

Conceptos Basicos de Estadistica

makita

10/14/2021

Conceptos básicos de estadística

Variables: género, edad, peso, talla, estrato, localidad, fecha, lugar de nacimiento..

Población (universo o colectivo)

Parámetros: IMC

Es el conjunto total de ELEMENTOS de la misma naturaleza cualquiera que sea, que son de interés para un problema dado

Clasificación de variables

- N = Representación de el tamaño de la población

Cualitativas (categóricas)

muestra

Cuantitativas

Variable aleatoria:

Los valores de las observaciones son numéricas y en consecuencia, ordenables.

Son fenómenos o características de los elementos de la población.

Discreta

Función de valor real que tiene como dominio el espacio muestral de un experimento aleatorio.

Recorridos finitos numerables sin tomar valores intermedios e.g. conteos.

Variables sobre las cuales tenemos un grado de incertidumbre respecto a los valores que puede tomar

Continua

Datos

Recorridos infinitos no numerables e.g. la distribución normal

Son los resultados observados de las variables aleatorias (Cuando se hace una medición)

Escalas de medición

Parámetro

Cualitativas

Es la medición global de cualesquier característica de los elementos de la población.

Nominales: Clasificación de objetos o fenómenos mediante símbolos o signos (No hay orden o dirección). e.g.

Es un valor teórico asociado a la población.

Ejemplos

- Nombre
- Número de la cédula
- Tipo de sangre
- Color de los ojos
- Número de camiseta de los jugadores

Población: Los niños y niñas de 0 a 5 años de edad localizados en Bogotá

Los números en la lista anterior no pueden ser sometidos a operaciones matemáticas

$$k = 1 + 3.322 \log(n) \quad (1)$$

Ordinales

k = intervalos de clase

Categorías ordenadas (Rangos, órdenes, escalamientos)

Para la longitud de los intervalos:

$$L = \frac{\text{Dato mayor} - \text{Dato menor}}{n} \quad (2)$$

- Sabor de un yogurt

- A menudo es prueba y error

Cuantitativas

Tipos de frecuencias

Intervalo

Los datos medidos en una escala ordinal para los cuales pueden clasificarse las distancias entre valores pero no existe un cero absoluto o no exista ausencia total de la característica

- Absoluta: Conteo de observaciones que cae en cada intervalo.
- Relativa: $\frac{\text{Absoluta}}{n}$.
- Acumulada: Suma de las frecuencias absolutas
- Relativa acumulada: Suma de las frecuencias relativas.

- Temperatura: a 0°C no deja de existir la temperatura
- Notas: se corre la escala e inicia desde 3.

Características a revisar de las distribuciones

Razón

Tiene todas las características de un intervalo, y además tiene un cero absoluto

- Distribucion
- Localizacion (sesgo)
- Dispersion

Resumen y descripción de datos de una variable

Medidas de localización

Datos en bruto en forma de listas (o bases no son fáciles de usar para tomar decisiones)

Media aritmética:

- Se necesita algún tipo de organización

Si $x_1, x_2, x_3, \dots, x_n$ es una muestra de una población de tamaño N entonces la media es N

Para esto podemos utilizar gráficos de barras, graficos de torta, o tablas de frecuencias.

Media poblacional

$$\mu = \frac{\sum_{i=1}^n x_i}{N} \quad (3)$$

Como agrupar los datos: Sturges

Estimador muestral

Si n no es demasiado grande, intervalos = \sqrt{n}

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad (4)$$

En caso contrario:

Características:

- Es facil de obtener
- Medida no robusta: Afectada por valores extremos o datos atípicos.

