

Workshop Vector distribution modeling in Geographic Information System (GIS)

Quito, 19th july 2016





	19 th july 2016
H 1	TerrSet introduction (map composition) Species presence points
H2 – H3	Climatological time series analysis (deseason, PCA, STA)
H4 – H5	Maxent application
H6	Result analysis for decision support



In order to locate risk zones for disease vectors, the use of species distribution models is becoming widely recognized. This workshop is aimed to give a practical example of how to prepare input data for the model, how to set modeling parameters and how to interpret the results obtained.

- The first part of the workshop will apply geospatial mathematical operations to climatological time series (deseason, PCA, STA) for use in the Maxent model.
- After the structuration of presence points of dengue vector in the GIS, the second part of the workshop will focus on the Maxent model, searching to optimize its configuration parameters in function of validation results.
- The workshop will conclude applying GIS operators to the relative probability maps for use as decision support by the public health authorities.

Workshop introduction

- **Purpose of the exercise** : using the example of the mosquito *Aedes aegypti*, illustrate the preparation of environmental input data and the application of the maximum entropy method to model the spatial distribution of the vector in a Geographic Information System (GIS)
- Mathematical concepts used :
 - Time series analysis : Deseasoning, Principal Component Analysis, T-mode, S-mode
 - Maximum entropy

Plataforma EpiSIG

The EpiSIG core utility is in charge of:

Compiling, structuring, **analyzing**, **modeling** and presenting results related to public, particularly to provide information and tolos to the health authorities (MSP), in order to support decision taking.



INSPI-Clark Labs agreement

• INSPI-EpiSIG = one of the 19 Resource Centers of Clark Labs in the world

https://clarklabs.org/

- Benefits :
 - INSPI is provided with a TerrSet server license
 - Impulse to be the reference in Ecuador for spatiotemporal modeling using TerrSet, mainly apply to health related problems





TerrSet characteristics

- More dedicated to raster images processing
- Basic modules (horizontal)
- Integrated modules (vertical)
- Allow adding our own modules

Note : in Regional Settings, make sure to use the point as decimal symbol



Product:	
Version:	

TerrSet 18.21

IDRISI GIS Analysis IDRISI Im	age Processing Land Change M	lod
Database Query	Earth Trends Modeler Clin	nate
Mathematical Operators	•	
Distance Operators	, 💾 🖲 🚥 💿 🕀 🗶 🥆	с₽
Context Operators	•	
Statistics	•	
Decision Support	•	
Change / Time Series	•	
Surface Analysis	•	
Model Deployment Tools	•	
INSPI-EpiSIG (Ecuador, 22 may 2016)	 Gestión integrada de la salud 	
	Importar imágenes	
	Cadena epidemiológica	- 1
	Espacialización de casos	
	Conversión de fechas	
		_

0	Use	r Preferences	- • ×
Sy	stem settings Display sett	ings 3rd-Party Software Settings	
	GDAL Gdal binaries folder :		
	Gdal data folder :		
	Gdal plugin folder :		
Γ	MaxEnt		
	MaxEnt program folder :	C:\SIG\MxnUva	









Spatial **raster** entity is based on the **cell** or **pixel** unit.



Reference system used in Ecuador



Familiarization with TerrSet

<u>Exercise</u> : prepare a map composition of the elevation of Ecuador





Type of collected data	Available method
None	Multi-criteria evaluation (MCE)
Presence (1)	Mahalanobis Typicality Weighted Mahalanobis Maximum entropy
Presence / Absence (1 / 0)	Logistic regresion (Multi-Layer Perceptron)
Abundance (0,1, , n)	Multiple regression

