

# Capítulo 8

## Técnicas de validación estadística

### 8.1. Introducción

Dado un conjunto de observaciones obtenidas a partir de una recolección de datos puede ser de interés determinar cuál es la distribución de estos datos. También puede haberse realizado una simulación que modele un sistema real, y luego se desea testear si los datos simulados tienen la misma distribución que los datos que se observaron en la realidad.

Una forma de determinar si un conjunto de observaciones proviene de una distribución dada es a través de las pruebas de **bondad de ajuste**. Una prueba de bondad de ajuste es un test de hipótesis, en la cual la hipótesis nula,  $H_0$ , afirma que los datos provienen de una determinada distribución  $F$ . La hipótesis alternativa,  $H_1$ , es la negación de  $H_0$ .

Según cuál sea la hipótesis se define un determinado **estadístico muestral**  $T = T(X_1, X_2, \dots, X_n)$ . El estadístico es una variable aleatoria, que, bajo la hipótesis nula, tiene una distribución conocida o de la cual se saben algunas propiedades. Esto es, se conoce algo de  $P_{H_0}(T \leq t)$ , donde el subíndice  $H_0$  indica que la distribución de  $T$  está basada en que las observaciones satisfacen la hipótesis  $H_0$ .

Así, dada una muestra de datos  $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ , se evalúa el estadístico en esta muestra y se toma una decisión de rechazar o no rechazar la hipótesis nula según "cuán probable" es haber obtenido ese valor. En algunos casos se trata de medir si el valor obtenido es demasiado alto. Así por ejemplo, si el valor obtenido es  $T = t$ , y

$$p = P_{H_0}(T \geq t) \leq \alpha \quad (8.1)$$

entonces se **rechaza** la hipótesis nula con un nivel de rechazo  $\alpha$ . En otros casos se analiza si el valor obtenido no es muy bajo ni muy alto, y en esos casos se considera:

$$p = \min\{P_{H_0}(T \geq t), P_{H_0}(T \leq t)\}. \quad (8.2)$$

Las cantidades  $p$  en (8.1) y (8.2) se denominan también  $p$ -valor. Un  $p$ -valor pequeño es un indicador de rechazo de la hipótesis nula.

## 8.2. Pruebas de bondad de ajuste

### 8.2.1. Datos discretos - Test chi-cuadrado de Pearson

Denotamos con  $Y_1, Y_2, \dots, Y_n$  una muestra de observaciones independientes, que toman alguno de los valores en el conjunto  $\{1, 2, \dots, k\}$ .

#### Parámetros especificados

Supongamos que se desea testear si los datos observados provienen de una determinada distribución teórica  $F$  conocida, y sea  $X$  una variable aleatoria con distribución  $F$ . Llamamos

$$p_i = P(X = i), \quad N_i = \#\{j \mid Y_j = i, 1 \leq j \leq n\}.$$

Esto es,  $p_i$  es la probabilidad que una variable con distribución  $F$  tome el valor  $i$ , y  $N_i$  es la frecuencia con la que el valor  $i$  aparece en la muestra, es decir, la **frecuencia observada**.

Si los datos provienen realmente de la distribución  $F$ , es de esperar que  $N_i$  sea próximo a  $np_i$ , por lo cual podría considerarse

$$(N_i - np_i)^2$$

como una medida de cuán próximos están los datos de la distribución teórica. Ahora bien, notemos que si por ejemplo  $(N_i - np_i)^2 = 1$ , este valor 1 será mucho más significativo si  $np_i = 2$  que si  $np_i = 100$ . Por ello, es más adecuado considerar para cada  $i$  el valor

$$\frac{(N_i - np_i)^2}{np_i}$$

como una medida de distancia entre la distribución empírica de los datos y la distribución  $F$ . En particular, el estadístico para el **test chi cuadrado de Pearson** está dado por:

$$T = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i}.$$

Si el valor de  $T$  es grande, se considera que hay evidencias que la muestra no proviene de la distribución  $F$ : **se rechaza la hipótesis nula**. Por el contrario, si el valor de  $T$  es pequeño, no hay evidencias suficientes para rechazar la hipótesis. Si la hipótesis nula es cierta y  $n$  es grande, entonces el estadístico  $T$  tiene una distribución  $\chi$ -cuadrado con  $k - 1$  grados de libertad:  $\chi_{k-1}^2$ .

Si el nivel de rechazo es del 5 %, entonces un  $p$ -valor menor que 0.05 es indicativo que la muestra no proviene de la distribución  $F$ , es decir, que debe rechazarse la hipótesis nula. Por el contrario, valores de  $p$  más grandes no dan evidencia que se deba rechazar la hipótesis.

Si en cambio se rechaza al 1 % entonces un  $p$ -valor menor a 0.01 indicará que se debe rechazar la hipótesis nula.

**Ejemplo 8.1.** Supongamos que se ha tomado una muestra de 100 datos (o se han simulado 100 valores), que toman alguno de los valores entre 0 y 7. Las frecuencias observadas  $N_i$  son las siguientes:

$$2, 7, 20, 22, 24, 23, 0, 2.$$

Se quiere testear la hipótesis que estos valores provienen de una distribución binomial  $Bi(7, 0.5)$  con un nivel de rechazo del 5 %.

El **diseño** de este test implica establecer la hipótesis nula, el estadístico a utilizar y la fórmula del  $p$ -valor que será utilizada para decidir el rechazo o no de la hipótesis nula. En este caso la hipótesis nula es:

**H0)** La muestra proviene de una variable  $Y$  tal que  $P(Y = i) = P(Bin(Y, 0.5) = i)$ .

El **estadístico** a utilizar para esta prueba de hipótesis es

$$T = \sum_{i=0}^7 \frac{(N_i - np_i)^2}{np_i},$$

donde  $N_i$  son las frecuencias observadas en la muestra,  $n$  es el tamaño de la muestra,  $k = 8$  es el número de agrupamientos de valores considerados y las probabilidades teóricas  $p_i$  están dadas por:

$$[0.0078125, 0.0546875, 0.1640625, 0.2734375, 0.2734375, 0.1640625, 0.0546875, 0.0078125].$$

Si  $T = t$  es el valor del estadístico, el  $p$ -valor es:

$$P_{H_0}(T \geq t) = P(\chi_{k-1}^2 \geq t).$$

Así, para esta muestra en particular las frecuencias esperadas son:

$$[0.78125, 5.46875, 16.40625, 27.34375, 27.34375, 16.40625, 5.46875, 0.78125].$$

El valor del estadístico  $T$  es:

$$T = \frac{(2 - 0.78125)^2}{0.78125} + \frac{(7 - 5.46875)^2}{.46875} + \dots + \frac{(0 - 5.46875)^2}{5.46875} + \frac{(2 - 0.78125)^2}{0.78125} = \textcolor{red}{14.59}.$$

Como hemos considerado 8 términos en el estadístico  $T$ , entonces debemos testear con una distribución  $\chi_{8-1}^2$ . En particular,

$$P(\chi_7^2 \geq 14.59) \approx 0.04.$$

Para un nivel de confianza del 95 % la hipótesis nula se rechaza porque el  $p$ -valor es menor al 5 %. Para un nivel de confianza del 99 % la hipótesis nula no se rechaza.

**Simulación del  $p$ -valor**

Si el  $p$ -valor obtenido es muy próximo al nivel de rechazo, puede existir la duda si conviene o no rechazar la hipótesis. Esto es por ejemplo, si el nivel de rechazo es  $\alpha = 0.05$ , el valor del estadístico obtenido es  $T = t_0$  y se tiene que

$$P_{H_0}(T \geq t_0) \sim 0.05.$$

Una forma de ayudar a tomar esta decisión es simular muestras de tamaño  $n$  de la distribución  $F$  y para cada una de ellas calcular el estadístico  $T$ . Para un número de simulaciones suficientemente grande, la proporción de valores de  $T$  que exceden al valor  $T = t_0$  tomado en la muestra original es una buena estimación del  $p$ -valor.

En caso que  $n$  sea muy grande en relación a  $k$ , es conveniente generar directamente las frecuencias observadas. Esto es, supongamos que la variable aleatoria  $X$  toma valores  $1, 2, \dots, k$ , con  $p_i = P(X = i)$ , entonces se simulan los valores  $N_1, N_2, \dots, N_k$ . Como  $N_1$  es la cantidad de datos iguales a 1 en una muestra de tamaño  $n$ , entonces  $N_1$  tiene distribución binomial  $\text{Bin}(n, p_1)$ .

