

# Examen 1, versión A. Regresión lineal simple

Gonzalo Pérez, Leonardo de la Cruz, José Ángel Román y David González

Semestre 2023-1

El examen se deberá subir al classroom antes de las 11:59 PM del 30 de octubre de 2022. Todas las preguntas tienen un valor de 1.5 puntos.

Favor de argumentar con detalle las respuestas.

NOTA. En caso de que se identifiquen respuestas iguales en otros exámenes, se procederá a la anulación de los exámenes involucrados.

NOTA. Incluir el(los) nombre(s) completo(s) de la(s) persona(s) que está(n) resolviendo los ejercicios. Equipos de máximo tres integrantes.

Usar una confianza de 95 % o una significancia de .05 en los casos en donde no se requiera otro nivel de forma explícita. En el caso de realizar alguna transformación de las variables, se tiene que hacer explícita la variable que se usa y la interpretación en las pruebas de hipótesis o intervalos de confianza.

## 1. Regresión a través del origen.

Ocasionalmente, un modelo en donde el valor del intercepto es conocido a priori y es igual a cero puede ser apropiado. Supongamos que además se considera el posible uso de una regresión ponderada, es decir, el modelo está dado por:

$$y_i = \beta x_i + \varepsilon_i, \quad i = 1, \dots, n,$$

donde  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  son variables independientes tal que  $\varepsilon_i \sim N\left(0, \frac{\sigma^2}{w_i}\right) \quad \forall \quad i = 1, \dots, n$ .

En general  $\sigma^2$  es desconocida, pero en lo que sigue suponga que es conocida. Además suponga que

$$w_i = \frac{1}{x_i}, \quad i = 1, \dots, n.$$

- i) Encuentre el estimador de  $\beta$  obtenido por el método de máxima verosimilitud,  $\hat{\beta}$ .
- ii) Encuentre la expresión de la varianza de  $\hat{\beta}$ .
- iii) Demuestre que  $\hat{\beta}$  es el UMVUE de  $\beta$ , es decir, que es el mejor estimador insesgado de  $\beta$ .

## 2.

Considere el modelo de regresión

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i,$$

donde  $E(\varepsilon_i) = 0$ ,  $V(\varepsilon_i) = \sigma^2$  y  $Cov(\varepsilon_i, \varepsilon_j) = 0 \quad \forall \quad i \neq j; \quad i, j = 1, \dots, n$ .

Calcular  $Cov(e_i, \hat{\beta}_0)$ , donde  $e_i = y_i - \hat{y}_i$  y  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ , con  $\hat{\beta}_0$  y  $\hat{\beta}_1$  los estimadores de los parámetros del modelo.

Hint: Recordar que  $\hat{y}_i$ ,  $\hat{\beta}_0$  y  $\hat{\beta}_1$  se pueden escribir como una combinación lineal de las  $y_{i's}$ .

**3. Expresión alternativa para  $R^2$** 

Considere el coeficiente de correlación muestral o de Pearson para dos variables  $X$  y  $Y$ :

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{(\sum_{i=1}^n (x_i - \bar{X})^2 \sum_{i=1}^n (y_i - \bar{Y})^2)^{1/2}}.$$

Considere el modelo de regresión

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

a. Demuestre que:

$$R^2 = r_{xy}^2.$$

Hint: Puede usar lo encontrado en la expresión (77) de las notas.

b. Demuestre que  $t^* = t$ , donde  $t$  es la estadística usada para contrastar " $H_0 : \beta_1 = 0$  vs  $H_1 : \beta_1 \neq 0$ ":

$$t = \frac{\hat{\beta}_1}{\sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{X})^2}}}.$$

Por otra parte,  $t^* = \frac{r_{xy}\sqrt{n-2}}{\sqrt{1-r_{xy}^2}}$  es la estadística usada para contrastar " $H_0 : \rho = 0$  vs  $H_a : \rho \neq 0$ "

cuando  $(X, Y)$  sigue una distribución normal bivariada con coeficiente de correlación  $\rho = \rho_{xy}$ .

Hint: Puede usar la relación en la expresión (68) de las notas.

