



Universidad Nacional Autónoma de México

Facultad de Estudios Superiores Acatlán

Apuntes
Estadística 2

Autor:

17 de enero de 2025

Índice

1. 29/01/2024	3
1.1. Objetivo general	3
1.2. Temario	3
1.3. Bibliografía	3
1.4. Evaluación	3
1.5. Contacto	3
1.6. Examen diagnostico	3
2. 02/02/2024	4
2.1. Solución del examen diagnóstico	4
3. 07/02/2024	6
3.1. Unidad 1: Pruebas no parametricas	6
3.1.1. Introducción	6
3.1.2. Pruebas de bondad de ajuste	6
4. 12/02/2024	7
4.1. Función de Distribución Empírica	7
4.2. Estadístico de la Prueba en la Prueba de Kolmogorov-Smirnov	7
5. 14/02/2024	10
5.1. Pruebas de Bondad de Ajuste	10
5.1.1. Para Variables Continuas	10
5.1.2. Para Variables Discretas	10
5.2. Pruebas Basadas en Rangos	10
5.2.1. Para 2 Poblaciones	10
5.2.2. Para Más de 2 Poblaciones	11
5.3. Prueba de Mann-Whitney	11
6. 16/02/2024	14
6.1. Prueba de Wilcoxon	14
7. 19/02/2024	17
7.1. Prueba de Kruskal-Wallis	17
8. 21/02/2024	20
8.1. Pruebas basadas en corridas	20
9. 23/02/2024	22
9.1. Pruebas de Independencia	22
9.2. Coeficiente de Correlación de Spearman	22
10.26/02/2024	24
10.1. Tau de Kendall	24
11.01/03/2024	26
11.1. Tablas de contingencia	26
11.1.1. Ejemplo 9	27
11.2. Pruebas de una cola (Prueba de Mann-Whitney y Prueba de Wilcoxon)	29
12.11/03/2024	30

13.13/03/2024	30
13.1. Análisis de varianza	31
13.2. Diseño completamente aleatorizado	31
13.3. Modelo	31
13.4. Modelo completo	31
13.5. Modelo reducido	32
14.15/03/2024	33
15.01/04/2024	34
15.1. Diseño de bloques completamente aleatorizados	34
16.5/04/2024	35
16.1. Tabla ANOVA	35
17.8/04/2024	36
18.10/04/2024	37
18.1. Unidad 3 - Análisis de regresión	37
18.2. 3.1 - Modelo de regresión lineal simple	37
18.3. Estimacion de los parametros por minimos cuadrados	38
18.4. Estimación de los parámetros por mínimos cuadrados	39
19.17 de abril del 2024	40
20.03/05/2024	43
21.13/05/2024 Clase	44
22.13/05/2024 ChatGPT	46
23.17/05/2024	47
23.1. Ejemplo 6	47
23.2. Ejemplo 7	48

1. 29/01/2024

1.1. Objetivo general

El estudiante aplicará pruebas no paramétricas, análisis de varianza, estadística bayesiana, análisis de regresión a la solución de problemas dentro de diversos campos de conocimiento.

1.2. Temario

1. Pruebas no paramétricas
2. Análisis de varianza y diseño de experimentos
3. Análisis de regresión
4. Inferencia bayesiana

1.3. Bibliografía

- Canavos, G. (1987). *Probabilidad y estadística: Aplicación y métodos*.
- Lee, Peter M. (2012). *Bayesian Statistics: An Introduction*.
- Montgomery, D. (2004) *Diseño y análisis de experimentos*.
-
- Montgomery, D. (2002) *Introducción al análisis de Regresión*.

1.4. Evaluación

- Exámenes parciales: 40 % (4 exámenes, uno por unidad)
- Prácticas: 30 %
- Tareas: 20 %
- Cuestionario: 10 % (Sobre el examen)

1.5. Contacto

1.6. Examen diagnóstico

- Medidas descriptivas
- Propiedades de los estimadores puntuales
- Intervalos de confianza
- Prueba de hipótesis

2. 02/02/2024

2.1. Solución del examen diagnóstico

Ejercicio 1

Consideremos el conjunto $s = \{4, 2, 0, 9, 4, 2, -1, 1, -4, 2\}$.

A continuación, calculamos las medidas de tendencia central y de dispersión para el conjunto s :

- **Media:** La media aritmética se calcula como el promedio de los valores del conjunto:

$$\text{Media} = \frac{1}{n} \sum_{i=1}^n s_i$$

donde n es el número de elementos en el conjunto.

- **Moda:** La moda es el valor o valores que aparecen con mayor frecuencia en el conjunto. En caso de que todos los elementos aparezcan con la misma frecuencia, el conjunto se considera amodal.
- **Mediana:** La mediana es el valor que ocupa la posición central del conjunto una vez que este ha sido ordenado. Si el número de elementos es par, la mediana se calcula como el promedio de los dos valores centrales. Para nuestro conjunto ordenado, debemos calcularlo correctamente.
- **Varianza:** La varianza mide la dispersión de los valores del conjunto respecto a la media. Se calcula con la fórmula:

$$\text{Varianza} = \frac{1}{n} \sum_{i=1}^n (s_i - \bar{x})^2$$

donde \bar{x} es la media del conjunto.

- **Rango:** El rango es la diferencia entre el valor máximo y el valor mínimo del conjunto. Se calcula como:

$$\text{Rango} = X_{\max} - X_{\min}$$

Ejercicio 2

a)

$$\hat{\theta}_1 = \frac{3x_1 + 4x_4}{10} - \mu$$

$$\hat{\theta}_2 = \frac{5x_2 + 3x_3 + 4x_4}{10}$$

$$\hat{\theta}_3 = \frac{x_1 + x_2 + x_3 + \sigma^2}{4}$$

$$\hat{\theta}_4 = \frac{2x_1 + 4x_2 + 4\mu + 2x_4}{12}$$

$$\hat{\theta}_5 = X_n$$

$$\hat{\theta}_6 = \bar{X}$$

b)

Sesgo

$$\beta(\hat{\theta}) = E(\hat{\theta}) - \theta$$

$$-\beta(\hat{\theta}_2) = E(\hat{\theta}_2) - \mu$$

$$-\beta(\hat{\theta}_2) = E\left(\frac{5x_2 + 3x_3 + 4x_4}{10}\right) - \mu$$

$$\frac{1}{10}[5E(x_2) + 3E(x_3) + 4E(x_4)] - \mu$$

$$\frac{5\mu + 3\mu + 4\mu}{10} - \mu$$

$$1,2\mu - \mu = 0,2\mu$$

$$-\beta(\hat{\theta}_5) = E(X_n) - \mu$$

$$-\beta(\hat{\theta}_5) = \mu - \mu = 0$$

$$-\beta(\hat{\theta}_6) = E(\bar{X}) - \mu$$

$$-\beta(\hat{\theta}_6) = \mu - \mu = 0$$

Ejercicio 3

Pruebas t : $X_1, \dots, X_n \sim N(\mu, \sigma^2), \mu = ?, \sigma^2 = ?$

Pruebas basadas en la distribución normal: $X_1, \dots, X_n \sim N(\mu, \sigma^2), \mu = ?$

Estandarizar una variable: $\frac{X - E(x)}{\sqrt{Var(x)}}$

En el ejercicio tenemos:

$X_1, \dots, X_n \sim N(\mu, \sigma^2), \mu = ?, \sigma = 30$

Intervalo de confianza:

$$\bar{X} \pm z_{\frac{\alpha}{2}} \sqrt{Var(\hat{\mu})}$$

$$\bar{X} \pm z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}$$

$$1550 \pm 1,96 \left(\frac{30}{\sqrt{50}} \right) = (1541,684, 1558,316)$$

Ejercicio 4:

$$X_1 \sim N(\mu, \sigma^2 = 1)$$

$$H_0 : \mu = 0 \text{ vs } H_1 : \mu = 2$$

Región de rechazo = $\{x | x_1 \geq 1,5\}$

$$X_1 = 2 \geq 1,5 \rightarrow X_1 \sim N(2, 1)$$

Probabilidad del cometer el error 1: Error 1 = Rechazar H_0 cuando es cierta

$$P(X_1 \geq 1,5 | \mu = 0) = 0,0668072$$

Probabilidad del cometer el error 2: Error 2 = No rechazar H_0 cuando es falsa

$$P(X_1 < 1,5 | \mu = 2) = 0,3085375$$

c) $1 - \beta = 0.6914625$

3. 07/02/2024

3.1. Unidad 1: Pruebas no paramétricas

3.1.1. Introducción

En este tipo de pruebas no se tiene el supuesto de que la muestra pertenece a una familia paramétrica, como es el caso de las pruebas F y T que requieren que los datos se distribuyan como una normal. En el caso paramétrico se tiene que $X \sim f_x(x|\theta)$ donde $f_x(x|\theta)$ es una función de distribución conocida y θ es desconocida.

Mientras que, en el caso no paramétrico $X \sim f_x(\cdot)$ con $f_x(\cdot) \in \mathbb{F}$ es el conjunto de todas las funciones de distribución

3.1.2. Pruebas de bondad de ajuste

Consisten en comparar las observaciones de una muestra aleatoria con aquellas que se esperan observar si la hipótesis nula es correcta.

El objetivo de las pruebas de bondad de ajuste es encontrar las hipótesis, $H_0 : F = F_0$ contra $H_1 : F \neq F_0$, donde F_0 es una distribución completa o parcialmente conocida.

Prueba Ji-Cuadrada

Se utiliza para variables discretas que tienen un número finito de categorías y se puede adaptar a variables continuas pero no es recomendable.

En esta prueba se comparan las frecuencias observadas en cada categoría respecto a las frecuencias esperadas bajo el supuesto de la hipótesis nula

Estadístico de la prueba

$$X^2 = \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i} \sim X_{(k-1)}^2$$

donde N_i es la frecuencia observada en la i -ésima categoría y np_i es la frecuencia esperada bajo H_0 , además n es el tamaño de la muestra y k es el total de categorías.

Se rechaza H_0 con un nivel de significancia α , si el Estadístico $x^2 > q_{X_{(k-1)}^2}^{(1-\alpha)}$ es decir.

$$R = \left\{ \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i} > q_{X_{(k-1)}^2}^{(1-\alpha)} \right\}$$

p-valor = $\mathbb{P}(X_{(k-1)}^2 > \sum_{i=1}^k \frac{(N_i - np_i)^2}{np_i})$
si p-valor $< \alpha$ entonces se rechaza H_0

Interpretación del p-valor

El nivel mínimo de probabilidad de obtener los valores observados o más extremos si la hipótesis nula es correcta.

Es el nivel de significancia más pequeño para el que la muestra obtenida obligada a rechazar la hipótesis nula.

4. 12/02/2024

4.1. Función de Distribución Empírica

La función de distribución empírica es una herramienta esencial en estadística para estimar la distribución de probabilidad de un conjunto de datos. Se construye a partir de una muestra y ofrece una estimación de la función de distribución acumulativa (CDF, por sus siglas en inglés) de la población original. Esta función se emplea en diversas áreas como el análisis exploratorio de datos, pruebas de hipótesis y el ajuste de distribuciones teóricas a conjuntos de datos empíricos.

La definición matemática de la función de distribución empírica, dada una muestra ordenada x_1, x_2, \dots, x_n , es:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x) \quad (1)$$

donde:

- $F_n(x)$ representa la función de distribución empírica evaluada en x ,
- n es el tamaño de la muestra,
- $I(x_i \leq x)$ es una función indicadora que adopta el valor de 1 si $x_i \leq x$ y 0 en caso contrario.

