

ANÁLISIS EXPLORATORIO DE DATOS ESPACIALES ESTADÍSTICA ESPACIAL

DEPARTAMENTO DE GEOGRAFÍA
FACULTAD DE HUMANIDADES



Prof. Silvia Stela Ferreyra

Estadística:

serie de datos numéricos presentados de manera ordenada y sistemática,

pero cuando comenzamos a adentrarnos en la investigación, empezamos a percibir a la estadística como una herramienta para obtener resultados confiables en una Investigación.

La Investigación puede tener como objetivos:

- El análisis exploratorio de datos que va más allá de una descripción de los mismos sino que es parte de la actitud que define al buen investigador, es decir aprender a observar que expresan esos datos, como ser encontrar patrones de comportamiento o anomalías en las distribuciones.
- El análisis confirmatorio de datos relacionado con la estadística inferencial, que comienza con la formulación de hipótesis que luego serán confrontadas con los datos empíricos por medio de metodologías estadísticas adecuadas.

“El Análisis de Datos” o Estadística descriptiva comprende:

Serie de técnicas orientadas a extraer información de un conjunto de datos, mediante medidas que resumen los principales rasgos de la distribución de frecuencias de ese conjunto de datos.

Las aplicaremos para el análisis de las diferencias espaciales entre unidades territoriales.



Cuando se dispone de una cantidad considerable de datos es útil presentarlos en:

TABLAS ESTADÍSTICAS

- Permiten conocer el número de veces que se repite un fenómeno. En su estado más desagregado son listados de valores, que muestran la distribución de los valores de la variable.
- Con el fin de reducir el listado de valores, se agrupan los valores en clases y se presenta la frecuencia de las mismas, a estas tablas se denominan **Tablas de Frecuencias**.
- Es común presentar estas frecuencias absolutas como frecuencias relativas o porcentuales, y también puede presentarse en forma acumulada.
- Para establecer las clases o grupos que estará en función del número de observaciones, algunos autores proponen que una distribución debe tener un número de clases no inferior a 6 y no superior a 20. Según Huntensberge lo estima con $K = 1 + 3,3 \cdot \log n$, siendo n el núm. total de observaciones.
- Rango = $V_{\max} - V_{\min}$
- $C(\text{ancho de clase}) = \text{Rango} / \text{núm. de clases}$

Para datos sin agrupar

X	Frecuencia Absoluta	Frecuencia %	Frecuencia Acumulada	Frecuencia Acumulada %
	fa	fr	f %	FA%
1	3	0,06	6	0,06
2	9	0,18	18	0,24
4	13	0,26	26	0,5
6	8	0,16	16	0,66
7	8	0,16	16	0,82
8	4	0,08	8	0,9
10	5	0,1	10	1
Total	50	1	100	

SERIES TEMPORALES

Distribución de una variable en el tiempo:

Para la descripción de una serie temporal es frecuente la elaboración de números índices (resultado de una comparación por cociente), y tomando el valor de un año base como el 100%, se analiza la evolución por comparación.

$$I(i/t) = (X_i/X_t) * 100$$

>100 % crecimiento

<100% decrecimiento

**2. NÚMEROS ÍNDICES SIMPLES:
CONCEPTO Y PROPIEDADES. ÍNDICES EN CADENA.
NÚMERO ÍNDICE SIMPLE**

Ej.: Estudiar la evolución del número de turistas (en miles) que entraron en España provenientes de Francia y Canadá entre los años 1992 y 1997, determinando los correspondientes índices.

Años	Turistas Franceses	Turistas Canadienses
1992	1035	44
1993	1230	56
1994	1351	65
1995	1230	85
1996	1750	116
1997	2554	153

Años	$I_{t/92}^F$	$I_{t/92}^C$
1992	1	1
1993	1,1884	1,2727
1994	1,3053	1,4772
1995	1,1884	1,9318
1996	1,6908	2,6363
1997	2,4676	3,4772

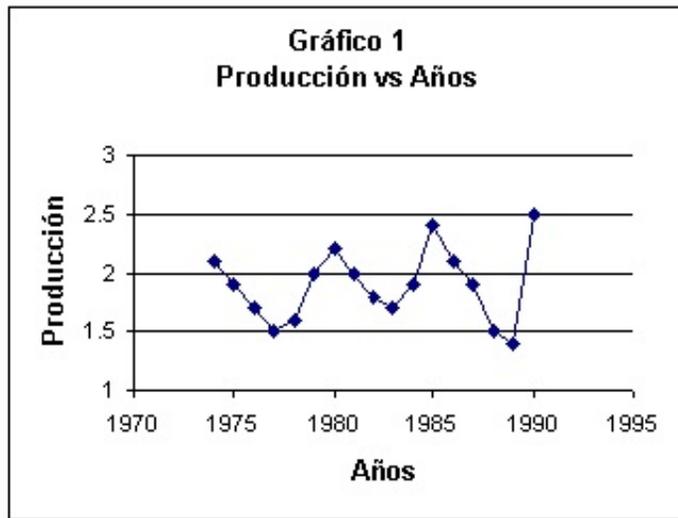
© Antonio Pajares Ruiz

6

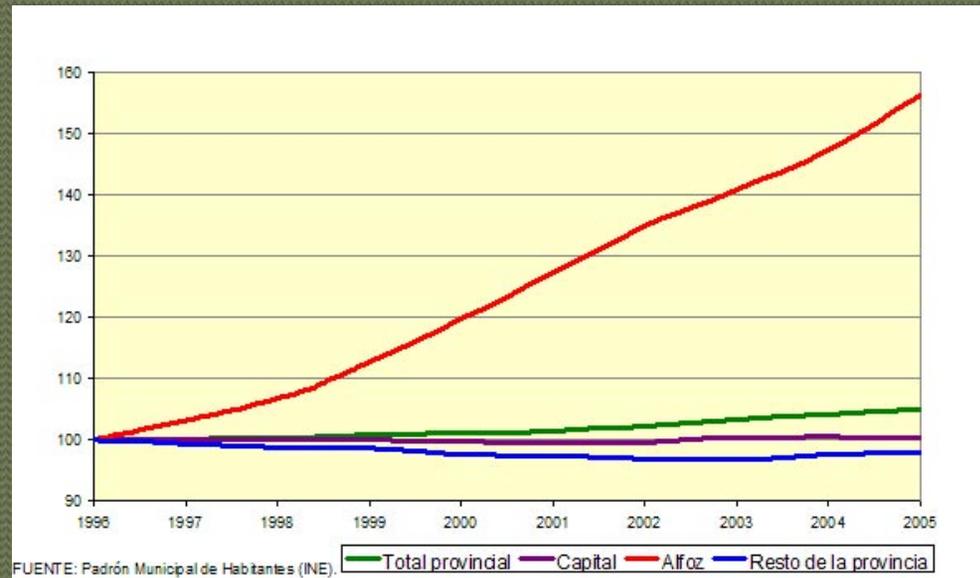
Evolución del número de turistas que entraron en España provenientes de Francia entre 1992 y 1997, respecto a 1992.



Los gráficos de series temporales, representan en el eje horizontal la variable temporal, mientras en el eje vertical los valores de la variable o la conversión en índices.).



