

TEMA 5

REDUCCION DE DATOS , MEDIDAS DESCRIPTIVAS BASICAS Y REPRESENTACIONES GRAFICAS

1-INTRODUCCION DE LA ORGANIZACION A LA DESCRIPCION DE DATOS.

Una vez que hemos depurado los datos y comprobado que la tabla de datos está bien, pasaremos a analizar esos datos.

Si para la depuración de datos hemos usado las frecuencias ya estará hecho el primer paso, sino deberemos realizar esta estadística descriptiva (la distribución de frecuencias)

La variabilidad es el estudio de la dispersión de los valores, contribuyendo a explicar lo grande que es la muestra y la naturaleza de los valores (o sea las medidas de dispersión nos muestra la variabilidad). Ahora el libro nos ofrece la descripción de conceptos ya vistos. Cuando usamos la estadística para contrastar hipótesis buscamos la relación posible entre variables

Variable es lo que puede cambiar de un sujeto a otro

Constante es un valor común para todos los sujetos algo que no cambia

La estadística descriptiva es el procedimiento de organizar, clasificar y resumir el conjunto de datos, o sea crear una tabla donde organizar y clasificar y después crear frecuencias y gráficos donde se resumen. Pero también se usarán fórmulas que nos ayudan a resumir los datos como las **medidas de tendencia central** y **las medidas de dispersión o variabilidad** (medidas centralizadas son la media aritmética, la mediana y la moda), (las medidas de dispersión son la desviación con respecto a la media, la desviación media, la varianza y la desviación típica) Cuando los datos son muy diferentes los valores, ej 150 valores como en el caso del CI, los reducimos a intervalos. Primero se hacen las medidas centralizadas y luego las de dispersión (lógico porque para hacer las de dispersión necesitamos la media aritmética).

5.2 MEDIDAS DE TENDENCIA CENTRAL : MEDIA , MEDIANA Y MODA. USOS Y SU INTERPRETACION

Busca cuál es el valor que mejor representa al grupo.

La media aritmética se halla sumando todos los valores dados y dividiéndolos entre el total. Se representa

$$\bar{x} = \frac{X_1 + X_2 + X_3 + \dots + X_n}{N}$$

con el signo \bar{x} y la fórmula es $\bar{x} = \frac{X_1 + X_2 + X_3 + \dots + X_n}{N}$ o sea la suma de todos los datos entre la

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{N}$$

población o en reducido $\bar{x} = \frac{\sum_{i=1}^n x_i}{N}$ que quiere decir exactamente lo mismo la suma de todos los valores, pero esta forma sería un poco difícil de llevar a cabo si tenemos muchos valores, imaginar una encuesta de 100 personas sobre el número de hijos, o sea que si nos pusiesemos a sumar los hijos de cada una de las cien personas nos perderíamos antes o después para muestra un botón

11111222222333333444444222224342323131, si no te pierdes sumando todo eso sería impresionante,

Por eso los estadísticos han hecho una trampita es la creación de una columna nueva que ya sabéis visto antes la de $X_i \times f_i$ en esta multiplicamos el valor (x_i) por el número de veces que aparece (f_i) y sumando la multiplicación de todos los valores $\sum X_i$ nos daría lo mismo que sumar todos los valores. Ya veréis que lo entendéis mejor con un ejemplo :

x_i	f_i	$X_i * f_i$
1	5	5
2	2	4
3	3	9
Σ	10	18

Hemos preguntado a 10 amigos por los hijos que tienen , y el resultado es el siguiente hay 5 amigos que tienen 1 hijo , 2 amigos que tienen 2 hijos y 3 amigos que tienen 3 hijos . O sea 1,1,1,1,1,2,2,3,3,3, la media la podríamos hacerla sumando todos los datos $1+1+1+1+1+2+2+3+3+3 = 18$, pero mirar que coincidencia la suma de la columna $X_i * f_i$ es también 18. Porque si sumamos $1+1+1+1+1=5$ y si multiplicamos 1 número de hijos por 5(f_i) también nos da 5 , y así sucesivamente. Por eso la suma de esas columnas es igual que la suma entre el total de valores. Pero si os cuesta entenderlo no lo penseis simplemente tomarlo como una fórmula más , si sumas todas las columnas de $X_i * f_i$ y lo divides entre el total de la muestra nos da la media aritmética.

Con esta medida podemos comparar los grupos miraremos si es inferior o superior al resultado... pero todo esto como pedagogo no es estadística. La media solo puede usarse en medidas de intervalo y de razón, o dicotómicas.