Esta función es escalonada, incrementándose en $1/n$ en cada uno de los datos x_i de la muestra. La función $F_n(x)$ refleja la proporción de observaciones en la muestra que no superan el valor x .

Propiedades destacadas de la función de distribución empírica incluyen:

1. **No decreciente:** $F_n(x)$ es una función que solo puede mantenerse constante o incrementar.
2. **Límites:** $F_n(x) = 0$ para todo x menor que el mínimo de la muestra, y $F_n(x) = 1$ para x superior al máximo de la muestra.
3. **Convergencia:** Bajo ciertas condiciones, $F_n(x)$ converge hacia la verdadera CDF de la población, $F(x)$, a medida que el tamaño de muestra n se incrementa indefinidamente. Este comportamiento se deriva del teorema del límite central, entre otros resultados fundamentales en teoría de la probabilidad.

La construcción y visualización de la función de distribución empírica son sencillas, facilitando el análisis descriptivo y comparativo de distribuciones de datos. Además, su utilidad se extiende a métodos no paramétricos, como la prueba de Kolmogorov-Smirnov, que compara la distribución empírica de una muestra con una distribución teórica o con la de otra muestra.

4.2. Estadístico de la Prueba en la Prueba de Kolmogorov-Smirnov

El estadístico de la prueba de Kolmogorov-Smirnov (K-S) se emplea para comparar una distribución empírica con una distribución teórica, o dos distribuciones empíricas entre sí. Este cuantifica la máxima discrepancia entre las funciones de distribución acumulativa (CDF) implicadas. Existen dos variantes de la prueba K-S: unidireccional, que compara una muestra con una distribución teórica, y bidireccional, que compara dos muestras.

Prueba K-S Unidireccional

En la variante unidireccional, el estadístico de la prueba D se define como:

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

donde:

- D_n es el estadístico de la prueba.
- \sup_x representa el supremo (máximo absoluto) sobre todos los valores de x .
- $F_n(x)$ es la función de distribución empírica de la muestra.
- $F(x)$ es la CDF teórica para comparación.

Prueba K-S Bidireccional

Para la comparación bidireccional de dos muestras empíricas, el estadístico se calcula como:

$$D_{n,m} = \sup_x |F_n(x) - G_m(x)| \quad (3)$$

donde:

- $D_{n,m}$ es el estadístico de la prueba.
- $F_n(x)$ y $G_m(x)$ son las funciones de distribución empíricas de las dos muestras, de tamaños n y m , respectivamente.

Interpretación

El valor de D indica la mayor diferencia vertical entre las CDF comparadas. Un valor bajo sugiere que las distribuciones son similares, mientras que uno alto indica diferencias notables. La significancia estadística de D se evalúa comparándolo con un valor crítico de referencia, dependiendo del nivel de significancia α y, en algunos casos, del tamaño de la muestra. Si D supera este valor crítico, se rechaza la hipótesis nula de igualdad de distribuciones.

Ejemplo 2

Verifica que los siguientes datos se distribuyen como una normal estándar con un nivel de significancia de 0.01.

$[-1,26, -1,12, -0,99, -0,72, -0,15, 0,07, 0,18, 0,29, 0,39, 0,45, 0,55, 0,59, 0,84, 0,86, 2,30, 2,57]$

Paso 1: Entender la Hipótesis Nula y Alternativa

Hipótesis Nula (H0): Los datos siguen la distribución teórica propuesta (normal estándar).

Hipótesis Alternativa (H1): Los datos no siguen la distribución teórica propuesta.

Paso 2: Calcular el Estadístico de la Prueba

Cuadro 1: Verifica que los siguientes datos se distribuyen como una normal estándar con un nivel de significancia $\alpha = 0,01$

X	$F_n(x)$	$F(x)$	$ F_n(x) - F(x) $
-1.26	1/16	0.1038	0.0413
-1.12	2/16	0.1314	0.0064
-0.99	3/16	0.1611	0.0264
-0.72	4/16	0.2358	0.0142
-0.15	5/16	0.4404	0.1279
0.07	6/16	0.5279	0.1529
0.18	7/16	0.5714	0.1339
0.29	8/16	0.6141	0.1141
0.39	9/16	0.6517	0.0892
0.45	10/16	0.6736	0.0486
0.55	11/16	0.7088	0.0213
0.59	12/16	0.7224	0.0276
0.84	13/16	0.7995	0.0130
0.86	14/16	0.8051	0.0699
2.30	15/16	0.9893	0.0518
2.57	16/16	0.9949	0.0051

$$D_{n,m} = \sup_x |F_n(x) - G_m(x)| = 0,1529 \quad (4)$$

$$q_n(1 - \alpha), n = 16, \alpha = 0,01 \quad (5)$$

Buscar en la tabla el cuantil

$$0,1529 < q_n(1 - \alpha) \quad (6)$$

$$0,1529 < 0,392 \quad (7)$$

Por lo tanto no se rechaza H_0 y $X \sim N(0, 1)$

5. 14/02/2024

5.1. Pruebas de Bondad de Ajuste

Estas pruebas determinan si una muestra observada sigue una distribución específica. La elección de la prueba depende del tipo de variable:

5.1.1. Para Variables Continuas

■ Prueba de Kolmogorov-Smirnov

- *Descripción:* Compara la distribución acumulativa de una muestra con una distribución teórica, evaluando la máxima diferencia.
- *Supuestos:* Muestra independiente y aleatoria.
- *Uso recomendado:* Comparar si la distribución de una muestra continua se ajusta a una distribución teórica específica.

5.1.2. Para Variables Discretas

■ Prueba de Ji-Cuadrada (χ^2)

- *Descripción:* Evalúa la discrepancia entre las frecuencias observadas y las esperadas bajo una hipótesis nula.
- *Supuestos:* Muestras grandes, expectativas de frecuencia en cada categoría mayores a 5.
- *Uso recomendado:* Verificar ajuste en datos categóricos o independencia entre dos variables categóricas.

5.2. Pruebas Basadas en Rangos

Alternativa no paramétrica para comparación de medianas cuando no se cumplen las suposiciones de normalidad.

5.2.1. Para 2 Poblaciones

■ ¿Son Independientes? Sí

- *Mann-Whitney U*
 - *Descripción:* Compara medianas de dos grupos independientes.
 - *Supuestos:* Muestras independientes.
 - *Uso recomendado:* Comparar tendencias centrales de dos muestras no relacionadas.

■ ¿Son Independientes? No

- *Wilcoxon de Rangos con Signos*
 - *Descripción:* Compara diferencias medias dentro de pares de observaciones relacionadas.
 - *Supuestos:* Pares emparejados o medidas repetidas.
 - *Uso recomendado:* Datos emparejados o medidas repetidas en el mismo grupo.

5.2.2. Para Más de 2 Poblaciones

■ ¿Son Independientes? Sí

- *Kruskal-Wallis*
 - *Descripción*: Generalización de Mann-Whitney para más de dos grupos, comparando medianas.
 - *Supuestos*: Muestras independientes entre sí.
 - *Uso recomendado*: Comparar medianas de tres o más grupos independientes.

■ ¿Son Independientes? No

- *Friedman*
 - *Descripción*: Comparación de tres o más grupos relacionados, basada en rangos.
 - *Supuestos*: Muestras relacionadas, como medidas repetidas.
 - *Uso recomendado*: Para datos relacionados con tres o más condiciones o tratamientos.

En las pruebas basadas en rango no se tiene que indicar la distribución de los datos, basta con que la distribución de la muestra sea continua

Definición (Rangos)

Sea x_1, x_2, \dots, x_n una muestra aleatoria sin empates, es decir, $x_i \neq x_j, \forall i \neq j$. Se define el rango de x_i , denotado por $R(x_i)$, como la posición que ocupa esta observación cuando los datos se ordenan de menor a mayor.

5.3. Prueba de Mann-Whitney

Sea $x_1, x_2, \dots, x_n \sim F_x$ y $y_1, y_2, \dots, y_m \sim F_y$ dos muestras aleatorias independientes de dos poblaciones. La prueba de Mann-Whitney U contrasta las siguientes hipótesis en relación a las distribuciones, las medias esperadas y las medianas de las dos poblaciones:

- Para las distribuciones:
 - $H_0 : F_x = F_y$ contra $H_1 : F_x \neq F_y$
- Para las medias esperadas:
 - $H_0 : E(x) = E(y)$ contra $H_1 : E(x) \neq E(y)$
- Para las medianas:
 - $H_0 : \text{Med}(x) = \text{Med}(y)$ contra $H_1 : \text{Med}(x) \neq \text{Med}(y)$

Estadístico de la prueba

El estadístico de la prueba Mann-Whitney u se calcula como el menor de los dos valores u_1 y u_2 , donde:

$$u_1 = nm + \frac{n(n+1)}{2} - T_1$$
$$u_2 = nm + \frac{m(m+1)}{2} - T_2$$

Donde T_1 y T_2 son las sumas de los rangos de las variables X y Y , respectivamente.

$$u = \min(u_1, u_2)$$

La esperanza matemática y la varianza del estadístico u son:

$$E(u) = \frac{nm}{2}$$
$$Var(u) = \frac{nm(n+m+1)}{12}$$

Bajo la hipótesis nula, el estadístico normalizado Z sigue una distribución normal estándar $N(0, 1)$:

$$Z = \frac{u - E(u)}{\sqrt{Var(u)}}$$

Se rechaza la hipótesis nula H_0 con un nivel de significancia α , si el valor absoluto de Z es mayor que el cuantil crítico $q_z(1 - \frac{\alpha}{2})$ de la distribución normal estándar:

$$|Z| > q_z\left(1 - \frac{\alpha}{2}\right)$$

El p -valor se calcula como la probabilidad de obtener un valor tan extremo o más que el valor observado de Z bajo la hipótesis nula:

$$p\text{-valor} = 2\mathbb{P}(Z > |z|)$$

Si $p\text{-valor} < \alpha$, se rechaza H_0

Ejemplo 3

Investigadores reunieron datos sobre el número de admisiones hospitalarias resultantes de choques de vehículos en diferentes días. A continuación, se presentan los resultados obtenidos para los viernes 6 y los viernes 13. Se utiliza un nivel de significancia de 0.05 para probar la aserción de que los viernes 13, el número de admisiones hospitalarias por choques de vehículos no se ve afectado.