Producción de motocicletas en millones

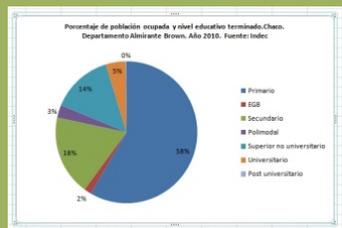


Evolución de la población de Valladolid, de su alfoz (entorno) y del resto de la provincia 1996-2005. Números índice (1996=100)

REPRESENTACIONES GRAFICAS



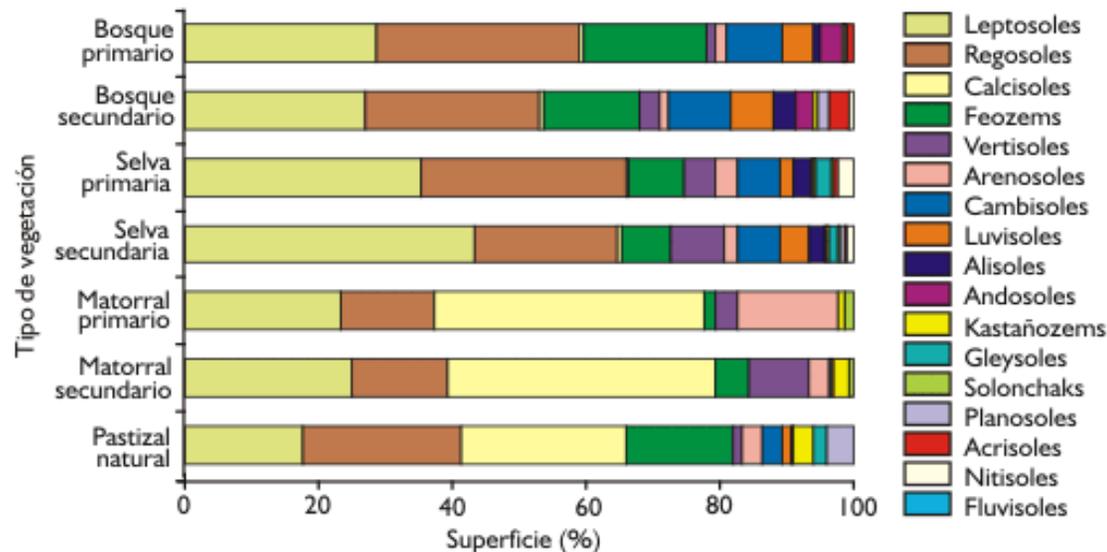
Año	2005	2006	2007	2008	2009	2010
Sexo Masculino	100	160	200	120	230	140
Sexo Femenino	50	90	80	90	100	600
Total	150	250	280	310	330	200



Las tablas estadísticas pueden ser representadas a partir de una numerosa serie de gráficos, que permiten visualizar de forma rápida la información geográfica.

- Diagrama de barras: permite representar la importancia cuantitativa de las diferentes categorías de una variable cualitativa o cuantitativa, por ejemplo: distribución de usos de suelos en un determinado espacio.

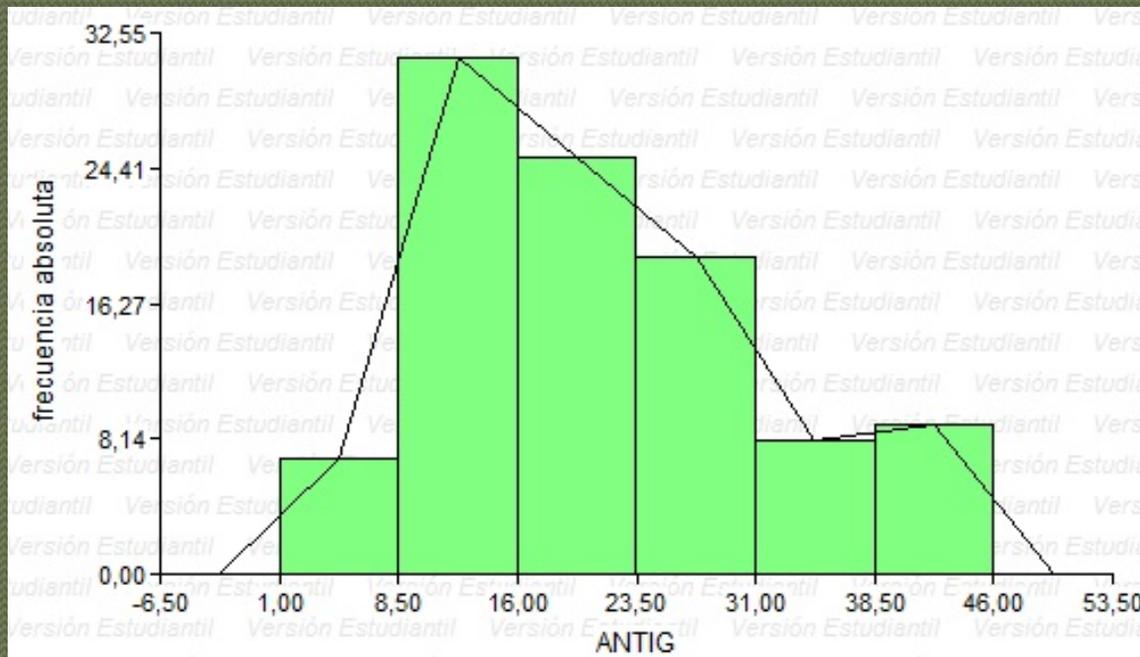
Figura 3.2 Suelos asociados a los diferentes tipos de vegetación presentes en México



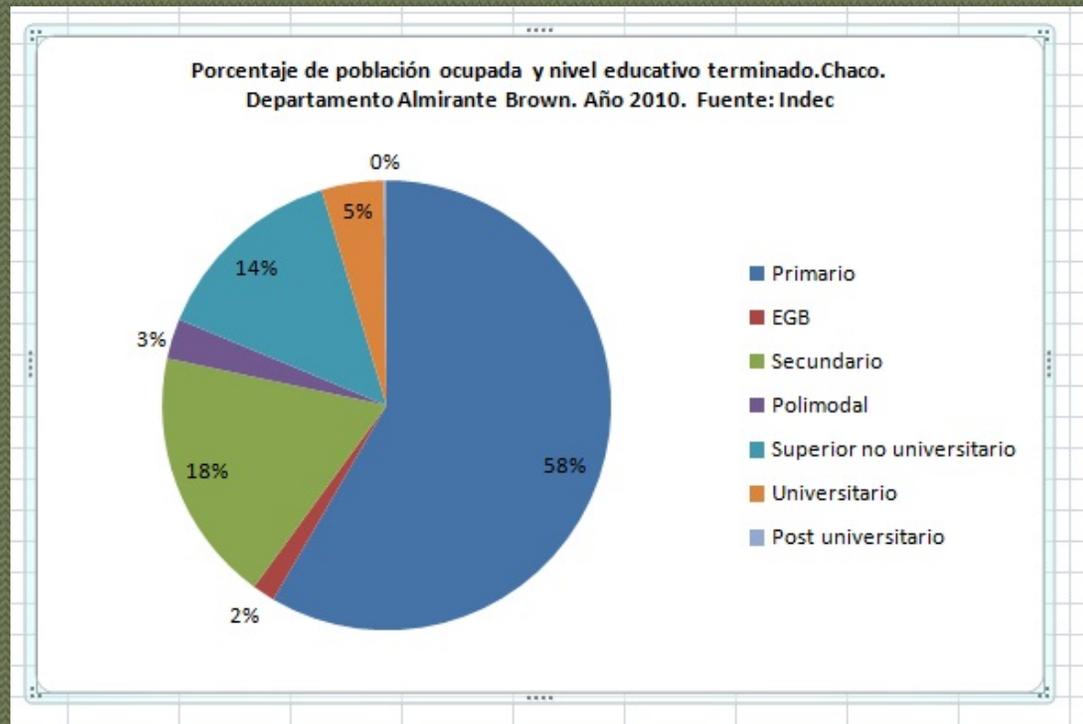
Fuente:

Elaboración propia con datos de: Semarnat. *Inventario Nacional de Suelos 2002*. México. 2002.

- El histograma de frecuencias se construye sobre un eje vertical, que representa la frecuencia absoluta o relativa y en el eje horizontal van las clases definidas por su intervalo y punto medio o marca de clase, que sirven para trazar el polígono de frecuencias.



- Diagrama de sectores puede usarse indistintamente para variables cualitativas o cuantitativas, hace equivaler la frecuencia de la categoría a un área proporcional del círculo.



MEDIDAS DESCRIPTIVAS

Medidas de tendencia central:

- 1) Media Aritmética: la más conocida, sumatoria de todos los valores de la variable dividido por el número de datos.

$$X_{\text{media}} = \sum X_i / N, \text{ para datos sin agrupar}$$

$$X_{\text{media}} = \sum f_i \cdot X_i / \sum f_i, \text{ para datos agrupados, } X_i \text{ es la marca de clase.}$$

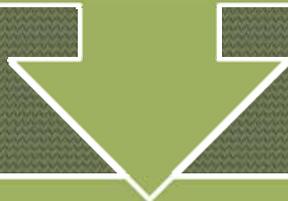
$$X_{\text{media ponderada}} = \sum X_i \cdot W_i / \sum W_i,$$

si los datos están agrupados X_i es la marca de clase.

