Big Data & Data Science

Manejando información espacio-temporal

José Ramón Ríos Viqueira

Centro Singular de Investigación en Tecnoloxías da Información

UNIVERSIDADE DE SANTIAGO DE COMPOSTELA

citius uscles





Guion

- Introducción
- Datos espacio-temporales
- Captura de datos
- Análisis de datos espaciales
- Tecnologías
- Conclusiones



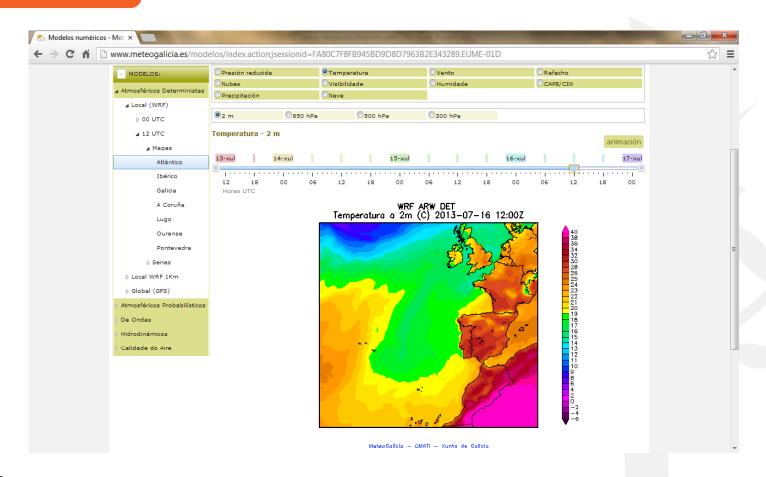
Guion

- Introducción
 - Aplicaciones
 - Motivación
 - Objetivo de la charla
- Datos espacio-temporales
- Captura de datos
- Análisis de datos espaciales
- Tecnologías
- Conclusiones