Datos

Viernes 6: 9, 6, 11, 11, 3, 5

Viernes 13: 13, 12, 14, 10, 4, 12

Hipótesis

Las hipótesis a contrastar son:

$$H_0 : F_x = F_y \quad \text{contra} \quad H_1 : F_x \neq F_y$$

$$H_0 : E(x) = E(y) \quad \text{contra} \quad H_1 : E(x) \neq E(y)$$

Ordenación de los Datos y Cálculo de Rangos

Al ordenar todos los datos de menor a mayor, obtenemos la siguiente secuencia:

$$\{3_1, 4_2, 5_3, 6_4, 9_5, 10_6, 11_7, 11_8, 12_9, 12_{10}, 13_{11}, 14_{12}\}$$

- 3: 1
- 4: 2
- 5: 3
- 6: 4
- 9: 5
- 10: 6
- 11: 7.5
- 12: 9.5
- 13: 11
- 14: 12

suma de los rangos de viernes 6

$$T_1 = 5 + 4 + 7,5 + 7,5 + 1 + 3 = 28$$

suma de los rangos de viernes 13

$$T_2 = 11 + 9,5 + 12 + 6 + 2 + 9,5 = 50$$

$$n = 6 \quad m = 6$$

entonces

$$u_1 = nm + \frac{n(n+1)}{2} - T_1$$

$$u_2 = nm + \frac{m(m+1)}{2} - T_2$$

$$u_1 = 6 * 6 + \frac{6(6+1)}{2} - 28 = 29$$

$$u_2 = 6 * 6 + \frac{6(6+1)}{2} - 50 = 7$$

$$u = \min(u_1, u_2) = 7$$

6. 16/02/2024

Continuacion de la clase pasada

$$E(u) = \frac{nm}{2} = 18$$

$$Var(u) = 39$$

$$z = \frac{u - E(u)}{\sqrt{Var(u)}} = -1,76140$$

$$|z| = 1,76140$$

$$u = 7$$

$$q_z(1 - \frac{\alpha}{2}) = q_z(0,975) = 1,96$$

Se rechaza H_0 si

$$|z| > q_z(1 - \frac{\alpha}{2})$$

$$H_0 : F_x = F_y$$

$$H_1 : F_x \neq F_y$$

$$\text{p-valor} = 2\mathbb{P}(z > 1,76)$$

$$\text{p-valor} = 2[1 - \mathbb{P}(z < 1,76)]$$

$$\text{p-valor} = 2[0,0392]$$

$$\text{p-valor} = 0,0392$$

Se rechaza H_0 si

$$\text{p-valor} < \alpha$$

$$0,0784 < 0,05$$

6.1. Prueba de Wilcoxon

La prueba de Wilcoxon para muestras relacionadas es un test no paramétrico utilizado para comparar dos muestras relacionadas, emparejadas o medidas repetidas en un solo sujeto, con el objetivo de evaluar si sus distribuciones de población difieren.

Hipótesis

Sea X y Y dos variables aleatorias representando dos tratamientos o condiciones con pares de observaciones dependientes $(X_1, Y_1), \dots, (X_n, Y_n)$. Se desean contrastar las siguientes hipótesis:

$$H_0 : F_x = F_y$$

contra

$$H_1 : F_x \neq F_y$$

Estadístico de la Prueba

El estadístico de la prueba, T , se define como el mínimo entre T_+ y T_- , donde:

T_+ = suma de los valores absolutos de los rangos positivos

T_- = suma de los valores absolutos de los rangos negativos

Para cada par (X_i, Y_i) , calculamos la diferencia $D_i = X_i - Y_i$. Los rangos se asignan a los valores absolutos de estas diferencias, ignorando las diferencias que sean cero. T_+ es la suma de los rangos para las diferencias positivas, y T_- es la suma de los rangos para las diferencias negativas.

Aproximación Normal

Bajo la hipótesis nula y para un tamaño de muestra grande, el estadístico T puede aproximarse por una distribución normal con las siguientes expectativas y varianza:

$$E(T) = \frac{n(n+1)}{4}$$

$$Var(T) = \frac{n(n+1)(2n+1)}{24}$$

El valor z se calcula entonces como:

$$z = \frac{T - E(T)}{\sqrt{Var(T)}}$$

Este valor z se compara con los valores críticos de la distribución normal estándar para determinar si se rechaza o no la hipótesis nula.

Se rechaza H_0 con un nivel de significancia α , si

$$|z| > q_z\left(1 - \frac{\alpha}{2}\right)$$

o

$$\text{p-valor} = 2\mathbb{P}(z > |z|) < \alpha$$

Ejemplo 4

Los números de fusibles eléctricos defectuosos producidos por las líneas A y B se registraron a diario durante 10 días, con los siguientes resultados:

Cuadro 2: Verifica que los siguientes datos se distribuyen como una normal estándar con un nivel de significancia $\alpha = 0,05$

Línea A	Línea B	$ A - B $	Rango	R con signo
170	201	31	10	-10
164	179	15	7	-7
140	159	19	8	-8
184	195	11	5	-5
174	177	3	1	-1
142	170	28	9	-9
191	183	8	2	2
169	179	10	4	-4
161	170	9	3	-3
200	212	12	6	-6

Datos

El número 2 positivo en la tabla es positivo porque la diferencia de $A - B$ de sus respectivas líneas es positiva

$$T_+ = 2$$

$$T_- = 53$$

Estas T son la suma de los rangos

$$T = \min(T_+, T_-) = 2$$

$$E(T) = \frac{n(n+1)}{4} = \frac{10(10+1)}{4} = 27,5$$

$$Var(T) = \frac{n(n+1)(2n+1)}{24} = 96,25$$

$$z = \frac{T - E(T)}{\sqrt{Var(T)}} = -2,5993$$

$$\alpha = 0,05$$

$$q_z(1 - \frac{\alpha}{2}) = 1,96$$

Se rechaza H_0 si

$$|z| > q_z(1 - \frac{\alpha}{2})$$

$$2,5993 > 1,96$$

Por lo tanto las funciones de densidad son diferentes

7. 19/02/2024

7.1. Prueba de Kruskal-Wallis

Consideremos una característica X presente en k poblaciones independientes. Supongamos que estamos interesados en determinar si la distribución de X es la misma en todas las k poblaciones.

Sea F_i la distribución de X en la i -ésima población. Las hipótesis que se plantean son:

- Hipótesis nula (H_0): Todas las distribuciones son iguales, es decir,

$$H_0 : F_1 = F_2 = \dots = F_k.$$

- Hipótesis alternativa (H_1): Al menos una de las distribuciones F_i es diferente de las demás, es decir,

$$H_1 : \text{Al menos una } F_i \text{ es diferente de las demás.}$$

Estadístico de la prueba

El estadístico de la prueba de Kruskal-Wallis, H , se calcula como:

$$H = \frac{12}{n(n+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(n+1) \sim \chi_{k-1}^2,$$

donde n es el tamaño total de la muestra (la suma de los tamaños de todas las muestras), n_i es el tamaño de la muestra para la i -ésima población, y R_i es la suma de los rangos de las observaciones en la i -ésima muestra. Este estadístico sigue aproximadamente una distribución chi-cuadrado (χ^2) con $k-1$ grados de libertad, bajo la hipótesis nula.

Reglas de decisión

Para tomar una decisión sobre la hipótesis nula H_0 con un nivel de significancia α , se utilizan las siguientes reglas:

- Se rechaza H_0 si el estadístico de la prueba H es mayor que el valor crítico de la distribución chi-cuadrado con $k-1$ grados de libertad, es decir,

$$H > q_{\chi_{k-1}^2}(1-\alpha),$$

donde $q_{\chi_{k-1}^2}(1-\alpha)$ es el cuantil $(1-\alpha)$ -ésimo de la distribución chi-cuadrado con $k-1$ grados de libertad.

- Alternativamente, se rechaza H_0 si el p-valor asociado al estadístico de la prueba es menor que el nivel de significancia α , es decir,

$$\text{p-valor} = \mathbb{P}(\chi_{k-1}^2 > H) < \alpha.$$

Esto significa que si el estadístico de la prueba calculado H es suficientemente grande o si el p-valor es suficientemente pequeño, tenemos evidencia estadística para rechazar la hipótesis nula de que todas las poblaciones tienen distribuciones idénticas.

Ejemplo 5

Se tomaron muestras aleatorias independientes de casas recientemente vendidas en 4 zonas residenciales de una ciudad. El objetivo es determinar si existen diferencias significativas entre las zonas con respecto al cociente entre el precio de venta y el valor catastral de las propiedades.

Datos

Los datos son los siguientes:

Cuadro 3: **Cocientes entre los precios de venta y el valor catastral**

Zona 1	Zona 2	Zona 3	Zona 4
1.19	1.08	0.98	1.12
1.05	1.23	1.19	1.14
1.14	1.26	1.08	1.31
1.25	1.10	0.93	1.12
1.29	1.18	1.23	1.19
	1.14	1.18	

Asignación de Rangos

Cuadro 4: Asigancion de rangos (Ordenamiento del menor al mayor asignandoles un numero)

Zona 1	Zona 2	Zona 3	Zona 4
1,19 ₁₆	1,08 ₄	0,98 ₂	1,12 ₇
1,05 ₃	1,23 ₁₈	1,19 ₁₅	1,14 ₁₁
1,14 ₉	1,26 ₂₀	1,08 ₅	1,31 ₂₂
1,25 ₁₉	1,10 ₆	0,93 ₁	1,12 ₈
1,29 ₂₁	1,18 ₁₂	1,23 ₁₇	1,19 ₁₄
	1,14 ₁₀	1,18 ₁₃	

- 1.08 se repite en 4 y 5 $4 + 5/2 = 4,5$
- 1.12 se repite en 7 y 8 $7 + 8/2 = 7,5$
- 1.14 se repite en 9, 10 y 11, $9 + 10 + 11/3 = 10$
- 1.18 se repite en 12 y 13 $12 + 13/2 = 12,5$
- 1.19 se repite en 14 y 15 y 16 $14 + 15 + 16/2 = 15$
- 1.23 se repite en 17 y 18 $17 + 18/2 = 17,5$

Por lo tanto los rangos se asignaron de la siguiente manera:

Cuadro 5: **Asignación de rangos y su promedio en caso de repetición**

Zona 1	Zona 2	Zona 3	Zona 4
1,19 ₁₅	1,08 _{4,5}	0,98 ₂	1,12 _{7,5}
1,05 ₃	1,23 _{17,5}	1,19 ₁₅	1,14 ₁₀
1,14 ₁₀	1,26 ₂₀	1,08 _{4,5}	1,31 ₂₂
1,25 ₁₉	1,10 ₆	0,93 ₁	1,12 _{7,5}
1,29 ₂₁	1,18 _{12,5}	1,23 _{17,5}	1,19 ₁₅
	1,14 ₁₀	1,18 _{12,5}	

Cuadro 6: **Número de observaciones y suma de rangos**

Zona 1	Zona 2	Zona 3	Zona 4
$n_1 = 5$	$n_2 = 6$	$n_3 = 5$	$n_4 = 4$
$R_1 = 68$	$R_2 = 70,5$	$R_3 = 52,5$	$R_4 = 62$

Cálculo de la Estadística de Prueba

El número de observaciones (n) y la suma de rangos (R) para cada zona son:
La estadística de prueba H se calcula como:

$$H = \frac{12}{n(n+1)} \left(\sum \frac{R_i^2}{n_i} \right) - 3(N+1)$$

Donde N es el número total de observaciones. Para este caso:

$$H = \frac{12}{22 \times 23} \left(\frac{68^2}{5} + \frac{70,5^2}{6} + \frac{52,5^2}{5} + \frac{62^2}{4} \right) - 3(22+1) = 1,70395$$

Conclusión

Se rechaza H_0 si

$$H > q_{X_{k-1}^2}(1-\alpha)$$

$$q_{X_{k-1}^2}(0,95) = 7,815$$

$$1,70395 > 7,815$$

No rechazamos H_0 . Por lo tanto, concluimos que no hay evidencia suficiente para afirmar que existen diferencias significativas entre las zonas en cuanto al cociente entre el precio de venta y el valor catastral de las propiedades.

La expresión para el valor crítico de la distribución chi-cuadrado con 3 grados de libertad es:

$$X_3^2 = q_{X_3^2}(1-\alpha)$$

El cálculo del p-valor se realiza como sigue:

$$\text{p-valor} = \mathbb{P}(X_3^2 > 1,70395)$$

Para estimar el p-valor, se puede utilizar la aproximación a la distribución normal estándar, de manera que:

$$\mathbb{P}(X_3^2 > 1,70395) \approx \mathbb{P}\left(\frac{X - E(X)}{\sqrt{\text{Var}(X)}} > \frac{1,70395 - 3}{\sqrt{6}}\right)$$

Simplificando la expresión, obtenemos:

$$= \mathbb{P}(Z > -0,53)$$

Finalmente, utilizando la tabla de la distribución normal estándar, calculamos el p-valor como:

$$= 1 - \mathbb{P}(Z < -0,53) = 1 - 0,2981 = 0,7019$$

8. 21/02/2024

8.1. Pruebas basadas en corridas

Las pruebas basadas en corridas son utilizadas en estadística para determinar si una secuencia de elementos es aleatoria. Esto es particularmente útil para analizar patrones dentro de series de datos donde se espera una distribución aleatoria. Las hipótesis que se contrastan en este tipo de prueba son:

- H_0 : Los datos siguen una secuencia aleatoria.
- H_1 : Los datos no siguen una secuencia aleatoria.