Xmedia Armónica = $1 / \sum 1 / X_i$, cuando se trabaja con variables referidas a tasas de cambio o movimiento

La media armónica se define como el recíproco de la media aritmética de los recíprocos:

Este valor se emplea para promediar variaciones con respecto al tiempo.



Xmedia Geométrica = $\sum \log X_i / N$, cuando la variable crece de forma geométrica

Existen dos usos principales de la media geométrica:

Para promediar porcentajes, índices y cifras relativas y

Para determinar el incremento porcentual promedio en ventas, producción u otras actividades o series económicas de un periodo a otro. La media geométrica no se ve tan afectada por valores extremos.

MEDIDAS DE TENDENCIA CENTRAL



2) Moda: es el valor más observado o que más se repite, puede haber más de un a moda



$$\text{Moda} = L_i + c \cdot \left(\frac{d_1}{d_1 + d_2} \right), \text{ para datos agrupados.}$$

MEDIDAS DE TENDENCIA CENTRAL

3) Mediana: es el valor tal que ordenados los datos de menor a mayor, deja por debajo y por arriba el 50% de los datos de la distribución.

$$\text{Mediana} = Li + c \cdot \left(\frac{N}{2} - \text{facum_interv. Anteriores} \right) / \text{fintMna}$$
, para datos agrupados.

Li: límite real inferior del intervalo que contiene a la mediana. C amplitud del intervalo

MEDIDAS DE POSICIÓN NO CENTRALES

1) Cuartiles: se calculan ordenando los datos de menor a mayor y se buscan los tres valores que dividen al conjunto de datos en cuatro partes de igual cantidad de datos. Entre C_1 y C_3 se encuentra el 50% de los datos.

C_1 : el 25% de los datos son menores al valor de C_1 .

$C_1 = Li + c \cdot (N/4 - \text{facum_interv. Anteriores} / \text{fint} C_1)$, para datos agrupados.

Li : límite real inferior del intervalo que contiene al primer cuartil. C amplitud del intervalo

C_2 : el 50% de los datos son menores al valor de C_2 .

C_3 : el 75% de los datos son menores al valor de C_3 .

MEDIDAS DE POSICIÓN NO CENTRALES

2) Quintiles: se calculan ordenando los datos de menor a mayor y se buscan los cuatro valores que dividen al conjunto de datos en cinco partes de igual cantidad de datos.

Q₁: el 20% de los datos son menores al valor de Q₁.

$Q_1 = Li + c \cdot (N/5 - \text{facum_interv. Anteriores} / \text{fint} Q_1)$, para datos agrupados.

Li: límite real inferior del intervalo que contiene al primer quintil. C amplitud del intervalo

Q₂: el 40% de los datos son menores al valor de Q₂.

Q₃: el 60% de los datos son menores al valor de Q₃.

Q₄: el 80% de los datos son menores al valor de Q₄.

MEDIDAS DE POSICIÓN NO CENTRALES

3) Deciles: se calculan ordenando los datos de menor a mayor y se buscan los nueve valores que dividen al conjunto de datos en diez partes de igual cantidad de datos, $(D_1, D_2, D_3, \dots, D_9)$.

4) Percentiles: se calculan ordenando los datos de menor a mayor y se buscan los noventa y nueve valores que dividen al conjunto de datos en cien partes de igual cantidad de datos, $(P_1, P_2, P_3, \dots, P_{99})$.

$P_1 = L_i + c \cdot (N/100 - \text{facum_interv. Anteriores} / \text{fint} P_1)$, para datos agrupados.

L_i : límite real inferior del intervalo que contiene al primer Percentil. C amplitud del intervalo

MEDIDAS DE DISPERSIÓN

Las medidas de centralización proporcionan una información parcial de la distribución de la variable. El rango y los cuantiles proporcionan cierta información complementaria, sin embargo es necesario completar esa información con medidas relativas a las desviaciones de los datos respecto a las medidas de centralización.

1) La Varianza =

$$v = \frac{\sum (X_i - X_m)^2}{N}$$

$$v = \frac{\sum f_i \cdot (X_i - X_m)^2}{N}$$

para datos agrupados, X_i es la marca de clase.

2) La Desviación Típica

$$\sigma = \sqrt{\frac{\sum (X_i - X_m)^2}{N}}$$

para datos agrupados, X_i es la marca de clase.

$$\sigma = \sqrt{\frac{\sum f_i \cdot (X_i - X_m)^2}{N}}$$

3) Respecto de la mediana: Desviación intercuartílica y rango intercuartílico:

$$C = (Q_3 - Q_1) / 2$$

Cualquier distribución pueden ser caracterizados a partir de la media y la desviación típica. De esta forma se pueden comparar las diferencias de las distribuciones de una variable en dos espacios diferentes o la evolución de una misma variable en distintos momentos temporales.

* Si las distribuciones de dos series tienen igual media la que presente menor dispersión será la más conveniente(Ej.: se estudian las precipitaciones en dos zonas para escoger la más apropiada para un cultivo,

$$X_{m1} = 43\text{mm}, S_{x1} = 12.59\text{mm}$$

$$X_{m2} = 43\text{mm}, S_{x2} = 6.75\text{mm}$$

Coeficiente de variación de Pearson:

Expresa la homogeneidad de los datos, más o menos concentrados alrededor del promedio

1. Cuando se quiere comparar la variabilidad de datos que tienen distintas unidades
2. Si las distribuciones de dos series tienen distinta media la comparación lo haremos por el coeficiente de variación.

$$CV = \frac{\sum \sigma}{\bar{X}} \cdot 100$$

MEDIDAS DE FORMA

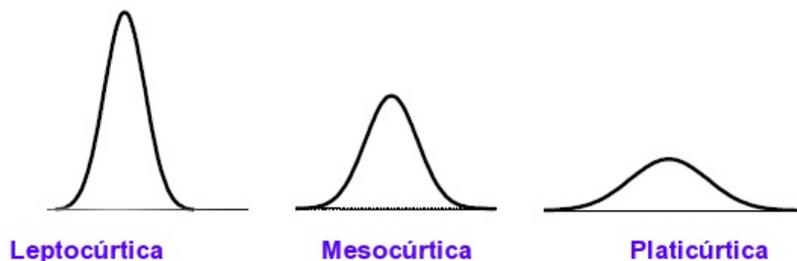
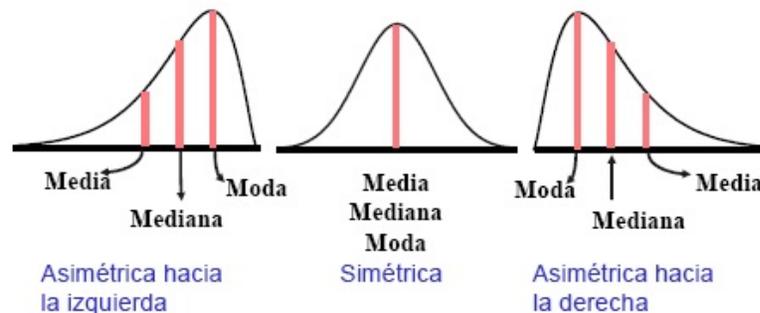
1) Para medir la asimetría de una distribución usamos unos parámetros que llamamos sesgo, cuyo valor nos da idea de la concentración de datos a un lado u otro de la media.

$$S = (x_m - M_o) / s$$

$S > 0$ sesgo derecho

$S = 0$ simétrica

$S < 0$ sesgo negativo



2) Para medir el grado de apuntamiento de la curva, si es muy aguda hacia arriba o muy aplanada.