Una vez generado  $N_1$ , se genera  $N_2$  que es la cantidad de datos restantes ( $n - N_1$ ) iguales a 2. Dado que ya se han contado los datos iguales a 1, cada uno de los datos restantes tomará el valor 2 con probabilidad:

$$\lambda_2 = P(X = 2 \mid X \neq 1) = \frac{P(X = 2)}{P(X \neq 1)} = \frac{p_2}{1 - p_1}.$$

Por lo tanto  $N_2$  tiene distribución binomial  $\text{Bin}(n - N_1, \frac{p_2}{1 - p_1})$ .

Los siguientes  $n - N_1 - N_2$  datos tomarán el valor 3 con probabilidad:

$$\lambda_3 = P(X = 3 \mid X \neq 1, X \neq 2) = \frac{p_3}{1 - p_1 - p_2}.$$

Así siguiendo, las variables  $N_j$  condicionadas a los valores obtenidos previamente tienen distribución binomial:

$$N_j \sim \text{Bin}(n - (N_1 + N_2 + \dots + N_{j-1}), \lambda_j), \quad \lambda_j = \frac{p_j}{1 - P(X < j)}.$$

Notar que  $\lambda_k = 1$ , por lo cual  $N_k = n - \sum_{j=1}^{k-1} N_j$ .

Así, para estimar el valor  $p$  se generan directamente los valores  $N_1, N_2, \dots, N_k$  y se calcula el estadístico  $T$ . Repitiendo este procedimiento una cierta cantidad de veces, el  $p$ -valor se calcula como la proporción de valores que superan el valor  $T = t$  en la muestra original.

**Parámetros no especificados**

Las pruebas de bondad de ajuste también pueden aplicarse si los parámetros de la distribución  $F$  no son todos conocidos. Por ejemplo, se podría testear si los datos provienen de una distribución de Poisson  $\mathcal{P}(\lambda)$ , desconociendo  $\lambda$ .

En este caso, se estima el o los parámetros no especificados. Esto determinará una cierta distribución  $\hat{F}$ . Por ejemplo, si se estima  $\lambda$  en la Poisson, se tendrá una distribución  $\hat{F} = \mathcal{P}(\hat{\lambda})$ .

Sea  $\hat{p}_i = P_{\hat{F}}(X = i)$ , donde el subíndice indica que  $X$  tiene distribución  $\hat{F}$ . El estadístico es en este caso:

$$T = \sum_{j=1}^k \frac{(N_i - n\hat{p}_i)^2}{n\hat{p}_i}.$$

Sea  $m$  el número de parámetros que se utilizan para el cálculo de  $p_i$  y que deben ser estimados. Es decir, no están especificados. Se puede demostrar que, para  $n$  suficientemente grande, el estadístico  $T$  tiene aproximadamente una distribución chi-cuadrado con  $k - 1 - m$  grados de libertad. En particular, el  $p$ -valor puede estimarse como

$$p - \text{valor} \approx P(\chi_{k-1-m}^2 \geq t).$$

En caso de utilizar simulaciones para estimar el  $p$ -valor, el procedimiento es como sigue:

1. Supongamos que la hipótesis nula  $H_0$  es que los datos  $Y_1, \dots, Y_n$  provienen de una distribución  $F$ , y asumamos que existen  $m$  parámetros de esta distribución que son desconocidos:  $\theta_1, \dots, \theta_m$ .
2. A partir de la muestra de datos, se estiman los parámetros obteniendo valores  $\hat{\theta}_1, \dots, \hat{\theta}_m$ . Esto determina una distribución  $\hat{F}$  y una probabilidad  $\hat{p}_i$  para cada valor  $i$  de la variable aleatoria. A partir de estas estimaciones se calcula el estadístico

$$T = \sum_{i=1}^k \frac{(N_i - n\hat{p}_i)^2}{n\hat{p}_i}.$$

Llamamos  $t$  al valor obtenido.

3. En cada simulación, se generan  $n$  datos a partir de la distribución  $\hat{F}$ . Luego se vuelven a estimar los parámetros  $\theta_1, \dots, \theta_m$  obteniendo estimaciones  $\theta_1(\text{sim}), \dots, \theta_m(\text{sim})$  a partir de la muestra simulada. Estos parámetros determinan una distribución  $F_{\text{sim}}$ . Con estas estimaciones se calculan las probabilidades  $p_i(\text{sim})$ , es decir,  $p_i(\text{sim}) = P_{F_{\text{sim}}}(X = i)$  si  $X$  tiene distribución  $F_{\text{sim}}$ . Luego se calcula el estadístico utilizando las probabilidades  $p_i(\text{sim})$ :

$$T_{\text{sim}} = \sum_{i=1}^k \frac{(N_i - n\hat{p}_i(\text{sim}))^2}{n\hat{p}_i(\text{sim})}.$$

4. El  $p$ -valor se estima como la proporción de  $T_{\text{sim}}$  mayores o iguales a  $t$ .

**Ejemplo 8.2.** Supongamos que a lo largo de 30 días han habido

- 6 días en que no ocurrió ningún accidente,
- 2 días en los que ocurrió 1 accidente,
- 1 en el que ocurrieron 2 accidentes,
- 9 días en que ocurrieron 3 accidentes,
- 7 en que ocurrieron 4 accidentes,
- 4 que ocurrieron 5, y
- 1 en que ocurrieron 8 accidentes.

Se quiere testear la hipótesis que los datos provienen de una distribución de Poisson  $\mathcal{P}(\lambda)$ . Es decir, que el número de accidentes por día tiene distribución Poisson, pero la tasa  $\lambda$  es desconocida y por lo tanto se estimará a partir de los datos. Como  $\lambda$  representa la media o valor esperado, entonces puede estimarse con la media muestral:

$$\hat{\lambda} = \frac{6 \cdot 0 + 2 \cdot 1 + 1 \cdot 2 + 9 \cdot 3 + 7 \cdot 4 + 4 \cdot 5 + 1 \cdot 8}{30} = \frac{87}{30} = 2.9.$$

Por otra parte, como una variable Poisson toma infinitos valores, decidiremos agruparlos en una cantidad finita. Podemos agrupar los datos en 6 grupos: los que toman el valor 0, 1, 2, 3, 4 y los mayores o iguales a 5. Así tendremos frecuencias observadas:

$$N_0 = 6, \quad N_1 = 2, \quad N_2 = 1, \quad N_3 = 9, \quad N_4 = 7, \quad N_5 = 5.$$

De esta manera, en el diseño del test la hipótesis nula es:

**H0:** Los datos provienen de una distribución  $Y$  tal que:

$$\hat{p}_i = P(Y = i) = e^{-\hat{\lambda}} \frac{\hat{\lambda}^i}{i!}, \quad 0 \leq i \leq 4, \quad \hat{p}_5 = P(Y = 5) = 1 - \sum_{j=1}^4 P(Y = j),$$

con  $\hat{\lambda}$  igual a la media de los datos observados. El estadístico a utilizar será:

$$T = \sum_{j=0}^4 \frac{(N_j - np_j)^2}{np_j},$$

y el  $p$ -valor se determinará a partir del valor  $T = t$  observado en la muestra como:

$$p = P_{H_0}(T \geq t) = P(\chi_{5-1-1}^2 \geq t) = P(\chi_3^2 \geq t).$$

Luego tenemos que:

$$\hat{p}_0 = 0.05, \quad \hat{p}_1 = 0.1596, \quad \hat{p}_2 = 0.2312, \quad \hat{p}_3 = 0.2237, \quad \hat{p}_4 = 0.1622, \quad \hat{p}_5 = 0.1682.$$

El valor del estadístico  $T$  está dado por:

$$\begin{aligned} T &= \frac{(6 - 30\hat{p}_0)^2}{30\hat{p}_0} + \frac{(2 - 30\hat{p}_1)^2}{30\hat{p}_1} + \frac{(1 - 30\hat{p}_2)^2}{30\hat{p}_2} + \frac{(9 - 30\hat{p}_3)^2}{30\hat{p}_3} + \frac{(7 - 30\hat{p}_4)^2}{30\hat{p}_4} + \frac{(5 - 30\hat{p}_5)^2}{30\hat{p}_5} \\ &= 19.887. \end{aligned}$$

Dado que el estadístico tiene  $k = 6$  sumandos y se ha estimado  $m = 1$  parámetro, el  $p$ -valor se estima con una chi cuadrado de  $k - 1 - m = 4$  grados de libertad::

$$p - \text{valor} \approx P(\chi_4^2 \geq 19.887) \sim 0.0005.$$

Para un nivel de rechazo  $\alpha = 0.01$  la hipótesis nula es rechazada.