Estadística de la prueba

La estadística de la prueba se basa en el número total de corridas (R_t), donde una *corrida* se define como una secuencia ininterrumpida de elementos similares (por ejemplo, una serie de números crecientes o decrecientes).

$$R_t = \text{Número total de corridas}$$

Aproximación normal

Bajo la hipótesis nula de aleatoriedad, el valor esperado y la varianza de R_t pueden aproximarse mediante las siguientes fórmulas, donde n_1 y n_2 son el número de elementos en cada uno de los dos grupos definidos (por ejemplo, números por encima y por debajo de la mediana, presencia o ausencia de una característica):

$$E(R_t) = \frac{2n_1n_2}{n} + 1$$

$$\text{Var}(R_t) = \frac{2n_1n_2(2n_1n_2 - n)}{n^2(n - 1)}$$

Una vez calculados el valor esperado y la varianza, se puede utilizar la siguiente fórmula para determinar el valor z , el cual indica cuántas desviaciones estándar se encuentra el número observado de corridas del esperado bajo la hipótesis nula:

$$z = \frac{R_t - E(R_t)}{\sqrt{\text{Var}(R_t)}}$$

Este valor de z se compara entonces con los valores críticos de la distribución normal estándar para determinar si se rechaza o no la hipótesis nula de aleatoriedad en los datos.

Ejemplo 6

Consideremos la siguiente secuencia de datos para determinar si es aleatoria:

$$1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 0, 1, 0$$

Para analizar la aleatoriedad, agrupamos los datos en corridas según su valor. Una *corrida* se define como una secuencia consecutiva de números iguales. La secuencia se divide en los siguientes grupos (corridas):

- Grupo 1: 1
- Grupo 2: 0
- Grupo 3: 0, 0, 0
- Grupo 4: 1, 1
- Grupo 5: 0
- Grupo 6: 1
- Grupo 7: 0
- Grupo 8: 1
- Grupo 9: 0
- Grupo 10: 1
- Grupo 11: 0

Cada grupo representa una corrida, por lo que el número total de corridas (R_t) en esta secuencia es igual al número de grupos:

$$R_t = 11$$

Para calcular el valor esperado ($E(R_t)$) y la varianza ($Var(R_t)$) de las corridas bajo la hipótesis de aleatoriedad, primero identificamos n_1 y n_2 , que son las cantidades de 1s y 0s respectivamente:

$$n_1 = 8 \quad (\text{número de 1s})$$

$$n_2 = 7 \quad (\text{número de 0s})$$

El total de observaciones (n) es la suma de n_1 y n_2 :

$$n = n_1 + n_2 = 15$$

Con n_1 , n_2 , y n definidos, podemos proceder a calcular $E(R_t)$ y $Var(R_t)$ utilizando las fórmulas para la aproximación normal de la prueba de corridas, y luego determinar si la secuencia es aleatoria comparando el valor observado de R_t con el valor esperado bajo la hipótesis nula de aleatoriedad.

9. 23/02/2024

9.1. Pruebas de Independencia

Definición

Una medida numérica de asociación para las variables aleatorias continuas X y Y , denotada por $\mu_{x,y}$, es considerada una medida de dependencia si cumple con las siguientes propiedades:

1. $0 \leq \mu_{x,y} \leq 1$, donde los límites inferior y superior representan la independencia total y la dependencia funcional perfecta, respectivamente.
2. $\mu_{x,y} = \mu_{y,x}$, lo que indica que la medida de dependencia es simétrica respecto a X y Y .
3. X y Y son independientes si y solo si $\mu_{x,y} = 0$. Esto establece un criterio numérico para la independencia de X y Y .
4. $\mu_{x,y} = 1$ si y solo si X y Y son casi seguramente funciones estrictamente monótonas una de la otra. Esto describe una relación de dependencia funcional completa.
5. Si α y β son funciones estrictamente monótonas con probabilidad 1, entonces $\mu_{\alpha(x),\beta(y)} = \mu_{x,y}$. Esta propiedad asegura que la medida de dependencia es invariante bajo transformaciones monótonas.
6. Si $\{X_n, Y_n\}$ es una secuencia de variables aleatorias continuas con cópulas subyacentes C_n y si $\{C_n\}$ converge a una cópula C , entonces

$$\lim_{n \rightarrow \infty} \mu_{C_n} = \mu_C.$$

Esto establece la continuidad de la medida de dependencia en términos de convergencia de cópulas.

Estas propiedades definen formalmente cómo se debe medir la dependencia entre dos variables aleatorias continuas y establecen criterios claros para su evaluación.

9.2. Coeficiente de Correlación de Spearman

El coeficiente de correlación de Spearman, denotado como r_s , es una medida no paramétrica de la correlación de rango, que evalúa la dependencia entre dos variables cuantificando la dirección y la intensidad de su asociación monótona.

La fórmula para calcular el coeficiente de Spearman es:

$$r_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)},$$

donde d_i es la diferencia entre los rangos correspondientes de las dos variables, y n es el número de observaciones.

Prueba del Coeficiente de Spearman

Para evaluar la significancia estadística de la correlación de Spearman entre dos variables, se utiliza la siguiente prueba de hipótesis:

- Hipótesis nula (H_0): No existe dependencia monótona entre las dos variables. Esto se formaliza como $F(x, y) = F(x)F(y)$, donde $F(x, y)$ es la distribución conjunta de las variables y $F(x)$, $F(y)$ son sus distribuciones marginales.

- Hipótesis alternativa (H_1): Existe una dependencia monótona entre las dos variables, lo que implica que $F(x, y) \neq F(x)F(y)$.

Para realizar el contraste se sigue el procedimiento:

1. Ordenar los datos de cada variable y asignar rangos.
2. Calcular las diferencias de rango d_i para cada par de observaciones.
3. Insertar los valores de d_i en la fórmula de r_s para obtener el coeficiente de Spearman.
4. Utilizar el valor calculado de r_s y el número de observaciones n para determinar el valor p asociado a través de tablas o software estadístico.
5. Comparar el valor p con un nivel de significancia preestablecido (usualmente 0.05 o 0.01) para decidir si se rechaza H_0 .

Este procedimiento permite evaluar si la correlación observada entre las variables es estadísticamente significativa, indicando la presencia de una relación monótona entre ellas.

Ejemplo 7

Ocho profesores de ciencias han sido clasificados por un juez de acuerdo a su capacidad de enseñanza y todos han tomado la prueba. ¿Cuál es la correlación entre la calificación del juez y la calificación de la prueba?

Cuadro 7: Calificaciones y rangos de los profesores

Profesor	Calf. Juez	Calf. Prueba	$R(x_i)$	$R(y_i)$	d_i^2
1	7	44	7	1	36
2	4	72	4	5	1
3	2	69	2	3	1
4	6	70	6	4	4
5	1	93	1	8	49
6	3	82	3	7	16
7	8	67	8	2	36
8	5	80	5	6	1

La diferencia de rangos (d_i) se calcula restando el rango de la calificación de la prueba ($R(y_i)$) del rango de la calificación del juez ($R(x_i)$) para cada profesor. Luego, elevamos al cuadrado esta diferencia para obtener d_i^2 , lo que nos ayuda a medir la discrepancia entre los rangos de las dos variables sin importar la dirección de la diferencia.

$$d_i^2 = (R(X_i) - R(Y_i))^2$$

Con la suma de d_i^2 igual a 144, calculamos la correlación de Spearman entre las calificaciones del juez y las calificaciones de la prueba utilizando la fórmula:

$$r_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}$$

$$r_s = 1 - \frac{6 * 144}{8(8^2 - 1)}$$

donde n es el número de observaciones (profesores), que en este caso es 8. La correlación de Spearman calculada es $-0,714$, lo que indica una fuerte correlación negativa entre las dos variables. Esto significa que a medida que una variable aumenta, la otra tiende a disminuir.

10. 26/02/2024

Análisis de la Correlación de Spearman y Prueba de Hipótesis

Una vez calculada la correlación de Spearman, es crucial determinar si esta correlación es estadísticamente significativa. Bajo la hipótesis nula $H_0 : \rho = 0$, no esperamos encontrar correlación entre las calificaciones del juez y las de la prueba. La correlación de Spearman, r_s , debería ser 0 si H_0 es cierta. Utilizamos la siguiente estadística de prueba que se distribuye aproximadamente como una normal estándar $N(0, 1)$ para muestras de tamaño considerable:

$$\frac{r_s - E(r_s)}{\sqrt{\text{var}(r_s)}}$$

Donde $E(r_s) = 0$ bajo H_0 , y $\text{var}(r_s) = \frac{1}{n-1}$, siendo n el número de observaciones. Para nuestro caso, con $n = 8$ y $r_s = -0,714$, la estadística de prueba se calcula como:

$$\frac{-0,714 - 0}{\sqrt{\frac{1}{8-1}}} = -1,889$$

Este resultado se compara con los valores críticos de la distribución $N(0, 1)$. Para un nivel de significancia $\alpha = 0,05$ en una prueba de dos colas, los valores críticos son $\pm 1,96$. Al no ser -1.889 menor que -1.96 ni mayor que 1.96, no rechazamos la hipótesis nula al nivel de significancia del 5 %. Esto indica que no hay evidencia suficiente para afirmar que existe una correlación significativa entre las dos variables evaluadas.

Para formalizar el proceso de toma de decisiones, se rechaza H_0 si:

$$|z| > q_z\left(1 - \frac{\alpha}{2}\right)$$

Utilizando el valor absoluto de z , calculamos:

$$|z| = \left| \frac{r_s}{\sqrt{\frac{1}{n-1}}} \right| = |\sqrt{n-1} r_s|$$

Sustituyendo los valores conocidos:

$$|z| = |\sqrt{8-1} - 0,714| = |-1,889| \approx 1,889$$

Comparando este valor con $q_z(1 - \frac{\alpha}{2}) = q_z(0,975)$, ($q_z(0,975)$ se calcula en phyton con `norm.ppf(1 - alpha / 2)`) aproximadamente 1.96, concluimos:

Dado que $|z| = 1,889 > 1,96$ no se cumple, entonces no rechazamos la hipótesis nula H_0 al nivel de significancia de 0,05. Por lo tanto, la calificación del juez y la calificación de la prueba son independientes

10.1. Tau de Kendall

La Tau de Kendall es una medida de correlación basada en el orden de los rangos de los datos, que evalúa la similitud entre las ordenaciones de los datos en dos variables. A diferencia de la correlación de Pearson, que evalúa la relación lineal entre variables cuantitativas, la Tau de Kendall se centra en la relación ordinal, lo que la hace útil en situaciones donde la relación no es necesariamente lineal.

Definición

Consideremos un conjunto de observaciones de un vector aleatorio (X, Y) , donde X e Y son variables aleatorias continuas. Sea (X_i, Y_i) y (X_j, Y_j) dos observaciones distintas de este vector, con $i \neq j$.

Definimos estas observaciones como **concordantes** si los pares se mueven en la misma dirección, es decir, si:

$$(x_i < x_j \text{ y } y_i < y_j) \quad \text{o} \quad (x_i > x_j \text{ y } y_i > y_j),$$

lo que indica que tanto X como Y aumentan o disminuyen juntos.

