

Estadística No Paramétrica

Clase 4: Test de Kolgomorov-Smirnov

Joaquin Cavieres G.

Ingeniería en Estadística

Facultad de Ciencias, Universidad de Valparaíso



Test de Kolgomorov-Smirnov

Dentro de la estadística no paramétrica también podemos encontrar el test de Kolgomorov-Smirnov. Este test tiene la finalidad comparar entre dos distribuciones de probabilidad continuas y unidimensionales en función de una muestra aleatoria con una función de de distribución probabilidad como referencia.

Test de Kolgomorov-Smirnov

El test de Kolgomorov-Smirnov es uno de los test más populares en la estadística no paramétrica. El test **cuantifica la distancia entre la función de distribución empírica de la muestra F_n y la función de distribución acumulada de una función de distribución de probabilidad de referencia $F(x)$, esto es:**

Definición

La función de distribución empírica F_n para n observaciones *i.i.d* de una muestra aleatoria X_i esta definida como:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{(0,\infty)}(X_i) \quad (1)$$

donde $I_{(0,\infty)}(X_i)$ es la función indicadora igual a 1 si $X_i \leq x$ e igual a 0 en otro caso, por lo tanto:

$$D_n = \sup_x |F_n(x) - F(x)| \quad (2)$$

Test de Kolgomorov-Smirnov

Supongamos que tenemos X_1, \dots, X_n una v.a con una función de distribución la cual no conocemos y que denotamos como F_X . A nosotros nos interesaría probar que F_X^0 es la verdadera distribución. Dado lo anterior planteamos lo siguiente hipótesis:

$$H_0 : F_X(x) = F_X^0(x)$$

$$H_1 : F_X(x) \neq F_X^0(x)$$

Test de Kolgomorov-Smirnov

Supongamos que tenemos X_1, \dots, X_n una v.a con una función de distribución la cual no conocemos y que denotamos como F_X . A nosotros nos interesaría probar que F_X^0 es la verdadera distribución. Dado lo anterior planteamos lo siguiente hipótesis:

$$H_0 : F_X(x) = F_X^0(x)$$

$$H_1 : F_X(x) \neq F_X^0(x)$$

¿Le parece conocida esta misma hipótesis para funciones de distribución?

Test de Kolgomorov-Smirnov

Como en la prueba X^2 se contrastan hipótesis frente a clases "arbitrarias", el test de Kolgomorov-Smirnov evita esta categorización. Considerando la definición de la función de distribución empírica entonces:

Si asumimos que X_i es una v.a, entonces $F_n(x)$ también lo es, así:

$$\mathbb{E}(F_n(x)) = \mathbb{E}\left(\frac{\sum_{i=1}^n I_{(0,\infty)}(X_i)}{n}\right) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(I_{(0,\infty)}(X_i)) \quad (3)$$

es un estimador insesgado y consistente para $F_X(x)$

Ya que $I_{(0,\infty)}(X_i) \sim \text{Ber}(\mathbb{P}(X_i \leq x)) = \text{Ber}(F_X(x))$, entonces $\mathbb{E}(I_{(0,\infty)}(X_i)) = F_X(x)$ para todo i , así.

Estimador insesgado

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}(I_{(0,\infty)}(X_i)) = \frac{1}{n} n \mathbb{E}(I_{(0,\infty)}(X_i)) = F_X(x) \quad (4)$$

Consistencia del estimador

$$\begin{aligned} \text{VAR}(F_n(x)) &= \text{VAR}\left(\frac{\sum_{i=1}^n I_{(0,\infty)}(X_i)}{n}\right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{VAR}(I_{(0,\infty)}(X_i)) = \frac{F_X(x)(1 - F_X(x))}{n} \end{aligned} \quad (5)$$

Cuando suponemos observada la muestra aleatoria, F_n puede ser visualizada y se espera F_n se aproxime a la verdadera función de distribución de donde proviene la muestra.

Ejemplo 1

Ver código en R simulación de una v.a Normal, construcción de una función de distribución empírica y ajuste a los datos.

Parece ser que la información que nos entrega la función de distribución empírica es razonable para la construcción de un estadístico de prueba para así verificar la bondad de ajuste. Mediante la prueba anterior de insesgamiento y consistencia entonces podemos probar que $F_n(x)$ converge casi seguramente a $F_X(x)$ para cada x .