Aplicaciones

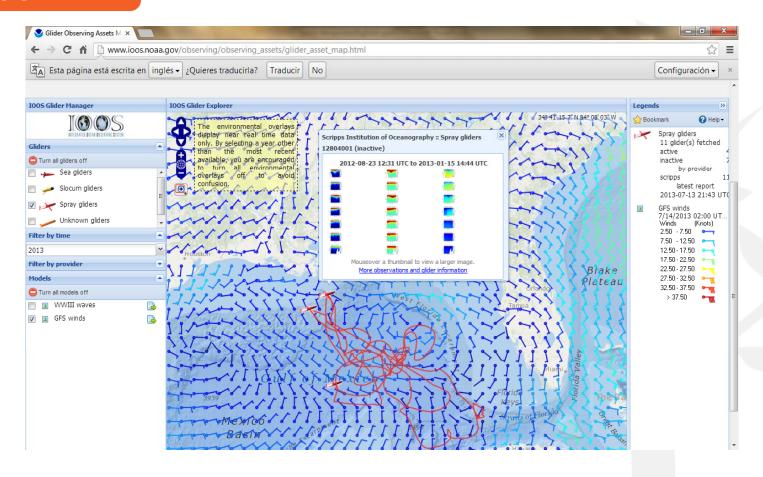
METEOROLOGÍA





Aplicaciones

OCEANOGRAFÍA

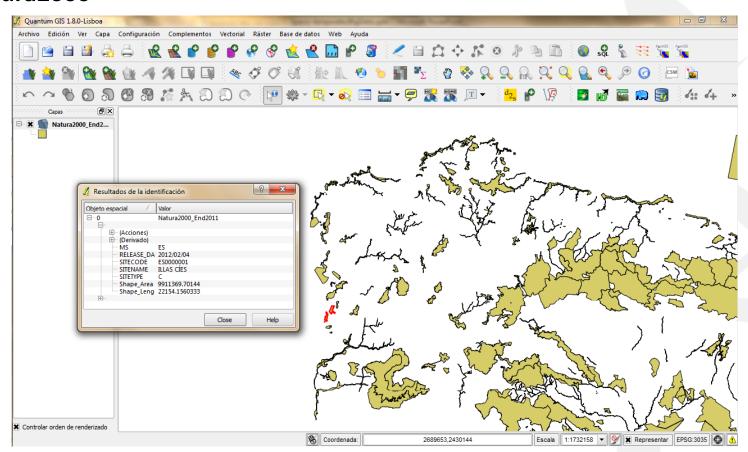




Aplicaciones

GESTIÓN MEDIOAMBIENTAL

Red Natura2000

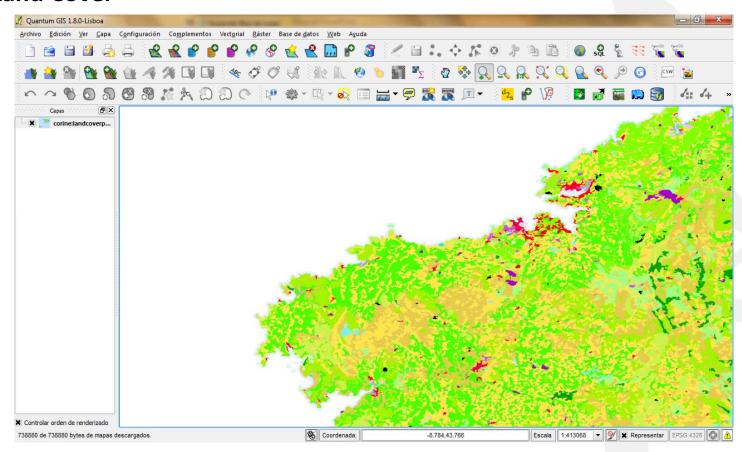




Aplicaciones

GESTIÓN MEDIOAMBIENTAL

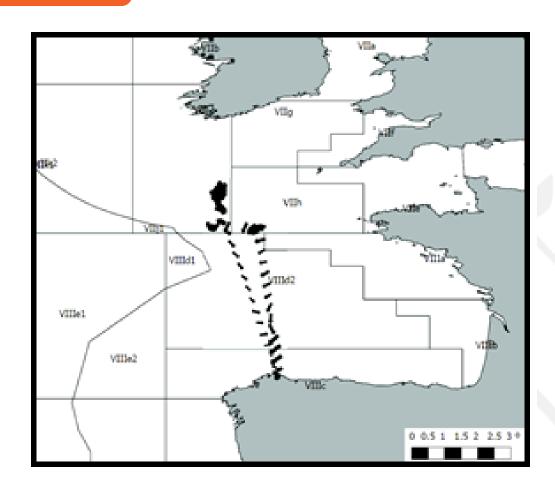
CORINE Land Cover





Aplicaciones

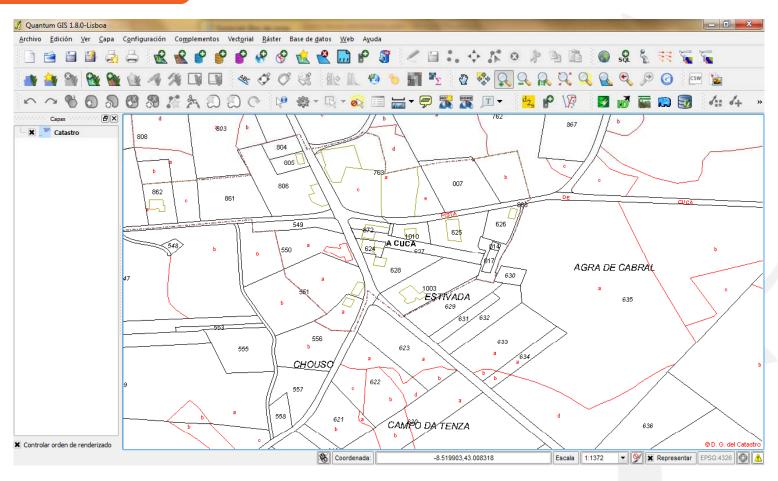
GESTIÓN DE FLOTAS





Aplicaciones

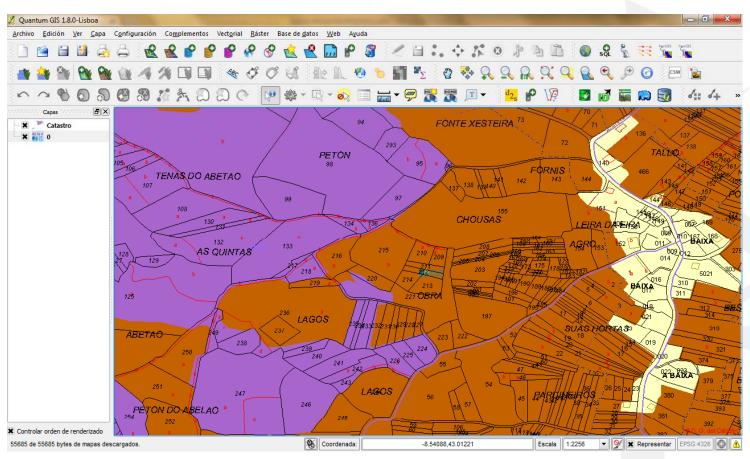
GESTIÓN CATASTRAL





Aplicaciones

AGRICULTURA

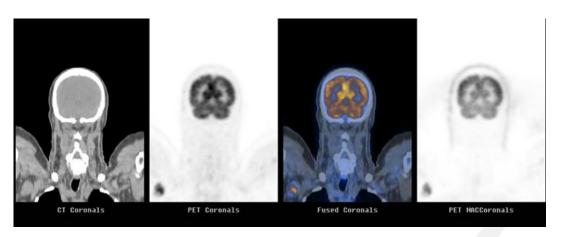


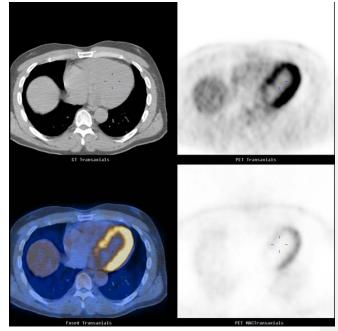


Aplicaciones

MEDICINA









Aplicaciones

ASTRONOMÍA File Edit Colors Tools Zoom Replot Help Graph coordinates: (1h48m35.44s, -0:56:14.42) Physical pixel: (1084, 825) Pixel value: 1020 (DNs) fpC-001752-u1-0042.fit_0 RA (deg) 1h48m50s 1h48m40s -1h48m30s 1h48m20s -0:57 -1:03 -1:00 -0:54 -0:51 DEC (deg) Cimus 12

Motivación

SISTEMAS DE INFORMACIÓN CONVENCIONALES

Información estructurada: Éxito rotundo del modelo relacional Información no estructurada: Modelos de recuperación de información

Estacion

ld_est	temp	dirViento	velViento	mun
est1	15	3	46	15078
est2	22	1,23	39	15078
est4	12	0,36	25	15086

Municipio

cod_mun	nombre	pob
15078	Santiago de Compostela	96000
15086	Trazo	6000

```
SELECT m.nombre AS municipio, AVG(e.temp) AS temperatura
FROM Estacion AS e, Municipio AS m
WHERE e.mun = m.cod_mun
GROUP BY m.cod_mun, m.nombre
HAVING COUNT(e.id_est) > 1
```



Motivación

DATOS TEMPORALES

Estacion

ld_est	mun
est1	15078
est2	15078
est4	15086

Medidas

ld_est	tiempo	temp	dirViento	velViento
est1	20/11/1999 15:00	15,30	1,45	22
est1	20/11/1999 15:10	15,45	1,45	22
est1	20/11/1999 15:20	16	1,56	24
			•••	••••

