

Control de Calidad de los Datos

Rémi COUSIN, DL
IRI



Objetivos

- Asegurarse de que las variables climáticas y sanitarias, y los dominios temporales y espaciales estén bien definidos
- **Validar** la calidad de las bases de datos al aplicar análisis preliminares utilizando una notación de las retículas apropiada a la Data Library del IRI

Plan

- Introducción
- Cargar un archivo Excel
- Añadir Metadatos (variables independientes y dependientes)
- Diagnosticar Problemas Comunes en los Datos (coherencia y uniformidad)
- Resumen

Datos de Malaria

MALARIA MONIT

Year	Month	Date	1-7	8-14	15-21	22-30	Total
1998	MAY	15-21/98	2	4	3	9	
1999	JANU	1-7/1999	2	4	2	8	
1998		8-14/98	2	2	1	5	
1999		1-7/99	4	3	2	9	
1998	JULY	8-14/98	4	3	3	10	
1999	MAR	1-7/99	3	2	2	7	
1998	AUGU	8-14/98	4	2	2	8	
1999	APRIL	1-7/99	5	2	2	9	
1998	SEPT.	8-14/98	4	5	5	14	
1999	MAY	1-7/99	3	4	7	14	
1998	OCTO	8-14/98	10	2	5	17	
1999	JUNE	1-7/99	2	2	2	6	
1998	NOVE	8-14/98	2	2	2	6	
1999	JULY	1-7/99	4	3	1	8	
1998	DECE	8-14/98	3	3	1	7	
1999	AUGU	1-7/99	4	2	2	8	

Cargamento

- Localizar archivo Excel o shapefile
- Explorar y exponer la estructura de la tabla
- Corregir/mejorar las descripciones de los datos

Cargar un Archivo Excel

<http://iridl.ideo.columbia.edu/rdfconfigs/>

The screenshot shows a web browser window titled "Rdfconfigs: index - Mozilla Firefox" with the address bar displaying "http://iridl.ideo.columbia.edu/rdfconfigs/". The page content is as follows:

IRI Data Table Upload

IRI Data Library username (for Tabular files) :

Table File Location(with **.xls**, **.csv**, **.tsv** . extension):

Acceptable Tabular file formats:
Excel (.xls, not .xlsx)
Tab separated (.tsv)
Comma separated (.csv)

General Guidelines for all files (Tabular and GIS):
-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 in column names and character data
-do not start a column name with a numeric value, or give it the name 'name'
-decimal numbers should follow the decimal point convention, not the decimal comma convention.
-missing values should be filled as -999
-file names should not contain any blanks
-longitude and latitudes should be in decimal degrees

IRI GIS File Upload

IRI Data Library username (for GIS files) :

GIS File Location(**.dbf** extension) :

GIS File Location(**.shp** extension) :

GIS File Location(**.shx** extension) :