La moda es el valor que más aparece, no necesita ningún cálculo . Mira la columna de la frecuencia absoluta (f_i) y el número más alto nos indicará cuál es el valor x_i de la moda.

Si hay dos valores con la misma frecuencia absoluta más alta lo llamaremos bimodal, o sea moda de dos valores.(importante la moda siempre es x_i o sea el valor a estudiar no se pondría el número de veces. Por ejemplo la moda de la tabla anterior será 1 hijo. Si son tres o más valores los altos será, plurimodales. La mediana es el valor que está en medio , el del centro cuando las puntuaciones están en tabla o ordenados es fácil , es el que ocupa la posición central , coincide con el percentil 50 (o sea percentil viene de por ciento y el percentil 50 es el que está en el 50% o sea la mitad, porque la mitad de 100 es 5'.

Para calcularlo miraremos la columna de frecuencia absoluta acumulada es decir F_i y marcaremos n (el número de gente) y lo dividiremos entre dos y cuando la frecuencia absoluta acumulada pase ese número de $n/2$ ese valor o sea x_i será la mediana. Veamos esta tabla de 15 personas entonces $n/2=15/2= 7,5$

x_i	f_i	F_i
1	7	7
2	4	11
3	5	15

15

Donde estará el valor 7,5 pues como con 1 hijo se queda en 7,7,5, estarán en dos hijos. Cuando la división $n/2$ da un número entero la mediana necesitará la media aritmética de esos dos valores centrales por ejemplo en la tabla $14/2=n/2=7$

xi	fi	Fi
1	5	5
2	2	7
3	7	14

Pero como es un numero entero no nos serviria porque la mediana seria un intervalo que va de 7 a 8 por eso su mediana sera dos hijos y tres hijos porque entre ellos esta 7,5.

La medida mas precisa de tendencia central es la media aritmetica, pero solo se puede hacer en medidas de intervalo o de razon (es decir las que tienen numeros) Es la suma de todos los valores entre el total de la muestra por eso se hace una columna nueva, donde se multiplica la fi por xi y al sumarlo nos seria la suma de todos los valores. En esta tabla seria $30/14=2,14$

xi	fi	Fi	Fi·xi
1	5	5	5
2	2	7	4
3	7	14	21
Σ	14		30

Si hay muchas diferencias entre la media y la moda nos indica un dato muy interesante, nos indica que hay muchos valores bajos , lo entenderéis mejor con un ejemplo de las notas de clase la media es 5,2 pero la moda es 7 eso significa que la nota que mas se repite es notable siete pero hay un grupo grande de miembro de la clase con notas muy muy bajas que nos bajan la media dos puntos y deberiamos preguntarnos porque (el metodo no se adapta bien a todos los alumnos por ejemplo) por eso la media es la mas importante porque al verse afectada por la puntuaciones extremas da una vision mas global no solo mira lo normal sino que ve un poco mas alla .La media se ve arrastada por las puntuaciones extremas. La distribucion de frecuencias normales cuando la media aritmetica , la mediana y la moda son el mismo valor

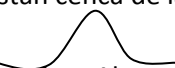

5.3 MEDIDAS DE VARIABILIDAD O MEDIDAS DE DISPERSION

Tambien necesitamos saber si hay muchos valores que se alejan de la media aritmetica , o si la mayoría estan cerca de la media para ello se usa la medida de dispersion

5.3.1 LA DESVIACION MEDIA .

Es la media aritmetica de la desviacion del valor con respecto a la media . Traducido es la media aritmetica de las desviaciones respecto a la media , que son la resta a cada valor de la media aritmetica, se suma todos los resultados de las restas y se dividen entre el numero total de la muestra. Por eso es necesario tener la desviacion respecto a la media para lo que se hara una nueva columnita que sera $xi - \bar{x}$. Es $\bar{x} = 18/10 = 1,8$

xi	fi	Fi	Fi·xi	$xi - \bar{x}$
1	5	5	5	$1-1,8=0,8$
2	2	7	4	$2-1,8=0,2$
3	3	10	9	$3-1,8=1,2$
Σ	10		18	2,2