Propiedades de la media aritmetica:

Si $x_1, x_2, x_3, \dots, x_n$ es una muestra de una poblacion de tamaño N entonces la media es \bar{x} , entonces

$$\sum_{i=1}^n x_i = x_1 + x_2 + x_3 + \dots + x_n \quad (5)$$

Si $x_i = c$ y a su vez c es constante, entonces

$$\sum_{i=1}^n x_i = \sum_{i=1}^n c = c + c + c + \dots \quad (6)$$

Entonces

$$\sum_{i=1}^n x_i = nc \quad (7)$$

- *Ejemplo:*

Si c es una constante que multiplica las observaciones:

$$\sum_{i=1}^n cx_i = c \sum_{i=1}^n x_i$$

Si $x_1, x_2, x_3, \dots, x_n$ y $y_1, y_2, y_3, \dots, y_n$ son sucesiones de numeros;

$$\sum_{i=1}^n (x_i + y_i) = \sum_{i=1}^n x_i + \sum_{i=1}^n y_i$$

Si $x_1, x_2, x_3, \dots, x_n$ y $y_1, y_2, y_3, \dots, y_n$ son sucesiones de numeros;

$$\sum_{i=1}^n (x_i - y_i) = \sum_{i=1}^n x_i - \sum_{i=1}^n y_i$$

5.

$$\sum_{i=1}^n (x_i - \bar{x}) = 0$$

$$\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x}$$

6.

promedio de y en funcion de promedio de x en regresion lineal simple

Si A

La mediana

Es el valor central (es el dato de la variable que esta en el centro de la misma). Deja por encima y por debajo mitad y mitad de las observaciones.

Calculo de la mediana

Depende si el conjunto es par o impar:

Si $x_1, x_2, x_3, \dots, x_n$ Son los valores ordenados en una muestra de una poblacion de tamaño N:

$$\hat{x} = \frac{x_{n/2} + x_{n+1/2}}{2} \text{ si n es par}$$

$$\hat{x} = x_{n+1/2} \text{ si n es impar}$$

Es un estimador robusto, no se ve afectado por valores extremos

Ejemplo

Edad de ninos

```
x1 <- c(6,7,8,9,10)
```

n es impar, entonces $\hat{x} = x_{n+1/2} = x_{6/2} = x_3 = 8$

De la muestra analizada la mitad de los ninos tienen entre 6 y 8 años, y la otra mitad entre 8 y 10 años.

Moda

- El valor que más se repite
- Usada para valores numéricos o categóricos

e.g. Cual es el color más frecuente en los ojos.

Medidas de dispersión o variación

Varianza

- Varianza poblacional:

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N} \quad (8)$$

- Varianza muestral:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \quad (9)$$

$$s^2 = \sum_{i=1}^n \frac{(X_i^2 - 2x_i\bar{x} + \bar{x}^2)}{n-1} \quad (10)$$

$$s^2 = \frac{\sum_{i=1}^n X_i^2 - 2\bar{x} \sum_{i=1}^n x_i + \bar{x}^2}{n-1} \quad (11)$$

Coefficientes de asimetría de Fisher

Permite interpretar la forma de la distribución, respecto a ser o no asimétrica

Coefficiente de curtosis

Mide el grado de aplastamiento o apuntamiento de la gráfica de la distribución.

Coefficiente de variación

$$CV = \frac{s}{\bar{x}} \cdot 100\%$$

Si CV es igual o menor a 5% hay homogeneidad

si esta entre 5% y 20% los datos son medianamente homogéneos

Si CV mayor a 20% hay heterogeneidad ## Rango

medida no robusta, si hay datos atípicos se ve muy afectado

Otras medidas de centralización

Desviación media absoluta

$$DM = \frac{\sum_{i=1}^n |x_i - \bar{x}|}{n} \quad (12)$$

Rango intercuartílico

boxplot

Cuartiles

Se divide en cuatro partes porcentuales el conjunto de observaciones

Se calcula de la siguiente manera

$$Q_k = k \cdot \frac{n}{4} \quad k = 1, 2, 3$$

Deciles

Se divide en diez partes porcentualmente iguales

$$D_k = \frac{n}{10} \quad 1, 2, 3, \dots, 9$$

Percentiles

es más detallado, nos da más acceso a distintos puntos de la distribución

$$P_k = \frac{n}{100} \quad 1, 2, 3, \dots, 99$$