Obtener una estimación de la temperatura a las 15:15 del 20/11/1999 en est1



Motivación

DATOS ESPACIALES



 Id_est
 temp
 dirViento
 velViento
 utmx
 utmy

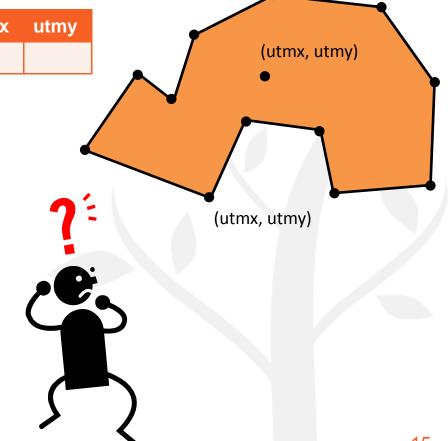
Municipio

cod_mun	nombre	pob

geo_municipio

cod_mun	num_punto	utmx	utmy

Obtener la media de temperatura en cada municipio que tenga más de dos estaciones





Motivación

DATOS ESPACIO-TEMPORALES

Estacion

ld_est	utmx	utmy

Municipio

cod_mun	nombre	pob

Medidas

ld_est	tiempo	temp	dirViento	velViento

geo_municipio

cod_mun	num_punto	utmx	utmy

Utiliza el método de interpolación espacial IDW para obtener medidas de temperatura en cada punto de la provincia con una resolución espacial de 25 metros y temporal de 10 minutos

utmx	utmy	tiempo	temp

Utiliza este resultado para obtener la evolución en el tiempo de la temperatura máxima en cada municipio





Objetivo de la charla

DATOS ESPACIO-TEMPORALES

Breve introducción a los datos y tipos de datos espaciales y espaciotemporales

CAPTURA Y ANÁLISIS DE DATOS

Breve introducción a los métodos de captura de datos y a las capacidades de análisis de las herramientas y aplicaciones

TECNOLOGÍAS

Breve introducción a las tecnologías de gestión de datos espacio-temporales existentes en el mercado. Ejemplos de algunas herramientas.



Guion

- Introducción
- Datos espacio-temporales
 - Datos espaciales
 - Datos temporales
 - Objetos móviles
- Captura de datos
- Análisis de datos espaciales
- Tecnologías
- Conclusiones

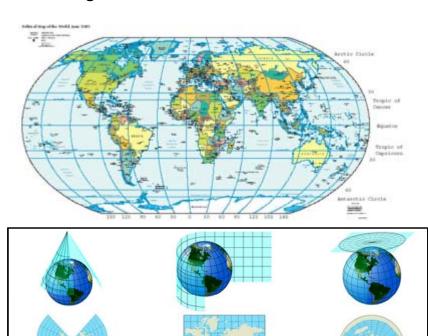


Datos Espaciales

REPRESENTACIÓN DEL ESPACIO

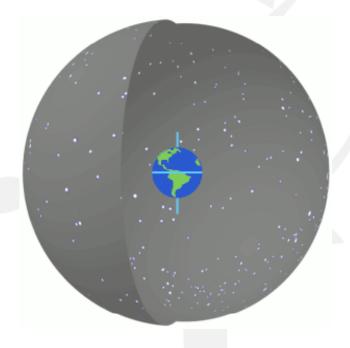
Asignar coordenadas numéricas a cada punto del espacio

Proyección Polar



Proyección Cilíndrica







Proyección Cónica

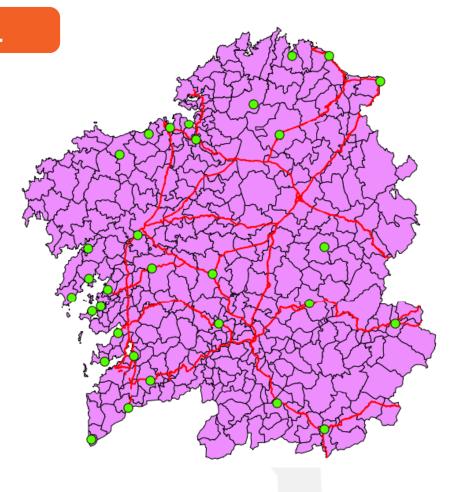
Datos Espaciales

TIPOS DE INFORMACIÓN ESPACIAL

Entidades Espaciales

- Propiedades convencionales
- Propiedades geométricas
 - Punto
 - Línea
 - Superficie
 - etc.