En una simulación para el cálculo del  $p$ -valor se generarán  $N$  muestras de tamaño 30 de una  $X \sim \mathcal{P}(2.9)$ . Por ejemplo, si en la simulación  $j$  se obtienen los valores:

$$3 \ 3 \ 3 \ 1 \ 6 \ 3 \ 4 \ 4 \ 1 \ 6 \ 5 \ 6 \ 2 \ 1 \ 5 \ 3 \ 8 \ 4 \ 1 \ 4 \ 1 \ 2 \ 7 \ 1 \ 2 \ 2 \ 2 \ 3 \ 4 \ 2,$$

entonces

$$N_0 = 0, \quad N_1 = 6, \quad N_2 = 6, \quad N_3 = 6, \quad N_4 = 5, \quad N_5 = 7,$$

y  $\lambda(\text{sim}) = \frac{99}{30} = 3.3$ . Para este valor de  $\lambda(\text{sim})$  se calculan las correspondientes probabilidades  $p(\text{sim})$ :

$$\hat{p}_i(\text{sim}) = e^{-3.3} \cdot \frac{3.3^i}{i!}, \quad 1 \leq i \leq 4, \quad \hat{p}_5(\text{sim}) = 1 - \sum_{i=0}^4 \hat{p}_i(\text{sim})$$

y se calcula el valor del estadístico:

$$T_{\text{sim}} = \sum_{i=0}^5 \frac{(N_i - n\hat{p}_i(\text{sim}))^2}{n\hat{p}_i(\text{sim})} = 2.7165.$$

El siguiente segmento de código estima el  $p$ -valor con 10000 simulaciones.

---

```

datos = np.zeros(30, int)    #muestras
N = np.zeros(6, int)        #frecuencias observadas
p = np.zeros(6, float)      #p(sim)
pvalor = 0
for _ in range(10000):
    for j in range(30):
        datos[j] = Poisson(2.9)
    N *= 0

```

---

```

for observacion in datos:
    if observacion < 5:
        N[observacion] += 1
    else:
        N[5] += 1
lamda = sum(datos)/len(datos)
for i in range(5):
    p[i] = exp(-lamda)*lamda**i/factorial(i)
p[5] = 1 - sum(p[:5])
T = 0
for i in range(6):
    T += (N[i]-30*p[i])**2/(30*p[i])
if T >= 19.887012:
    pvalor += 1
print('p valor: ', pvalor/10000)

```

---

### 8.2.2. Datos continuos - Test de Kolmogorov-Smirnov

Si las observaciones provienen de datos de tipo continuo, puede aplicarse también el test  $\chi$ -cuadrado realizando una discretización. Esto es, pueden agruparse los datos en  $k$  intervalos consecutivos:

$$(-\infty, y_1], (y_1, y_2], (y_2, y_3], \dots, (y_{k-1}, \infty),$$

y considerar  $N_i$  como el número de observaciones en el intervalo  $i$ , y  $p_i$  la probabilidad dada por la distribución  $F$  de que la variable esté en el  $i$ -ésimo intervalo:

$$\begin{aligned}
 p_0 &= F(y_1), \\
 p_i &= F(y_i) - F(y_{i-1}), \quad 1 < i < k, \\
 p_k &= 1 - F(y_k).
 \end{aligned}$$

Este método, si bien puede utilizarse, tiene la desventaja de agrupar los datos en intervalos y no considerar cómo se distribuyen estos datos dentro del intervalo. El test de Kolmogorov-Smirnov resulta más adecuado para datos de tipo continuo.

#### Con parámetros especificados

Consideramos al igual que antes una muestra  $Y_1, Y_2, \dots, Y_n$  de datos que se suponen independientes, y la hipótesis nula está dada por:

$$H_0 : \text{los datos provienen de la distribución continua } F.$$



En primer lugar se ordenan los datos de menor a mayor. Con  $Y_{(j)}$  denotamos al dato que ocupa el  $j$ -ésimo lugar luego del ordenamiento. Se considera luego la distribución empírica de los datos,  $F_e$ , donde

$$F_e(x) = \frac{\#\{j \mid Y_j \leq x\}}{n}.$$

En particular, si se asumen todos los datos distintos, se tiene que

$$F_e(x) = \begin{cases} 0 & x < Y_{(1)} \\ \frac{j}{n} & Y_{(j)} \leq x < Y_{(j+1)}, \quad 1 \leq j < n \\ 1 & x \geq Y_{(n)}. \end{cases}$$

El test de Kolmogorov-Smirnov esencialmente compara la distribución empírica de los datos con la distribución  $F$ , estimando la distancia máxima entre los dos gráficos. Así, el **estadístico de Kolmogorov-Smirnov** está dado por:

$$D = \sup_{x \in \mathbb{R}} |F_e(x) - F(x)| \quad (8.3)$$

$$= \sup_{x \in \mathbb{R}} \left\{ \sup_{x \in \mathbb{R}} \{F_e(x) - F(x)\}, \sup_{x \in \mathbb{R}} \{F(x) - F_e(x)\} \right\}. \quad (8.4)$$

Dado que  $|F_e(x) - F(x)|$  no es una función continua en todos los reales, no podemos asegurar que alcance un máximo propiamente. Sin embargo, por tomar valores en un subconjunto acotado de  $\mathbb{R}$  podemos garantizar la existencia de un supremo.

Como  $F_e(Y_{(n)}) = 1$ , y  $F(x) \leq 1$  para cualquier  $x$ , entonces  $\sup_x \{F_e(x) - F(x)\}$  es no negativo. Además, como  $F$  es monótona creciente en el intervalo donde no vale 0 ni 1, entonces  $F_e(x) - F(x)$  es decreciente en los intervalos donde  $F_e$  es constante. En particular,  $F_e(x) - F(x)$  alcanza el máximo en alguno de los  $n$  puntos  $Y_{(j)}$ . Luego

$$\sup_{x \in \mathbb{R}} \{F_e(x) - F(x)\} = \max_{1 \leq j \leq n} \left\{ \frac{j}{n} - F(Y_{(j)}) \right\}.$$

Por otra parte,  $F(x) \geq 0$  para todo  $x$  y  $F_e(x) = 0$  si  $x < Y_{(1)}$ , por lo tanto el máximo de  $F(x) - F_e(x)$  es no negativo. Por otra parte, como  $F$  es creciente en los intervalos donde  $F_e$  es constante, entonces  $F(x) - F_e(x)$  tiene una discontinuidad de salto en cada  $Y_{(j)}$ , y podría decirse que el supremo se alcanza justo antes de un valor  $Y_{(j)}$ . Este supremo es igual a  $F(Y_{(j)}) - F_e(Y_{(j-1)})$ . Luego:

$$\sup_{x \in \mathbb{R}} (F(x) - F_e(x)) = \max_{1 \leq j \leq n} \left\{ F(Y_{(j)}) - \frac{j-1}{n} \right\}.$$

Finalmente, el estadístico  $D$  en (8.3) puede escribirse como:

$$D = \max_{1 \leq j \leq n} \left\{ \frac{j}{n} - F(Y_{(j)}), F(Y_{(j)}) - \frac{j-1}{n} \right\}.$$

Bajo la hipótesis  $H_0$ , el estadístico  $D$  sigue una distribución de Kolmogorov, con una expresión bastante compleja y de la cual se tienen algunos valores tabulados para  $\lim_{n \rightarrow \infty} \sqrt{n}D^1$ . Por lo tanto, en la práctica es conveniente calcular el  $p$ -valor con simulaciones. Esto es, realizar  $k$  simulaciones de muestras de tamaño  $n$  de una variable con distribución  $F$ , calcular el correspondiente valor del estadístico  $D = d_i$  para cada muestra,  $1 \leq i \leq k$ . Finalmente, se estima el  $p$ -valor a la proporción de valores  $d_i$  que exceden al valor  $d$ :

$$p - \text{valor} = \frac{\#\{i \mid d_i > d\}}{k}.$$

Una simplificación de este paso es el hecho que la distribución de  $D$  es independiente de la distribución  $F$ . Esto es, si  $Y_1, Y_2, \dots, Y_n$  denota una muestra de tamaño  $n$  de una variable con distribución  $F$ , y  $X_1, X_2, \dots, X_n$  denota una muestra de tamaño  $n$  de una variable con distribución  $G$ , y definimos los estadísticos  $D_F$  y  $D_G$  por:

$$D_F = \max_{1 \leq j \leq n} \left\{ \frac{j}{n} - F(Y_{(j)}), F(Y_{(j)}) - \frac{j-1}{n} \right\},$$

$$D_G = \max_{1 \leq j \leq n} \left\{ \frac{j}{n} - G(X_{(j)}), G(X_{(j)}) - \frac{j-1}{n} \right\},$$

entonces para cualquier  $d$  se cumple que

$$P_F(D_F \geq d) = P_G(D_G \geq d).$$

Resumimos este hecho en el siguiente teorema:

**Teorema 8.1.** La probabilidad  $P_F(D \geq d)$  es la misma para cualquier distribución continua  $F$ .

*Demostración.* Denotamos con  $D$  al estadístico de Kolmogorov-Smirnov, y sea  $d$  el  $p$ -valor obtenido a partir de  $n$  observaciones independientes de una distribución continua  $F$ :

$$P_F(D \geq d) = P_F \left( \sup_{x \in \mathbb{R}} \left| \frac{\#\{i \mid Y_i \leq x\}}{n} - F(x) \right| \geq d \right).$$

Aquí el subíndice  $F$  hace referencia a la hipótesis nula. Esto es,  $P_F(D \geq d)$  indica la probabilidad de que el supremo de  $|F_e(x) - F(x)|$  sea mayor a  $d$  dado que las observaciones provienen de la distribución  $F$ .