Training data character None • Presence • Presence / Absence • Abundance Modeling approach : • MCE • MLP • MAXENT • Mahalanobis Typicality • Logistic Regression • Weighted Mahalanobis • Multiple Regression • Training site file type • XYZ • Text × XYZ • CSV Input training data file : …
O None • Presence • Presence / Absence • Abundance Modeling approach : • MCE • MLP • MAXENT • Mahalanobis Typicality O MCE • MLP • MAXENT • Mahalanobis Typicality O Logistic Regression • Weighted Mahalanobis • Multiple Regression Training site file type • Vector • Raster • XYZ • Text • XYZ • CSV Input training data file : • • •
Modeling approach : O MCE O MLP Image: MAXENT O Mahalanobis Typicality C Logistic Regression O Weighted Mahalanobis O Multiple Regression Training site file type Image: C Vector Image: Regression O Vector Image: Regression Image: C XYZ - Text Image: XYZ - CSV Input training data file : Image: Image
O MCE O MLP Image: MAXENT O Mahalanobis Typicality O Logistic Regression O Weighted Mahalanobis O Multiple Regression Training site file type Image: Comparison O XYZ - Text O XYZ - CSV O Vector Image: Regression Image: Comparison Image: Comparison Input training data file : Image: Comparison Image: Comparison
C Logistic Regression C Weighted Mahalanobis C Multiple Regression Training site file type C Vector Raster XYZ - Text XYZ - CSV ZXY - CSV Input training data file :
Training site file type O Vector • Raster O XYZ - Text O XYZ - CSV Input training data file : …
C Vector Raster C XYZ - Text C XYZ - CSV C ZXY - CSV Input training data file :
Input training data file :
Input training data file :
Retrieve settings
Environmental variables :
N. variables 1 💼 Continuous 💌
Insert laver group



Explicative variable	Measure unit	Time	Time unit	Spatial resolution	Source
Elevation SRTM	m (above sea level)	2000	1 mission	1" (~30 m)	http://earth explorer.usgs.gov
Precipitation TRMM	Intensity mm/h → mm	2010- 2015	Monthly	0.25° ~25 km	ftp://disc2.nascom.nasa.gov /ftp/data/s4pa//TRMM_L3/ TRMM_3B43/
Temperature LST MODIS	Daily mean 50×(°C+273.15)→°C	2010- 2015	Monthly	0.05° 5 km	http://e4ftl01.cr.usgs.gov/M OLT/MOD11C3.005/
Land cover MODIS	Vegetation index VI × 10000 \rightarrow [-1 1]	2010- 2015	8-days	250 m	http://e4ftl01.cr.usgs.gov/M OLT/MOD13Q1.005/
Population INEC	Density hab/km ²	2010	1 census	Census tract	http://www.ecuadorencigras .Gob.ec/banco-de- información/
Overcrowding INEC	Person / bedroom	2010	1 census	Census tract	http://www.ecuadorencigras .Gob.ec/banco-de- información/

Presence point preparation

				Idrisi [Database	Workshop			
File	Edit	Que	y Help						
	🚰 I		S 14		1 🔺 🔽] 📼 📼 💶	♦ 🖪 🗳	•	
aa_	bln	×	1 J	1					
		1	-90.405394	-0.649405				Filter	
		1	-90.338116	-0.700753					
		1	-89.60651	-0.898587	Select :	×			
		1	-89.442936	-0.901795	F				
		1	-80.906043	-1.067143	From :	Aa			
		1	-80.765051	-0.972743	Where :	[aa bln] = 1		~	
		1	-80.738526	-1.35239					
		1	-80.656382	-1.34862					
		1	-80.664702	-1.050026			. –		
		1	-80.7492	-0.9531	Order Pu	Aa1oot II		convert	Aa1ppt Llipt 🚽
		1	-80.7421	-0.9447	Older by		12		
		1	-80.63689	-2.410278					
							_ L/	project	Aa1pnt -
Note	: thi	s sa	mple of p	oints are o	only for				
exer	cise i	our	oose (inco	mplete se	et that				
nigh	t ha		ome erro	r in locatio	n)	Elv	⊢ ₩	pointras	Aa1pnt
		vc J	onic chio		5117.				
							7		
						└ <mark>→</mark> / initial		Ecd_0	
									-





Method	Purpose	Apply to deseasoned serie
Series Trend Analysis	Interannual trend	Should
STA (Seasonal Trend Analysis)	Annual seasonal progression	Νο
PCA (Principal Components Analysis) / EOF	Irregular but recurrent patterns in space/time	Should
EOT (Empirical Orthogonal Teleconnections)	Related patterns of variation between widely separated areas of the globe	Should Time-consuming
CCA (Canonical Correlation Analysis)	Requires 2 series	
Fourier PCA Spectral Analysis	Cyclical components	Should not
Linear Modeling	Requires 2 series	Should

Spatio-temporal cube preparation

Module developped by INSPI-EpiSIG (running under TerrSet)