Convergencia en forma uniforme

Teorema de Glivenko-Cantelli

Sea X_1, \dots, X_n una muestra aleatoria de $F_X(x)$ y sea $F_n(x)$ su función de distribución empírica, entonces:

$$D_n = \sup_x |F_n(x) - F_X(x)| \rightarrow 0 \quad (6)$$

Convergencia en forma uniforme

Teorema de Glivenko-Cantelli

Sea X_1, \dots, X_n una muestra aleatoria de $F_X(x)$ y sea $F_n(x)$ su función de distribución empírica, entonces:

$$D_n = \sup_x |F_n(x) - F_X(x)| \rightarrow 0 \quad (6)$$

Esto quiere decir que cuando existe un mayor número de muestra entonces F_n reproduce casi totalmente la verdadera función de distribución.

Test de Kolgomorov-Smirnov

Convergencia en forma uniforme

Teorema de Glivenko-Cantelli

Sea X_1, \dots, X_n una muestra aleatoria de $F_X(x)$ y sea $F_n(x)$ su función de distribución empírica, entonces:

$$D_n = \sup_x |F_n(x) - F_X(x)| \rightarrow 0 \quad (6)$$

Esto quiere decir que cuando existe un mayor número de muestra entonces F_n reproduce casi totalmente la verdadera función de distribución.

La clave en el test de Kolgomorov-Smirnov es que la distribución del supremo no depende de la distribución desconocida de F_X^0 de la muestra, siempre y cuando F_X^0 sea una distribución continua.

¿Como evaluar el test de Kolgomorov-Smirnov?

¿Como evaluar el test de Kolgomorov-Smirnov?

Si $X_1, \dots, X_n \sim F$, entonces $F_n(x)$ es un estimador no paramétrico de F (recordar que F es una función de distribución).

Entonces:

- Una medida de proximidad entre F_n (dado por la muestra) y F_0 (especificado por H_0) debería indicar la "distancia" entre F_n y F_0 para determinar si H_0 es falsa.

Si consideramos que:

$$H_0 : F = F_0$$

$$H_1 : F \neq F_0$$

Si consideramos que:

$$H_0 : F = F_0$$

$$H_1 : F \neq F_0$$

entonces

Estadístico D_n

Si $H_0 : F = F_0$ se cumple, entonces D_n tiende a ser pequeño. De lo contrario, si D_n tiende a ser grande, entonces $H_1 : F \neq F_0$ y rechazamos H_0 .

Calculo "a mano" de D_n

$$D_n = \max(D_n^+, D_n^-), \quad (6.2)$$

$$D_n^+ := \sqrt{n} \max_{1 \leq i \leq n} \left\{ \frac{i}{n} - U_{(i)} \right\},$$

$$D_n^- := \sqrt{n} \max_{1 \leq i \leq n} \left\{ U_{(i)} - \frac{i-1}{n} \right\},$$

donde U_j representa para el j -ésimo sorteado $U_i := F_0(X_i), i = 1, \dots, n$

Distribución bajo H_0

Si H_0 se cumple y F_0 es continua, entonces D_n tiene una función cumulativa asintótica dada por la función K :

$$\lim_{x \rightarrow \infty} \mathbb{P}[D_n \leq x] = K(x) := 1 - 2 \sum_{j=1}^{\infty} (-1)^{j-1} e^{-2j^2 x^2}$$

Test de Kolgomorov-Smirnov

Consideraciones importantes:

- Bajo H_0 el test Kolgomorov-Smirnov no depende de F_0 pero solamente si F_0 es continua y la muestra X_1, \dots, X_n también es continua. Si se suplen estos supuestos entonces una muestra $X_1, \dots, X_n \sim F_0$ i.i.d genera muestras $U_1, \dots, U_n \sim \text{Uniforme}(0, 1)$.
- Como consecuencia al punto anterior, D_n no depende de F_0 .
- Si F_0 no es continua entonces K no es la verdadera función de distribución asintótica.

Desarrollo en \mathbb{R}

Determinar si:

$$H_0 : F = F_0$$

$$H_1 : F \neq F_0$$

Test de Kolgomorov-Smirnov

Ventajas

- La distribución del estadístico de la prueba Kolgomorov-Smirnov en sí no depende de la función de distribución acumulativa subyacente que se esté probando.
- Es una prueba exacta (la prueba de bondad de ajuste de X^2 depende de un tamaño de muestra adecuado para que las aproximaciones sean válidas).

Desventajas

- Solo se aplica a distribuciones continuas.
- Tiende a ser más sensible cerca del centro de la distribución que en las colas.
- Si los parámetros de ubicación, escala y forma se estiman a partir de los datos, la región crítica de la prueba Kolgomorov-Smirnov ya no es válida. Por lo general, debe determinarse mediante simulación.