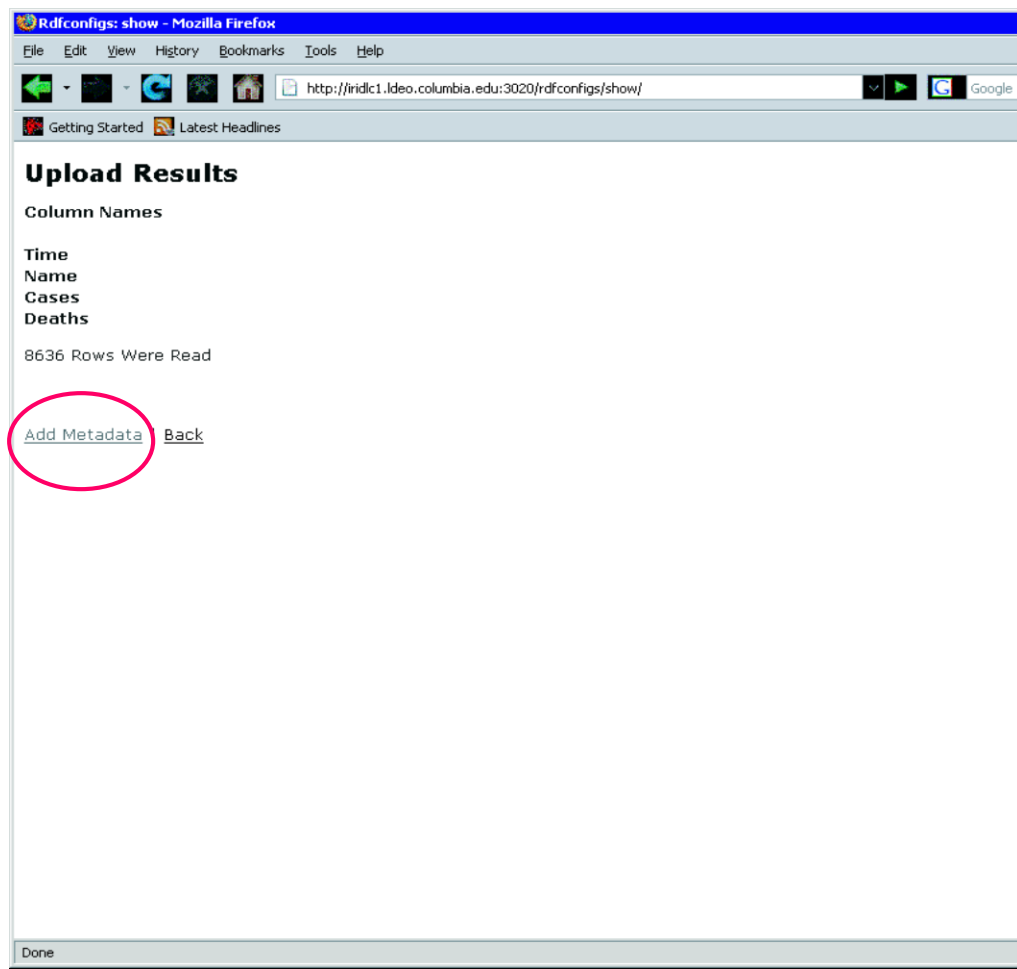
Specific Guidelines for Tabular files:
-only one line for column names
-no blank lines after the column names until the end of data
-each row should have a spatial location and a time
-one or more columns with spahal identifiers (administrative id or name, or latitude and longitude) that must correspond 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 format, or the excel Date format. For weekly data, 2 temporal identifiers will be needed

Done

The Windows taskbar at the bottom shows the Start button, several application icons, and the system tray with the time "2:50 PM" and the date "Wednesday".

Respuesta al Cargar del Archivo Excel

El botón 'Add Metadata' permite de continuar hasta la página siguiente donde se añade información a propósito de la base de datos.



Añadir más Metadatos

Adding Metadata to Uploaded File

IRI Data Library username :

Name of Table (one word, lowercase, and no dashes (-)) :

Short Description of Table:

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 'months since 1961-01-01' for monthly or seasonal data.
These changeunits are specific unit conversions used in the Data Library

Variable: Time
datatype: units: missing_value: changeunits:

Variable: District
datatype: units: missing_value: changeunits:

Variable: Cases
datatype: units: missing_value: changeunits:

Variable: Deaths
datatype: units: missing_value: changeunits:

Done

Página para Metadatos

Esta página permite añadir más información. La parte superior permite añadir una descripción de la nueva base de datos. Información para cada columna puede ser añadida también: asegurarse de que las columnas estén reconocidas como fechas o valores numerales es particularmente útil.

Variables Independientes: Tiempo y Espacio

Al final los datos serán considerados como dependientes del tiempo y del espacio (a lo menos). Seguramente el archivo Excel tiene información a propósito del tiempo y un tipo de indicación espacial (como distrito o estado).

Tiempo

Idealmente, una de las columnas indica el tiempo, con un formato de fecha Excel estándar o ISO 8601. En este caso, la columna está indicada como una variable independiente, y la herramienta de cargamento extrae una retícula de tiempo ordenada. Columnas múltiples para describir el tiempo es aceptable también (año, mes, día)

Describir el Tiempo Precisamente

- Empiezo y fin, es decir que **January 2011** indica el mes entero, o que **1 January 2011** indica el día entero
- Intervalo de tiempo (días, semanas, meses, años)
- No es siempre posible de definir precisamente
- Ejemplo: ‘weekly sea surface temperature data’

Espacio

Idealmente, una de las columnas del archivo Excel indica una entidad espacial, es decir que cada valor de los datos tiene un identificador espacial único. En este caso, la columna está indicada como una variable independiente, y la herramienta de cargamento extrae una retícula para el dominio espacial.

