

Estudio estadístico de dos muestras de datos independientes entre sí.

GRUPO 11

José Ignacio López Paz	C.I:4491310-6
Germán Ruiz Raviales	C.I:4317743-2
Federico Mujica Cazenave	C.I:4786543-9

INDICE GENERAL

1_3-4 Estudio de Aleatoriedad

1.1_3 Test de rachas de ascensos y descensos

1.2_4 Test de correlación de rangos de Spearman

2_5-7 Tests de Ajuste

2.1_5-6 Test de Lilliefors de ajuste para
exponenciales

2.2_6-7 Test de Kolmogorov-Smirnov para una
muestra

3_8 Comparación de las distribuciones

3.1_8 Test de Kolmogorov-Smirnov para dos
muestras

4_9 Tests Paramétricos de comparación

4.1_9 Intervalos de confianza para la media

4.2_10 Test de comparación

5_11 Conclusión final

Estudio de Aleatoriedad

Primero veremos si las muestras son aleatorias realizando 2 tests. El test de Rachas de ascensos y descensos y el test de correlación de rangos de Spearman.

Para ambas muestras en estos tests la hipótesis nula H_0 corresponde a que los datos son iid (independientes idénticamente distribuidos) y la hipótesis alternativa H_1 a que no lo son.

Test de Rachas:

Realizando los respectivos cálculos para la primera muestra llegamos a que el total de rachas es igual a 64 y para la segunda muestra es 68. Como las muestras son de tamaño mayor a 25, entonces se usa la "aproximación normal" del TCL (teorema central del límite) para hallar el p-valor.

¹ El p-valor de la primera muestra es igual a 0.3304003 y el p-valor de la segunda muestra es igual a 0.3900298. Tomando el error de tipo I igual al 10% ($\alpha=0.10$), no se rechaza la hipótesis nula para ninguna de las muestras ya que sus respectivos p-valores son mayores que α .

Como las dos muestras aceptan H_0 para este test podemos decir que los datos no tienen una tendencia sistemática (muy pocas rachas) y que no reflejan una periodicidad (demasiadas rachas).

¹ Para ver los cálculos realizados para los estadísticos y p-valores de los tests ver en los anexos el script de TestDeAleatoriedad

Test de Spearman:

El estadístico de Spearman es igual a:

$R_s = 1 - [6(\sum_{1 \leq i \leq n} (R(X_i) - i)^2) / (n(n^2 - 1))]$ siendo n la cantidad de datos en la muestra y $R(X_i)$ el rango del dato i .

Haciendo los cálculos llegamos a que el estadístico² de Spearman $R_{s1} = -0.0240264$ para la muestra1, y para la muestra2 $R_{s2} = 0.110003$.

Usando la “aproximación normal” del TCL (teorema central del límite) llegamos a que el p-valor de la primera muestra nos dio 0.4055296 y el de la segunda dio 0.1368644. Para las dos muestras el p-valor es mayor que el error de tipo1 (α) por lo cual no se rechaza la hipótesis nula del test para ambas muestras.

Al haber aceptado H_0 las dos muestras podemos decir que los datos no tienen ni una tendencia creciente ni decreciente.

Conclusión: Luego de realizar estos dos tests conjuntamente y como las dos muestras aceptaron llegamos a que los datos de ambas muestras se pueden considerar independientes idénticamente distribuidos.

² Para ver los cálculos realizados para los estadísticos y p-valores de los tests ver en los anexos el script de TestDeAleatoriedad

Tests de Ajuste

Primero realizamos graficas e histogramas para estimar cual distribución se adaptaba más a cada muestra.³

Luego de observar las graficas notamos que las dos podrían corresponder a una distribución exponencial, para verificar esta idea nos propusimos realizar algunos tests empezando con el test de Lilliefors de ajuste a exponenciales.

Para los tests presentados a continuación la hipótesis nula H_0 corresponde a que los datos de la muestra tienen una distribución $F = \exp(\lambda) = F_0$ y la hipótesis alternativa que no.

Test de Lilliefors para exponenciales:

A partir de la muestra estimamos por el método de máxima verosimilitud (EMV) el parámetro λ de dicha distribución. El estadístico de esta prueba es igual a $D_n = \sup_{t \in \mathbb{R}} |F_n(t) - F_0(t)|$ siendo la función de distribución empirica de la muestra

$F_n(t) = \# \{i: 1 \leq i \leq n: X_i \leq t\} / n$ (n es la cantidad de datos de la muestra).