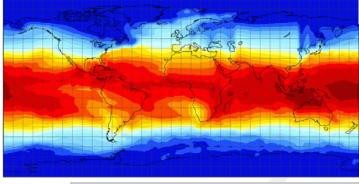
Datos Espaciales

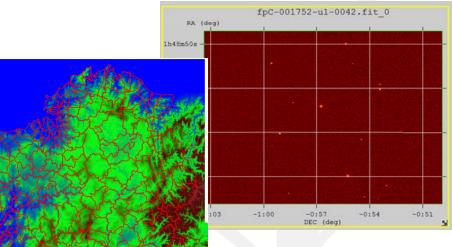
TIPOS DE INFORMACIÓN ESPACIAL

Coberturas Espaciales

- Conjuntos de funciones
- Dominio Espacial
- Rangos Convencionales
- Tipos
 - Cambio continuo
 - Cambio discreto





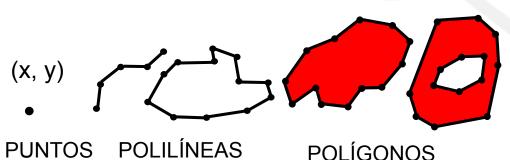




Datos Espaciales

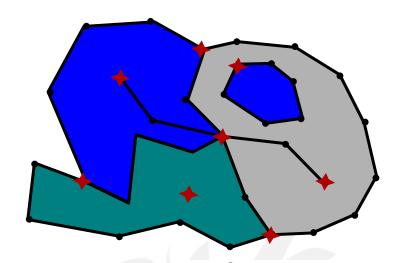
REPRESENTACIONES FINITAS

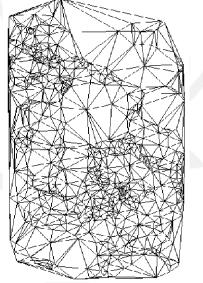
- Representaciones Vectoriales
 - **Vectorial Geométrica**
 - **Puntos**
 - Polilíneas
 - Polígonos
 - **Vectorial Topológica**
 - **Nodos**
 - Arcos (geometría polilínea)
 - Caras
 - Caso especial
 - Red Irregular de Triángulos (TIN)





POLÍGONOS

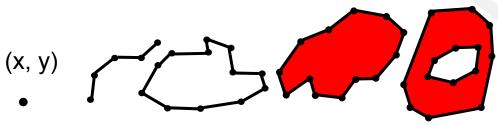




Datos Espaciales

REPRESENTACIONES FINITAS

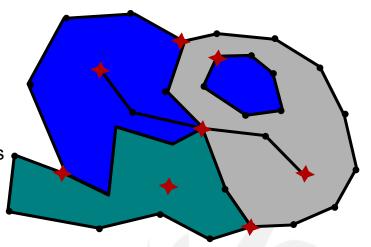
- Representaciones Vectoriales
 - Vectorial Geométrica
 - Ventaja: Representación compacta de geometrías
 - Inconvenientes
 - Redundancia
 - Análisis topológico poco eficiente
 - Vectorial Topológica
 - Ventaja: Sin redundancia. Análisis topológico eficiente
 - Inconvenientes
 - Reconstrucción topológica en inserción y borrado
 - Reconstrucción de objetos para visualización