**4. Problema Anova. Equivalencia con la prueba t para comparar dos poblaciones normales.**

Sea  $X_1, \dots, X_n$  una m.a. de la distribución  $N(\mu_x, \sigma^2)$  y  $Y_1, \dots, Y_m$  una m.a. de la distribución  $N(\mu_y, \sigma^2)$ , ambas muestras aleatorias son independientes entre sí. La prueba t se usa bajo este contexto para contrastar, por ejemplo:

$$H_0 : \mu_x = \mu_y \quad vs \quad H_a : \mu_x \neq \mu_y.$$

Sea  $t$  la estadística asociada a la prueba t antes mencionada, es decir:

$$t = \frac{\bar{X} - \bar{Y}}{s_p \sqrt{\frac{1}{n} + \frac{1}{m}}}, \quad (1)$$

donde

$$s_p^2 = \frac{(n-1)s_x^2 + (m-1)s_y^2}{m+n-2}, \quad s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n-1}, \quad s_y^2 = \frac{\sum_{i=1}^m (y_i - \bar{Y})^2}{m-1}.$$

i. Considere una variable  $Z$  tal que:  $Z = 1$  si la observación es de la población con distribución  $N(\mu_x, \sigma^2)$  y  $Z = 0$  si la observación es de la población con distribución  $N(\mu_y, \sigma^2)$ . Considere el modelo de regresión lineal simple:

$$w_j = \beta_0 + \beta_1 z_j + \varepsilon_j,$$

donde  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_{n+m}$  son variables independientes tal que  $\varepsilon_j \sim N(0, \sigma^2) \quad \forall \quad j = 1, \dots, n+m$ . En este modelo los valores de las variables  $X$  y  $Y$  componen la variable  $W$ , asumiendo que las primeras  $n$  observaciones son las que tienen valor  $Z = 1$  y el resto son las que tienen valor  $Z = 0$ . Indique cuál es la distribución de  $W$  para cada valor de la variable  $Z$ , haciendo énfasis en indicar la relación que esto implica entre los parámetros  $\mu_x$  y  $\mu_y$  con  $\beta_0$  y  $\beta_1$ .

- II. En términos de los parámetros del modelo de regresión lineal simple en I), indique cómo se deben escribir las hipótesis

$$H_0 : \mu_x = \mu_y \quad vs \quad H_a : \mu_x \neq \mu_y.$$

Además dé la expresión de la estadística asociada a la prueba que se usaría para contrastar estas hipótesis en el contexto del modelo de regresión lineal simple.

- III. Demuestre que la estadística dada en II) es algebraicamente equivalente en valor absoluto a la estadística  $t$  asociada a la prueba  $t$  dada en (1).

Hint: Puede usar todas las expresiones ya encontradas en clase para los estimadores y pruebas de hipótesis. Encuentre una expresión para  $\hat{\beta}_1$  en términos de  $x_i$  y  $y_i$  tomando ventaja de que  $z_j$  sólo toma el valor 0 o 1; además use una expresión para  $\hat{w}_j$  que sólo dependa de  $\hat{\beta}_1$ , en particular identificar  $\hat{E}(W|Z = 1)$  y  $\hat{E}(W|Z = 0)$  en términos de  $x_i$  y  $y_i$ .

## 5. Problema ANOVA. Medicamentos

Suponga que una empresa farmacéutica está ofreciendo al gobierno un nuevo medicamento para tratar a pacientes con la enfermedad Covid-19. El costo del medicamento es considerable y para tomar una buena decisión se han acercado a usted para analizar los datos que ha compartido la empresa farmacéutica. El archivo Ejercicio5A.csv contiene la información siguiente:  $Y$  es un índice de carga viral y  $Med$  es una variable con dos niveles dependiendo si se aplicó o no el nuevo medicamento. Se sabe que tener una menor carga viral está relacionado con una menor probabilidad de desarrollar una versión grave de la enfermedad y la empresa afirma que eso se logra al aplicar el medicamento, pues los pacientes que recibieron el medicamento tienen menor carga viral que los que sólo recibieron placebo.