El estadístico⁴ del test para la primer muestra nos dio $D_{n1} = 0.0585$ y de la segunda nos dio $D_{n2} = 0.0395$. Como la cantidad de datos es mayor a 30 usamos la “aproximación normal” para hallar el valor critico

³ Ver en anexos las graficas Histogramas_Densidades, EstimacionMuestra1 y EstimacionMuestra2

⁴ Para los cálculos ver sript de TestDeAjuste de los anexos

asintótico que corresponde al error de tipo I $\alpha=0.1$, que lo llamamos C_n . Llegamos a que $C_n=0.096$.

La región crítica de la primera muestra es igual a $R=\{D_{n1}>C_n\}$ y para la segunda es $R=\{D_{n2}>C_n\}$.

El estadístico de ambas muestras es menor a C_n por lo tanto no caen en la región entonces no se rechaza que los datos de las muestras tengan distribución exponencial.

A continuación realizamos el test de Kolmogorov-Smirnov para ambas muestras.

Test de Kolmogorov-Smirnov:

El estadístico y la región crítica de este test es el mismo que en el de Lilliefors. La diferencia es que el parámetro λ se estima con la mitad de la muestra y luego se usa la otra parte para aplicar el test. Decidir que parte de la muestra se usara para estimar los parámetros y que parte para aplicar el test, es una arbitrariedad. Para evitarla, lo hacemos en dos casos. En el primero aplicamos el procedimiento estimando los parámetros con la primera mitad de la muestra y aplicamos la prueba de ajuste con la segunda y en el segundo hacemos al revés, la segunda mitad para estimar y la primera para aplicar el test. De esta manera aseguramos más el rechazo o la aceptación de H_0 .

Para el primer caso llegamos a que el estadístico para la primera muestra es $D_{n1}=0.1195$ y para la segunda es $D_{n2}=0.1684$. En este caso el valor de

⁵ Para los cálculos ver script de TestDeAjuste en los anexos

$C_n = 0.1725341$. D_{n1} es claramente menor que C_n por lo tanto no se rechaza H_0 en este caso para la primer muestra. D_{n2} también es menor que C_n pero por muy poca diferencia, por lo tanto hay que probar con el segundo caso planteado para corroborar que los datos de la segunda muestra no rechazan H_0 .

En el segundo caso ⁶ $D_{n1} = 0.0777$ que también es claramente menor que C_n por lo tanto para la primer muestra no se rechaza H_0 . Luego $D_{n2} = 0.1413$. En este caso dio bastante menor que en el anterior por lo tanto nos aseguro que la segunda muestra no rechaza H_0 .

Conclusión: Luego de aplicar estos tests llegamos a la conclusión de que las dos muestras se pueden ajustar a distribuciones exponenciales ya que no se rechazo la hipótesis nula planteada en ninguno de ellos.

⁶ Para los cálculos ver sript de TestDeAjuste en los anexos

Comparación de las distribuciones

En la parte anterior concluimos que los datos de las dos muestras tienen el mismo tipo de distribución por lo tanto en esta sección haremos el test de Kolmogorov-Smirnov (KS) para dos muestras para ver si tienen la misma distribución.

Para este test la hipótesis nula corresponde a que la distribución de la primer muestra F_1 es igual a la de la segunda F_2 , la hipótesis alternativa es que son distintas.⁷

Test K-S para dos muestras:

El estadístico de este test es

$D_{nm} = \sup_{t \in \mathbb{R}} |F_1(t) - F_2(t)|$ con n y m la cantidad de datos de la primer y segunda muestra respectivamente. Como en este caso las dos muestras son de igual tamaño lo llamaremos simplemente D .

Luego llegamos a que $^8D=0.22$ y el C_n para $\alpha=0.1$ es $C_n= 0.5456006$. La región crítica de este test para rechazar H_0 es $R=\{nmD>C_n\}$ (nmD el producto entre ellos y el m y n el tamaño de las muestras). Luego $nmD=545.60$ que es claramente mayor que C_n por lo tanto se rechaza H_0 para las muestras.