Seleccionar también Año a partir de 6 con longitud 4

Entradas	Extensio	ón : Carpeta co	n 199 archivo	s a importar	:
TRMM mes	▼ .hdf	▼ H:\SIG\Fr	it\TRMM\Mo	nth\hdf\	
Primer archivo : 3843.19	9980101.7.HDF		Selecci precipi	onar capa : tation	•
Salida				_	_
Mes 🗾 🔽	Extraer a par	rtir del carácter	10 🌩 de	longitud 2	ŧ
Agregar : 📀 Prefijo 🛛 🕞	ProMes				
No volver a importar		🔲 MODIS tile	es		
Borrar resultados inte	ermedios				
Metadatos (de origen)		Temporalidad :		alor "sin dati	n'' -
Unidad : mm/h	•	nes 🗸	[
Sistema de georeferencia	a: "		L L L L L L L L L L L L L L L L L L L		
lationg		Min.	Max.	Cal 1440	
Resolución :	×	50	,	Lin 400	▼
0.25					-
Opciones de salida			100.100 5		1
ransponer Girar 90	a la izquierda	▼ 0-360 a	1-180+180	Lienar fail	ante
Convertir 🔲 Unidad :		T Valor	"sin dato" :		
🔲 Temporalio	dad :				
🗸 Adecuar a zona de e	studio : Alt		[.rst]		
🔽 Cambiar proyección	,			Bilineal	•
utm-17s .	ef .	Min.	Max.		
Resolución :	X: 47900		59000	Col. 680	‡
1000 👻	Y : [94300	100	174000	Lin. /44	÷
Acenta	ar	Cerrar	Ayud	a	
(Teepte					

Seasonal Trend Analysis



Two stage analytical process

- Harmonic Analysis of each year in the series to determine the best fit mean value (Amplitude 0), annual cycle (Amplitude 1 and Phase 1) and semi-annual cycle (Amplitude 2 and Phase 2)
- 2. Theil-Sen Median Slope operator to determine trends in these five parameters

Advantages

- 1. No need to identify seasonal/phenological events models each year as a totality
- 2. Removes short-term variability with less than a 6-month frequency
- 3. Removes interannual up to 30% of the length of the series.





Removal of seasonality in a time series : calculate the deviation from the mean; standardized anomalies includes the division by the standard deviation.



J O 13 A

O 14 A

203.78

173.78

143.78

113.78

83.78

53.78

23.78





Principal Component Analysis



orthogonal transformation of n-dimensional image data that produces a new set of images (components) that are uncorrelated with one another and ordered with respect to the amount of variation (information) they represent from the original image set



- With T-mode analysis, the image is the component and the graph expresses the loading

 the correlation between the component image pattern and each image in the series.
- With S-mode (Empirical Orthogonal Function analysis), the graph is the component and the image contains the loadings.

MaxEnt methodology

Let *X* be any geographic region of interest.

X:

- is a set of discrete cells (a finite set of cells)
- $x_1, \ldots, x_m \in X$, is a set of point.

It represents the locations where the species has been observed y recorded.

$$x_1, \dots, x_m \in X$$

selected independently of *X* according to an unknown probability distribution π and our objective is to estimate π .

(π coincides with the biological concept of potential distribution or fundamental niche)

MaxEnt methodology

Now the problem is one of "density estimation"

The set of points $x_1, ..., x_m$ has been chosen independently of an unknown distribution π and the purpose is to build a distribution $\hat{\pi}$ that approximates π .

For the construction of $\widehat{\pi}$

Let f_1, \dots, f_n be a set of functions where

$$f_j: X \to \mathbb{R}$$

f: represents the vector of the n functions or environmental variables

$$\forall f, \qquad f: X \to \mathbb{R}$$

 $\pi[f]$ is defined by way of

$$\widetilde{\pi}(x) = \frac{|\{1 \le i \le m: x_i = x\}|}{m}$$
, empirical distribution
 $\forall f, \widetilde{\pi}[f]$ is empirical average of f

which will be as near as possible to its true value $\pi[f]$), thus an approximation of π must be find under which

 $\hat{\pi}[f] \approx \tilde{\pi}[f]$



There are several distributions which satisfy these constraints.

The one with maximum entropy is chosen, it is the function which is closest to the uniform distribution.

Observe that the entropy of any distribution p in X is defined by:

$$H(p) = -\sum_{x \in X} p(x) \ln(p(x))$$

So the problem is:

Estimate π by $\hat{\pi}$ distribution of maximum entropy subject to the condition that

$$\forall f_j, \hat{\pi}[f_j] \approx \tilde{\pi}[f_j]$$

To solve this problem, a family of distributions functions named "Distributions Gibbs" and defined as follows is used:

$$q_{\lambda}(x) = \frac{e^{\lambda f(x)}}{Z_{\lambda}}$$



Where,

$$Z_{\lambda} = \sum_{x \in X} e^{\lambda . f(x)}$$

is a normalization constant and $\lambda \in \mathbb{R}^n$.

