Carga de Datos en la Data Library del IRI

John del Corral jdcorral@iri.columbia.edu traducido del inglés por : Rémi Cousin remic@iri.columbia.edu

Septiembre 30, 2011

<u>Resumen:</u> Muchos usuarios de la *Data Library* son profesionales del sector público (salud, agricultura y manejo de agua) que tienen datos relacionados con sus sectores para sus regiones o países y desean correlacionar estos datos con datos del clima y del medio ambiente que están disponibles en la *Data Library*. Dentro de la Data Library existe una herramienta para cargar datos de diferentes sectores, lo único que se requiere es , que estos datos tengan un formato a nivel espacial y temporal, que se ajuste a la plataforma de *Data Library*.

<u>Objetivos :</u> Guiar al usuario por los diferentes pasos que son necesarios para cargar un archivo local desde su computador hasta la *Data Library* del IRI usando una interface en internet.

<u>Perspectiva general sobre la Data Library del IRI :</u> La Data Library es un servicio en la red del internet que sirve para visualizar, analizar y bajar datos geofísicos de forma gratuita . La biblioteca contiene más de 400 bases de datos con aproximadamente 300 terabytes de datos. Hay varias opciones de visualización, incluyendo la superposición de diferentes capas incluyendo distribuciones geográficas o administrativas de cada sector. La Data Library tiene diferentes funciones de análisis estadísticos. La capacidad de cargar datos dentro de esta herramienta permite correlacionar datos de usuarios con datos climáticos . Visíte la pagina de la Data Library aquí : http://iridl.ldeo.columbia.edu.

La carga de datos en la *Data Library* del IRI es un proceso complejo. No se trata simplemente de abrir un archivo con su programa preferido. El archivo de datos esta integrado, almacenado, en una base de datos, y los metadatos están almacenados en otra base de datos. Ambas bases de datos están leídas por el programa de la *Data Library* para que sean accesibles en la misma . **Se sugiere empezar con una muestra pequeña de sus datos, hasta estar más familiarizado con el proceso de carga de datos y ahí cargar todas las series de datos que requieran.**

Etapas para Cargar Datos en la Data Library :

- I. Formatos del Archivo del Usuario
- II. Formatos de los Identificadores Espaciales y Temporales
- III. Formulario de Internet para la Carga
- IV. Adicionar Metadatos
- V. Definir Dimensiones Espaciales y Temporales
- VI. Control de Calidad
- VII. Uso en la Data Library

I. Formatos del Archivo del Usuario

La *Data Library* puede cargar varios formatos generalizados de archivos de datos (tabular y SIG). El archivo que se cargue está almacenado en una tabla de base de datos. Pues, todas las herramientas de visualización y de análisis de la *Data Library* pueden ser utilizadas con esta nueva base de datos. La base de datos aparece en el directorio 'home' del usuario en la lista de bases de datos de la *Data Library*, y puede ser protegida por una contraseña.

Formatos de archivos tabularios que se aceptan para cargar en la Data Library

Excel (.xls, no .xlsx) Tab separated (.tsv) Comma separated (.csv)

Guías para incluir archivos tabularios

-una línea sola para los nombres de las columnas

-ninguna línea vacía desde la línea de los nombres de las columnas hasta el fin de los datos

-cada línea debe tener un referente espacial y temporal

-una columna o más con identificadores espaciales (código o nombre administrativo, o latitud y longitud) los cuales tienen que corresponder con los identificadores espaciales en los archivos SIG del usuario, o los archivos SIG de la Data Library

-una columna o más con identificadores temporales (año, mes, semana, o día) en formato texto, ISO 8601 estandar o Fecha Excel.