$$K = Q / (P_{90} - P_{10}), \text{ con } Q = (Q_3 - Q_1) / 2$$

$K > 3$ Leptocúrtica

$K = 3$ Mesocúrtica

$K < 3$ Platicúrtica

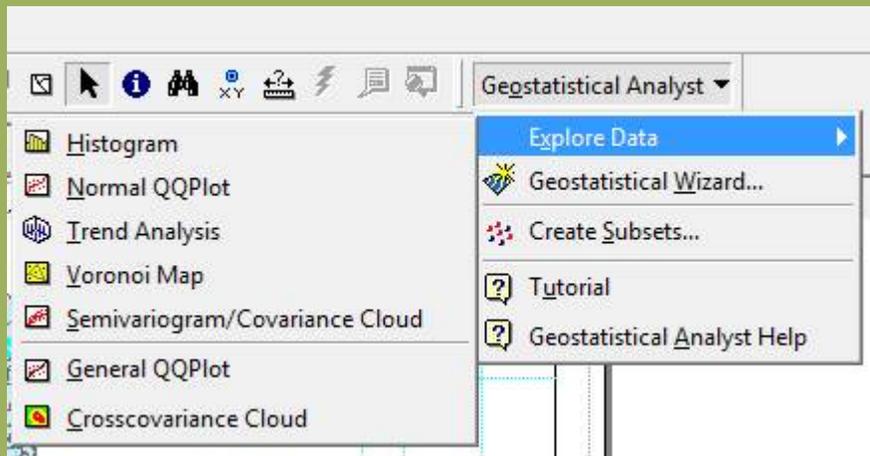
- El análisis exploratorio de datos espaciales lo podemos realizar con la extensión Geostatistical Analyst (Analista Geoestadística), del programa de Arc Gis.

Los objetivos del análisis exploratorio en general son:

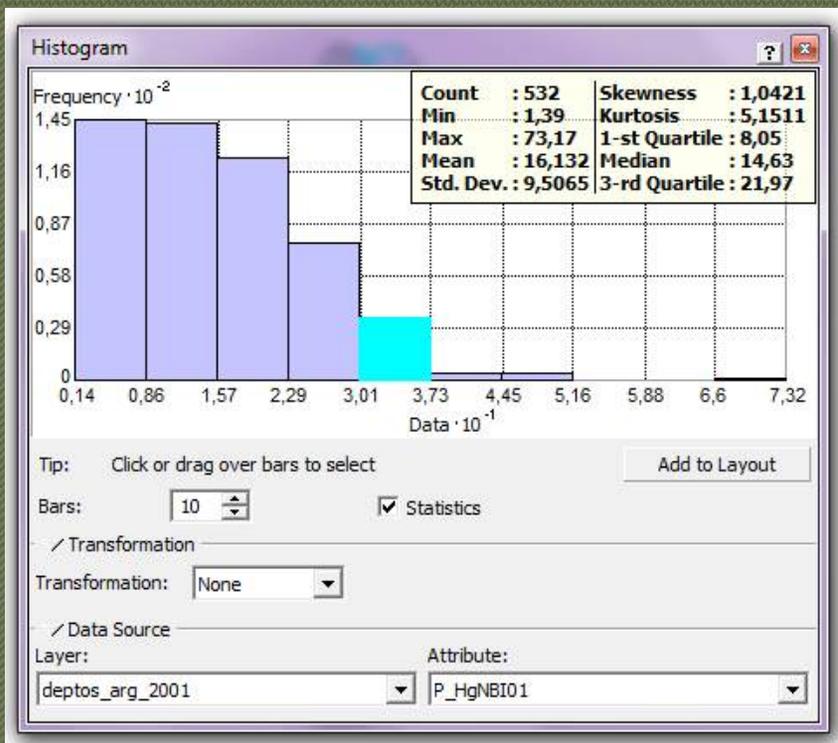
- ❖ Examinar la distribución de datos → Distribución de frecuencias
- ❖ Detectar valores extremos (outliers) → Diagrama Box –Plot
- ❖ Examinar la correlación espacial entre variables espaciales → Diagramas de Dispersión .
- ❖ Examinar la autocorrelación espacial.
- ❖ Entender la covariación entre múltiples variables.
- ❖ Realizar transformaciones de los datos si son requeridas.

HISTOGRAMA

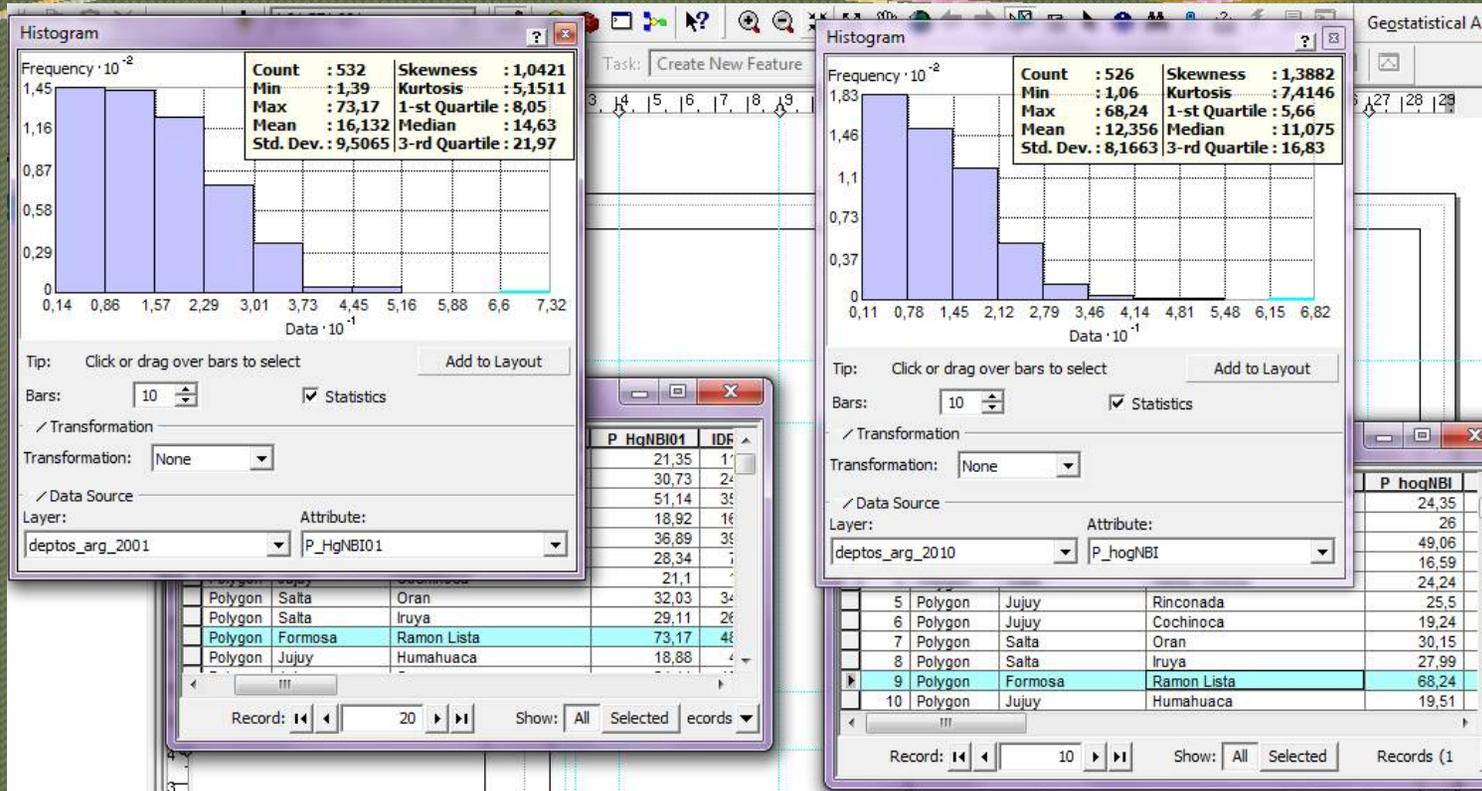
- Su obtención lo hacemos mediante el ícono que lo caracteriza.



- Permite una descripción univariada de los datos.



