = 0$ vs. $H_1: \beta_1 \neq 0$
- ▶ Consideremos el cociente

$$F = \frac{\text{SCM}/1}{\text{SCR}/(n-2)} = \frac{\text{SCM}}{s_R^2}$$

- ▶ Bajo H_0 , F sigue una distribución $F_{1,n-2}$
- ▶ Contraste a nivel α : rechazar H_0 si $F > F_{1,n-2;\alpha}$
- ▶ ¿Cuál es la relación con el contraste basado en la t de Student que vimos en el Tema 4? **Son equivalentes**

5.2: Ej. ANOVA

DatosContReg12.xlsx - Microsoft Excel

Archivo Inicio Insertar Diseño de página Fórmulas Datos Revisar Vista

Desde Access Desde web Desde texto De otras fuentes Conexiones existentes Actualizar todo Editar vínculos Conexiones Ordenar Filtro Avanzadas Ordenar y filtrar Borrar Volver a aplicar Avanzadas Texto en columnas Quitar duplicados Herramientas de datos Validación de datos Consolidar Análisis Y si Agrupar Desagrupar Subtotal Esquema Análisis de datos Solver Análisis

M1 Resumen

	M	N	O	P	Q	R	S	T	U	V
1	Resumen									
2										
3	<i>Estadísticas de la regresión</i>									
4	Coefficiente de correlación múltiple	0,684396733								
5	Coefficiente de determinación R ²	0,468398888								
6	R ² ajustado	0,432958814								
7	Error típico	0,038453336								
8	Observaciones	17								
9										
10	ANÁLISIS DE VARIANZA									
11		<i>Grados de libertad</i>	<i>Suma de cuadrados</i>	<i>Promedio de los cuadrados</i>	<i>F</i>	<i>Valor crítico de F</i>				
12	Regresión	1	0,019542912	0,019542912	13,21664527	0,002441485				
13	Residuos	15	0,022179885	0,001478659						
14	Total	16	0,041722797							
15										
16		<i>Coefficientes</i>	<i>Error típico</i>	<i>Estadístico t</i>	<i>Probabilidad</i>	<i>Inferior 95%</i>	<i>Superior 95%</i>	<i>Inferior 95,0%</i>	<i>Superior 95,0%</i>	
17	Intercepción	0,28515256	0,045407156	6,279903505	1,47518E-05	0,188369497	0,381935623	0,188369497	0,381935623	
18	Variable X 1	0,838324888	0,230595986	3,635470434	0,002441485	0,346821179	1,329828596	0,346821179	1,329828596	
19										

SR²

5.3: Relaciones no lineales y linealización

- ▶ Supongamos que la parte determinista $f(x_i; a, b)$ de la respuesta en el modelo

$$y_i = f(x_i; a, b) + u_i, \quad i = 1, \dots, n$$

es una **función no lineal** de x que depende de dos parámetros a y b (ej: $f(x; a, b) = ab^x$)

- ▶ En algunos casos podemos aplicar **transformaciones** a los datos para **linearizarlos**, y así poder aplicar los métodos de regresión lineal
- ▶ A partir de los datos (x_i, y_i) originales, obtenemos los datos transformados (x'_i, y'_i)
- ▶ Los parámetros β_0 y β_1 de la relación lineal entre las x'_i y las y'_i se obtienen como transformaciones de los parámetros a y b

5.3: Transformaciones para linealización

- ▶ Ejemplos de transformaciones para linealización:
 - ▶ Si $y = f(x; a, b) = ax^b$ entonces $\log y = \log a + b \log x$: tomamos $y' = \log y$, $x' = \log x$, $\beta_0 = \log a$, $\beta_1 = b$
 - ▶ Si $y = f(x; a, b) = ab^x$ entonces $\log y = \log a + (\log b)x$: tomamos $y' = \log y$, $x' = x$, $\beta_0 = \log a$, $\beta_1 = \log b$
 - ▶ Si $y = f(x; a, b) = 1/(a + bx)$ entonces $1/y = a + bx$: tomamos $y' = 1/y$, $x' = x$, $\beta_0 = a$, $\beta_1 = b$
 - ▶ Si $y = f(x; a, b) = \ln(ax^b)$ entonces $y = \ln a + b \ln x$: tomamos $y' = y$, $x' = \ln x$, $\beta_0 = \ln a$, $\beta_1 = b$