- I. Realice un análisis descriptivo y/o la visualización de los datos
- II. Escriba la prueba asociada para argumentar en favor o no de la afirmación de la compañía. Para esto deberá indicar qué modelo podría usar y cuáles son los supuestos de éste.
- III. Lleve a cabo la prueba de hipótesis, justificando que los supuestos del modelo que está usando son válidos. Dé la interpretación de los resultados.
- IV. Suponga ahora que dado que el costo del medicamento es considerable, le han vuelto a preguntar si los resultados en el inciso III) son contundentes. Para esto, usted ha decidido analizar más el proceso de generación de los datos y ha platicado con los empleados de la farmacéutica, logrando que le compartan una nueva variable  $Edad$ . Realice un análisis descriptivo y/o visualización de los datos incluyendo esta nueva información. Comente lo que observe analizando si las conclusiones en III) se pueden **atribuir** sólo al medicamento.
- V. Dependiendo de lo observado en IV) y si considera necesario, repita los incisos II) y III) con un conjunto de datos donde el efecto se pueda **atribuir** sólo al medicamento y concluya.

Hint: Recuerde que para poder **atribuir** un efecto a algún factor se deben comparar poblaciones homogéneas, es decir, que no exista otro factor oculto que pudiera estar asociado con las diferencias o no diferencias que se observen.

## 6. Uso del modelo de regresión lineal simple

Los *pingüinos Macaroni* ponen nidadas de dos huevos de tamaño diferente. El peso en gramos de los huevos de 11 nidadas se presenta en la tabla de abajo.

- I. Ajuste la recta de regresión del peso del huevo mayor ( $y$ ) dado el peso del huevo menor ( $x$ ). Comente sobre el ajuste del modelo, es decir, si parece correcto y si se cumplen los supuestos.

- II. Los investigadores tienen la sospecha de que en promedio se puede decir que la diferencia entre el peso mayor y el peso menor es constante (es decir, no depende del peso del huevo menor observado). Usando el modelo en I) realice una prueba de hipótesis para responder la pregunta de los investigadores.
- III. Posteriormente se observa el peso de los huevos de una nueva nidada, observándose un peso de 75 y 130 gramos. Usando un intervalo adecuado, comente sobre la sospecha de que la nidada de huevos sí proviene de pingüinos *Macaroni*.

```
x=c(79, 93, 100, 105, 101, 96, 96, 109, 70, 71, 87)
y=c(133, 148, 164, 171, 165, 159, 162, 170, 127, 133, 148 )
Datos6=data.frame(cbind(x,y))
kable(t(Datos6)) %>%
  kable_styling(bootstrap_options = "striped", full_width = F)
```

x	79	93	100	105	101	96	96	109	70	71	87
y	133	148	164	171	165	159	162	170	127	133	148

## 7.

Considere los datos en la base *performance.csv* y las variables:  $y$  = academic performance of the school (api00) y  $x$  = percentage of English language learners (ell). Estos datos corresponden a una muestra aleatoria de 400 escuelas primarias en California, en donde por escuela se realizaron mediciones que tienen que ver con su desempeño en 2000.

- I. Ajustar un modelo de regresión lineal simple. Verificar los supuestos a partir de este modelo. Deberá indicar para cada supuesto qué gráfica o prueba sirve para argumentar el cumplimiento o no del supuesto.
- II. En caso de que alguno de los supuestos no se satisfaga en I), realizar modificaciones a las variables para encontrar un modelo en donde sí se satisfagan los supuestos:
  - a. Para transformar la variable Y, probar con transformaciones Box-Cox u otras conocidas como  $\log()$  o  $\exp()$ .
  - b. Para transformar la variable X, probar con transformaciones Box-Tidwell u otras conocidas como  $\log()$  o  $\exp()$ .
  - c. Recuerde que siempre se puede sumar una constante (e.g. +1) para hacer positivas a las variables.

Al finalizar, deberá indicar el modelo de regresión lineal simple que se ajustará, haciendo explícito qué variables fueron transformadas y cómo. También deberá indicar para cada supuesto del modelo de regresión qué gráfica o prueba sirve para argumentar su cumplimiento.

- III. En una misma gráfica incluir los puntos en escala original, la recta de regresión del modelo en I) y la curva del modelo en II).
- IV. Interpretar  $R^2$  y la prueba anova del modelo en II).
- v. Con el modelo final ayude a un investigador a argumentar a favor o en contra de la hipótesis: “A mayor porcentaje de estudiantes que requieren cursos de recuperación de inglés es menor el desempeño de la escuela”.