⁷ Ver en los anexos grafica de Funciones DeDistribución

⁸ Ver la comparación de las muestras en los anexos en el script de Comparación para más detalles

Conclusión: Cómo se rechaza H_0 concluimos con mucha seguridad que no tienen distribución idéntica.

Test Paramétricos de comparación

En esta sección del informe ya suponemos que las muestras tienen una distribución exponencial de parámetro λ_1 y λ_2 estimados por el método de máxima verosimilitud (EMV), que son iguales a $1/X_n$ (X_n =Promedio de los datos de la muestra1), y $1/Y_n$ (Y_n =Promedio de los datos de la muestra2) respectivamente.

Tenemos los datos de la muestra1 $X = X_1, \dots, X_n$, y de la muestra2 $Y = Y_1, \dots, Y_n$ que los suponemos iid y que tienen una distribución que no es la normal. La cantidad de datos es grande entonces si $EX_i = \mu_1$, $EY_i = \mu_2$, $\text{Var}(X_i) = \sigma_1^2 > 0$ y $\text{Var}(Y_i) = \sigma_2^2$ vale el TCL entonces podemos calcular los intervalos de confianza para la media de dichas muestras.

Intervalos de confianza para la media:

El intervalo de confianza al 95% para μ_1 es igual a ${}^9IC_1 = [1.046226, 2.925395]$ y para μ_2 es igual a $IC_2 = [1.437317, 4.629810]$.

A primera vista se puede observar que el IC_1 es más chico que el IC_2 . El valor mínimo del IC_1 es más chico del IC_2 por poco pero en cambio el máximo es claramente menor al máximo del IC_2 . Para poder

⁹ Ver en los anexos los cálculos de los intervalos en el script de TestParametricos

verificar que μ_1 es efectivamente menor que μ_2 aplicamos unos tests de comparación.

Test de comparación de las medias:

Tomamos la hipótesis nula como $H_0: \mu_1 = \mu_2$ y la hipótesis alternativa como $H_1: \mu_1 \neq \mu_2$. La región crítica del test es $R = \{|E| \geq Z_{\alpha/2}\}$ con el estadístico $E = (X_n - Y_n) / \sqrt{((S_n^2(X) + S_n^2(Y)) / n)}$.

Haciendo los respectivos cálculos llegamos a que ¹⁰ $E = -2.912908$ y para 5% ($\alpha=0.05$) $Z_{\alpha/2} = 1.96$.

Haciendo el valor absoluto de E resulta que se cumple la condición de la región crítica por lo tanto se rechaza H_0 para las muestras.

Luego dejamos $H_0: \mu_1 = \mu_2$ y cambiamos a $H_1: \mu_1 \leq \mu_2$. El estadístico E es el mismo y la región crítica en este caso es $R = \{E \leq -Z_{\alpha}\}$.

Para $\alpha=0.05$ en este caso $Z_{\alpha} = 1.65$ entonces nuevamente se cumple la condición de R por lo tanto se rechaza H_0 .

Conclusión: Luego de aplicar estos dos tests pudimos confirmar lo que se había observado en los intervalos de confianza. Efectivamente la media de la primer muestra μ_1 es menor que la de la segunda muestra μ_2 .

¹⁰ Ver en anexos los cálculos reaalizados en el script de TestParametricos

Conclusión Final

Para este estudio los datos de las dos muestras corresponden al tiempo necesario para realizar la misma tarea en dos computadoras diferentes.

Luego de realizar todo este estudio estadístico pudimos concluir que los datos de ambas muestras podrían ser independientes idénticamente distribuidos para luego poder ajustar los datos a una distribución. Concluimos que ambas siguen una distribución exponencial de parámetro λ estimado por máxima verosimilitud. También concluimos que no tienen distribución idéntica. Las muestras se suponen independientes entre sí por lo tanto pudimos aplicar tests paramétricos para comparar las medias que serían la esperanza de cada dato. La media de los datos de la primera muestra resultó menor que el de la segunda.

Con todo esto en consideración pudimos concluir que los tiempos de ejecución del programa generados por la primera computadora tienden a ser menores que los generados por la segunda.

Como conclusión final nosotros consideraríamos que la primera computadora es más eficiente que la segunda realizando la tarea en cuestión.