Si el valor de esta desviacion media es muy alta significa que la mayoría de los datos sus valores estan muy lejos de la media aritmetica, vamos para pedagogia que hay un problema. Pero si es pequena significa que los valores estan cerca de la media , significa que todo esta bien. Ej grafico todo bien es una campana de guas normal  todo mal  Si la desviacion media es 0 significa que todos los valores son iguales . Ahora veremos como se halla la desviacion media tomaremos los valores sin

importarnos el signo que sea negativo o positivo. Es por eso que una desviación media 0 significa que todos los valores son iguales porque si tuviésemos estos valores.

x	fi	xi · fi	$Xi - \bar{x}$
80	2	160	-24
128	2	256	24

$$\bar{X} = 416/4 = 104$$

Ahora haremos la desviación media que es una media aritmética de la desviación respecto a la media, y como lo que es la desviación media es una media, lo que debemos hacer es sumarla desviación respecto a la media de todos los valores que han dado los encuestados. Por eso como ocurre al hacer la media aritmética debemos hacer una nueva columna que será $|xi - \bar{x}| \cdot fi$


xi	fi	Xi · fi	$xi - \bar{x}$	$ Xi - \bar{X} \cdot fi$
1	5	5	-0,8	$0,8 \cdot 5 = 4$
2	2	4	0,2	$0,2 \cdot 2 = 0,04$
3	3	9	1,2	$1,2 \cdot 3 = 3,6$
Σ	10	18		7,64

$$\bar{X} = 1,8$$

Otra cosita si os fijáis en la desviación respecto a la media del valor 1 tendría que restar a 1- 1.8 y por tanto me da un producto negativo, pero en estadística para la elaboración de tablas no vamos a tener en cuenta el signo negativo (ese es el significado de esa especie de parentesis en forma de líneas que tienen algunos valores). Una vez hecha esta nueva columna nos iremos a la fila de los sumatorios y sumaremos todos sus valores y lo dividiremos entre el total de la muestra diez en este caso y esta será la desviación media

$$Dm = 7,64/10 = 0,76$$

5.3.2 La desviación típica

Otro índice conocido es la desviación típica representada por una letra S para estadístico (basado en la muestra), y $\bar{\sigma}$ para parámetro (para población). Además esta la varianza que se representa con una  para parámetros y una s al cuadrado para muestra. Se realiza elevando al cuadrado cada desviación respecto a la media de todos los valores o sea tendríamos que poner al cuadrado (multiplicar por sí mismo) las desviaciones respecto a la media de todos y cada uno de los sujetos, por tanto haremos una nueva columna que multiplicaría la frecuencia absoluta por el cuadrado de la desviación respecto a la media. La suma de estos cuadrados se ha de fi y así no tendremos que sumarlo uno a uno, además debemos una vez hecha la suma dividirlo entre el número de la muestra.

Vamos paso a paso

- 1) Restamos la media aritmética a los valores
- 2) Ponemos al cuadrado la desviación respecto a la media (paso 1)
- 3) Multiplicamos por el número de veces que aparece cada valor fi.
- 4) Sumamos el resultado de toda esta nueva columna
- 5) Dividimos el resultado de la suma entre el total de la muestra

6) Al resultado le sacamos su raíz cuadrada

Aquí da igual que utilicemos el número de la desviación respecto a la media en negativo o en positivo porque al ponerlo al cuadrado el negativo desaparecería porque menos por menos es más y al hacer el cuadrado desaparece por tanto el negativo

xi	fi	Xi·fi	xi- \bar{x}	Xi - \bar{X} al cuadrado	Xi - \bar{X} al cuadrado por fi
1	5	5	-0,8	0,64	3,2
2	2	4	0,2	0,04	0,08
3	3	9	1,2	1,44	4,32
Σ	10		8		7,6

Media aritmética es = a 1,8

Desviación media 0,76

Varianza = 0,76

La desviación típica es = 0,87

$$\sigma = \sqrt{\frac{\sum f_i (M_i - \bar{x})^2}{n}}$$

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

La fórmula es

La desviación típica es difícil de interpretar y se utiliza para comprobar la dispersión entre grupos distintos. Y por tanto eso a veces se hace la misma desviación típica y la máxima

Para la mínima desviación típica todos los valores deben ser el mismo valor y para la máxima, la mitad de los sujetos deberían obtener la puntuación máxima de la escala y la otra mitad la mínima. Esto solo se ve en variables dicotómicas o dicotomizadas y cuanto más grande sea la distancia entre valores más improbable es que se de esta situación... Si una desviación típica es más baja en un grupo que en otro significa que los sujetos están más igualados entre sí. La desviación típica en una población será más elevada que en una muestra, al haber más sujetos, serán más probables encontrar mayores diferencias interindividuales.