-identificadores espaciales deben ser iguales durante toda la serie temporal de datos que se cargue

- no se debe incluir ningún comentario en el archivo

-no se debe incluir ningúna suma o totales en el archivo

-no se debe duplicar ninguna línea

-No deben existir espacios en blanco , sihayun dato faltante debe ser llenados con identificadores espaciales y temporales, como -999 como valores de omision -Archivos Excel tienen un límite de ~60000 líneas. Bases de datos más grandes deben ser archivadas en un formato .tsv o .csv

-La carga de datos en el sistema solo admite archivos Excel con una tabla unica

Formatos de archivos SIG

.dbf .shp .shx

Los archivos SIG están generalmente representados por 3-5 archivos. Los requisitos mínimos son una serie de 3 archivos (.shp, shx y .dbf). La proyección necesaria para cargar en la *Data Library* es la proyección geográfica, es decir que la coordenadas de longitud y de latitud esten en grados decimales, no en metros. Es posible convertir una proyección que no es geográfica, como *Universal Transverse Mercator (UTM)*, en una proyección geográfica utilizando ArcView o ArcMap.

La serie de archivos SIG tiene una tabla de información adicional sobre cada rasgo geográfico. Esta información puede ser: nombre, populación, superficie, suma de niños de edades de 5 o menos, etc.

Guías Generales para todos archivos

- no se debe incluir ningún carácter especial como ', `, @, \$, ^, (,), -, el nombre del archivo o cualquiera entrada del archivo

-el uso del carácter guión bajo ('_') está permitido

-Nunca empiece un nombre de columna con una valor numérico, o con un nombre

si usa -números decimales ,tienen que seguir la convención del punto decimal, no la de la coma decimal.

-valores vacíos deben ser llenados con un indicador de valor vacío, como por ejemplo -999

-los nombres de los archivos no deben tener ningún blanco -longitud y latitud tienen que ser en grados decimales

La fuente de los datos tiene que ser identificada y reconocida. Por favor indique cualquier restricción del uso de los datos, y cualquier proceso de autorización necesario para el acceso de los datos.

Un *username* de *Data Library* es necesario para cargar datos en la *Data Library*. Por favor contacte <u>help@iri.columbia.edu</u> o jefft@iri.columbia.edu, si no tiene un *username*. Si ya tiene un *username* del IRI, utilícelo.

II. Formatos de los Identificadores Espaciales y Temporales

Los identificadores espaciales en archivos tabulares tienen que corresponder con los identificadores en los archivos SIG del usuario o en las bases de datos SIG de la *Data Library* del IRI. Las bases de datos SIG de la *Data Library* se pueden encontrar en : <u>http://iridl.ldeo.columbia.edu/SOURCES/.Features/</u>. Una base de datos SIG de la *Data Library* tiene generalmente identificadores espaciales numéricos (código) y alfabéticos (nombres). Este identificador de la *Data Library* debe ser añadido al archivo tabular del usuario, antes de que se realice cualquier correlación con variables climáticas .

Como ejemplos de identificadores espaciales, se tiene la división administrativa de países, regiones climáticas, distritos o municipios o áreas especificas de ciertos cultivos los cuales tienen referentes espaciales distintos.

Ejemplos de identificador espacial por jerarquia

Nombres de Divisiones Administrativas de un País

-Primer Nivel Administrativo (estados, provincias, o regiones)
-Segundo Nivel Administrativo (distritos, zonas, departamentos, o municipios)
-Tercer Nivel Administrativo (sub-distritos, sub-zonas, o municipios)
-Cuarto Nivel Administrativo (poco utilizado, pero puede incluir ciudades, o pueblos o barrios o veredas)

Codigos de Identificación asociados con una entidad

-Código de Provincia (por ejemplo 01)
-Código de Distrito (por ejemplo 0101 -- incluyendo el código de provincia, o 0001-sin incluir el código de provincia)
-Código de Sub-Distrito (por ejemplo 010101 -- incluyendo los códigos de provincia y de distrito, o 000001 -- sin incluir los códigos de provincia o de distrito)

Los identificadores temporales que recomendamos son los que siguen el estándar Internacional ISO 8601 o el formato Fecha Excel.