Dado que  $F$  es no decreciente, entonces  $Y_i \leq x$  si y sólo si  $F(Y_i) \leq F(x)$ . Por otra parte,  $F(x)$  toma todo el rango de valores en  $[0, 1]$  si  $x$  toma cualquier valor en  $\mathbb{R}$ . Luego podemos sustituir la variable  $F(x)$  para  $x \in \mathbb{R}$  por una variable  $y$ , con  $0 \leq y \leq 1$ . Así resulta:

$$P_F(D \geq d) = P \left( \sup_{0 \leq y \leq 1} \left| \frac{\#\{i \mid F(Y_i) \leq y\}}{n} - y \right| \geq d \right).$$

<sup>1</sup>Smirnov N (1948). *Table for estimating the goodness of fit of empirical distributions*. Annals of Mathematical Statistics. 19: 279–281. doi:10.1214/aoms/1177730256.

Por otra parte, ya hemos visto que si  $Y$  tiene distribución  $F$ , entonces  $F(Y)$  tiene distribución uniforme en  $(0, 1)$ . Por lo tanto, si  $Y_1, Y_2, \dots, Y_n$  son observaciones independientes de una variable con distribución  $F$ , entonces  $F(Y_1), F(Y_2), \dots, F(Y_n)$  son observaciones independientes de una variable  $U$  con distribución uniforme. Entonces, si  $G$  es la función de distribución uniforme en  $(0, 1)$  tenemos que:

$$\begin{aligned} P_F(D \geq d) &= P \left( \sup_{0 \leq y \leq 1} \left| \frac{\#\{i \mid F(Y_i) \leq y\}}{n} - y \right| \geq d \right) \\ &= P \left( \sup_{0 \leq y \leq 1} \underbrace{\left| \frac{\#\{i \mid U_i \leq y\}}{n} - y \right|}_{G_e(y) - G(y)} \geq d \right) \\ &= P_G(D \geq d). \end{aligned}$$

Luego, el  $p$ -valor obtenido bajo la hipótesis nula que los datos provienen de la distribución continua  $F$  es el mismo que si los datos provienen de la distribución uniforme. Como  $F$  es cualquier distribución continua, concluimos que el  $p$ -valor es independiente de  $F$ .  $\square$

Volviendo a la estimación del  $p$ -valor a través de simulaciones, el Teorema 8.1 permite utilizar muestras a partir de la distribución uniforme en lugar de muestras de distribución  $F$ : Esto es, una vez observado el estadístico  $D = d$  con la muestra  $Y_1, Y_2, \dots, Y_n$ , se realizan  $k$  simulaciones de muestras de tamaño  $n$  de una variable uniforme  $U \sim \mathcal{U}(0, 1)$ . Para cada una de estas muestras simuladas, se calcula el correspondiente estadístico  $d_i$ ,  $1 \leq i \leq k$ . Notemos que para la distribución uniforme,  $G(u) = u$  para  $u \in (0, 1)$ . Luego el estadístico  $D$  está dado por:

$$D = \max_{1 \leq j \leq n} \left\{ \frac{j}{n} - U_{(j)}, U_{(j)} - \frac{j-1}{n} \right\}.$$

De esta manera, para estimar el  $p$ -valor a través de simulaciones es suficiente con generar muestras de tamaño  $n$  de variables aleatorias uniformes en  $(0, 1)$  y calcular la proporción de valores  $d_i$  que exceden a  $d$ .

**Ejemplo 8.3.** Se quiere testear la hipótesis que una determinada muestra proviene de una distribución exponencial con media 100:

$$F(x) = 1 - e^{-x/100}.$$

Los valores ordenados para una muestra de tamaño 10 para esta distribución son:

55, 72, 81, 94, 112, 116, 124, 140, 145, 155,

¿qué conclusión puede obtenerse con un nivel de rechazo del 5 %?

La siguiente tabla resume los valores que deben analizarse para determinar el estadístico  $D$ . Para hacer más clara la lectura se han considerado sólo tres decimales. En particular, se observa que el valor máximo de  $|F_e(x) - F(x)|$  se alcanza en  $x = 55$ :

$j$	$Y_{(j)}$	$F(Y_{(j)})$	$\frac{j}{n} - F(Y_{(j)})$	$F(Y_{(j)}) - \frac{j-1}{n}$
1	55	0.423	-0.323	<b>0.423</b>
2	72	0.513	-0.313	0.413
3	81	0.555	-0.255	0.355
4	94	0.609	-0.209	0.309
5	112	0.674	-0.174	0.274
6	116	0.687	-0.087	0.187
7	124	0.711	-0.011	0.111
8	140	0.753	0.047	0.053
9	145	0.765	0.135	-0.035
10	155	0.788	0.212	-0.112
$d = 0.423$				

La estimación del  $p$ -valor para este caso se hará generando  $k$  (por ejemplo  $k = 10000$ ) muestras de tamaño 10 de una distribución uniforme. Para la  $i$ -ésima muestra, se calcula el valor  $d_i$  del estadístico  $D$ :

$$d_i = \max_{1 \leq j \leq 10} \left\{ \frac{j}{10} - U_{(j)}, U_{(j)} - \frac{j-1}{10} \right\}.$$

---

```

d_KS = 0.423 #estadístico
pvalor = 0.
n = 10
Nsim = 10000
for _ in range(Nsim):
    uniformes = random.uniform(0,1,n)
    uniformes.sort()
    d_j = 0
    for j in range(n):
        u_j = uniformes[j]
        d_j = max(d_j, (j+1)/n-u_j, u_j-j/n)
    if d_j >= d_KS:
        pvalor += 1
print(pvalor/Nsim)

```

---

El  $p$ -valor calculado es:

$$p - \text{valor} \approx \frac{\#\{i \mid d_i \geq 0.423\}}{10000} = 0.0407.$$

Para un  $\alpha = 0.05$ , (confianza del 95 %), la hipótesis nula es rechazada.

**Parámetros no especificados.**

Si se quiere testear la hipótesis que las observaciones  $Y_1, Y_2, \dots, Y_n$  provienen de una distribución  $F$  con ciertos parámetros desconocidos, entonces en primer lugar se estiman los parámetros  $\theta_1, \theta_2, \dots, \theta_m$ , y luego se calcula el estadístico de Kolmogorov-Smirnov:

$$D = \sup_{x \in \mathbb{R}} |F_e(x) - F_{\hat{\theta}}(x)|,$$

donde  $F_e$  es la distribución empírica de los datos, y  $F_{\hat{\theta}}$  es la función de distribución obtenida con la estimación de los parámetros  $\theta$ .

Si el valor de  $D$  que se obtiene es  $d$ , entonces se puede aproximar el  $p$ -valor como en el caso de parámetros especificados:

$$P_{F_{\hat{\theta}}}(D \geq d) = P_U(D \geq d),$$

donde  $U \sim \mathcal{U}(0, 1)$ .

En caso que el  $p$ -valor resultara en el área de rechazo ( $< 0.05$  por ejemplo), es conveniente realizar una segunda simulación más certera. Más específicamente:

1. Se generan  $N$  simulaciones de muestras de tamaño  $n$ , generadas a partir del  $F_{\hat{\theta}}$ .
2. Para cada una de estas muestras  $sim$ ,  $1 \leq sim \leq N$ :

$$X_{1,sim}, X_{2,sim}, \dots, X_{n,sim},$$

se vuelven a estimar los parámetros. Llamamos  $\hat{\theta}_1(sim), \hat{\theta}_2(sim), \dots, \hat{\theta}_m(sim)$ . Con estas estimaciones se calcula el estadístico de Kolmogorov Smirnov a partir de la distribución empírica de la muestra simulada, y la distribución  $F_{\hat{\theta}(sim)}$ :

$$d_{sim} = \sup_{x \in \mathbb{R}} |F_{e,sim}(x) - F_{\hat{\theta}(sim)}(x)|.$$

3. La proporción de valores  $d_{sim}$  que superen el valor  $d$  de la muestra original será la estimación del  $p$ -valor.