Por eso se ve que cuanto más baja, los valores son más iguales. Estas que hemos visto son desviaciones típicas sesgadas que acentúan el efecto de las grandes desviaciones y es siempre superior a la desviación media.

La desviación típica y varianza insesgadas miden lo mismo que la sesgada aunque tienen diferentes propiedades, que es en vez de para la muestra para la población. Y es más elevada que en una muestra (al haber más sujetos es más probable encontrar más diferencias por lo que ponemos el 1 para disminuir el denominador, el resultado de la desviación será por tanto mayor. La insesgada se codifica para muchos $\sigma-1$ y la sesgada σ

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

5.3.3 La amplitud o recorrido

Otra medida de dispersión son la amplitud o recorrido de una variable que es la puntuación mayor menos la menor más 1.

$$A = X_1 - X_{\text{ultimo}} + 1$$

Y se usa para organizar intervalos y para los gráficos. Y afecta a la dispersión y las otras medidas porque no es lo mismo una varianza de 5 para 10 personas que para 1000. Y se utiliza como medida de dispersión única cuando no se pueden usar otras, o para complementar la moda. Solo se hace con dos puntuaciones por lo que puede dar resultados que confundan sino se toma con cuidado porque puede haber valores extremos (outliers)

5.3.4 La desviación semi intercuartílica

Se denomina desviación semiintercuartílica Q que muestra la dispersión en el 50% central de la distribución. Útil para nivel de medida ordinal y como complemento a la mediana. Así se eliminan las puntuaciones extremas al prescindir del 25% superior e inferior de las puntuaciones. Para hallarla se ha de saber el 1º cuartil y el 3º cuartil recuerda $1Q = 1 \cdot N \div 4$. El $3Q = 3 \cdot N \div 4$

$$\text{desv. int cuartil.} = \frac{Q_3 - Q_1}{2}$$

5.3.5 El coeficiente de variación

Se representa con V o CV. Permite comparar variables (valores) de medida de datos diferentes, por ejemplo el coeficiente intelectual de amplitud 50 y actitud hacia el estudio de amplitud 10.

$$C.V. = \frac{S}{\bar{X}} * 100$$

Este índice tiene límites fijos dependiendo de valor de la media. Suele ser menos que 1. Se interpreta en porcentaje depende de la desviación típica o estándar. Vamos es la relación entre esta desviación típica y su media.

5.4 Media y desviación típica para variables dicotómicas

Las variables dicotómicas es un tipo de variable X_i que solo tiene dos posibles valores. Valores a los que se puede hacer la media aritmética para cada uno. La suma de ambas medias debe dar siempre 1

$$P + q = 1$$

La varianza también es muy fácil de calcular

$$\sigma^2 = p \cdot q$$

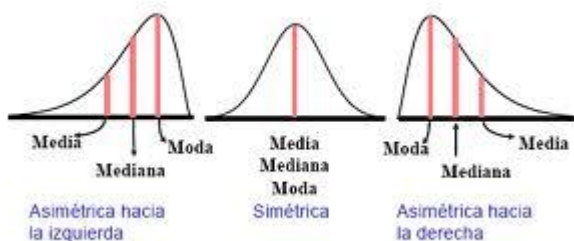
$$S = \sqrt{p \cdot q}$$

5.5 Asimetría y apuntamiento : relation con la curva normal

La representación gráfica de las variables, sobre todo cuando trabajamos con muestras grandes, tienden a ser curvas que por su grado de asimetría, pueden asemejarse a una de estas tres.

Una curva normal carece de asimetría, es simétrica o tiene índice de asimetría igual a 0

Por su grado de asimetría puede asemejarse a una de estas tres



La asimetría positiva indica que la mayoría de los sujetos está en la parte baja de la puntuación de la distribución de frecuencia. La cola de la distribución está a la derecha. Lo cual no indica que los sujetos tengan puntuaciones bajas, todo depende de donde se encuentra el grupo más grande de puntuaciones que puede ser más alta.

La asimetría negativa indica lo contrario, que los sujetos tienden a agruparse en torno a las puntuaciones altas de la distribución, lo que como antes no significa necesariamente que los sujetos tengan puntuaciones altas.