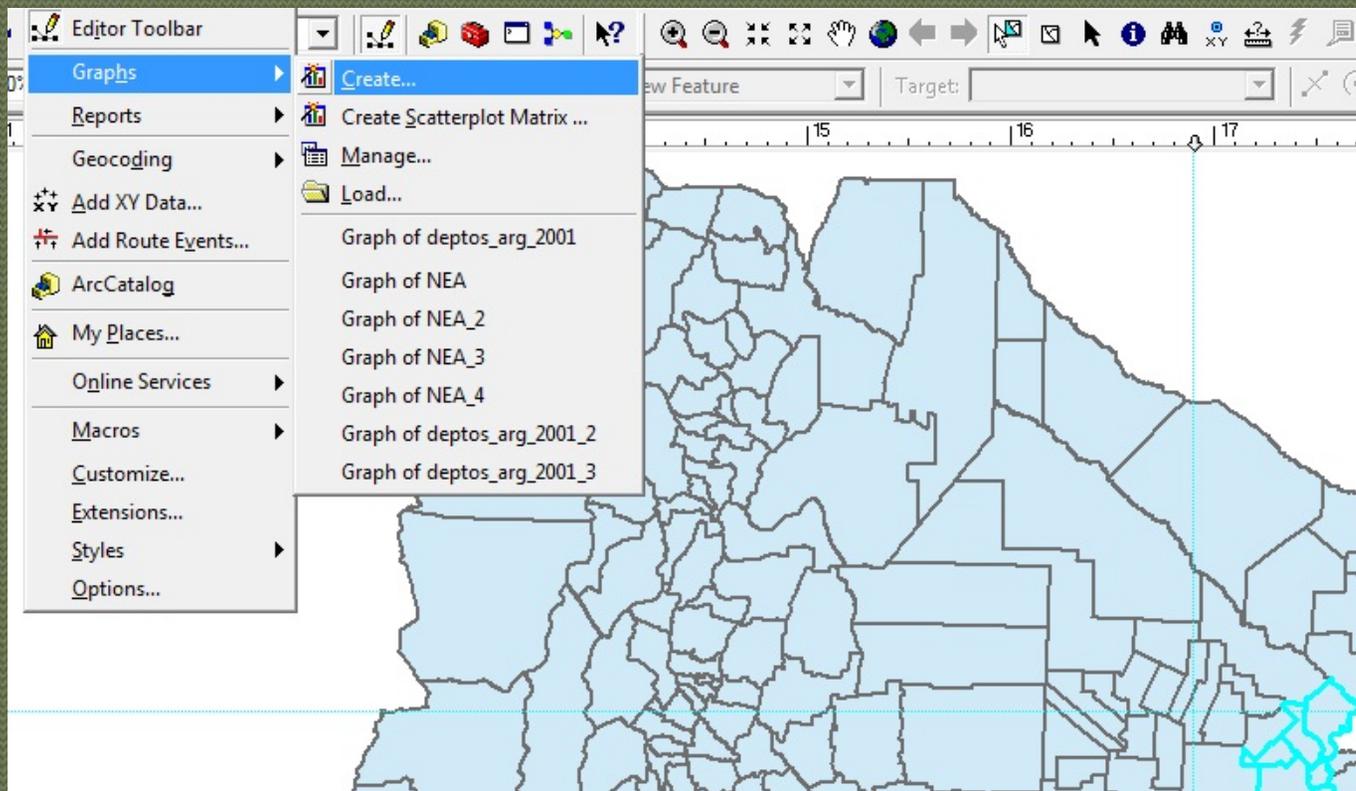
- Presenta una distribución de frecuencias y un cuadro resumen de algunos estadísticos descriptivos de la centralidad de los datos, la dispersión y la forma.
- La distribución de frecuencias se presenta mediante una serie de barras que agrupan los valores en un número de clases o intervalos. La altura de cada barra representa la frecuencia absoluta o relativa, en el ejemplo observamos la distribución de frecuencias del NBI de los departamentos de la Argentina.

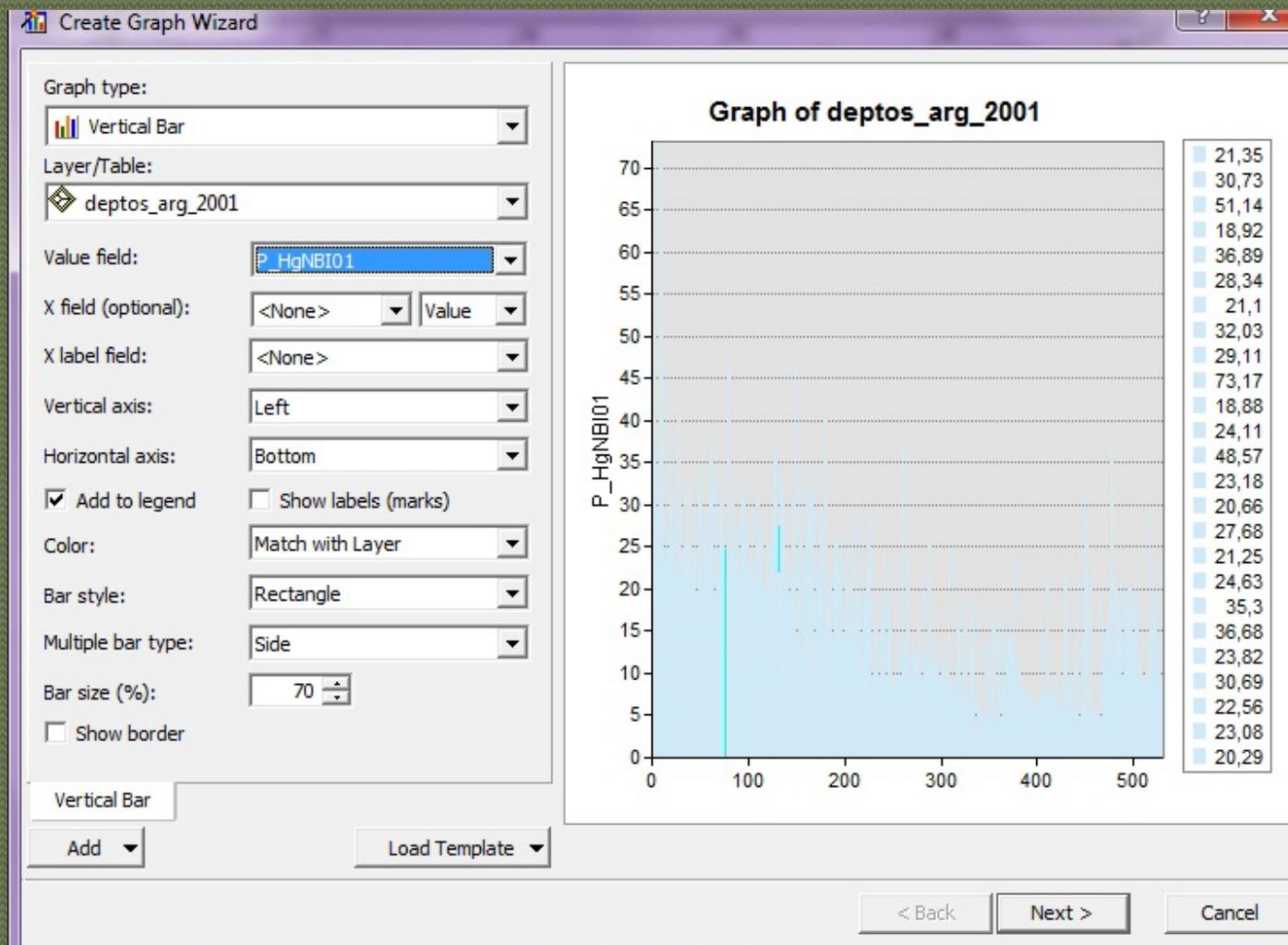


- Las medidas de centralidad nos dan idea de donde está el centro y otras partes de la distribución.
- La media es el promedio aritmético de los datos, da idea del centro de la distribución.
- La mediana representa una proporción acumulada del 0,50, el 50% de los datos está por debajo de ella y el otro 50% por arriba, también da idea del centro de la distribución.
- El primer y tercer cuartil corresponden a una proporción acumulada del 0,25 y 0,75 respectivamente.

- Como los datos están ordenados en forma creciente el 25% de los datos están por debajo del primer cuartil y el 25% de los datos están por encima del tercer cuartil.
- Las medidas de dispersión indican la mayor o menor agrupación de los datos alrededor del valor medio.
- La desviación típica que es la raíz cuadrada de la varianza (media del cuadrado de las desviaciones), describe la dispersión de los datos respecto a la media en las mismas unidades que los datos originales.
- El hecho de que la media y mediana no coincidan y el coeficiente de asimetría esté lejos del 0, evidencian la falta de normalidad de la variable. La cola derecha de la distribución indica la presencia de un punto con un NBI elevado.

