

Estadística Descriptiva

Permite

- Conocer reglas y métodos usados en el tratamiento de datos
- Evaluar y cuantificar la importancia de los resultados estadísticos
- Entender fenómenos de la realidad (sociales, biológicos, ...)
- Dar visión clara de la información proveniente de distintas fuentes

Estadística

↳ Manejo de Información Cuantificable

- Descripción de datos
- Inferencia = concluir

↳ Estudia fenómenos aleatorios y su comportamiento (variabilidad)

↳ Funciones

- Resumir = Gráficamente o Numericamente
- Simplificar
- Comparar
- Relacionar
- Proyectar

↳ Tareas

1. Delimitar estudio
 - Establecer unidades
 - Definir variables
2. Observar
 - Censo
 - Muestreo
 - Diseño experimental
3. Recolección y registro de datos
4. Depurar información
5. Resumir = Graf. Numéric.
6. Interpretar

Temas

- Diseño experimental
- Estadística Descriptiva (tareas)
- Inferencia
- Estadística no paramétrica (modelos sin supuestos previos)
- Elementos de regresión (variables respuesta y explicativas)

TIPOS DE VARIABLES

↳ Características que varían

- Estatura
- Presión sanguínea

1. Continuas = Mediciones, Estatura de un grupo
2. Discretas = Conteos, # autos en un semáforo, horas
3. Categoricals = clasificación, Sexo, estrato

NIVELES DE MEDICION

- ① Nominal = representan categorias, Estado civil = O, S, M no cumple 1<2<3
No medidas basicas de resumen
- ② Ordinal = representan jerarquias, Quemaduras = 1, 2, 3 se cumple 1<2<3
No medidas basicas de resumen
- ③ Intervalo = cuantitativas
- ④ Razon = Cuantitativa con escalas, Estatura, Presion sanguinea

DESCRIPCION DE DATOS

Distribuciones de frecuencia

• Tabla de frecuencia

- Encuentre Minimo y Maximo de los valores registrados
- Defina Intervalos de clase

- Sturges \times

$$K = 1 + 3.33 \log_{10}(n)$$

K = Numero de Clases

- Velleman n Pequeño

$$K = 2 \sqrt{n}$$

n = Numero de datos

- Dixon y Kronmal n Grande

$$K = 10 \cdot \log_{10}(n)$$

Notas:

- * Entre 5 y 25 clases es lo adecuado
- * Pocas clases \Rightarrow Gran perdida de informacion
- * Pocas datos \Rightarrow No se evidencian comportamientos
- * Muchas clases \Rightarrow No se evidencian comportamientos
- * Dentro de las clases es importante organizar (menor a mayor)

- Cuente las observaciones por cada subintervalo (Frecuencia de clase)

- Calcule la frecuencia relativa

$$FR = \frac{\text{Frecuencia de clase}}{\# \text{ total de observaciones}}$$

Ejemplo: tabla dada:

Masa Kg	Estatura cm	Edad días	Genero H, M	Estivato 1 a 5	Costo \$	Longitud Alcaren cm	Tipo colega Pu opri	Horas a estudiar horas
Tiempo en llegar a la U	Longitud eres prom	Longitud Manos prom						
Mm	cm	cm						

• Tabla de frecuencia para genero

	Frecuencia	Porcentaje
Hombre	46	68,66
Mujer	21	31,34

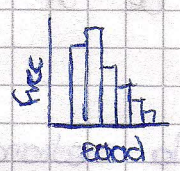
67 DATOS

• En la Edad el grafico de barras no es dicente (se hizo intervalos de clase) (8 clases)

Categoría	Frecuencia	Frecuencia Relativa
≥ 6300	20	0,299
(6300, 6900)	24	0,358
(6900, 7500)	11	0,164
(7500, 8100)	4	0,060
(8100, 8700)	3	0,045
(8700, 9300)	2	0,030
(9300, 10000)	2	0,030
> 10000	1	0,015

Se aplica la formula

* Nuevo Grafico frecuencia vs edad



• Para la estatura definieron clases

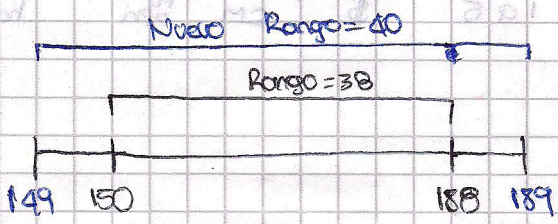
Sturges $\rightarrow K = 1 + 3,33 \log_{10}(67) = 7,08 \approx \boxed{8}$