La asimetría positiva a negativa se debe solo al signo del cálculo del índice

El índice de Pearson es la forma de calcularlo más fácil

El coeficiente de correlación de Pearson, pensado para variables cuantitativas (escala

mínima de intervalo), es un índice que mide el grado de covariación entre distintas

variables relacionadas. El coeficiente de correlación de Pearson es un índice de fácil ejecución e, igualmente, de

fácil interpretación. Digamos, en primera instancia, que sus valores absolutos oscilan

entre 0 y 1. Hemos especificado los términos "valores absolutos" ya que en realidad si se contempla

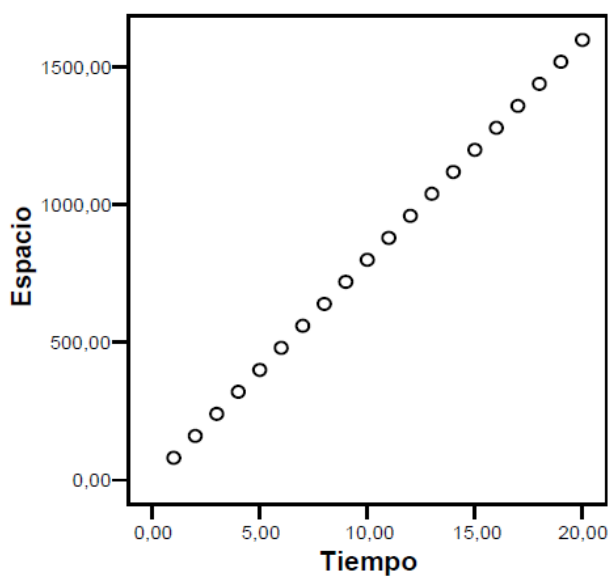
el signo el coeficiente de correlación de Pearson oscila entre -1 y +1.

En este sentido, tan fuerte

es una relación de +1 como de -1. En el primer caso la relación es perfecta positiva y en el segundo perfecta negativa.

Decimos que la correlación entre dos variables X e Y es perfecta positiva cuando exactamente en la medida que aumenta una de ellas aumenta la otra. Esto sucede cuando la relación entre ambas variables es funcionalmente exacta.

Por ejemplo, la relación entre espacio y tiempo para un móvil que se desplaza a velocidad constante. Gráficamente la relación ser del tipo:



Vamos resumiendo cuanto más intensa sea la concordancia (en sentido directo o inverso) de las posiciones relativas de los datos en las dos variables, el producto del numerador toma mayor valor (en sentido absoluto). Si la concordancia es exacta, el numerador es igual a N (o a -N), y el índice toma un valor igual a 1 o -1.

La formula es

$$A = \frac{\overline{X} - M_o}{S}$$

Esa si porque la media , a diferencia de la moda y la mediana se ve atraido por portuaciones extremas, luego su posicion relativa se vera desplazada hacia puntuaciones extremas.

El apuntamiento o la curiosi indica el grado en que la distribucion es mas o menos picudo , es decir si la distribucion es mas abierta o dispersa respecto a la media.

Si las puntuaciones estan poco concentradas respecto a la media y por tanto mas chata o aplastada (platicurtica) o por el contrario mas apuntada y por tanto mas estrecha o con las distribuciones mas concentradas en torno a la media (leptocurtica). Una curiosi igual a cero viene representada por la distribucion normal. Curtosis superior a cero nos indica una distribucion leptocurtica maestra que si es inferior a cero , esta sera platicurtica. Par calcular la curiosi puede utilizarse la siguiente formula.

Hemos comentado que el concepto de asimetría se refiere a si la curva que forman los valores de la serie presenta la misma forma a izquierda y derecha de un valor central (media aritmética)

$$g_2 = \frac{1}{N} \cdot \frac{\sum_{i=1}^n (X_i - \bar{X})^4 \cdot f_i}{\sigma^4} - 3$$

5.6 Representaciones graficas

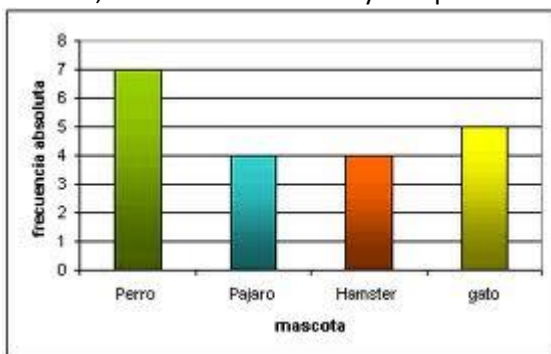
La representacion grafica de los datos como histogramas , poligono de frecuencias, ciclo grama... son buen complemento de los indices numericos y ayudan a comprender rapidamente la information descriptiva. A partir de una distribucion de frecuencias es muy facil realizar una representacion grafica, como hemos visto en el capitulo anterior. Y es recomendable adaptar el tipo de grafico al nivel de medida de la variables

5.6.1 Grafica de sectores

Tb llamado ciclo grama o grafico de torta representa variables con nivel de medida nominal. Tiene forma circular y dividido en porciones cada porcion representa la presencia proporcional de cada una de los niveles de la variable.