GRAFICOS DE BARRA





Graph of deptos_arg_2001

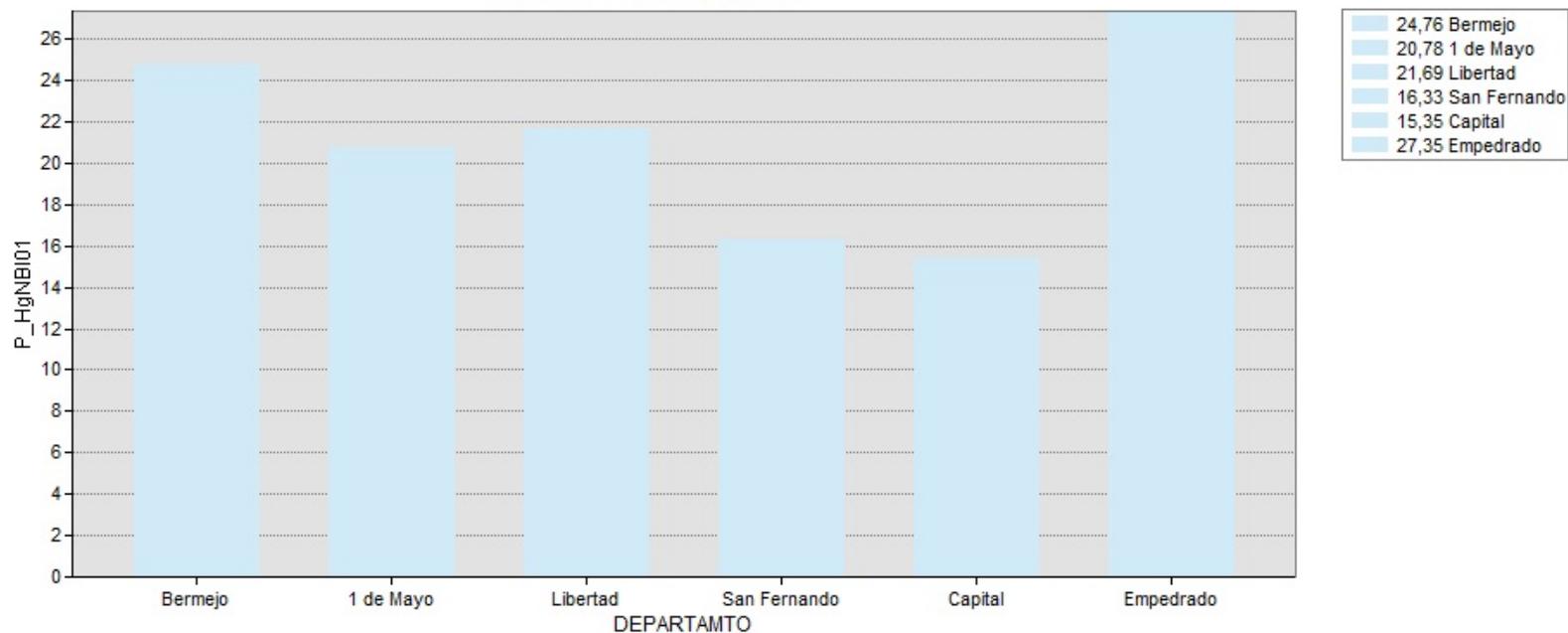
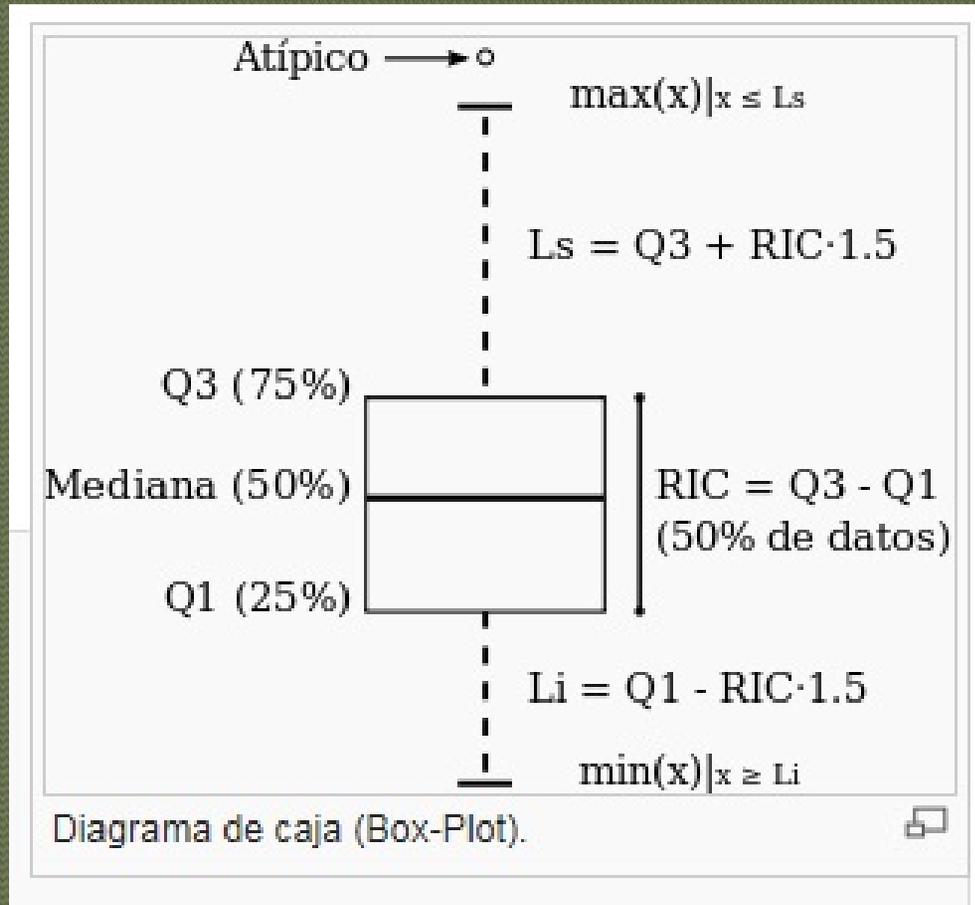


DIAGRAMA DE CAJAS O BOX PLOT



- Un Diagrama de caja es un gráfico, basado en cuartiles, mediante el cual se visualiza un conjunto de datos. Está compuesto por un rectángulo, la "caja", y dos brazos, los "bigotes".
- Es un gráfico que suministra información sobre los valores mínimo y máximo, los cuartiles Q_1 , Q_2 o mediana y Q_3 , y sobre la existencia de valores atípicos y la simetría de la distribución.
- Los bigotes, las líneas que se extienden desde la caja, se extienden hasta los valores máximo y mínimo de la serie o hasta 1.5 veces el RIC.
- Cuando los datos se extienden más allá de esto, significa que hay valores atípicos en la serie y entonces hay que calcular los límites superior e inferior, Li y Ls .
- Para ello, se consideran atípicos los valores inferiores a $Q_1 - 1.5 * RIC = 5.66 - 1.5 * (16.83 - 5.66) = -11.095$
- o superiores a $Q_3 + 1.5 * RIC = 16.83 + 1.5(16.83 - 5.66) = 33.585$
- Después se buscan los últimos valores que NO son atípicos, que serán los extremos de los bigotes. En el ejemplo NBI 2010: 1.06 y 33.33
- Marcar como atípicos todos los datos que están fuera del intervalo (Li , Ls).
- En el ejemplo: 1.06 y 33.33
- Además, se pueden considerar valores extremadamente atípicos aquellos que exceden $Q_1 - 3 * RIC = 5.66 - 3 * (16.83 - 5.66) = -27.87$
- o $Q_3 + 3 * RIC = 50.34$