Por otro lado, los pares se consideran **discordantes** si se mueven en direcciones opuestas, es decir, si:

$$(x_i < x_j \text{ y } y_i > y_j) \quad \text{o} \quad (x_i > x_j \text{ y } y_i < y_j),$$

indicando que uno aumenta mientras el otro disminuye.

La Tau de Kendall, τ , se calcula como la diferencia entre la proporción de pares concordantes y la proporción de pares discordantes, sobre el total de pares. Esto proporciona una medida de la correlación entre las variables, donde un valor de $\tau = 1$ indica una correlación perfecta, $\tau = -1$ indica una correlación inversa perfecta, y $\tau = 0$ sugiere que no hay correlación.

Si $(x_i - x_j)(y_i - y_j)$ son concordantes

$$(x_i - x_j)(y_i - y_j) \geq 0$$

Si $(x_i - x_j)(y_i - y_j)$ son discordantes

$$(x_i - x_j)(y_i - y_j) \leq 0$$

Cálculo de la Tau de Kendall

Para calcular la Tau de Kendall, se siguen los siguientes pasos:

1. Ordenar las observaciones de acuerdo con los valores de X , manteniendo el emparejamiento con sus correspondientes valores de Y .
2. Para cada observación (X_i, Y_i) , contar el número de pares concordantes y discordantes comparándola con las demás observaciones (X_j, Y_j) donde $j > i$.
3. Calcular τ utilizando la fórmula:

$$\tau = \frac{\text{Número de pares concordantes} - \text{Número de pares discordantes}}{\binom{n}{2}},$$

donde $\binom{n}{2}$ es el número total de combinaciones de pares únicas posibles entre las observaciones.

La Tau de Kendall ofrece una medida robusta y fiable de la correlación ordinal, útil en análisis estadísticos donde las suposiciones de linealidad y normalidad no se cumplen, o cuando se trabaja con datos ordinales.

11. 01/03/3034

11.1. Tablas de contingencia

1.5.3

Version 1

Dados los siguientes datos:

- r : número de renglones
- c : número de columnas

Las hipótesis son:

- H_0 : $F(x, y) = F(x)F(y)$
- H_1 : $F(x, y) \neq F(x)F(y)$

El estadístico de la prueba es:

$$\chi^2 = \sum \left(\frac{(Fo - Fe)^2}{Fe} \right) \sim \chi^2_{(r-1)(c-1)}$$

Donde:

- Fo : frecuencia observada
- Fe : frecuencia esperada

Se rechaza H_0 si:

$$\chi^2 > \chi^2_{\alpha, (r-1)(c-1)}$$

Se rechaza H_0 si el p-valor α :

$$p\text{-valor} = P(\chi^2_{(r-1)(c-1)} > \chi^2_{\text{estadístico}})$$

Version 2

ados dos eventos X e Y , las hipótesis son:

- H_0 : $F(x, y) = F(x)F(y)$
- H_1 : $F(x, y) \neq F(x)F(y)$

El estadístico de la prueba es:

$$X^2 = \sum \frac{(Fo - FE)^2}{FE} \sim \chi^2_{(r-1)(c-1)}$$

donde:

- r es el número de renglones
- c es el número de columnas

Se rechaza H_0 si:

$$X^2 > \chi^2_{\alpha, (r-1)(c-1)}$$

Se rechaza H_0 si p-valor $< \alpha$

$$p\text{-valor} = P(X^2_{(r-1)(c-1)} > \text{estadístico})$$

11.1.1. Ejemplo 9

Un total de $n = 309$ defectos en muebles fueron registrados y los defectos fueron clasificados en 4 tipos: A, B, C y D. Al mismo tiempo, cada pieza de mueble fue identificada por la línea de producción en la que se elaboró. Verifica que los defectos son independientes de la línea de producción.

Version Mike

Frecuencias observadas: F_0

Tipos	1	2	3	suma
A	15	26	33	74
B	21	31	17	69
C	45	34	49	128
D	13	5	20	38
Suma	94	96	119	309

Frecuencias esperadas: F_E

Tipos	1	2	3	suma
A	$\frac{94(74)}{309}$	$\frac{96(74)}{309}$...	74
B	$\frac{94(69)}{309}$	$\frac{96(69)}{309}$...	69
C	$\frac{94(128)}{309}$	$\frac{96(128)}{309}$...	128
D	$\frac{94(38)}{309}$	$\frac{96(38)}{309}$...	38
Suma	94	96	119	

Frecuencias esperadas: F_E

Tipos	1	2	3
A	22.51	22.99	28.49
B	20.99	21.43	26.57
C	38.94	39.76	49.29
D	11.56	11.80	14.63

$$X^2 = \frac{(15 - 22,51)^2}{22,51} + \frac{(26 - 22,99)^2}{22,99} + \dots + \frac{(20 - 14,63)^2}{14,63} \quad (8)$$

$$X^2 = 2,50 + 0,39 + 0,71 + \dots = 19,19 \quad (9)$$

$$q_{(X^2)_{(r-1)(c-1)}^{1-\alpha}} = q_{(X^2_6)^{(0,95)}} = 12,59$$

Se rechaza H_0 si

$$X^2 > q_{(X^2_6)^{0,95}}$$

$$19,19 > 12,59$$

Se rechaza H_0 , Por lo tanto, los tipos de defectos dependen de las líneas de producción. El p-valor se calcula como sigue:

$$\begin{aligned}
\text{p-valor} &= P(\chi^2 > 19,19) \\
&= P\left(\frac{\chi^2 - E(\chi^2)}{\sqrt{\text{Var}(\chi^2)}} > \frac{19,19 - E(\chi^2)}{\sqrt{\text{Var}(\chi^2)}}\right) \\
&= P\left(Z > \frac{19,19 - 6}{\sqrt{12}}\right) \\
&= P(Z > 3,81) \\
&= 1 - P(Z \leq 3,81) \\
&= 1 - 0,9999 \approx 0
\end{aligned}$$

Dado que el p-valor $\leq \alpha = 0,05$, entonces se acepta H_1 .

Version ChatGPT

Enunciado Para ello, se observaron las siguientes frecuencias de defectos:

Tipo de Defecto	Línea 1	Línea 2	Línea 3	Total
A	15	26	33	74
B	21	31	17	69
C	45	34	49	128
D	13	5	20	38
Total	94	96	119	309

Cuadro 8: Frecuencias observadas (F_O)

Y se calcularon las frecuencias esperadas (F_E) bajo la hipótesis nula de independencia:

Tipo de Defecto	Línea 1	Línea 2	Línea 3
A	$\frac{94 \times 74}{309}$	$\frac{96 \times 74}{309}$...
B	$\frac{94 \times 69}{309}$	$\frac{96 \times 69}{309}$...
C	$\frac{94 \times 128}{309}$	$\frac{96 \times 128}{309}$...
D	$\frac{94 \times 38}{309}$	$\frac{96 \times 38}{309}$...

Cuadro 9: Frecuencias esperadas (F_E)

Luego, se procedió a calcular el estadístico de prueba X^2 :

$$X^2 = \sum \frac{(F_0 - F_E)^2}{F_E} = \frac{(15 - 22,51)^2}{22,51} + \frac{(26 - 22,99)^2}{22,99} + \dots + \frac{(20 - 14,63)^2}{14,63} = 19,19 \quad (10)$$

Se compara el valor calculado con el valor crítico del estadístico χ^2 para el nivel de significancia deseado y los grados de libertad correspondientes $((r - 1)(c - 1))$:

$$\chi_{\text{crítico}, 0,95, 6}^2 = 12,59$$

La hipótesis nula H_0 , que postula la independencia, se rechazará si el estadístico de prueba es mayor que el valor crítico:

$$19,19 > 12,59$$

Dado que esto se cumple, se concluye que hay suficiente evidencia para rechazar la hipótesis nula, indicando que los defectos no son independientes de la línea de producción

11.2. Pruebas de una cola (Prueba de Mann-Whitney y Prueba de Wilcoxon)

Las hipótesis para la Prueba de Mann-Whitney son:

$$\begin{aligned} H_0 : E(X) = E(Y) & \quad \text{vs} \quad H_1 : E(X) < E(Y) \\ H_0 : E(X) \geq E(Y) & \quad \text{vs} \quad H_1 : E(X) < E(Y) \end{aligned}$$

Se rechaza H_0 si:

$$\frac{T - E(T)}{\sqrt{\text{Var}(T)}} \leq z_\alpha$$

Para la Prueba Wilcoxon, las hipótesis son:

$$\begin{aligned} H_0 : E(X) = E(Y) & \quad \text{vs} \quad H_1 : E(X) > E(Y) \\ H_0 : E(X) \leq E(Y) & \quad \text{vs} \quad H_1 : E(X) > E(Y) \end{aligned}$$

Se rechaza H_0 si:

$$\frac{T - E(T)}{\sqrt{\text{Var}(T)}} \geq z_\alpha$$

En ambos casos, T es el estadístico de prueba y z_α es el valor crítico de la distribución normal estándar para un nivel de significancia α .

Fin Unidad 1

12. 11/03/2024

Unidad 2 - Análisis de varianza y diseño de experimentos

Introducción

Un experimento es una prueba o una serie de pruebas en las que se hacen cambios deliberados en las variables de entrada de un sistema o un proceso para identificar las razones de los cambios que pueden registrarse en la respuesta de salida.

El diseño de experimentos se refiere al proceso para planificar el experimento de tal forma que se recaben los datos adecuados que puedan analizarse con métodos estadísticos que lleven a conclusiones válidas y objetivas.

Los tres principios del diseño de experimentos son:

- Realización de réplicas
- Aleatorización
- Formación de bloques

Lineamientos para diseñar experimentos

- Identificar el problema.
- Elección de factores, niveles y tratamientos.
- Definir la unidad experimental.
- Seleccionar las variables de respuesta.
- Elección del diseño experimental.
- Determinar el número de réplicas.
- Ejecutar el experimento.
- Analizar los datos.
- Conclusiones y recomendaciones.

13. 13/03/2024

Factor controlable

Es la característica que controlamos y cuyo efecto se desea estudiar.

Nivel

Es la categoría estudiada del factor.

Tratamientos

Combinaciones de los factores estudiados.

Unidad experimental

Es la subdivisión del material que puede recibir un tratamiento de forma independiente.

Variable respuesta

Es la característica que se va a medir en cada unidad experimental.

Diseño experimental

Es la forma de asignar los tratamientos a las unidades experimentales.

Error experimental

Describe la variación entre las unidades experimentales.

Ejemplo 1

Supóngase que de un grupo se seleccionaron al azar 20 hombres y 20 mujeres. Después, estos grupos se dividen al azar en grupo experimental y grupo de control. Al grupo de control se le aplica una prueba después de desayunar; al grupo experimental se le aplica la misma prueba en ayuno. Identifica los factores, niveles, tratamientos de este experimento.

13.1. Análisis de varianza

El método de análisis de varianza consiste en probar si varias medidas poblacionales son iguales. Cuando las varianzas son iguales y desconocidas, de manera que las hipótesis que se contrastan son:

$$H_0 : \mu_1 = \cdots = \mu_k$$

$$H_1 : \mu_i \neq \mu_j, \text{ para } i \neq j$$

para k poblaciones normales independientes y con varianza σ^2 .

13.2. Diseño completamente aleatorizado

Es un diseño completamente aleatorio, la asignación de los tratamientos a cada unidad experimental es de forma aleatoria. Además, todas las unidades experimentales son homogéneas.