* el valor maximo es ~~100~~ 188cm y el minimo es el rango es igual a 38

La Amplitud de cada intervalo estar dado por

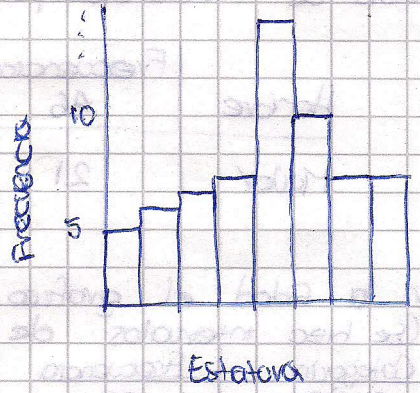
$A = \frac{\text{Rango}}{K} = \frac{38}{8} = 4,75 \approx \boxed{5}$

como hay un excedente de datos se reparte equifreentemente



La tabla quedar de esta

Intervalo	Frecuencia	FR	Moved
149, 154	5	0,075	151,5
154, 159	6	0,090	156,5
159, 164	7	0,104	161,5
164, 169	8	0,119	166,5
169, 174	17	0,254	171,5
174, 179	10	0,149	176,5
176, 184	7	0,104	181,5
184, 189	7	0,104	186,5

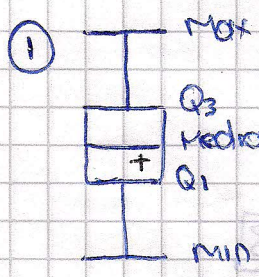


NOTA: Las tablas de frecuencia deben contar con

- Uniformidad = Clases de igual amplitud o variable que dependan del tipo de datos
- Claridad = Clases no traslapadas
- Completez = Cada dato pertenecer a una y solo una clase

Box-Plot (Cajas y Bigotes)

Lo señalan las características importantes de un conjunto de datos

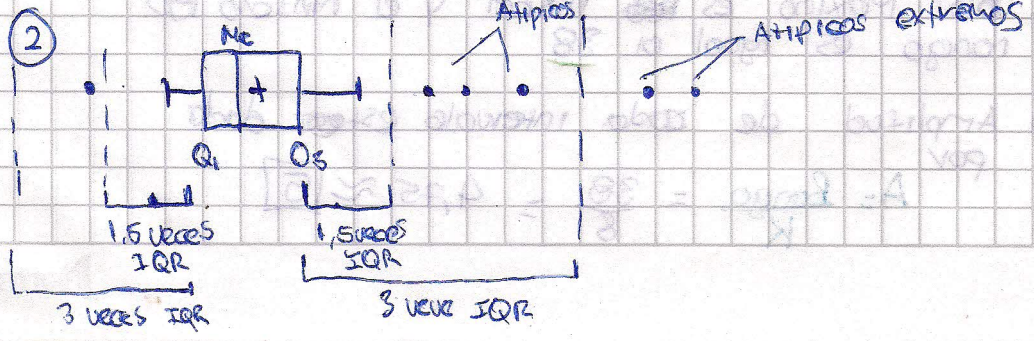


Q = cuartil

Nota = la mediana y la moda no siempre estan en el mismo lugar

→ 50% de los datos mas centrales

+ Medida muestral



ANÁLISIS DESCRIPTIVO DE UN CONJUNTO DE N DATOS

- Datos agrupados = Crear clases y hacer frecuencia
- Datos sin agrupar = Manipulación de la información como se fomo

Datos Agrupados
Medida Muestral

Datos Sin agrupar
Medida Muestral

sensible a valores extremos

$$\bar{X} = \frac{\sum_{i=1}^{\# \text{clases}} \text{Marca de clase} \cdot \text{Frecuencia de clase}}{\text{Total de frecuencia}}$$

\bar{X} = Promedio

$$= \frac{\sum \text{de datos}}{\# \text{ de datos}}$$

Ejemplo

$$\bar{X} = 170,38$$

Moda

↳ Marca de clase con la Mayor frecuencia

Ejemplo

$$\text{Moda} = 171,5$$

Moda

↳ Dato con mayor frecuencia (Dato que mas se repite)

Ordenar datos

Percentiles (Xp)

↳ Valores abajo y arriba de los cuales se encuentra una cierta proporción de datos del conjunto

$$P_{50} \Rightarrow p = 0,5$$

$$P_{25} \Rightarrow p = 0,25$$

$$X_p = L + \frac{(P-a) \times h}{f}$$

Percentiles

↳ Homologo de los Cuantiles (%)

~~FR~~ Frecuencia

F = FR de la clase que contiene el percentil

L = Limite inferior de la clase que contiene el percentil

P = Percentil a calcular (P50 ⇒ p=0,5)

a = FR Acumulada del intervalo anterior al del percentil

h = longitud de clase del percentil

Ejemplo P50

$$P_{50} = 169 + \frac{(0,5 - 0,38) \times 5}{0,254}$$

$$P_{50} = 171,2$$

Mediana

↳ Es igual al P50

Varianza

$$S_a^2 = \frac{\sum_{i=1}^k (m_i - \bar{x}_a)^2 \times f_i}{n-1}$$

Rango intercuartil

↳ diferencia entre el percentil 75 y el 25 (Datos mas centrales)