DIAGRAMA DE CAJAS O BOX PLOT

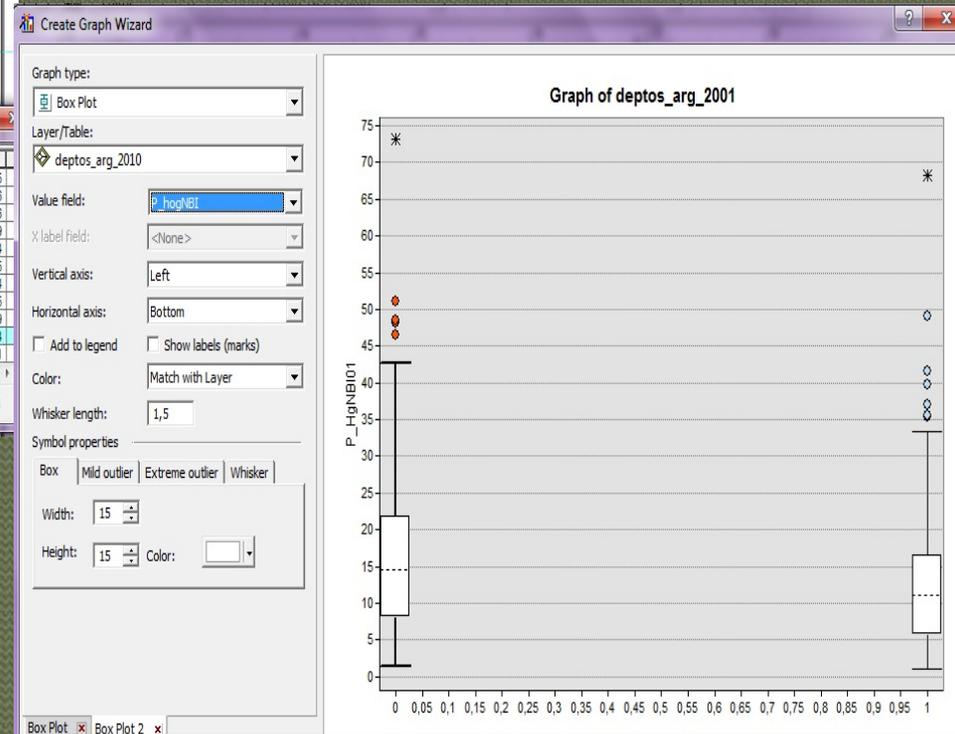
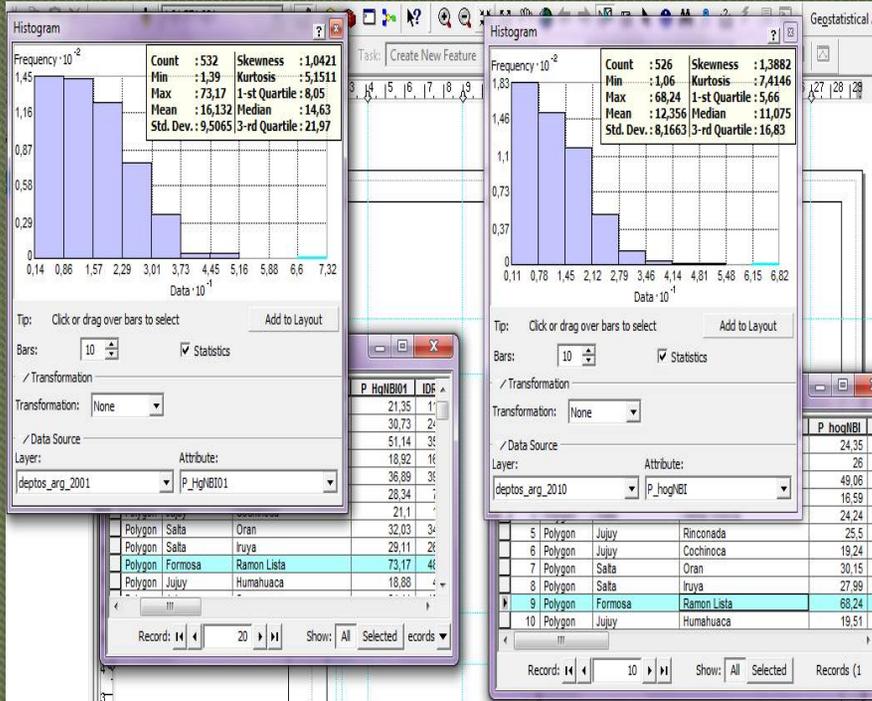