**Ejemplo 8.4.** En el Ejemplo 8.3 podría validarse la hipótesis de que los datos provienen de una distribución exponencial, pero estimar la media a partir de la muestra. Tendríamos entonces

$$\hat{\theta} = \bar{X}(10) = 109.4.$$

El estadístico de Kolmogorov Smirnov para este valor del parámetro resulta  $d_{KS} = 0.3951$ , y el  $p$ -valor obtenido usando uniformes es:

$$P(D \geq 0.3951) \sim 0.0658.$$

Con este  $p$ -valor no se rechazaría la hipótesis nula a un nivel de rechazo del 5 %. Sin embargo, al ser un valor cercano a 0.05 se podría realizar una segunda simulación que no utilice uniformes sino muestras de exponenciales con media 109.4. En cada simulación el parámetro  $\theta$  se estima nuevamente.

---

```

def F_exponencial(x, lamda):
    'Distr. acumulada de la exponencial'
    return 1-exp(-x*lamda)
def K_S(datos, theta):
    'Estadístico de Kolmogorov Smirnov'
    n = len(datos)
    d = 0
    for j in range(n):
        x = datos[j]
        d = max(d, (j+1)/n-F_exponencial(x, theta),
                F_exponencial(x, theta)-j/n)
    return d
....
d_KS = 0.3951 #estadístico
landa = 1./109.4
pvalor = 0.
Nsim = 10000
for _ in range(Nsim):
    muestra = []
    for _ in range(n):
        muestra.append(-log(1-random())/landa)
    muestra.sort()
    landa_est = n/sum(muestra)
    d_j = K_S(muestra, landa_est)
    if d_j >= d_KS:
        pvalor += 1
pvalor = pvalor/Nsim

```

---

El  $p$ -valor obtenido es ahora:

$$P(D \geq 0.3951) \sim 0.0064,$$

que fortalece la decisión de rechazar la hipótesis nula.

### 8.3. El problema de las dos muestras

En una simulación es posible generar valores de una variable aleatoria que sean útiles para el modelo, pero que no necesariamente se conozca su distribución. Por ejemplo: el tiempo total de permanencia de clientes en un servidor a lo largo de un día, o la cantidad de clientes que llegan en determinada franja horaria. Aún desconociendo la distribución, lo que sí es deseable es que

los datos que se simulan sean coherentes con las observaciones que se han obtenido del modelo real. Esto es, si se ha simulado una muestra de valores provenientes de una distribución:  $X_1, X_2, \dots, X_m$ , y se tiene una muestra observada  $Y_1, Y_2, \dots, Y_n$ , también de una distribución, interesa conocer si el conjunto de las  $n + m$  variables corresponden a observaciones independientes y si provienen de la misma distribución.

En general, el **problema de las dos muestras** considera dos muestras de observaciones que provienen de distribuciones  $F_1$  y  $F_2$  respectivamente, y se trata de validar la siguiente hipótesis:

$H_0$ : Las  $n + m$  variables  $X_1, X_2, \dots, X_n, Y_1, Y_2, \dots, Y_m$ , son independientes y provienen de una misma distribución  $F$ .

Para validar esta hipótesis, consideremos un ordenamiento de las  $n + m$  variables y supongamos que todos los elementos son distintos para asegurar que el ordenamiento es único. Si las  $n + m$  variables están igualmente distribuidas y son independientes, entonces todos los ordenamientos son equiprobables. Consideremos las variables  $X_1, X_2, \dots, X_n$  (Se podría elegir las otras  $m$  indistintamente.)

Denotaremos con  $R(X_i)$  a la posición que ocupa el elemento  $X_i$  luego del ordenamiento, y

$$R = \sum_{i=1}^n R(X_i).$$

Por ejemplo, si  $X_1 = 12, X_2 = 4, X_3 = 6, Y_1 = 9, Y_2 = 1$ , el ordenamiento resulta:

$$Y_2 = 1, \quad X_2 = 4, \quad X_3 = 6, \quad Y_1 = 9, \quad X_1 = 12,$$

y entonces

$$R = R(X_1) + R(X_2) + R(X_3) = 5 + 2 + 3 = 10.$$

Diremos que  $R(X_i)$  es el **rango del elemento**  $X_i$  y  $R$  es el **rango de la muestra** de tamaño  $n$ . Si se observa el valor  $R = r$ , y  $r$  es un valor muy grande, es indicativo que los valores  $X_i, 1 \leq i \leq n$  son en general mayores que los  $Y_j, 1 \leq j \leq m$ . Análogamente, si  $r$  es muy pequeño, esto indica que los valores de los  $Y_j$  son mayores que los  $X_i$ . Como estas dos situaciones dan razón para rechazar la hipótesis nula  $H_0$ , el  $p$ -valor estará asociado con las siguientes probabilidades:

$$P_{H_0}(R \geq r), \quad P_{H_0}(R \leq r).$$

Esto es, si alguna de estas probabilidades es muy pequeña se rechaza  $H_0$ . Así, el  $p$ -valor se define como

$$p - \text{valor} = 2 \cdot \min \{ P_{H_0}(R \geq r), P_{H_0}(R \leq r) \}.$$

Se toma  $2 \min \{ \}$  porque la región de confianza del  $100(1 - \alpha) \%$  se elige entre dos valores  $r_1$  y  $r_2$  tales que

$$P_{H_0}(R \leq r_1) = P_{H_0}(R \geq r_2) = \frac{\alpha}{2}.$$

Así, si el nivel de confianza es  $1 - \alpha = 0.90$ , entonces la hipótesis nula será rechazada si alguna de las dos probabilidades es menor a 0.05, o lo que es lo mismo, si dos veces el mínimo es menor a 0.1.

El test de hipótesis que utiliza este  $p$ -valor se denomina **test de suma de rangos**, o de **Wilcoxon** o de **Mann-Whitney**.

Resta ahora la tarea de calcular estas probabilidades. Para ello pueden usarse dos métodos diferentes, según si  $n$  y  $m$  son valores pequeños o grandes.

### 8.3.1. Test de suma de rangos para $n$ y $m$ pequeños

Si  $n$  y  $m$  no son valores grandes y los datos son todos distintos, puede utilizarse una fórmula recursiva para calcular el  $p$ -valor. Si  $n$  o  $m$  son grandes, esta fórmula es válida pero poco eficiente.

Usaremos la notación

$$P_{n,m}(r) := P(R \leq r),$$

donde los subíndices  $n, m$  indican que los datos provienen de dos muestras de tamaño  $n$  (primera muestra) y  $m$  (segunda muestra) respectivamente, y  $R$  es el rango de la muestra de tamaño  $n$ . La fórmula recursiva se obtiene del siguiente modo. El elemento más grande de los  $n + m$  valores pertenece a la primera o a la segunda muestra, y el rango de este elemento es obviamente  $n + m$ .

- Si este elemento pertenece a la primera muestra, entonces el rango de esta muestra es  $n + m$  más el rango de los  $n - 1$  elementos restantes. Luego, la probabilidad de que  $R$  sea menor o igual a  $r$  **dado que** que el mayor elemento esté en la primera muestra es igual a la probabilidad que el rango de los  $n - 1$  elementos restantes de la primera muestra sea menor o igual a  $n + m - r$ :

$$P(R \leq r \mid \text{mayor elemento en la 1ra. muestra}) = P_{n-1,m}(r - m - n).$$

- Si el elemento mayor pertenece a la segunda muestra, entonces  $R \leq r$  independientemente que este elemento esté en la segunda muestra o se lo excluya del conjunto total:

$$P(R \leq r \mid \text{mayor elemento en la 2da. muestra}) = P_{n,m-1}(r).$$

- Ahora, como el elemento mayor puede estar en la primera muestra con probabilidad  $\frac{n}{n+m}$  y en la segunda con probabilidad  $\frac{m}{n+m}$ , entonces:

$$P_{n,m}(r) = \frac{n}{n+m} P_{n-1,m}(r - m - n) + \frac{m}{n+m} P_{n,m-1}(r).$$



Por último, si  $n + m = 1$ , entonces  $R = 1$  en caso que  $n = 1$  y  $R = 0$  si  $m = 1$ . Así:

$$P_{1,0}(k) = P(R \leq k) = \begin{cases} 0 & k < 1 \\ 1 & k \geq 1 \end{cases}, \quad P_{0,1}(k) = P(R \leq k) = \begin{cases} 0 & k < 0 \\ 1 & k \geq 0 \end{cases}.$$

Por último, como el valor observado  $r$  es un número entero, entonces:

$$P_{H_0}(R \geq r) = 1 - P_{H_0}(R < r) = 1 - P_{H_0}(R \leq r - 1) = 1 - P_{n,m}(r - 1),$$

el  $p$ -valor puede obtenerse recursivamente.