Ejemplos de ISO 8601 (cuando se utiliza este formato, llame la columna 'time')

-Año YYYY (por ejemplo 1997) -Mes YYYY-MM (por ejemplo 1997-07) -Día YYYY-MM-DD (por ejemplo 1997-07-16) -Semana YYYY-MM-DD/YYYY-MM-DD (por ejemplo 1997-07-16/1997-07-22)

Ejemplos Excel

-Año YYYY (por ejemplo 1997) -Mes (por ejemplo 7/1997 o Jul-1997) -Día (por ejemplo 7/16/1997 o 16-Jul-1997)

Ejemplos aceptados, pero que no son estándares:

dia	mes	ano	
3	Jan	2001	

mes	ano		
01	2001		

Nota : los nombres de las columnas no son impuestos, pero hay cueros caracteres que no se puede usar , por ejemplo no se usan tilde en 'dia' y no n en 'ano'. Las valores alfabéticas de los meses tienen que corresponder a la ortografía inglesa (por ejemplo para enero : 'Jan' y no 'Ene').

Por favor considere los ejemplos siguientes de identificadores temporales para que no tenga problemas cuando cargue su archivo en la Dta Library.

Archivo Excel (.xls)

fecha	distrito_id	distrito_nombre	casos	precipitacion
1/1997	504	Colonia	5	30
2/1997	504	Colonia	3	24

Archivos Tab Separated (.tsv)

time	distrito	_id	distrito_nom	casos	precipitacion
1997-0	1	504	Colonia	5	30
1997-0	2	504	Colonia	3	24

Archivos Comma Separated (.csv)

time,distrito_id,distrito_nombre,casos,precipitacion 1997-01-01/1997-01-07,504,Colonia,2,5 1997-01-08/1997-01-14,504,Colonia,1,8

Archivos SIG

municipios.dbf municipios.shp municipios.shx

Ahora que ya ha verificado sus sus datos están , avance en la página internet para realizar la carga de datos.

III. Formulario de Internet para la carga de datos.

El formulario de internet está disponibel en la siguiente direccion URL, <u>http://iridl.ldeo.columbia.edu/rdfconfigs/</u>. No olvide incluir la última barra oblicua ('/'). Utilice el botón 'Browse' para localizar su archivo (en su computador). Para archivos Excel y tabulares, indique un solo archivo . Para archivos SIG, indique tres archivos (.shp, .shx, and .dbf) y su *username* de *Data Library*. Asegúrese de que el nombre del archivo no tiene valores en blanco. Haga click sobre el botón 'Upload'.

🕄 Diffconfine: inday - Mozilla Firafov	
Elle Edit View History Bookmarks Tools Help	
Service Servic	ि Google
🙍 Most Visited 🐢 Getting Started <u>ଲ</u> Latest Headlines	
IRI Data Table Upload	<u>^</u>
Table File Location(with xtls, .csv, .tsv . extension): Upload Excel, Comma Separated, or Tab Separated File Reset	[Browse_]
Acceptable Tabular file formats:	
Excel (.xls, not .xlsx)	File Upload
Tab separated (.tsv)	Look in: 🔁 ISO8601 🗾 🕤 🤣 📂 🖽 -
Comma separated (.csv)	L 10 10 10 10 10 10 10 10 10 10 10 10 10
IRI GIS File Upload	Image: My Beacht Image: My Beacht Document Image: My Beacht Document Image: My Beacht
IRI Data Library username (for GIS files) :	■ bf601.tsv~ 風 bf601tsv~
CIC File I and African (_ Desktop 🗒 bf_8601_wk
GIS File Location(.dbi extension):	
CIC File Location(.snp extension):	- Witsvo
OLS File Location(.six extension)	My Documents dbf_8601_wk_span.tsv~
Lipipad GIS Files	Timeconvert.rb
	The second secon
	My Computer
General Guidelines for all files (Tabular and GIS):	5mm
-no special characters like ', `, @, $, $, `, (,), -, should appear in the file name or any of the entries in the file	
-use of the underscore ('_') character is permitted	My Network File name: Open
-do not start a column name with a numeric value	Files of type: All Files Cancel
-decimal numbers should follow the decimal point convention, not the decimal comma convention.	
-missing values should be filled as -999	
-file names should not confam any blanks	
-iongnude and lanudes should be in decimal degrees	
Guidelines Specifically for Tabular files:	
-only one line for column names	
-no blank lines after the column names until the end of data	
-each row should have a spatial location and a time	
-one or more columns with spatial identifiers (administrative id or name, or latitude and longitude) that must c	orrespond with the spatial identifiers in user's GIS files, or GIS files in the IRI Data Library
-one or more columns with a temporal identifier (year, month, week, or day) in text, the ISO 8601 standard is	format, or the excel Date format. For weekly data, 2 temporal identifiers will be needed (begindate
and enddate)	
-snatial identitiers must remain the same for the entire time snan of the data Done	<u>_</u>
John Upload to the IRI D S Rdfconfigs: index - M	J 3:21 PM