It can be demonstrated that the distribution of maximum entropy (MaxEnt) is the same as the Gibbs distribution maximum likelihood, this is the distribution q_{λ} that minimizes

 $RE(\tilde{\pi} \| q_{\lambda})$

Where

$$RE(p||q) = \sum_{x \in X} p(x) \ln\left(\frac{p(x)}{q(x)}\right)$$

denotes the relative entropy or Kullback-Leibler divergence.

Reference : S. J. Phillips, R. P. Anderson, and R. E. Schapire, 2006. <u>Maximum entropy</u> <u>modeling of species geographic distributions</u>. Ecological Modelling, 190:231-259.



- 75% of the presence points for training, 25% for validation
- 10 replicates (Bootstrapping, for sensibility analysis)
- To improve the AUC, it was observed that it is necessary to have a number of *background* points twice more than presence points

Notes :

- explicative rasters must have Flag value (indicate -9999 if none) and Flag description (Background);
- it is highly suggested to copy all the ecplicative rasters in the same folder and create a raster group;
- MaxEnt requires Java program installed.

Setting parameters

Maximum Entropy Parameters		
Basic Advanced Experimental		
Advanced Experimental Random Seed Give visual warnings Ask before overwriting Skip if output exists Remove duplicate presence records Write clamp grid when projecting Do MESS analysis when projecting Random test percentage 25 Regularization multiplier 1 Max number of background points 10000 Replicated run type Bootstrap Test sample file	Maximum Entropy Parameters X Basic Advanced Experimental Add all samples to background Add all samples to background Add all samples to background Maximum Entropy Vrite plot data Maximum Entropy Do clamping Maximum Entropy Vrite output grids Maximum Entropy Vrite output grids Maximum Entropy Vrite output grids Vrite output grids Vrite plots Per species results Cache ascii files Write background predictions Maximum iterations Show exponent in response curves Convergence threshold Fade by clamping Verbose Use samples with some missing data Threads Lq to lqp threshold Log file Use samples with some missing data Threads Lq to lqp threshold Hinge threshold Beta threshold Beta threshold Beta threshold Beta threshold Beta threshold	py Parameters ×
	Beta hinge	-1



Basic

Give visual warnings No Random test percentage 25 Replicates 10 Replicated run type Bootstrap

Advanced

Write output grids **No** Apply threshold rule **Equal training sensitivity and specificity**

Maxent memory usage 1024mb (2048mb can give error)

Threshold for reclassification

Logistic threshold	Description
0.253	Equal training sensitivity and specificity

Type of file to reclass — 「Image	Classifica	ation type	
C Vector C Attribute values file	C Equa	l-interval reclass	
nput file :	Aa_avg		
Dutput file :	Aa_ f sk		
Reclass parameters			
Assign a new value of	To all values from	To just less than	
0	0	0.253	
1	0.253	99999	



- How to summarize time series images for MaxEnt?
- How to set general parameters of MaxEnt optimally?
- Period to consider for time series?
- Presence points very close : useful or noisy?

Precisión	Decimal de grado	Puntos de presencia	
100 km	0	27	
10 km	1	162	
1 km	2	486	
100 m	3	1305	
10 m	4	1943	
1 m	5	2002	



• E-mails

episig.inspi@gmail.com, episig@inspi.gob.ec

• Web site of EpiSIG

http://www.investigacionsalud.gob.ec/webs/episig/

Request form online

Plataforma integrada de epidemiología, gromática, bioinformática y bioestadostica. INSPI			
Email Contraseña	2		
Iniciar Sesión			
<u>;Regístrate!</u>	©2015 EpiSIG.		





El Instituto Nacional de Investigación en Salud Pública INSPI "Dr. Leopoldo Izquieta Pérez", a través de su Plataforma integrada de epidemiología, geomática, bioinformática y bioestadística EpiSIG, invita al:

CURSO TALLER DE GEOMÁTICA "APLICACIONES INTEGRADAS DE SIG UTILIZANDO TERRSET"

3-7 octubre 2016 / 8h00 - 17h30 / Quito, Ecuador

Cupo limitado a 45 personas

Se hará una previa selección basada en la formación o empleo actual y motivación de los candidatos.

Taller sin costo

Para más información ingresa a http://www.investigacionsalud.gob.ec/webs/episig/taller/#taller3