13.3. Modelo

$$y_{ij} = \mu_i + E_{ij}$$

Donde y_{ij} es la variable respuesta del i -ésimo tratamiento y E_{ij} es el error experimental. En este caso, hay t tratamientos y r repeticiones.

13.4. Modelo completo

$$y_{ij} = \mu_i + E_{ij}$$

Cada tratamiento tiene una media diferente, es decir, el modelo completo es la representación de H_1 .

13.5. Modelo reducido

$$y_{ij} = \mu + E_{ij}$$

Cada tratamiento tiene la misma media, es decir, el modelo reducido es la representación de H_0 .

El ojetivo es estimar los parametros de ambos modelos y determinar cual de los dos se ajusta mas a los datos observados.

Para la solucionar este problema, se utilizaran los estimadores de minimos cuadrados que se obtienen al minimizar la suma de los cuadrados de los errores experimentales

Estimador del modelo completo

$$y_{ij} = \mu_i + E_{ij}$$

$$E_{ij} = y_{ij} - \mu_i$$

Por lo anterior, la suma de los cuaddrados de los errores del modelo completo es:

$$SS_{Ec} = \sum_t \sum_r^{j=1} (y_{ij} - \mu_i)^2$$

para una i fija

$$\frac{d[\sum_r^{j=1} (y_{ij} - \mu_i)^2]}{d\mu_i} = \frac{d[\sum_r^{j=1} (y_{ij}^2 - 2y_{ij}\mu_i + \mu_i^2)]}{d\mu_i} = \sum_r^{j=1} (-2y_{ij} + 2\mu_i)$$

al igualar a cero se tiene que

$$\begin{aligned} -2\left(\sum_r^{j=1}\right) &= 0 \\ -2\sum_r^{j=1} -r\mu_i &= 0 \\ \mu_i &= \frac{\sum_r^{j=1} y_{ij}}{r} \end{aligned}$$

Estimador de modelo reducido

$$y_{ij} = \mu + E_{ij}$$

$$E_{ij} = y_{ij} - \mu$$

Por lo tanto la suma de los cuadrados de los errores del modelo reducido es

$$SS_{Er} = \sum_t \sum_r^{j=1} (y_{ij} - \mu)^2$$

A continuacion, encontramos el punto critico de la funcion SS_{Er}

14. 15/03/2024

Ejemplo 2

Un ingeniero de desarrollo de productos tiene interés en investigar la resistencia a la tensión de una fibra que se usará para hacer playeras. El ingeniero sabe por experiencia previa que la resistencia a la tensión se afecta por el peso porcentual de algodón utilizado en la mezcla de materiales de fibra. Además sospecha que al aumentar el contenido de algodón se incrementará la resistencia, al menos en un principio el ingeniero decide probar ejemplares en cinco niveles del peso porcentual del algodón: 15, 20, 25, 30, 35 %. También decide probar 5 ejemplares en cada nivel del contenido de algodón. Los resultados que obtuvo, se muestra en la siguiente tabla

Datos

Cuadro 10: Verifica que los siguientes datos se distribuyen como una normal estándar con un nivel de significancia $\alpha = 0,05$

15 %	20 %	25 %	30 %	35 %
7	12	14	19	7
7	17	18	25	10
15	12	18	22	11
11	18	19	19	15
9	18	19	23	11

- $\bar{Y}_1 = 9,8$
- $\bar{Y}_2 = 15,4$
- $\bar{Y}_3 = 17,6$
- $\bar{Y}_4 = 21,6$
- $\bar{Y}_5 = 10,8$
- $\bar{Y} = 15,04$

$$SS_{\text{Tratamiento}} = \sum_t \sum_{i=1}^r (\bar{Y}_i - \bar{Y})^2$$

$$SS_{\text{Tratamiento}} = r \sum_t (\bar{Y}_i - \bar{Y})^2$$

$$SS_{\text{Tratamiento}} = 5 [(9,8 - 15,04)^2 + (15,4 - 15,04)^2 + (17,6 - 15,04)^2 \dots = 475,76]$$

15. 01/04/2024

Ejemplo 3 (continuación)

$$T_\alpha = 5,38 \quad \text{Prueba de Tukey}$$

Región de rechazo

$$|\bar{y}_i - \bar{y}_j| > T_\alpha$$

$$\binom{5}{2} = 10 \text{ Comparaciones}$$

$$|\bar{y}_1 - \bar{y}_2| = 5,6 > T_\alpha$$

$$|\bar{y}_1 - \bar{y}_3| = 7,8 > T_\alpha$$

$$|\bar{y}_1 - \bar{y}_4| = 11,8 > T_\alpha$$

$$|\bar{y}_1 - \bar{y}_5| = 1 \not> T_\alpha$$

$$|\bar{y}_2 - \bar{y}_3| = 2,2 \not> T_\alpha$$

$$|\bar{y}_2 - \bar{y}_4| = 6,2 > T_\alpha$$

$$|\bar{y}_2 - \bar{y}_5| = 4,6 \not> T_\alpha$$

$$|\bar{y}_3 - \bar{y}_4| = 4 \not> T_\alpha$$

$$|\bar{y}_3 - \bar{y}_5| = 6,8 > T_\alpha$$

$$|\bar{y}_4 - \bar{y}_5| = 10,8 > T_\alpha$$

$$y_1 = 9,8$$

$$y_2 = 15,4$$

$$y_3 = 17,6$$

$$y_4 = 21,6$$

$$y_5 = 10,8$$

15.1. Diseño de bloques completamente aleatorizados

Un diseño de bloques completamente aleatorizados se utiliza para controlar los factores de ruido conocidos que tienen un efecto en la variable respuesta pero que no son de interés para el estudio.

Modelo:

$$y_{ij} = \mu + \tau_i + \beta_j + E_{ij}$$

donde:

$$\mu = \text{Media general}$$

$$\tau_i = \text{Efecto del } i\text{-ésimo tratamiento}$$

$$\beta_j = \text{Efecto del } j\text{-ésimo bloque}$$

$$E_{ij} = \text{Error experimental de la } j\text{-ésima unidad y del } i\text{-ésimo tratamiento}$$

En este modelo, hay t tratamientos y b bloques.

Los bloques se utilizan para agrupar experimentalmente las unidades de manera que se elimine o reduzca el efecto del factor de ruido en la comparación de las medidas de los tratamientos. Esto permite una estimación más precisa de los efectos de los tratamientos al controlar las variables de confusión.

Tabla ANOVA

Cuadro 11: Ejemplo de Tabla

Factor de variabilidad	g.l	SS	MS	Estadística de la prueba
Tratamientos	t-1	$b \sum_{i=1}^t (\bar{y}_i - \bar{Y})^2$	$SS_T/t - 1$	
Bloques	b-1	$t \sum_{j=1}^b (\bar{y}_j - \bar{Y})^2$	$SS_B/b - 1$	Fila 2, Col 5
Errores	(t-1)(b-1)	Fila 3, Col 3	Fila 3, Col 4	Fila 3, Col 5
Total	bt-1	Fila 4, Col 3	Fila 4, Col 4	Fila 4, Col 5
Fila 5, Col 1		Fila 5, Col 3	Fila 5, Col 4	Fila 5, Col 5

Ejemplo 4

Supon que se quiere verificar si cuatro tipos de punta producen o no lecturas diferentes en una maquina para probar la dureza

hay cuatro tipos de puntas y cuatro ejemplares de material.cada punta se prueba una vez en cada ejemplar. los resusltados se muestran a continuación

Cuadro 12: Material

Tipo de punta	1	2	3	4	media tratamientos
1	9.3	9.4	9.6	10	9.575
2	9.4	9.3	9.8	9.9	9.6
3	9.2	9.4	9.5	9.7	9.45
4	9.7	9.6	10	10.2	9.875
Media bloques	9.4	9.425	9.725	9.25	

16. 5/04/2024

16.1. Tabla ANOVA

Factor de variabilidad	gl	SS	MS	F0
A	$a - 1$	SS_A	$\frac{SS_A}{a-1}$	$F_{01} = \frac{MS_A}{MS_E}$
B	$b - 1$	SS_B	$\frac{SS_B}{b-1}$	$F_{02} = \frac{MS_B}{MS_E}$
AB	$(a - 1)(b - 1)$	SS_{AB}	$\frac{SS_{AB}}{(a-1)(b-1)}$	$F_{03} = \frac{MS_{AB}}{MS_E}$
Errores	$ab(n - 1)$	SS_E	$\frac{SS_E}{ab(n-1)}$	
Total	$abn - 1$	SS_{Total}		

Las hipotesis que se contrastan son:

$$\begin{cases} H_{01} : \tau_i = 0, \forall_i \\ H_{11} : \text{Al menos un } \tau_i \neq 0 \end{cases}$$

$$\begin{cases} H_{02} : \beta_j = 0, \forall_j \\ H_{12} : \text{Al menos un } \beta_j \neq 0 \end{cases}$$

$$\begin{cases} H_{01} : r_{ij} = 0, \forall_{ij} \\ H_{11} : \text{Al menos un } r_{ij} \neq 0 \end{cases}$$

Regiones de rechazo:

- Se rechaza H_{01} si $F_{01} > q_{F(a-1, ab(n-1))}^{(1-\alpha)}$
- Se rechaza H_{02} si $F_{02} > q_{F(b-1, ab(n-1))}^{(1-\alpha)}$
- Se rechaza H_{03} si $F_{03} > q_{F((a-1)(b-1)-1, ab(n-1))}^{(1-\alpha)}$

Ejemplo 5

Un ingeniero está diseñando una batería que usará en un dispositivo que se somete a variaciones de temperatura extremas. Se prueban cuatro baterías para cada combinación de materiales y temperatura. En la siguiente tabla se presenta la vida, en horas, de cada batería.

Tipo de material	15	70	90
1	130 155	34 40	20 70
	74 180	80 75	82 58
2	150 188	136 122	25 70
	159 126	106 115	58 45
3	138 110	174 120	96 104
	168 160	50 139	82 60

$$a = 3, b = 3, n = 4$$

verifica si el tipo de material afecta la vida en horas de las baterías

$$\begin{cases} H_0 : \tau_i = 0, \forall_i = 1, 2, 3 \\ H_1 : \text{Al menos un } \tau_i \neq 0 \end{cases}$$

Tipo de material	15	70	90	
1	$\bar{y}_{11} = 134,75$	$\bar{y}_{12} = 57,25$	$\bar{y}_{13} = 57,5$	$\bar{y}_{1..} = 83,16$
2	$\bar{y}_{21} = 155,75$	$\bar{y}_{22} = 119,75$	$\bar{y}_{23} = 49,5$	$\bar{y}_{2..} = 108,33$
3	$\bar{y}_{31} = 144$	$\bar{y}_{32} = 145,75$	$\bar{y}_{33} = 85,5$	$\bar{y}_{3..} = 125,08$
	$\bar{y}_{.,1} = 144,83$	$\bar{y}_{.,2} = 107,58$	$\bar{y}_{.,3} = 64,16$	$\bar{y} = 105,52$

F.V	gl	SS	MS	F0
Material	2	10,683.72		$F_{01} = \frac{MS_A}{MS_E}$
Temperatura	2	39,118.72		$F_{02} = \frac{MS_B}{MS_E}$
Mat * Temp	4	9,613.77		$F_{03} = \frac{MS_{AB}}{MS_E}$
Errores	27	18.230.75		
Total	35	77,646.97		

17. 8/04/2024

Verificación del cumplimiento de los supuestos del análisis de varianza

Normalidad

- Pruebas de bondad de ajuste (Anderson-Darling, saphiro-wilk)

- Los errores se distribuyen como una $N(0, \hat{\sigma}^2)$, $\hat{\sigma}^2 = MS_E$
 - Histograma
 - Q-plot

18. 10/04/2024

18.1. Unidad 3 - Análisis de regresión

La regresión busca establecer la relación entre variables, clasificándose en dos modelos principales:

Modelos $\begin{cases} \text{Deterministas: } y = f(x), \text{ donde la relación es directa y sin aleatoriedad.} \\ \text{Probabilísticos: } y = f(x) + E, \text{ incorporando un término de error } E \text{ que añade variabilidad.} \end{cases}$

18.2. 3.1 - Modelo de regresión lineal simple

En el modelo de regresión lineal simple, la relación entre dos variables se expresa como:

$$y = \beta_0 + \beta_1 x + E$$

donde:

- y es la variable respuesta o dependiente, que intentamos predecir.
- x es la variable predictora o independiente, utilizada para predecir y .
- β_0 es la intersección con el eje Y, representando el valor esperado de y cuando x es 0.
- β_1 es la pendiente de la línea de regresión, indicando el cambio esperado en y por cada unidad de cambio en x .
- E representa el error o residuo, que es la diferencia entre el valor observado y el valor ajustado por el modelo.