---

```
def rangos(n,m,r):
    if n == 1 and m == 0:
        if r < 1:
            p = 0.
        else:
            p = 1.
    elif n == 0 and m == 1:
        if r < 0:
            p = 0.
        else:
            p = 1.
    else:
        if n == 0:
            p = rangos(0,m-1,r)
        elif m == 0:
            p = rangos(n-1,0,r-n)
        else: # n>0, m>0
            p = (n*rangos(n-1,m,r-n-m) + m*rangos(n,m-1,r)) / (n+m)
    return p
```

---

La desventaja de este método es que puede implicar un gran número de recursiones. En particular, si elegimos  $r$  como el menor de los rangos entre la primera y la segunda muestra,  $r$  podría tomar un valor cercano a la mitad de la suma de todos los rangos:  $\frac{(n+m)(n+m+1)}{4}$ . Luego la cantidad de valores de  $P_{k,l}(d)$  que deberán efectuarse es del orden de

$$\frac{nm(n+m)(n+m+1)}{4},$$

que para  $n = m$  es  $O(n^4)$ .

### 8.3.2. Test de suma de rangos para $n$ y $m$ grandes

En el caso en que  $n$  y  $m$  son grandes se sigue el siguiente procedimiento. Recordemos que

$$R = \sum_{i=1}^n R(X_i).$$

Bajo la hipótesis  $H_0$ , se puede probar que  $R$  tiene una distribución aproximadamente normal:

$$R \sim N(E[R], \sqrt{\text{Var}(R)}), \quad \text{o bien} \quad \frac{R - E[R]}{\sqrt{\text{Var}(R)}} \sim N(0, 1).$$

**Ejemplo 8.5.** La Figura 8.1 muestra dos histogramas de rangos de 1000 muestras de tamaño  $n = 10$ , comparadas con muestras de tamaño  $m$ , para  $m = 10$  y  $m = 20$  respectivamente. de variables exponenciales  $X_i \sim \mathcal{E}(1)$ ,  $Y_j \sim \mathcal{E}(1)$ .

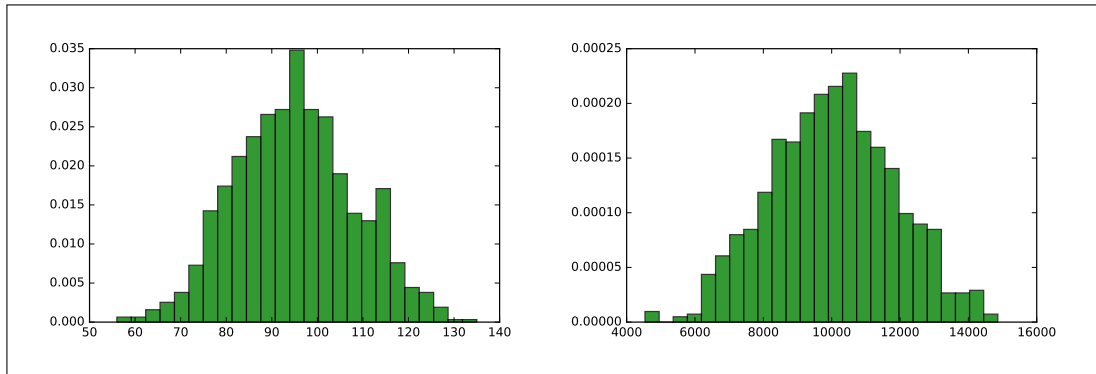


Figura 8.1: Suma de rangos:  $n = m = 10$  y  $n = 10$ ,  $m = 2000$

Tenemos que  $R(X_i)$  puede ser cualquier valor entre 1 y  $n + m$ , con igual probabilidad. Por lo tanto

$$E[R(X_i)] = \frac{n + m + 1}{2}, \quad E[R] = n \frac{n + m + 1}{2}.$$

Para el cálculo de la varianza, notemos que  $R(X_i)$  y  $R(X_j)$  no son independientes, en particular porque no pueden tomar simultáneamente el mismo valor. Puede probarse que

$$\text{Var}(R(X_i)) = \frac{(n + m + 1)(n + m - 1)}{12}$$

$$\text{cov}(R(X_i), R(X_j)) = -\frac{n + m + 1}{12},$$

y por lo tanto

$$\begin{aligned}\text{Var}(R) &= \sum_{i=1}^n \text{Var}(R(X_i)) + \sum_{i \neq j} \text{cov}(R(X_i), R(X_j)) \\ &= \frac{n(n+m-1)(n+m+1)}{12} - n(n-1) \frac{n+m+1}{12} \\ &= nm \frac{n+m+1}{12}.\end{aligned}$$

Así, bajo la hipótesis nula  $H_0$ , se tiene que

$$\frac{R - n(n+m+1)/2}{\sqrt{nm(n+m+1)/12}} \sim N(0, 1).$$

Luego, si  $Z \sim N(0, 1)$  y

$$r^* = \frac{r - n(n+m+1)/2}{\sqrt{nm(n+m+1)/12}},$$

entonces el  $p$ -valor puede calcularse como

$$2 \min \{P(Z \leq r^*), P(Z \geq r^*)\}.$$

Por la propiedad de simetría de  $Z$ , este mínimo es  $P(Z \leq r^*)$  si  $r^* \leq 0$  y es  $P(Z \geq r^*)$  en caso contrario. En términos de  $r$  esto es:

$$p - \text{valor} = \begin{cases} 2 P(Z \leq r^*) & \text{si } r \leq n \frac{n+m+1}{2} \\ 2 P(Z > r^*) & \text{en caso contrario.} \end{cases} \quad (8.5)$$

Si los  $n + m$  datos son todos distintos, entonces todos los ordenamientos son igualmente probables y equivalen a todos los ordenamientos del conjunto de números  $\{1, 2, 3, \dots, n + m\}$ . Por lo tanto, una vez observado el valor  $R = r$ , el  $p$ -valor puede determinarse simulando  $N$  permutaciones de los primeros  $n + m$  números naturales y calculando en cada simulación el valor  $R = R(1) + R(2) + \dots + R(n)$ . Finalmente,

$$p - \text{valor} = 2 \min \left\{ \frac{\#\{R \mid R \geq r\}}{N}, \frac{\#\{R \mid R \leq r\}}{N} \right\}.$$

Por último, si los  $n + m$  datos no son todos distintos entonces hay más de un ordenamiento posible y en consecuencia puede haber más de un rango para la muestra de tamaño  $n$ . En este caso, el rango  $R$  se define como el promedio de los rangos de cada ordenamiento. Por ejemplo, si los datos de las dos muestras son:

$$2 \quad 5 \quad 3, \quad \quad \quad 3 \quad 4 \quad 4,$$

los ordenamientos posibles y los correspondientes rangos de la primera muestra son:

$$\begin{array}{cccccc} \mathbf{2} & \mathbf{3} & 3 & 4 & 4 & \mathbf{5} & R = 9 \\ \mathbf{2} & 3 & \mathbf{3} & 4 & 4 & \mathbf{5} & R = 10 \end{array}$$

y en este caso se define el rango como  $R = 9.5$ . En este caso el  $p$ -valor se estima con la aproximación a la normal estándar, es decir, con la fórmula (8.5).

## 8.4. Test de rangos para varias muestras

En el caso que se quiera testear que varias muestras provienen de observaciones independientes de una misma distribución  $F$ , se aplica el **Test de rangos para varias muestras**. En este caso, se consideran  $m$  muestras provenientes de distribuciones  $F_1, F_2, \dots, F_m$ , respectivamente:

$$\begin{array}{c} X_1^{(1)}, X_2^{(1)}, \dots, X_{n_1}^{(1)}, \\ X_1^{(2)}, X_2^{(2)}, \dots, X_{n_2}^{(2)}, \\ \vdots \\ X_1^{(m)}, X_2^{(m)}, \dots, X_{n_m}^{(m)}, \end{array}$$

de tamaños  $n_1, n_2, \dots, n_m$  respectivamente. La hipótesis nula es:

- $H_0$ : Las  $n = n_1 + n_2 + \dots + n_m$  observaciones son independientes y provienen de una misma distribución  $F$ .

Asumimos en principio que todos los valores son distintos.