Si cargo exitosamente los datos , ahora encontrara una página 'Upload Results'. Esta página muestra los nombres de las columnas, y la suma de líneas leídas. Si un mensaje de error aparece, haga click sobre la flecha 'Atrás' para volver al formulario de internet, y buscar errores, o intentar con otro archivo.

🕲 Rdfconfigs: show - Mozilla Firefox	_ 0 ×
Elle Edit View Higtory Bookmarks Tools Help	0
🛃 👻 📓 - 🤮 🎆 🔞 http://ridi.ldeo.columbia.edu/rdfconfigs/show/	٩
🧱 Getting Started 🔯 Latest Headines	
Excel Upload Results	
Column Names	
T dane NASA GPCP V1DD prep test data	
375 Rows Were Read	
Add Metadata Back	
Lone	
🍯 Start 🗊 🕭 🖫 🛞 🐘 🛞 🖗 Adobe Illustrator - 🔯 Rdfconfigs: show - Mo	9:38 AM 💑 Tuesday

IV. Añadir Metadatos

Al final de la página 'Upload Results', hay un enlace que dice 'Add Metadata'. Haga click en este enlace para ir al formulario 'Adding Metadata'. Los tipos de metadatos necesarios están en *username*, un nombre corto para la base de datos, una descripción de la base de datos con la fuente de los datos, y el *datatype* para cada columna de datos. Indicar las unidades para cada columna de datos es algo opcional.

El *datatype* 'Excel date' es un *datatype* muy específico. Es un *datatype* utilizado por Excel para representar una fecha en un formato predefinido como '2-Jan-1997', 'Jan-1997', '1/2/1997', o '1/1997'. Un *datatype* 'Numeric' puede representar números decimales o enteros.

Un *datatype* 'Character' esta utilizado para nombres o texto, y para fechas ISO 8601. Las unidades de fechas ISO 8601 tienen que ser definidas ISO8601 también. Cuando se utiliza identificadores temporales que son estándares, utilice el *datatype* 'Character' para 'Jan' y 'Numeric' para '01' en el caso del mes de enero.