$$Q_{RANGE} = Q_3 - Q_1 = P_{75} - P_{25}$$

Ejemplo

$S_a^2 = 101,4$
 $S_a = 10,07$
 $P_{25} = 163,1$
 $P_{75} = 177,6$
 $Q_{RANGE} = 14,5$

Mediana

↳ Dato central

Impar = 1

Par = el promedio de los dos centrales

$$\tilde{x} = \frac{X_{(n/2)} + X_{(n/2)+1}}{2} \quad \text{Par}$$

$$\tilde{x} = X_{(\frac{n+1}{2})} \quad \text{Impar}$$

• 5 7 9

• 5 8 3 7
3 5 7 8

↳ $\frac{5+7}{2} = \boxed{6}$

Varianza

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

A mayor Varianza mayor dispersion

Desviacion Estandar

↳ Raiz de la Varianza

Rango intercuartil

$$Q_{RANGE} = Q_3 - Q_1 = P_{75} - P_{25}$$

Coficiente de Variacion

$$C.V. = \frac{\sigma}{\bar{x}} = \frac{\text{Desviacion estandar}}{\text{media}}$$

a mayor valor mayor dispersion

1. Importar Datos

%>% = Clic

Datos %>% Importar datos %>% Desde archivo Excel
Desde archivo de texto

2. Resumen y tipo de datos

Summary (df)
Str (df)

3. Modificar Datos

3.1 Conversion a factores

Datos %>% Modificar Variable del %>% Recodificador variables
conjunto de datos
activos

%>% Escoger variables

Para Proporciones

1 == "verdadero"
2 == "falso"

3.2 Nuevas Variables

Datos %>% Modificar variable del %>% Calcular nueva variable
conjunto de datos
activo

- factor (IFelse (Condicion, "Verdadero", "falso")) Prop
- Formula : Variable 1 * variable 2 =

3.3 Eliminar Factores no necesarios

Datos %>% Modificar variable del %>% Descartar niveles sin
conjunto de datos
activo
USO

Todos los factores

4. Agrupar Datos (filtrar)

Datos %>% Conjunto de datos activos %>% Filtrar el conjunto
de datos activos

Incluir todas las variables

Expresion de seleccion

- Para un factor : df == "factor 1"
- Para 2 o mas factores : df == "factor 1" | df == "factor 2"

5. Estadísticas de Interés (Pruebas No Normales o σ conocidas)

Estadísticas $\% > \%$ Resúmenes $\% > \%$ Resúmenes numéricos $\% > \%$

Resumir por grupo $\% > \%$ Seleccionar Variables $\% > \%$ Estadísticas de Interés

▣ Medida

▣ Desviación típica

6. Juego de hipótesis para medias

H_0 = Complemento

Cola derecha

H_a = Lo que queremos : Cola izquierda

Bilateral

6.1 Test de Normalidad

• Prueba de hipótesis

Estadísticas $\% > \%$ Resúmenes $\% > \%$ Test de Normalidad

$P\text{-Value} > \alpha$ $\% > \%$ Test por grupos
Son normales

• Q-Q Plot

Gráficas $\% > \%$ Gráficas de comparación de cuantiles

• Histograma

Gráficas $\% > \%$ Histograma

6.2 Homocedasticidad (Razon de varianzas)

Cuidado = "factor en orden alfabético"

Teniendo los dos factores agrupados

Estadísticas $\% > \%$ Varianzas $\% > \%$ Test F para dos varianzas $\% > \%$ Bilateral

Convertir var # en factor

$P\text{-Value} > \alpha$
Son iguales $\sigma_1^2 = \sigma_2^2$

6.3 P_h (T-Student) Para 1 media

Estadísticas $\% > \%$ Medias $\% > \%$ Test t para una muestra

- Ingresar H_a
 $\mu = x$

- Revisar α
- Analizar p-Value
- Concluir

6.4 Ph (T-Student) Para dos medias

Estadísticos $\% > \%$ Medias $\% > \%$ Test t para muestras $\% > \%$ independientes Seleccionar Variable

Revisar factores
Voltear H_a en
caso de estar
truncados

- Seleccionar H_a
- Revisar α
- Varianzas iguales
OSI
ONO

*
Teste para
diferencia

- Verificación Ph (Gráficamente)

Gráficos $\% > \%$ Gráficos de las medias

⊙ Intervalos de confianza

6.5 Ph para una proporción

Estadísticos $\% > \%$ Proporciones $\% > \%$ Test de proporción es para una muestra

$$H_0 = P_0$$

(Si los factores están truncados
 $H_0 = 1 - P_0$)

6.6 Ph (Multinomial) (Chi-Cuadrado)

H_0 = Todos las probabilidades son iguales a las propuestas

H_a = Al menos una de las probabilidades es diferente a la propuesta

Estadísticos > Resúmenes > Distribución de frecuencias

- Ingresar valores

Test Chi-Cuadrado de bondad de ajuste

- No hay evidencia en los datos con una significancia

$\alpha = 0.01$

- Correlation
- Analysis of Variance
- Regression

Estadísticas de los datos

Estadísticas de los datos y Test de hipótesis

- Intervalos de confianza
- Pruebas de hipótesis
- Pruebas de hipótesis
- Pruebas de hipótesis

Intervalos de confianza

Intervalos de confianza

Intervalos de confianza

Intervalos de confianza

Estadísticas de los datos y Test de hipótesis

$H_0 = \mu = 10$

Intervalos de confianza

Intervalos de confianza

Intervalos de confianza

Intervalos de confianza

Intervalos de confianza

Intervalos de confianza