Variables Dependientes

Las variables dependientes son los datos que analizar, y pueden ser dadas en columnas en el archivo Excel. En este caso, la página 'Add Metadata' lista las columnas, y información descriptiva adicional puede ser añadida.

Diagnosticar Problemas Comunes en los Datos

Los nombres de las entidades espaciales deben ser coherentes con sus mismos y con sus homólogos en los shapefiles.

Este tipo de problemas puede ser detectado fácilmente después de cargar una versión preliminar de los datos, y pues corregido.

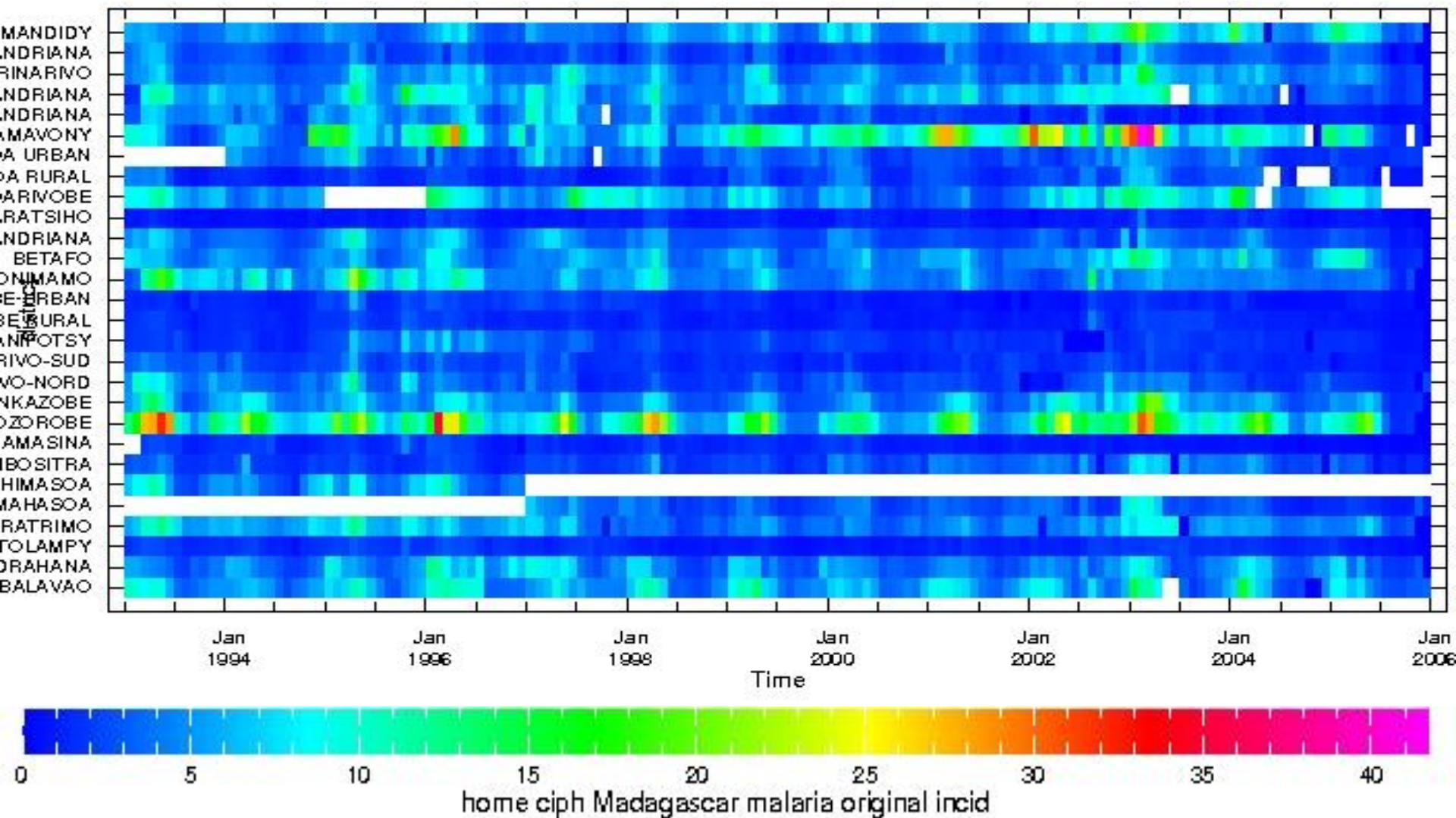
Examinación Inicial de los Datos

Por ejemplo, considere la incidencia de malaria en los Altiplanos de Madagascar para 1993-2005. Esta base de datos está accesible en la Data Library aquí:

expert

home .ciph .Madagascar .malaria .original

Examinación Gráfica Inicial de los Datos



Examinación Grafica Inicial de los Datos

El grafico muestra incidencia como color según el tiempo (eje X) y los distritos (eje Y). **Blanco** indica datos que faltan, y parece evidente que **Ambohima** y **Ambohimaso** son un mismo distrito disjunto.

Sugiere que ambos nombres representan el mismo lugar y que, a un momento, un nombre diferente fue utilizado. Al consultar el productor de estos datos, se verifica que es lo que pasó.

Examinación de la Revisión de los Datos

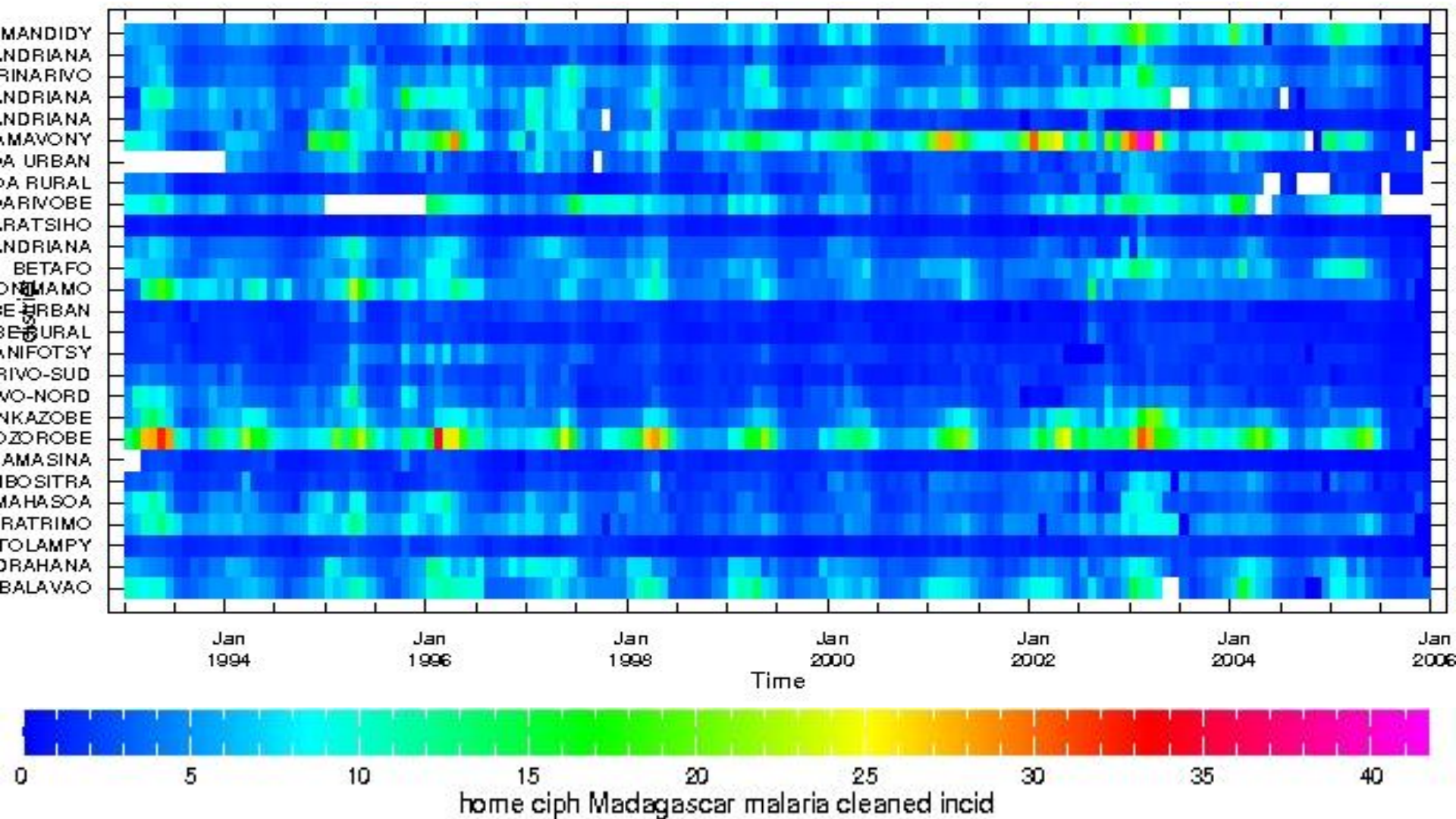
La revisión de los datos está cargada con un nuevo nombre:

expert

home .ciph .Madagascar .malaria .cleaned

Los datos que faltaban que eran sospechosos desaparecieron.

Examinación Gráfica de la Revisión de los Datos



Geo-localización

Ahora que esta coherente, se puede analizar la base de datos. Sin embargo, se quiere visualizar los datos de incidencia en su locación geográfica. Para geo-localizar los datos, explore las bases de datos de la Data Library en **SOURCES .Features.**

**SOURCES .Features .Political .Madagascar
.Districts**

use_as_grid

Una de las variables de esta base de datos es *District geometry* (the_geom), que son las geometrías de cada distrito.

En este base de datos, las variables dependen de 'codefiv', un código numeral para cada distrito. Una variable *FIV name* (nomfiv) corresponde a los nombres de los distritos de los datos cargados.

Se puede cambiar la variable the_geom tal que depende de 'nomfiv' y dar esta nueva retícula el mismo nombre que ella tiene en la base de datos cargados, es decir 'district', con la función **use_as_grid**:

**SOURCES .Features .Political .Madagascar .Districts
nomfiv /district use_as_grid**

add_variable

Por fin, hay que añadir la variable `the_geom` a la base de datos cargados utilizando la función `add_variable` para referenciarla. En **Expert Mode** se traduce:

```
home .ciph .Madagascar .malaria .cleaned  
SOURCES .Features .Political .Madagascar .Districts  
nomfiv /district use_as_grid .the_geom add_variable
```

Verificación de la Coherencia Espacial

Asegúrese que ninguno distrito falta. Para comparar los distritos de los datos de incidencia con los de la variable `the_geom`, utilice la función `SAMPLE_MISSING` que resulta en la variable `incid` restringida a los distritos que faltan en la variable `the_geom`.

`incid the_geom[district]SAMPLE_MISSING`

Resultados

Idealmente, ninguno distrito de la variable incid falta en la variable the_geom. Pero en este caso, si...
Quedan siete distritos que existen para la variable incid pero no para la variable the_geom:

**(ANTANANARIVO-NORD) (ANTANANARIVO-SUD)
(ANTSIRABE RURAL) (ANTSIRABE URBAN)
(FENOARIVOBE) (FIANARANTSOA RURAL)
(FIANARANTSOA URBAN)**

Consulte con el Productor de los Datos

Una nueva consultación con el productor de los datos nos indica la correspondencia entre los nombres que se quedan sin geometría para la variable incid y nombres de distritos en la base de datos de geometrías.

Malaria District	Feature District
Antananarivo-Nord	ANTANANARIVO-AVARADRANO
Antananarivo-Sud	ANTANANARIVO-ATSIMONDRANO
Antsirabe Rural	ANTSIRABE II
Antsirabe Urban	ANTSIRABE I
Fenoarivobe	FENOARIVO-AFOVOANY
Fianarantsoa Rural	FIANARANTSOA II
Fianarantsoa Urban	FIANARANTSOA I

Version Corregida

Una versión corregida de la base de datos está llamada **geolocated** y está verificada de nuevo con la función **SAMPLE_MISSING**

**home .ciph .Madagascar .malaria .geolocated
incid the_geom[district]SAMPLE_MISSING**

Ahora, la incidencia y la cuenta de casos dependen del tiempo y de los distritos, y las geometrías de los distritos fueron incluidas en la base de datos.

Resumen

Describir precisamente el tiempo y la referencia espacial de los datos:

- Simplifica los análisis que siguen
- Permite aplicar funciones más sofisticadas
- Permite comparar con otras bases de datos instantáneamente