Luego de haber ordenado los  $n = n_1 + n_2 + \dots + n_m$  valores, se calcula el rango de cada una de las muestras. Denotamos con  $R_i$  al rango de la  $i$ -ésima muestra, para  $1 \leq i \leq m$ . Si todas las muestras provienen de la misma distribución, entonces todos los ordenamientos son igualmente probables. Al igual que antes, el valor esperado de  $R_i$  está dado por:

$$E[R_i] = n_i \frac{n+1}{2}.$$

El **Test de rangos para múltiples muestras** o **Test de Kruskal-Wallis** se basa en el siguiente estadístico:

$$R = \frac{12}{n(n+1)} \sum_{i=1}^m \frac{(R_i - n_i \frac{n+1}{2})^2}{n_i}.$$

Notemos que bajo la hipótesis que todos los valores provienen de la misma distribución, es razonable suponer que los rangos  $R_i$  estén próximos a su valor esperado  $E[R_i]$ , en relación a su varianza. Es decir, es aceptable tener un valor *pequeño* de  $R$ . Por el contrario, si se observa un valor  $R = r$  *grande*, entonces se rechaza la hipótesis nula. Luego el  $p$ -valor se define en este caso como

$$p - \text{valor} = P_{H_0}(R \geq r),$$

y la hipótesis nula será rechazada si este valor es menor que  $\alpha$  determinado por un cierto grado de confianza  $1 - \alpha$ .

Bajo la hipótesis nula  $H_0$ ,  $R$  se distribuye aproximadamente como una variable aleatoria chi cuadrado con  $m - 1$  grados de libertad:

$$p - \text{valor} = P(\chi_{m-1}^2 \geq r). \quad (8.6)$$

Finalmente, en el caso en que las observaciones tomen valores repetidos, el rango  $R_i$  se define como el promedio de los rangos de la muestra  $i$  en todos los ordenamientos posibles. Para el  $p$ -valor se utiliza la misma aproximación que en (8.6).

## 8.5. Validación de Procesos de Poisson

Supongamos que se han observado o simulado los tiempos de arribo de clientes a un servidor a los largo de varios días, e interesa testear que el número de arribos constituye un proceso de Poisson no homogéneo. Para validar esta hipótesis, consideramos la siguiente hipótesis nula:

**Hipótesis 1:** Los procesos de arribos de cada día responden a procesos de Poisson no homogéneos, independientes y con una misma función de intensidad.

Sea  $m$  el número de días en que se han observado los tiempos de llegada, y denotamos  $N_1, N_2, \dots, N_m$  el número de arribos en cada día, respectivamente. Sea  $[0, T]$  el intervalo de tiempo correspondiente a un día.

Si los procesos observados son de Poisson, con la misma función de intensidad e independientes entre sí, entonces  $N_1, N_2, \dots, N_m$  es una muestra de una variable aleatoria Poisson  $\mathcal{P}(m_T)$ . Aquí  $m_T$  corresponde al valor medio de la función de intensidad en un día.

Testeamos entonces la siguiente hipótesis nula, más débil que la anterior:

**Hipótesis 2:** Las observaciones  $N_1, N_2, \dots, N_m$  son independientes y provienen de una misma distribución de Poisson.

Para testear la Hipótesis 2, puede utilizarse el test chi cuadrado estimando el parámetro no especificado,  $m_T$ . Este parámetro indica la tasa promedio de llegada en  $[0, T]$  y puede ser estimado como

$$\hat{m}_T = \frac{N_1 + N_2 + \dots + N_m}{m}.$$

Otra forma de analizarlo y que puede resultar más eficiente, es basarse en la propiedad que el valor esperado y la varianza de una variable aleatoria Poisson son iguales. Luego se consideran la media muestral  $\bar{N}$  y la varianza muestral  $S^2$  como estimadores de la media y la varianza:

$$\bar{N} = \frac{1}{m} \sum_{i=1}^m N_i, \quad S^2 = \frac{1}{m-1} \sum_{i=1}^m (N_i - \bar{N})^2.$$

Si la hipótesis nula es cierta, el estadístico

$$T = \frac{S^2}{\bar{N}} \quad (8.7)$$

no debería tomar valores ni muy pequeños ni muy grandes. Como siempre, el concepto de valor *pequeño* o *grande* es siempre en relación a los valores que toma  $T$  cuando la hipótesis nula es cierta.

Llamamos  $\hat{\lambda}$  el parámetro  $m_T$  estimado en la muestra. Entonces el  $p$ -valor para la observación  $T = t$  del estadístico, se define como:

$$p - \text{valor} = 2 \min\{P_{\hat{\lambda}}(T \leq t), P_{\hat{\lambda}}(T \geq t)\}, \quad (8.8)$$

donde el subíndice en  $P_{\hat{\lambda}}$  indica que la probabilidad es calculada asumiendo que las variables son de Poisson con media  $\hat{\lambda}$ . El  $p$ -valor en (8.8) se aproxima realizando  $M$  simulaciones. En la  $j$ -ésima simulación,  $1 \leq j \leq M$ ,

1. se simula una muestra de tamaño  $m$ ,  $X_1^{(j)}, X_2^{(j)}, \dots, X_m^{(j)}$  de una variable con distribución de Poisson. Esto es,  $X_i^{(j)} \sim \mathcal{P}(\hat{\lambda})$ .
2. Se evalúa el estadístico  $T$  dado en (8.7) en la muestra, obteniendo el valor  $T = t_j$ .

Finalmente, el  $p$ -valor se obtiene como

$$2 \min \left\{ \frac{\#\{j \mid t_j \leq t\}}{M}, \frac{\#\{j \mid t_j \geq t\}}{M} \right\}.$$

En el caso en que este  $p$ -valor no sea muy pequeño, la Hipótesis 2 no se rechaza y podemos asumir que *la cantidad de arribos* en los  $m$  días observados:  $N_1, N_2, \dots, N_m$ , son independientes y provienen de una distribución de Poisson.

Queda aún por testear la Hipótesis 1, que reescribimos como:

- a) cada día el proceso de arribos es un proceso de Poisson no homogéneo, y
- b) la intensidad del proceso de Poisson es la misma en todos los días.

Consideramos entonces los tiempos de arribos en cada uno de los días:

$$\begin{aligned} &X_1^{(1)}, X_2^{(1)}, \dots, X_{n_1}^{(1)}, \\ &X_1^{(2)}, X_2^{(2)}, \dots, X_{n_2}^{(2)}, \\ &\vdots \\ &X_1^{(m)}, X_2^{(m)}, \dots, X_{n_m}^{(m)}. \end{aligned}$$

Se puede demostrar que si cada una de estas  $m$  muestras corresponden a los tiempos de arribos de procesos de Poisson no homogéneos con la misma función de intensidad  $\lambda(t)$ , entonces las  $n = n_1 + n_2 + \dots + n_m$  observaciones son independientes y provienen de una misma distribución.

En base a este resultado, se puede aplicar el **test de rangos para múltiples muestras**, pero atendiendo a la siguiente diferencia. Este test presupone que cada una de las muestras efectivamente contiene observaciones independientes y de *alguna* distribución, y lo que se analiza es si en realidad las  $m$  muestras provienen todas de la *misma* distribución.

En cambio, en el caso que se han observado los tiempos de arribos en los  $m$  días, no se presupone que cada día correspondan a un proceso de Poisson no homogéneo, o provengan de alguna distribución, pues esto es precisamente lo que se quiere testear. Así, si bien se considera el mismo estadístico:

$$R = \frac{12}{n(n+1)} \sum_{i=1}^m \frac{(R_i - n_i \frac{n+1}{2})^2}{n_i},$$

no se esperan ni valores grandes ni tampoco pequeños. Esto es, si  $t$  es pequeño puede ser un indicativo que las observaciones no son independientes. Luego, el  $p$ -valor para la observación  $R = r$  se define por:

$$\begin{aligned} p\text{-valor} &= 2 \min \{P(R \leq r), P(R \geq r)\} \\ &= 2 \min \{P(\chi_{m-1}^2 \leq r), P(\chi_{m-1}^2 \geq r)\}. \end{aligned}$$

### 8.5.1. Validación de un proceso de Poisson homogéneo

Todo el análisis previo también es válido para un proceso de Poisson homogéneo, ya que en este caso se estaría considerando una función de intensidad constante. Sin embargo, un proceso de Poisson homogéneo posee propiedades que no se hacen extensivas a un proceso no homogéneo en general. Una de estas propiedades es la siguiente:

- Dado el número de arribos  $N(T)$ , los tiempos de arribos se distribuyen uniformemente en  $(0, T)$ .

Así, si las  $m$  muestras provienen de  $m$  procesos de Poisson homogéneos, todos con la misma tasa, entonces cada una de las muestras contiene observaciones independientes de una variable con distribución uniforme en  $(0, T)$ . Por lo tanto el conjunto de los  $n = n_1 + n_2 + \dots + n_m$  tiempos de arribo observados deberían también estar uniformemente distribuidos en  $(0, T)$ . Esta hipótesis puede testearse con Kolmogorov-Smirnov, o el test chi cuadrado.