El término de error E asume que los residuos siguen una distribución normal con media cero y varianza constante σ^2 , es decir, $E \sim N(0, \sigma^2)$. Esto implica que la variabilidad de y se mantiene constante a lo largo de los diferentes valores de x .

El objetivo es estimar β_0, β_1 y σ^2 a partir del conjunto de observaciones registradas

Supuestos del modelo de regresión lineal simple

El modelo de regresión lineal simple se basa en varios supuestos clave que deben cumplirse para que las estimaciones de los parámetros sean no sesgadas, mínimamente variables y confiables:

- **Linealidad:** La relación entre la variable dependiente y la variable independiente es lineal. Esto implica que el cambio esperado en la variable dependiente Y es una función lineal del cambio en la variable independiente X .

$$E(y|x) = \beta_0 + \beta_1 x$$

Aquí, $E(y|x)$ representa el valor esperado de Y dado X , con β_0 como la intersección y β_1 como la pendiente de la línea de regresión.

- **Homocedasticidad:** La varianza de los errores residuales (diferencia entre los valores observados y los valores predichos) es constante para todas las observaciones. Esto significa que la dispersión o variabilidad de Y alrededor de la línea de regresión es la misma para todos los valores de X .

$$\text{var}(Y|X) = \sigma^2$$

Donde σ^2 representa la varianza constante de los errores residuales.

- **Independencia:** Las observaciones son independientes entre sí. Esto indica que no existe correlación entre los términos de error de diferentes observaciones; el término de error de una observación no influye en el término de error de otra.

$$\text{cor}(y_i, y_j | x_i, x_j) = 0, \quad \forall i \neq j$$

Donde $\text{cor}(y_i, y_j | x_i, x_j)$ representa la correlación entre las respuestas y_i y y_j para diferentes observaciones i y j . Este supuesto es fundamental para la validez de las pruebas estadísticas asociadas al modelo.

18.3. Estimación de los parámetros por mínimos cuadrados

La estimación de los parámetros en un modelo de regresión lineal simple se realiza mediante el método de mínimos cuadrados. Este método busca encontrar los valores de β_0 y β_1 que minimizan la suma de los cuadrados de las diferencias entre los valores observados de la variable dependiente Y y los valores que el modelo predice. Matemáticamente, se busca minimizar la función de coste, que es la suma de los cuadrados de los residuos (errores).

Dado un conjunto de n observaciones (x_i, y_i) , la función de coste SS_E se define como:

$$SS_E = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

- y_i son los valores observados de la variable dependiente.
- x_i son los valores de la variable independiente.
- β_0 es el término de intersección con el eje Y .
- β_1 es el coeficiente de la variable independiente X , que indica cómo el cambio en X afecta a Y .

Para encontrar los valores de β_0 y β_1 que minimizan S , se toman las derivadas parciales de S con respecto a β_0 y β_1 , se igualan a cero y se resuelven las ecuaciones resultantes. Este proceso lleva a un sistema de ecuaciones lineales cuya solución da los estimadores de mínimos cuadrados de los parámetros:

$$\hat{\beta}_1 = \frac{n \sum (x_i y_i) - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\hat{\beta}_0 = \frac{\sum y_i - \hat{\beta}_1 \sum x_i}{n}$$

Estos estimadores, $\hat{\beta}_0$ y $\hat{\beta}_1$, proporcionan los coeficientes de la línea de regresión que mejor se ajusta a los datos en el sentido de mínimos cuadrados, minimizando así la suma total de los cuadrados de los residuos entre los datos observados y los predichos por el modelo.

Utilizar la regla de kramer para resolver el siguiente sistema de ecuaciones

$$\left(\begin{array}{cc|c} n & \sum x_i & \sum y_i \\ \sum x_i & \sum x_i^2 & \sum x_i y_i \end{array} \right)$$

El siguiente paso es demostrar que minimizan la SS_E . esto se hara mostrando que la matriz hessiana es definida

$$H = \begin{pmatrix} \frac{\partial^2 SS_E}{\partial \beta_0^2} & \frac{\partial^2 SS_E}{\partial \beta_0 \beta_1} \\ \frac{\partial^2 SS_E}{\partial \beta_1 \beta_0} & \frac{\partial^2 SS_E}{\partial \beta_1^2} \end{pmatrix}$$

Si $h_{11} > 0$ y $\det(H) > 0$, entonces H es definida positiva

CHATGP

18.4. Estimación de los parámetros por mínimos cuadrados

La estimación de los parámetros en un modelo de regresión lineal simple se realiza mediante el método de mínimos cuadrados. Este método busca encontrar los valores de β_0 y β_1 que minimizan la suma de los cuadrados de las diferencias entre los valores observados de la variable dependiente Y y los valores que el modelo predice. Matemáticamente, se busca minimizar la función de coste, que es la suma de los cuadrados de los residuos (errores).

Dado un conjunto de n observaciones (x_i, y_i) , la función de coste SS_E se define como:

$$SS_E = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

Para encontrar los valores de β_0 y β_1 que minimizan SS_E , se toman las derivadas parciales de SS_E con respecto a β_0 y β_1 , se igualan a cero y se resuelven las ecuaciones resultantes. Este proceso lleva a un sistema de ecuaciones lineales cuya solución da los estimadores de mínimos cuadrados de los parámetros:

$$\hat{\beta}_1 = \frac{n \sum (x_i y_i) - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\hat{\beta}_0 = \frac{\sum y_i - \hat{\beta}_1 \sum x_i}{n}$$

Estos estimadores, $\hat{\beta}_0$ y $\hat{\beta}_1$, proporcionan los coeficientes de la línea de regresión que mejor se ajusta a los datos en el sentido de mínimos cuadrados, minimizando así la suma total de los cuadrados de los residuos entre los datos observados y los predichos por el modelo.

A continuación, se utiliza la regla de Kramer para resolver el sistema de ecuaciones lineales. El sistema se puede representar como:

$$\left(\begin{array}{cc|c} n & \sum x_i & \sum y_i \\ \sum x_i & \sum x_i^2 & \sum x_i y_i \end{array} \right)$$

El siguiente paso es demostrar que los estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$ minimizan la función de coste SS_E . Esto se hace mostrando que la matriz Hessiana H asociada es definida positiva:

$$H = \begin{pmatrix} \frac{\partial^2 SS_E}{\partial \beta_0^2} & \frac{\partial^2 SS_E}{\partial \beta_0 \beta_1} \\ \frac{\partial^2 SS_E}{\partial \beta_1 \beta_0} & \frac{\partial^2 SS_E}{\partial \beta_1^2} \end{pmatrix}$$

Si los elementos diagonales de H son mayores que cero y el determinante de H es positivo, entonces la matriz Hessiana H es definida positiva, lo que confirma que la función de coste SS_E es mínima en los valores de $\hat{\beta}_0$ y $\hat{\beta}_1$.

19. 17 de abril del 2024

Ejemplo

Obs	Profundidad	Velocidad
1	0.34	0.636
2	0.29	0.319
3	0.28	0.734
4	0.42	1.327
5	0.29	0.487
6	0.41	0.924
7	0.76	7.350
8	0.73	5.890
9	0.46	1.979
10	0.40	1.124

Cuadro 13: Ejemplo de una tabla sencilla

a) Obtener

$$\hat{B}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

$$\hat{B}_0 = \bar{y} - \hat{B}_1 \bar{x}$$

$$\hat{\sigma}^2 = \frac{SS_{Residuos}}{n - 2} = \frac{SS_{Total} - \hat{B}_1 S_{xy}}{n - 2}$$

$$SS_{Total} = \sum_n^{i=1} (y_i - \bar{y})^2$$

```

# Librerías
import pandas as pd
import numpy as np

# Cargar archivo
df = pd.read_csv("/Ruta/...")

# Vista parcial
df.head()

# Crear columnas
df["Profundidad^2"] = df["Profundidad"]**2
df["Profundidad*Velocidad"] = df["Profundidad"] * df["Velocidad"]

# Estimadores
n = len(df)
n = df.shape[0]
x_bar = df["Profundidad"].mean()
y_bar = df["Velocidad"].mean()
S_xy = df["Profundidad*Velocidad"].sum() - n * x_bar * y_bar
S_xx = df["Profundidad^2"].sum() - x_bar**2
estimador_beta_1 = S_xy / S_xx
estimador_beta_0 = y_bar - estimador_beta_1 * x_bar
print(f"El estimador de beta_1 es {estimador_beta_1} y el estimador de beta_0 es {estimador_beta_0}")
SS_Total = (n-1) * df["velocidad"].var()
SS_Residuos = SS_Total - estimador_beta_1 * S_xy
estimador_sigma^2 = SS_Residuos/(n-2)

```

Se llevó a cabo un diseño completamente aleatorizado para comparar los efectos de cinco estímulos en los tiempos de reacción y se obtuvieron los siguientes resultados:

	A	B	C	D	E
	0,8	0,7	1,2	1,0	0,6
	0,6	0,8	1,0	0,9	0,4
	0,6	0,5	0,9	0,9	0,4
	0,5	0,5	1,2	1,1	0,7
	0,6	1,3	0,7	0,3	
	0,9	0,8			
	0,7				

- a) Realiza un análisis de varianza (ANOVA) y evalúa si existe una diferencia significativa en los tiempos medios de reacción entre los cinco estímulos.

Resultados del Análisis ANOVA

Medias y varianzas de los grupos:

- Grupo 1: Media = 0.67, Varianza = 0.019
- Grupo 2: Media = 0.77, Varianza = 0.087
- Grupo 3: Media = 1.00, Varianza = 0.045
- Grupo 4: Media = 0.84, Varianza = 0.098

- Grupo 5: Media = 0.52, Varianza = 0.022

Media global: 0.74

Sumas de cuadrados:

$$SST = 1,783 \quad (\text{Total})$$

$$SSB = 0,596 \quad (\text{Entre grupos})$$

$$SSW = 1,187 \quad (\text{Dentro de los grupos})$$

Grados de libertad:

- Entre grupos: 4
- Dentro de los grupos: 22

Resultado ANOVA:

- Valor F = 2.761
- p-valor = 0.0533

Conclusion

Como el pvalor resulto en 0,0533 esto indica que no hay suficiente evidencia para rechazar la hipótesis nula de igualdad de medias entre los grupos

- b) Compara los estímulos B y C para determinar si hay una disparidad en los tiempos medios de reacción.

Para realizar la compracion de estos estímulos usaremos: $f_{oneway}(B, C)$ que nos resultan en

- Valor F = 2.179
- p-valor = 0.174

Esto indica que no hay suficiente evidencia estadística para rechazar la hipótesis nula, sugiriendo que las diferencias observadas en los tiempos de reacción entre los estímulos B y C podrían ser debidas al azar.