5.4: Regresión lineal en forma matricial

- ▶ Recordemos el modelo de regresión lineal simple:

$$y_i = \beta_0 + \beta_1 x_i + u_i, \quad i = 1, \dots, n$$

- ▶ Escribiendo una ecuación para cada observación obtenemos

$$y_1 = \beta_0 + \beta_1 x_1 + u_1$$

$$y_2 = \beta_0 + \beta_1 x_2 + u_2$$

⋮

$$y_n = \beta_0 + \beta_1 x_n + u_n$$

5.4: Regresión lineal en forma matricial

- ▶ En forma matricial, podemos escribir

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_1 \\ \beta_0 + \beta_1 x_2 \\ \vdots \\ \beta_0 + \beta_1 x_n \end{bmatrix} + \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix},$$

o, separando los parámetros β_j de las x_i ,

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix},$$

5.4: Regresión lineal en forma matricial

- ▶ Escribimos la relación matricial

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix},$$

como

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$$

- ▶ \mathbf{y} : vector de respuestas; \mathbf{X} : matriz de variables explicativas (o del diseño experimental); $\boldsymbol{\beta}$: vector de parámetros; \mathbf{u} : vector de errores

5.4: La matriz de covarianzas de los errores

- ▶ Denotamos por $\text{Cov}(\mathbf{u})$ la matriz $n \times n$ de covarianzas de los errores; su elemento (i, j) es

$$\text{cov}(u_i, u_j) = \begin{cases} 0 & \text{si } i \neq j \\ \text{Var}[u_i] = \sigma^2 & \text{si } i = j \end{cases}$$

- ▶ $\text{Cov}(\mathbf{u})$ es la matriz identidad $\mathbf{I}_{n \times n}$ multiplicada por σ^2 :

$$\text{Cov}(\mathbf{u}) = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}$$

5.4: Estimación de mínimos cuadrados

- ▶ El vector estimado $\hat{\beta}$ de mínimos cuadrados es la solución única de la ecuación matricial 2×2 (comprueba las dimensiones)

$$(\mathbf{X}^T \mathbf{X}) \hat{\beta} = \mathbf{X}^T \mathbf{y},$$

es decir,

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- ▶ El vector $\hat{\mathbf{y}} = (\hat{y}_i)$ de respuestas estimadas es

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta}$$

y el vector de residuos es $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$

5.5: El modelo de regresión lineal múltiple

- ▶ Modelo de regresión lineal simple: predecir una respuesta y a partir de una variable explicativa x
- ▶ En numerosas aplicaciones, buscamos predecir la respuesta y a partir de múltiples variables explicativas x_1, \dots, x_k
- ▶ Ej: predecir el precio de una casa en función de su superficie, localización, planta, y número de baños
- ▶ Ej: predecir el tamaño de un parlamento en función de la población, su tasa de crecimiento, el número de partidos políticos con representación, etc.

5.5: El modelo de regresión lineal múltiple

- ▶ Modelo de regresión lineal múltiple: predecir una respuesta y a partir de múltiples variables explicativas x_1, \dots, x_k
- ▶ Tenemos n observaciones: para $i = 1, \dots, n$,

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + u_i$$

- ▶ Suponemos que las u_i son v.a. independientes con distribución Normal($0, \sigma^2$)

5.5: Ajuste de mínimos cuadrados

- ▶ Tenemos n observaciones: para $i = 1, \dots, n$,

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + u_i$$

- ▶ Buscamos ajustar a los datos $(x_{i1}, x_{i2}, \dots, x_{ik}, y_i)$ un hiperplano de la forma

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \dots + \hat{\beta}_k x_{ik}$$

- ▶ El **residuo** para la observación i es: $e_i = y_i - \hat{y}_i$
- ▶ Utilizamos la estimación de los parámetros $\hat{\beta}_j$ que minimiza la suma de los cuadrados de los residuos

5.5: Modelo en forma matricial

- ▶ Escribimos la relación matricial

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix},$$

como

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$$

- ▶ \mathbf{y} : vector de respuestas; \mathbf{X} : matriz de variables explicativas (o del diseño experimental); $\boldsymbol{\beta}$: vector de parámetros; \mathbf{u} : vector de errores

5.5: Estimación de mínimos cuadrados de β

- ▶ El vector estimado $\hat{\beta}$ de mínimos cuadrados es la solución única de la ecuación matricial $(k + 1) \times (k + 1)$ (comprueba las dimensiones)

$$(\mathbf{X}^T \mathbf{X}) \hat{\beta} = \mathbf{X}^T \mathbf{y},$$

como en el caso $k = 1$ visto anteriormente, es decir,

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- ▶ El vector $\hat{\mathbf{y}} = (\hat{y}_i)$ de respuestas estimadas es

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta}$$

y el vector de residuos es $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$