### 8.5.2. Estimación de la función de intensidad

Una vez que se ha validado que los procesos de arribos en los diferentes días responden a procesos de Poisson con una misma función de intensidad, puede interesar determinar cuál es esta función de intensidad.

Si se ha testado que corresponden a procesos homogéneos, entonces la tasa de llegada puede estimarse como el cociente entre el número total de llegadas y el tiempo total en los  $m$  días:

$$\hat{\lambda} = \frac{n_1 + n_2 + \cdots + n_m}{mT}.$$

Recalcamos que un día corresponde a un intervalo de tiempo  $[0, T]$ , es decir,  $T$  unidades de tiempo.

Si en cambio se ha rechazado la hipótesis que el proceso sea homogéneo, la función de intensidad puede reconstruirse del siguiente modo. Consideramos los  $n$  tiempos de arribos de los  $m$  días, ordenados de menor a mayor:  $Y_1, Y_2, Y_3, \dots, Y_n$ .

Definimos  $Y_0 = 0$ , y observamos que en el intervalo de tiempo

$$(Y_j, Y_{j+1}], \quad j \geq 0$$

se ha observado un único arribo en los  $m$  días, de donde podríamos estimar un promedio de  $1/m$  arribos diarios. Por lo tanto un estimador de  $\lambda(t)$  puede ser tal que:

$$N(Y_{j+1}) - N(Y_j) = \int_{Y_j}^{Y_{j+1}} \hat{\lambda}(t) dt = \frac{1}{m}.$$

Luego puede aproximarse la función de intensidad como constante entre dos observaciones consecutivas:

$$\hat{\lambda}(t) = \frac{1}{m(Y_{j+1} - Y_j)}, \quad Y_j < t \leq Y_{j+1},$$

## 8.6. Ejemplo

### 8.6.1. El problema

A lo largo de cuatro jornadas se han observado los siguientes tiempos de llegada en un intervalo de dos horas y media de tiempo:



Tiempos de arribo $X_{ji}$						$N_j$
1.626	1.767	1.771	1.955	2.335		5
0.927	1.593	1.832				3
0.291	0.442	0.769	1.512	2.109		5
0.679	1.110	2.085	2.340			4
0.265	2.177	2.225				3
1.664	1.945	2.472				3
0.123	1.991	2.106	2.156	2.267		5
0.201	0.775	1.002	1.217	2.077	2.449	6
0.202	1.321	1.347	1.754			4
1.854	2.400					2

Se desea validar la hipótesis que los arribos han ocurrido de acuerdo a un proceso de Poisson no homogéneo con la misma intensidad  $\lambda(t)$ .

### 8.6.2. Validación del número de arribos como variable aleatoria Poisson

En primer lugar, validamos la hipótesis:

**Hipótesis a)** El número de arribos es una variable aleatoria Poisson.

Para esto, estimamos el parámetro  $\lambda$  del proceso:

$$\hat{\lambda} = \frac{N_1 + N_2 + \cdots + N_{10}}{10} = 4.0,$$

y calculamos el estadístico para la prueba chi-cuadrado de bondad de ajuste. En este caso podemos tomar  $k = 6$ , el número máximo de arribos. Llamamos  $A_j$  a la frecuencia observada del valor  $j$ . Sea

$$p_j = e^{-\hat{\lambda}} \frac{\hat{\lambda}^j}{j!}, \quad 0 \leq j < k, \quad p_k = 1 - \sum_{j=0}^{k-1} p_j.$$

Así las probabilidades teóricas y el estadístico resultan:

$$[p_0, p_1, p_2, p_3, p_4, p_5, p_6] = [0.018, 0.073, 0.147, 0.195, 0.195, 0.156, 0.215]$$

$$T = \sum_{j=0}^k \frac{(A_j - 10p_j)^2}{10p_j} = 3.560.$$

Para estimar este  $p$  valor, realizamos 10000 simulaciones. En cada simulación se genera una muestra de tamaño 10 de variables Poisson con media 4.0. Por cada muestra se vuelve a estimar su media, se calcula el valor del estadístico y se compara con el de la muestra original:  $t_0 = 3.560$ . El  $p$ -valor obtenido con 10000 simulaciones es

$$P_{H_0}(T \geq 3.560) = P(\chi_{7-1-1}^2 \geq 3.560) \simeq 0.6123.$$

Un segundo test consiste en calcular el estadístico  $\tilde{T} = \frac{S^2(10)}{\bar{X}(10)}$ , que en este caso resulta con el valor

$$\tilde{T} \simeq \frac{0.527}{1.528} \simeq 0.345.$$

La misma simulación de 10000 muestras de tamaño 10 de variables aleatorias  $\mathcal{P}(4.0)$  arroja un  $p$ -valor

$$p = 2 \cdot \min\{P_{H_0}(\tilde{T} \geq 0.345), P_{H_0}(\tilde{T} \leq 0.345)\} = 0.108.$$

A un nivel de rechazo del 5 % la hipótesis nula no es rechazada.

### 8.6.3. Validación del proceso no homogéneo

Dado que no se rechaza la hipótesis de que el número de arribos diario sea una variable aleatoria Poisson, podemos avanzar sobre la segunda hipótesis:

**Hipótesis 2:** Las muestras provienen de observaciones independientes de una misma distribución.

Recordamos que si las observaciones provienen de muestras de Procesos de Poisson no homogéneos, todos con la misma intensidad, entonces los tiempos de arribos son muestras con observaciones independientes todas ellas de una misma distribución. Aplicamos entonces el Test de Rangos para múltiples muestras, con la observación que en este caso no hay una hipótesis de cuál es la distribución común.

Los rangos de las muestras son los siguientes:

$$[R_1, R_2, \dots, R_{10}] = [123, 50, 66, 85, 71, 84, 125, 102, 52, 62].$$

El total de observaciones es  $N = 40$ , por lo cual

$$\frac{12}{N(N+1)} \sum_{i=1}^{10} \frac{(R_i - A_i \frac{N+1}{2})^2}{A_i} \simeq 8.897.$$

Para estimar el  $p$ -valor realizamos 10000 simulaciones, cada una con 10 muestras de una misma distribución de tamaños  $A_1, A_2, \dots, A_{10}$  respectivamente. Eligiendo muestras de una distribución uniforme  $U \sim \mathcal{U}(0, 1)$ , y calculando el estadístico para el test de rangos de múltiples muestras obtenemos:

$$p - \text{valor} \sim P(\chi_9^2 \geq 8.897) = 0.4738.$$

Por lo tanto **no se rechaza la hipótesis nula** de que los datos provienen de 10 muestras de una misma distribución. En particular no hay evidencias para rechazar que provengan todas de observaciones de proceso de Poisson no homogéneo.

### 8.6.4. Validación del proceso homogéneo

En particular podría tratarse de un proceso de Poisson homogéneo. En este caso, los tiempos de arribo deberían estar uniformemente distribuidos en el intervalo  $(0, 2.5)$ . El estadístico de Kolmogorov-Smirnov arroja el valor:

$$D = d_{KS} = 0.2372.$$

Una simulación para el  $p$ -valor con 10000 muestras de tamaño 40 tomadas de una distribución  $U \sim \mathcal{U}(0, 1)$  estima:

$$p - \text{valor} = P(D \geq 0.2372) = 0.0176.$$

A un nivel del 5 % se rechaza la hipótesis nula.

### 8.6.5. Estimación de la función de intensidad

Para la estimación de la función de intensidad  $\lambda(t)$ , se ordenan las 40 observaciones de menor a mayor,

$$Y_1 < Y_2 < \dots < Y_{40},$$

se define  $Y_0 = 0$  y una estimación de  $\lambda(t)$  es:

$$\hat{\lambda}(t) = \frac{1}{10(Y_{j+1} - Y_j)}, \quad Y_j \leq t < Y_{j+1}, \quad 0 \leq j < 40. \quad (8.9)$$

Esto es, dado que  $N(Y_{j+1}) - N(Y_j) = 1$ , significa que en 10 jornadas ocurrió un único evento en el intervalo  $(Y_j, Y_{j+1})$ . Luego la tasa de llegada promedio es:

$$N(Y_{j+1}) - N(Y_j) = \int_{Y_j}^{Y_{j+1}} \lambda(t) dt = \frac{1}{10},$$

de donde se sigue que una aproximación de  $\lambda(t)$  como constante es como en (8.9). El gráfico de  $\hat{\lambda}(Y_j)$ ,  $0 \leq j \leq 40$  se muestra en la Figura 8.6.5. Como puede observarse, existen varios valores extremos o *outliers* producidos por valores muy cercanos en el tiempo: 0.201, 0.202 en comparación con otros valores vecinos.

Esto hace necesario la construcción de otros estimadores más *robustos* para la función de intensidad  $\lambda(t)$ , tema que está fuera del objetivo de este curso.