Dado que es de esperar que el tiempo medio de reacción pueda variar entre las personas, se podría haber llevado a cabo el experimento del ejercicio anterior de manera más eficiente utilizando un diseño de bloques aleatorizados, con las personas como bloques. Por lo tanto, se llevó a cabo un nuevo experimento con la participación de cuatro personas, y cada una de ellas fue expuesta a los cinco estímulos en un orden aleatorio. Los tiempos de reacción (en segundos) se presentan en la siguiente tabla.

Cuadro 14: Tiempo de reacción para cada sujeto

Sujeto	A	B	C	D	E
1	0.7	0.8	1.0	1.0	0.5
2	0.6	0.6	1.1	1.0	0.6
3	0.9	1.0	1.2	1.1	0.6
4	0.6	0.8	0.9	1.0	0.4

Realiza un análisis de varianza (ANOVA) para evaluar las diferencias en los tiempos medios de reacción para los cuatro estímulos.

20. 03/05/2024

Ejemplo 2

Para estudiar la relación de publicidad e inversión de capital con utilidades, los datos registrados en millones de pesos se recolectan para 5 empresas en el mismo año. La variable Y representa la utilidad, X_1 representa inversión de capital y X_2 representa gastos en publicidad

Y	X_1	X_2
15	25	4
16	1	5
2	6	3
3	30	1
12	29	2

```
# Librerías
import pandas as pd
import numpy as np
import scipy.stats as stats

#Crear diccionario con los datos de las empresas
empresa = {
    "utilidad": [15,16,2,3,12],
    "capital": [25,1,6,30,29],
    "publicidad": [4,5,3,1,2]
}

#Crear dataframe
df = pd.DataFrame(empresa)

#Crear columna de unos
df["intercept"] = 1

#Definir las matrices X y Y
X = df[["intercept", "capital", "publicidad"]]
```

21. 13/05/2024 Clase

Ejemplo 4

Encuentra la funcion de distribucion a posteriori para el modelo de Poisson
Funcion de distribucion a posteriori

$$f(\lambda) \propto ((\lambda)^\alpha) \exp[W_i(\lambda)\beta]$$

$$\begin{aligned} f(\lambda) &\propto ((\lambda)^\alpha) \exp[W_i(\lambda)\beta] \\ &= (e^{-\lambda})^\alpha \exp[Ln(\lambda)\beta] \\ &= e^{-\lambda\alpha} \lambda^\beta \end{aligned}$$

Quien es la familia conjugada del modelo de Poisson?

R: La distribucion gamma
si $\lambda \sim \text{Gamma}(\alpha, \beta)$

$$f(\lambda|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}$$

Asi que,

$$f(\lambda) \sim \text{Gamma}(\beta + 1, \alpha)$$

Funcion de distribucion a posteriori por otro lado

$$f(\lambda|x) \propto c(\lambda)^{\alpha+n} \exp[W_i(\beta + \sum t_i(x_j))]$$

Finalmente, obtenemos que

$$\begin{aligned} f(\lambda|x) &\sim \text{Gamma}(\beta + \sum x_i + 1, \alpha + n) \\ &\text{Gamma}(\alpha + \sum x_i, \beta + n) \end{aligned}$$

Ejemplo 5

Encuentra la funcion de distribucion a priori para el modelo exponencial

$$f(x|\lambda) = \lambda e^{-\lambda x}$$

- Mostrar que es familia exponencial
- Identificar la familia conjugada del modelo exponencial
- Encontrar la funcion de distribucion a priori
- Encontrar la funcion de distribucion a posteriori

Solucion

1. - si $f(x|\theta)$ es una familia exponencial, entonces

$$f(x|\theta) = h(x)c(\theta)\exp[W_i(\theta)T(x)]$$

Luego, si $X \sim \text{Exp}(\lambda)$, entonces

$$\begin{aligned} f(x|\lambda) &= \lambda e^{-\lambda x} I_{(0,\infty)}(x) \\ &= \lambda I_{(0,\infty)}(x) e^{-\lambda x} \end{aligned}$$

Con $h(x) = I_{(0,\infty)}(x)$, $c(\lambda) = \lambda$, $W_i(\lambda) = -\lambda$ y $T(x) = x$. Por lo tanto, la familia exponencial es $f(x|\lambda) = \lambda e^{-\lambda x}$. **Verificado**

2.- La familia conjugada para el modelo exponencial es la familia Gamma. **Verificado**

3.- La funcion de distribucion a priori para el modelo exponencial es:

$$\begin{aligned} f(\lambda) &\propto c(\lambda)^\alpha \exp[W_i(\lambda)\beta] \\ &= \lambda^\alpha \exp[-\lambda\beta] \\ &= \lambda^{\alpha+1-1} \exp[-\lambda\beta] \\ f(\lambda) &\sim \text{Gamma}(\alpha+1, \beta) \end{aligned}$$

4.- La funcion de distribucion a posteriori para el modelo exponencial es:

$$\begin{aligned} f(\lambda|x) &\propto c(\lambda)^{\alpha+n} \exp[W_i(\lambda)(\beta + \sum x_i)] \\ &= \lambda^{\alpha+n-1} \exp[-\lambda(\beta + \sum x_i)] \\ f(\lambda|x) &\sim \text{Gamma}(\alpha+n, \beta + \sum x_i) \end{aligned}$$

Ejemplo 6

Encuentra la funcion de distribucion a priori para el modelo normal donde μ es desconocida y σ^2 es conocida

- Mostrar que es familia exponencial
- Identificar la familia conjugada del modelo normal
- Encontrar la funcion de distribucion a priori
- Encontrar la funcion de distribucion a posteriori

1.- Si $f(x|\theta)$ es una familia exponencial, entonces

$$f(x|\theta) = h(x)c(\theta)\exp[W_i(\theta)T(x)]$$

luego, si $X \sim N(\mu, \sigma^2)$, entonces

$$f(x|\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

Desarrollando el binomio tenemos:

$$f(x|\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{x^2 - 2x\mu + \mu^2}{2\sigma^2}\right]$$

$$\begin{aligned}
&= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{x^2}{2\sigma^2} + \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2} \right] \\
c(\mu) &= \frac{-\mu^2}{\sqrt{2\pi\sigma^2}} \\
h(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{x^2}{2\sigma^2} \right] \\
W_i(\mu) &= \frac{\mu}{\sigma^2} \\
T(x) &= x
\end{aligned}$$

La familia conjugada para el modelo normal es la familia normal. **Verificado**

22. 13/05/2024 ChatGPT

Ejemplo 4: Distribución a Posteriori para el Modelo de Poisson

Supongamos que tenemos una variable aleatoria Y que sigue una distribución de Poisson con parámetro λ . Queremos determinar la distribución a posteriori de λ dado que observamos $Y = y$. Asumimos que el prior para λ es una distribución Gamma, comúnmente usada como prior conjugado en modelos de Poisson debido a sus propiedades analíticas.

Función de Distribución a Posteriori:

Si consideramos que λ sigue a priori una distribución Gamma con parámetros α y β , la función de densidad de probabilidad del prior es:

$$\pi(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}$$

El likelihood basado en el modelo de Poisson para una observación y es:

$$L(\lambda) = \frac{e^{-\lambda} \lambda^y}{y!}$$

Aplicando el teorema de Bayes para obtener el posterior, tenemos:

$$f(\lambda | y) \propto \pi(\lambda) L(\lambda) = \left(\frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \right) \left(\frac{e^{-\lambda} \lambda^y}{y!} \right)$$

Simplificando, llegamos a:

$$f(\lambda | y) \propto \lambda^{y+\alpha-1} e^{-(\beta+1)\lambda}$$

Esto indica que el posterior también es una distribución Gamma:

$$f(\lambda | y) \sim \text{Gamma}(y + \alpha, \beta + 1)$$

Familia Conjugada del Modelo de Poisson:

La familia conjugada para el parámetro de un modelo de Poisson es la familia de distribuciones Gamma. La conjugación se demuestra por el hecho de que el posterior, $\text{Gamma}(y + \alpha, \beta + 1)$, pertenece a la misma familia que el prior, $\text{Gamma}(\alpha, \beta)$.

Distribución a Posteriori Extendida:

Si consideramos un conjunto de datos $\{x_1, x_2, \dots, x_n\}$ donde cada x_i sigue independientemente una distribución de Poisson con el mismo λ , el likelihood conjunto es:

$$L(\lambda) \propto \lambda^{\sum x_i} e^{-n\lambda}$$

El posterior resultante, integrando este likelihood con el prior Gamma, se convierte en:

$$f(\lambda | x) \propto \lambda^{\alpha + \sum x_i - 1} e^{-(\beta + n)\lambda}$$

Por lo tanto:

$$f(\lambda | x) \sim \text{Gamma}(\alpha + \sum x_i, \beta + n)$$

Estas ecuaciones proporcionan un marco completo para entender cómo actualizar nuestras creencias sobre λ en un modelo de Poisson a medida que observamos nuevos datos.

23. 17/05/2024

23.1. Ejemplo 6

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

$$= \frac{\exp\left[-\frac{x^2}{2\sigma^2}\right]}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{\mu^2}{2\sigma^2}\right] \exp\left[\frac{\mu x}{\sigma^2}\right]$$

$$h(x) = \frac{\exp\left[-\frac{x^2}{2\sigma^2}\right]}{\sqrt{2\pi\sigma^2}}$$

$$c(\mu) = \exp\left[-\frac{\mu^2}{2\sigma^2}\right]$$

$$W_i(\mu) = \frac{\mu}{\sigma^2}$$

$$t_i(x) = x$$

Funcion de distribucion a priori

$$\begin{aligned} P(\mu) &\propto c(\mu)^\alpha \exp[W_i(\mu)\beta] \\ &= \exp\left[-\frac{\mu^2}{2\sigma^2}\right]^\alpha \exp\left[\frac{\mu\beta}{\sigma^2}\right] \\ &= \exp\left[-\frac{\mu^2}{2\sigma^2}\right]^\alpha \exp\left[\frac{\mu\beta}{\sigma^2}\right] \\ &= \exp\left[-\frac{\alpha}{2\sigma^2} \left(\mu^2 - \frac{2\mu\beta}{\alpha}\right)\right] \propto \exp\left[-\frac{\alpha}{2\sigma^2} \left(\mu^2 - \frac{2\mu\beta}{\alpha} + \left(\frac{\beta}{\alpha}\right)^2\right)\right] \\ &= \exp\left[\frac{-\alpha}{2\sigma^2} \left(\mu - \frac{\beta}{\alpha}\right)^2\right] \end{aligned}$$

Funcion de distribucion a posteriori

$$p(\mu|x) \sim N\left(\frac{\beta + \sum x_i}{\alpha + n}, \frac{\sigma^2}{\alpha + n}\right)$$

ESTO SE USA PARA LA TAREA

Si $p(\mu) \sim N(\mu_0, \sigma_0^2)$ y $p(\sigma^2)$

$$\mu_0 = \frac{\beta}{\alpha} \quad \sigma_0^2 = \frac{\sigma_1^2}{\alpha}$$

$$\beta = \mu_0 \alpha \quad \alpha = \frac{\sigma_1^2}{\sigma_0^2}$$

23.2. Ejemplo 7

La estatura de un grupo se distribuye como una $N(\mu, \sigma_1^2)$, donde μ_1 es desconocida y $\sigma_1^2 =$
Con informacion a priori, se obtiene que $p(\mu) \sim N(\mu_0, \sigma_0^2)$, $\mu = 1,70$