Dedfconfigs: edit - Mozilla Firefox	×
Ee Edit Yew Higtory Rookmarks Iools Help	
🕢 🕞 C 🗶 🏠 💿 http://ridl.ldeo.cokumbia.edu/rdfconfigs/1/edit 🏠 🕤 🖸 - 🖸 - 🔂 -	P
🗵 Most Visited 🏟 Getting Started 🔝 Latest Headines	
Adding Metadata to Uploaded File	<u> </u>
IRI Data Library usemame : jidcorral	
Name of Table (one word, lowercase, and no dashes (-)): [bf_example	
Short Description of Table: clean example of Malaria data from Burkino Faso with ISO8601 dates	
Tip: An 'Excel Date' will look like '2-Jan-1997', 'Jan-1997', '1/2/1997', or '1/1997' in the Excel File (e.g. 'Jan' is NOT an 'Excel Date', it is a 'Character' datatype) ISO 8601 dates are 'Character' type and must have the variable name 'time' and the 'units' of 'ISO8601'	
Missing Values will only apply to the dependent variables, not the independent variables for spatial and temporal identifiers. An example of a Missing Data value is -999	
Changeunits should be applied to temporal identifiers, when they are in ISO8601 format. The unit change will allow the dataset to be manipulated by more Data Library functions. Examples of allowed changeunits values are 'days since 1961-01-01' for daily, pentad, dekad, or weekly data, or 'monthe since 1961-01-01' for monthly or easeonal data. These changeunits are specific unit conversions used in the Data Library	
Variable: time	
datatype: Character 💌 units: ISO8601 missing_value: changeunits: deys since 1960-01-01	
Variable: name	
datatype: Character 🔽 units: missing_value: changeunits:	
Variable: cases	
datatype: Numeric 🔽 units: [cases missing_value: -999 changeunits:	
Variable: deaths	
datatype: Numeric 🔽 units: deaths missing_value: 1999 changeunits:	•
Done	
3 Start 3 Start 3 Start 3 Start 3 Start 3 Start 2 Start 2 Start 1 Start 1 Start	зу

Envíe los metadatos haciendo click sobre el botón 'Submit Metadata and add Dataset to Data Library'. Aparece en la página 'Success' si los metadatos fueron aceptados.

Esta página muestra un enlace hasta la nueva base de datos en la *Data Library*. La base de datos está localizada en la sección personal del usuario de la *Data Library* (home .username) o en la sección de los estudiantes *CIPH* en la *Data Library* (home .ciph .students .username). Si aparece un mensaje de error, haga click en la flecha 'Atrás' para volver al formulario internet, y buscar errores.



V. Definir las Dimensiones Espaciales y Temporales

La 'table diagnostics' aparece en la página de la nueva base de datos. Haga click sobre este enlace para seleccionar cuales columnas de la base de datos deben ser utilizadas para definir las dimensiones espaciales y temporales de esta base de datos. Se llaman variables independientes 'Independent Variables' de la base de datos. Seleccione una columna o más para describir la fecha y una columna o más para definir la localidad. Haga click sobre el botón 'submit changes' para aplicar estas definiciones.



Pie Edt yew Hejtory Bookmarks Tools Help Image: Pie Edt yew Hejtory Bookmarks Tools Too	4
Setting Started Market headines Trying to understand home jdcorral gridtableaddcol[test data t dane nasa gpcp vldd prcp] Independent Variables Describing time Which columns (taken together) describe time (more specific columns later)? I mone mone mone mone	
Control of the state of the st	I
Trying to understand home jdcorral gridtableaddcol[test data t dane nasa gpcp vldd prcp] Independent Variables Describing time Which columns (taken together) describe time (more specific columns later)? It none none none none specific columns later)? Undependent variables Which columns (taken together) describe locale (more specific columns later)? Idane none none none none none none none n	1
Independent Variables Describing time Which columns (taken together) describe time (more specific columns later)? I Oescribing locale Which columns (taken together) describe locale (more specific columns later)? Idane Inone Inon	
Describing time Which columns (taken together) describe time (more specific columns later)? t none none none none none none none non	
Which columns (taken together) describe time (more specific columns later)? t Describing locale Which columns (taken together) describe locale (more specific columns later)? dane none none none submit changes Preset	
t none Describing locale Which columns (taken together) describe locale (more specific columns later)? dane inone submit changes Preset	
Describing locale Which columns (taken together) describe locale (more specific columns later)? dene none none submit changes Reset	
Which columns (taken together) describe locale (more specific columns later)? dane none r none r submit changes Preset	
dane none v submit changes Reset	
submit changes Reset	
Done	//
<u>≹7</u> start] 🚮 🏟 🛱 😳 🗟	м

VI. Control de Calidad

Cuando las dimensiones espaciales y temporales están definidas, la *Data Library* evalúa si hay líneas duplicadas en la base de datos. El usuario puede ver dos resultados 1) un mensaje 'Passed duplicate line test', o 2) una lista de las líneas duplicadas. Cuando la *Data Library* hace la evaluación, el usuario ve una tabla de variables dependientes 'Dependent Variables' y sus dimensiones. La etapa siguiente es hacer click sobre el botón 'Mark dataset with dependent and independent variables'.