5.5: Estimación de la varianza σ^2

- ▶ Para el modelo de regresión lineal múltiple, estimamos la varianza σ^2 mediante la cuasi-varianza residual,

$$s_R^2 = \frac{\sum_{i=1}^n e_i^2}{n - k - 1},$$

que es un estimador insesgado (nótese que para regresión lineal simple el denominador vale $n - 2$)

5.5: Distribución muestral de $\hat{\beta}$

- ▶ Bajo los supuestos del modelo, el estimador de mínimos cuadrados $\hat{\beta}$ del vector de parámetros β sigue una distribución **normal multivariante**
- ▶ $E[\hat{\beta}] = \beta$ (i.e., es un estimador insesgado)
- ▶ La matriz de covarianzas de $\hat{\beta}$ es $\text{Cov}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$
- ▶ Estimamos $\text{Cov}(\hat{\beta})$ por $s_R^2 (\mathbf{X}^T \mathbf{X})^{-1}$
- ▶ La estimación de $\text{Cov}(\hat{\beta})$ nos da estimaciones $s^2(\hat{\beta}_j)$ de la varianza $\text{Var}(\hat{\beta}_j)$; $s(\hat{\beta}_j)$ es el **error estándar** del estimador $\hat{\beta}_j$
- ▶ Al tipificar $\hat{\beta}_j$ obtenemos:
$$\frac{\hat{\beta}_j - \beta_j}{s(\hat{\beta}_j)} \sim t_{n-k-1} \text{ (t de Student)}$$

5.5: Inferencia sobre los parámetros $\hat{\beta}_j$

- ▶ Intervalo de confianza a nivel $1 - \alpha$ para β_j :

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

- ▶ Contraste de hipótesis a nivel α para $H_0: \beta_j = 0$ vs. $H_1: \beta_j \neq 0$
- ▶ Rechazar H_0 si $|T| > t_{n-k-1; \alpha/2}$, donde $T = \hat{\beta}_j / s(\hat{\beta}_j)$ es el estadístico de contraste

5.5: La descomposición ANOVA

- ▶ ANOVA: *AN*alysis *Of* *VA*riance
- ▶ Al ajustar un modelo de regresión lineal múltiple

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \cdots + \hat{\beta}_k x_{ik}$$

a un conjunto de datos $(x_{i1}, \dots, x_{ik}, y_i)$, para $i = 1, \dots, n$, podemos distinguir tres **fuentes de variación** en las respuestas:

- **variación debida a la regresión:** $SCM = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$, donde las siglas “SC” se refieren a “suma de cuadrados”, y la “M” se refiere al “Modelo”
 - **variación residual:** $SCR = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2$
 - **variación total:** $SCT = \sum_{i=1}^n (y_i - \bar{y})^2$
- ▶ La descomposición ANOVA indica que $SCT = SCM + SCR$

5.5: El coeficiente de determinación R^2

- ▶ La descomposición ANOVA indica que $SCT = SCM + SCR$
- ▶ Notemos que: $y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$
- ▶ $SCM = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ mide la variación de las respuestas debida a la regresión (explicada por los valores predichos \hat{y}_i)
- ▶ Por lo tanto, el cociente SCR/SCT es la proporción de variación de la respuesta no explicada por la regresión
- ▶ El cociente $R^2 = SCM/SCT = 1 - SCR/SCT$ es la proporción de variación de las respuestas explicada por las variables explicativas; se conoce como **coeficiente de determinación múltiple**
- ▶ Resultado: $R^2 = r_{\hat{y}y}^2$ (coef. de correlación al cuadrado)
- ▶ Ej: si $R^2 = 0.85$, las variables x_1, \dots, x_k explican un 85% de la variación de la variable y

5.5: Tabla ANOVA

Fuente de variación	SC	G.L.	Media	Cociente F
Modelo	SCM	k	SCM/k	$(SCM/k)/s_R^2$
Residuos/Errores	SCR	$n - k - 1$	$\frac{SCR}{n-k-1} = s_R^2$	
Total	SCT	$n - 1$		

5.5: Contraste de hipótesis ANOVA

- ▶ Consideremos el contraste de hipótesis $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$ vs. $H_1: \beta_j \neq 0$ para algún $j = 1, \dots, k$
- ▶ H_0 : la respuesta no depende de las x_j
- ▶ Consideremos el cociente

$$F = \frac{\text{SCM}/k}{\text{SCR}/(n-k-1)} = \frac{\text{SCM}/k}{s_R^2}$$

- ▶ Bajo H_0 , F sigue una distribución $F_{k, n-k-1}$
- ▶ Contraste a nivel α : rechazar H_0 si $F > F_{k, n-k-1; \alpha}$