GRAFICO CIRCULAR

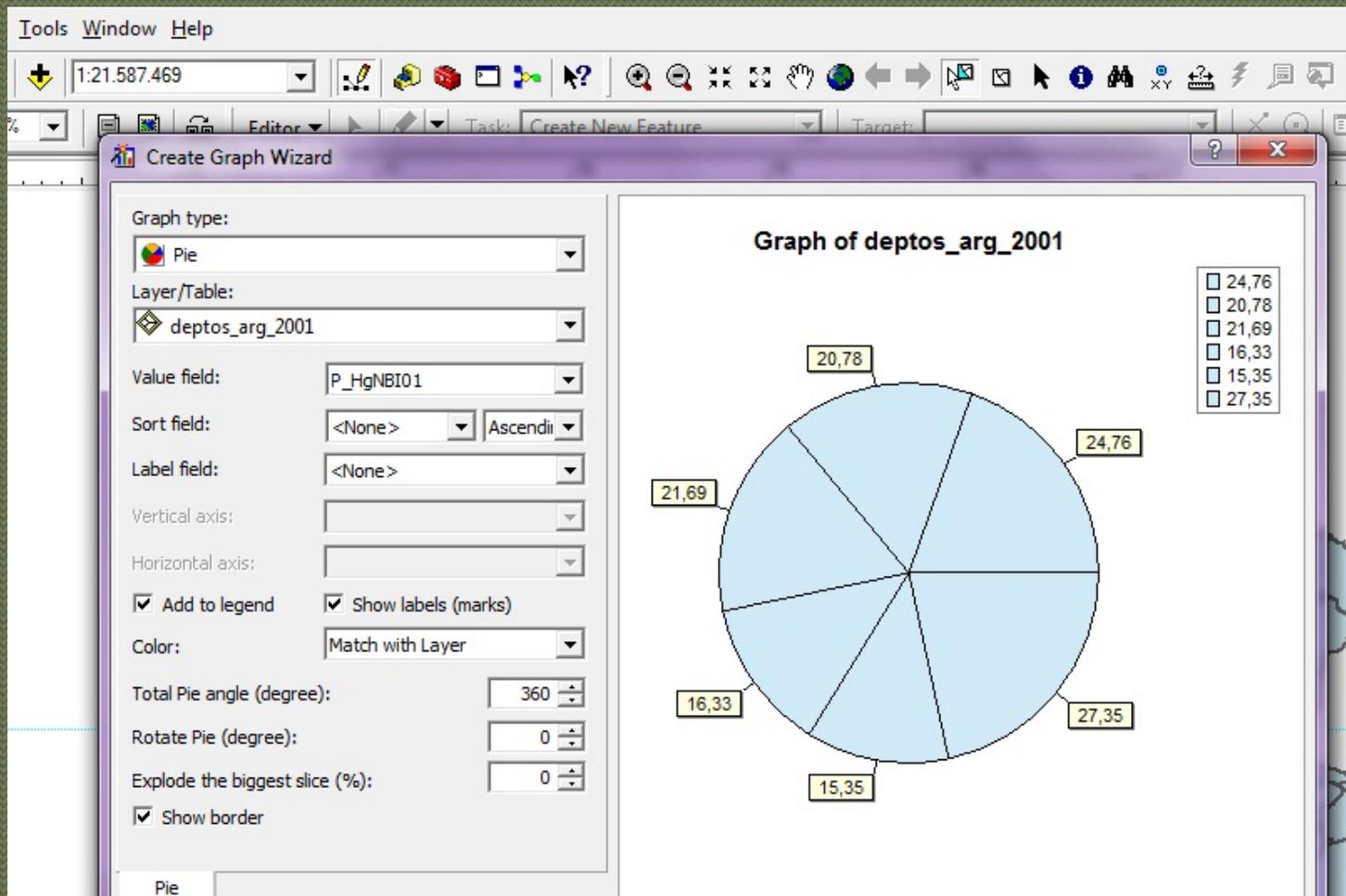
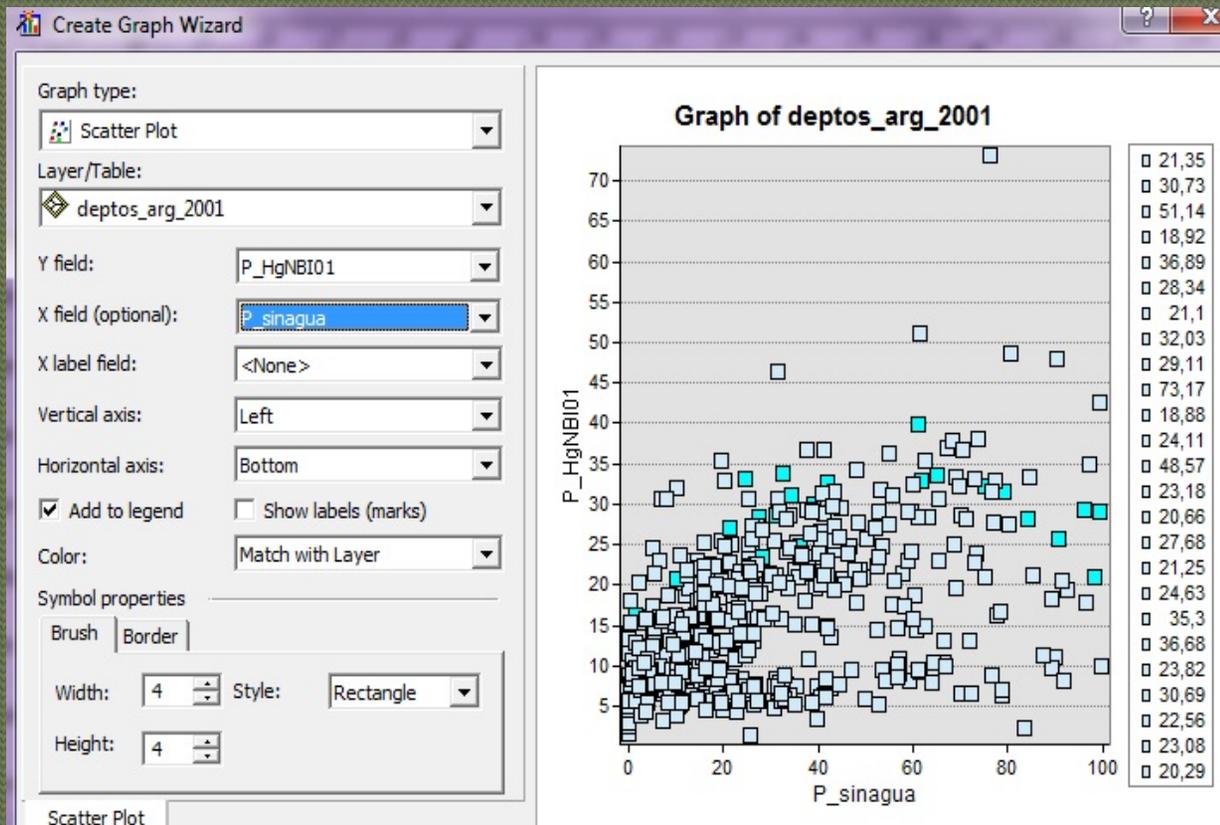
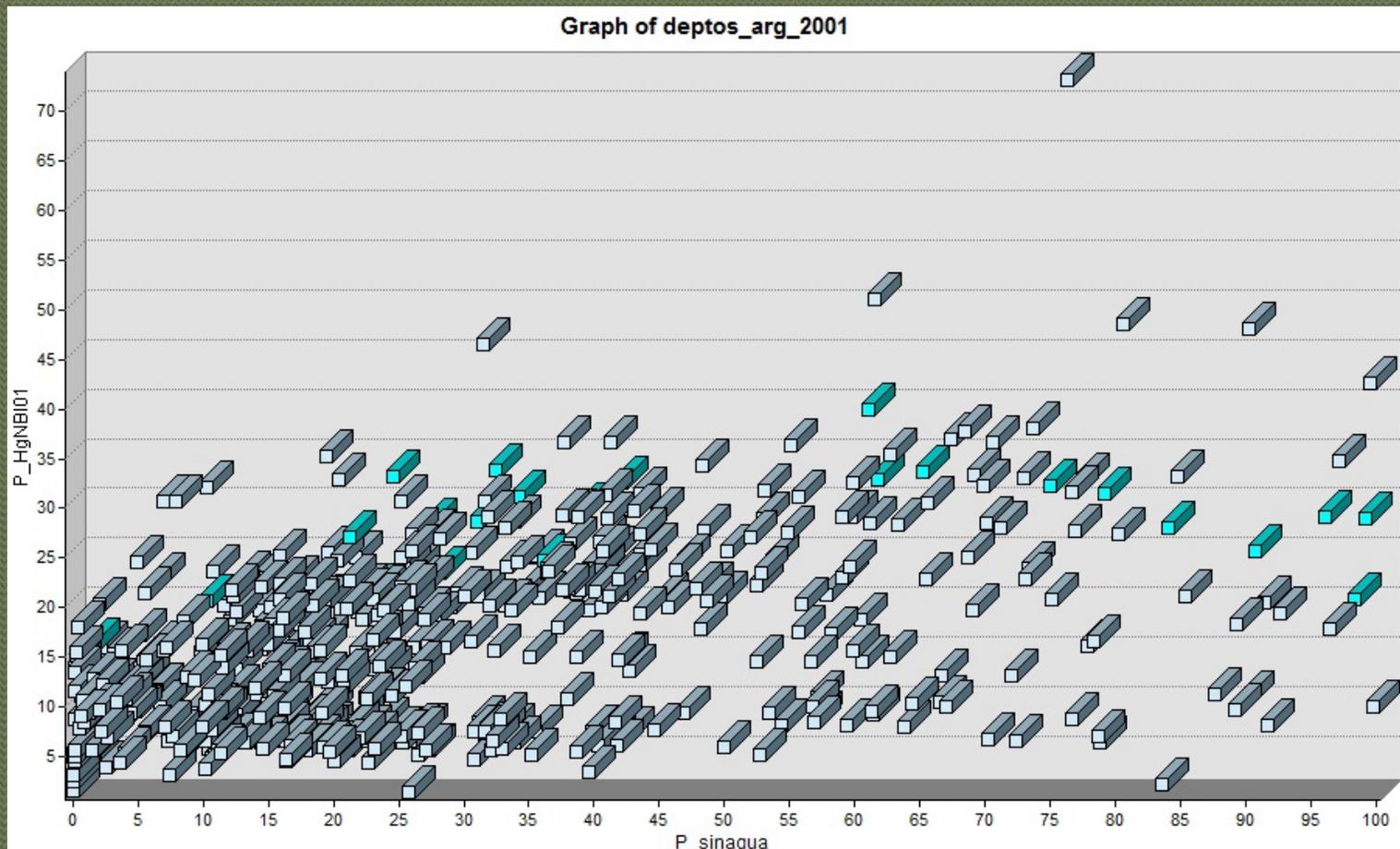


DIAGRAMA DE DISPERSIÓN

Diagrama de dispersión es la nube de puntos que resulta de la representación gráfica bidimensional de las variables analizadas, donde cada punto tiene como abscisa (X_i) el valor de la variable que se considera independiente y como ordenada (Y_i) el valor de la variable considerada dependiente.

$$r = \frac{\sum (X_i - \bar{X}) \cdot (Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2} \cdot \sqrt{\sum (Y_i - \bar{Y})^2}}$$





BIBLIOGRAFIA

- Módulos II. Exploración de datos y Estadística Espacial. Unidad 4. Conceptos de estadística. Curso e- learning de Cartografía temática avanzada. Ministerio de Fomento. Dirección General del Instituto Geográfico Nacional.
- Módulos III y IV del Curso de Posgrado en estadística aplicada a la Investigación de la UNC.
- Sistemas y Análisis de la Información Geográfica. Unidad 10.1. Análisis Exploratorio de Datos Espaciales. M.R. Cañada Torrecilla
- Nociones de Estadística y Probabilidad , Nélica Chica de Galassi –Analía Piccini – E. Rossi de Greco y Lilián Chica de Bonnet.