20 Mozilla Firefox	
Elle Edit Yjew Higtory Bookmarks Iools Help	0
🗲 - 💽 - 🥵 📓 🚮 💿 http://ridl.ldeo.columbia.edu/home//fullname/%28home%29def//name/%28home%29def//last_modified/1221164377/def/ 🔽 🕨 💽 Google	٩
🧱 Getting Started 🔂 Latest Headlines	
Trying to understand home jdcorral gridtableaddcol[dane nasa gpcp vldd prcp t test data]	
Independent Variables	
Describing time	
Which columns (taken together) describe time (more specific columns later)?	
t none v	
Describing locale	
Which columns (taken together) describe locale (more specific columns later)?	
dane none none	
submit changes Reset	
Diagnostics	
OK, lets see if I understand your dataset. Time is t. Locale is dane.	
show time show locale show duplicate lines show table for nasa_gpcp_v1dd_prcp 💌	
Passed duplicate line test Would you like a dataset with variables marked as functions of time and locale as appropriate?	I
Dependent Variables [(t)] [(dane)]	
nasa_gpcp_v1dd_prcp_X_X_X	I
test_data X X	I
Mark dataset with dependent and independent variables	
Done	
🚺 Start 🕜 🍘 🙄 🕑 🗉	 10:02 AM Tuesday

La página siguiente muestra un resumen de la base de datos. Si el resumen es satisfactorio, haga click sobre el botón 'Continue to dataset' para ver la base de datos en el directorio 'home' del usuario en la *Data Library*.

🔯 Mozilla Firefox	- O ×
Ele Edit View History Bookmanks Iools Help	0
🚰 🛩 🎆 🕤 🦉 🎆 👔 😰 http://indl.ldeo.columbia.edu/home//fullname/%28home%29def//name/%28home%29def//last_modified/1221164377/def/ 🔽 🕨 💽 Google	٩
🗱 Getting Started 🔝 Latest Headines	
Trying to understand home jdcorral gridtableaddcol[test data t dane nasa gpcp vldd prcp]	
Independent Variables	
Describing time	
Which columns (taken together) describe time (more specific columns later)?	
t none v	
Describing locale	
Which columns (taken together) describe locale (more specific columns later)?	
dane 💌 none 💌	
submit changes Reset	
Diagnostics	
OK, lets see if I understand your dataset. Time is t. Locale is dane.	
show time show locale show duplicate lines show table for nasa_gpcp_v1dd_prcp 💌	
Marking dataset	
Changes have been made <u>Continue to dataset</u>	
Done	1.
Image: Start I	10:03 AM Tuesday

VII. Uso en la Data Library

Si no utilizó fechas ISO8601 con conversión de unidades, es posible que la unidad de tiempo de la base de datos tenga que ser modificada. La mayoría de las bases de datos de la *Data Library* siguen una convención para la unidad de tiempo 'days since 1996-10-01 12:00:00' o 'months since 1960-01-01', o 'julian_day'. Si el usuario quiere correlacionar su base de datos con bases de datos en la *Data Library*, y utilizar las herramientas de visualización de la *Data Library*, la unidad de la dimensión temporal tiene que corresponder a la convención de la *Data Library*. Hay una función de la *Data Library* que se llama 'setunits'. Esta se puede utilizar para convertir la unidad temporal de la base de datos en una unidad que concuerda con la *Data Library*. Esta función puede ser utilizada después de hacer click sobre el enlace 'Expert Mode' en el esquina superior derecha de la página de la base de datos (vea el ejemplo siguiente).



Si tiene cualquier problema con la carga de datos, por favor contacte a John del Corral (jdcorral@iri.columbia.edu) o Rémi Cousin (remic@iri.columbia.edu).