5.6.2 Grafico de barras

Es comparativo , compara porcentajes f_i , en una muestra. Suele usarse para variables con nivel de medida ordinal , o cuando es nominal y comparaciones de variables clasificatorias o categoricas.

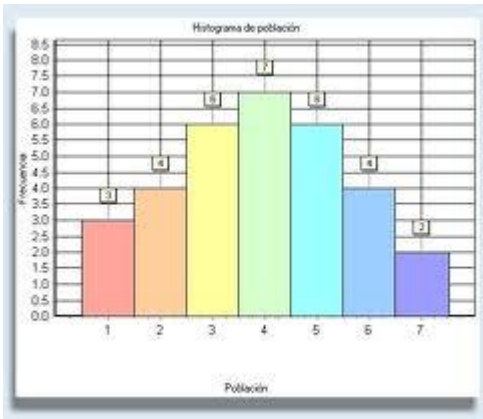


5.6.3 Histograma

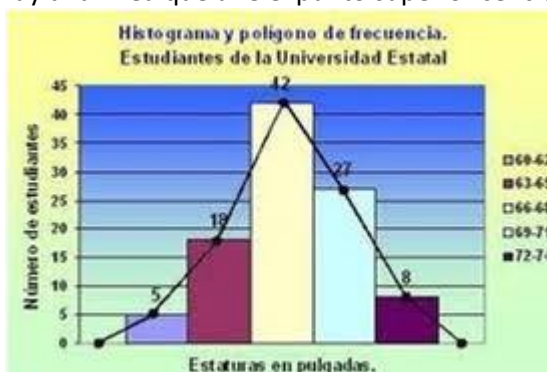
Es parecido al diagrama de barras, pero se usa en variables cuantitativas continuas con nivel de medida de intervalo o de razon . Las barras estan juntas. No aparecen todos los valores de la variable continua , son muchos los valores que pueden obtener los sujetos y seria muy farragoso con un grafico con muchas barras. Por eso se agrupan las puntuaciones en intervalo. Asi cada puntuacion tiene un numero determinado de puntuaciones, denomina valor intervalo .

Cuanto menor es el numero de intervalo ,, mas suavizados aparece la grafica. Los intervalos suelen usarse . Estimar por tanteo el numero de intervalos que se desea, segun lo que acabamos de indicar Iniciar el primer intervalo que incluya a puntuacion directa mas baja y que sea multiplo del valor del intervalo

Situar la marca de clase (punto central del intervalo) en el eje de abscisas en el centro de cada barra



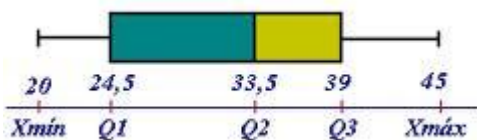
Hasta hace unas décadas, la agrupación por intervalos se utilizaba también para el cálculo de medidas de tendencia y variabilidad. Porque no se usaban programas informáticos. El polígono de frecuencias es cuando hay una línea que une el punto superior central de cada barra, abriendo y cerrando el eje de



abscisas.

5.6.4 Grafico de caja

También llamado como caja y pastillas o caja y bigotes. Es práctico porque permite hacerse una idea rápida de la distribución de las puntuaciones en la zona central desde el primer cuartil o percentil 25 hasta el tercer cuartil o percentil 75 y en los extremos. Las pastillas representan las puntuaciones hasta los extremos de la distribución o sea el último y el primer valor



5.6.5 Grafica de tallo y hojas

El gráfico de tallo y hojas combina la representación numérica y gráfica. Es una especie de histograma horizontal, construido con los números correspondientes a las puntuaciones. Las hojas es el último dígito de la puntuación y el tallo, es el resto de dígitos

15,16,21,23,23,26,26,30,32,41

Tallo	Hoja
1	5 6
2	1 3 3 6 6
3	0 2
4	1