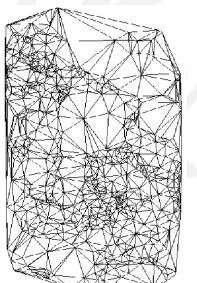




PUNTOS POLILÍNEAS

POLÍGONOS

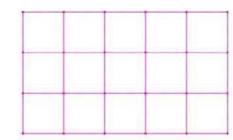


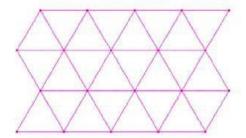


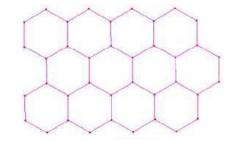
Datos Espaciales

REPRESENTACIONES FINITAS

- Mosaicos regulares
 - ▶ Triángulos, Hexágonos o Cuadrados
 - Raster
 - Celdas de forma cuadrada (píxeles)
- Usos comunes
 - Geometrías de entidades espaciales
 - Vectorial Geométrica
 - Coberturas discretas
 - Vectorial topológica
 - Coberturas continuas y discretas
 - Raster
 - Renderización 3D de coberturas continuas
 - TIN









Datos Temporales

TIEMPO DE VALIDEZ

- Tiempo durante el cual un hecho es cierto en la realidad que se está modelando
- Instante, período o conjunto de períodos
- Proporcionado por el usuario

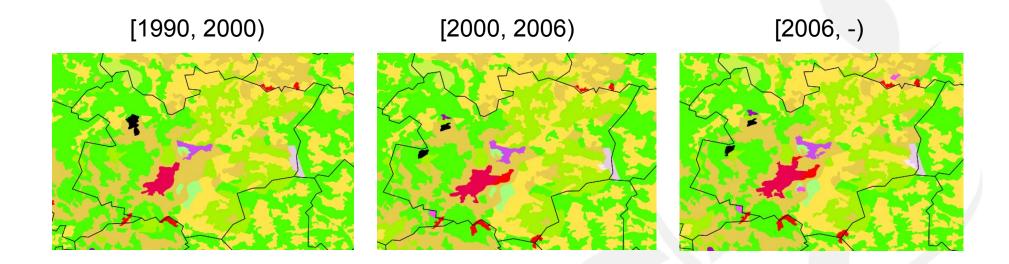
TIEMPO DE TRANSACCIÓN

- Tiempo durante el cual un hecho está almacenado en el sistema
- Período
- Actualizado por el sistema en inserciones, modificaciones y borrados



Datos Temporales

CAMBIO DISCRETO

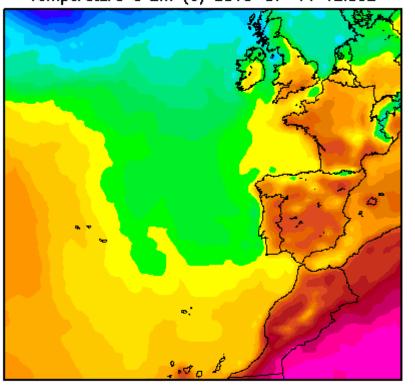


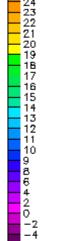


Datos Temporales

CAMBIO CONTINUO

WRF ARW DET Temperatura a 2m (C) 2013-07-14 12:00Z



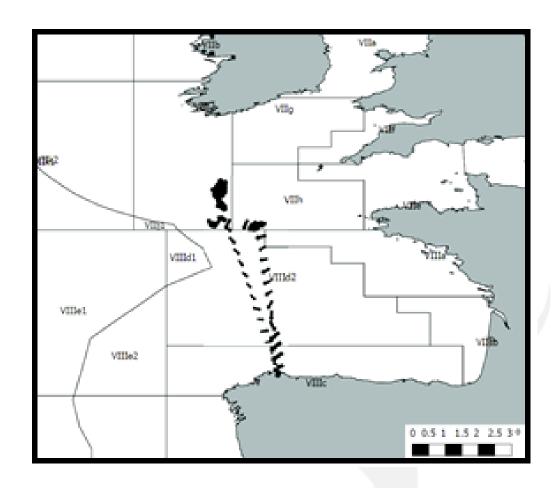






Objetos Móviles

- Cambio continuo de propiedades geométricas de entidades espaciales
- Tipos
 - Moving Points (muy comunes)
 - Moving Lines
 - Moving Regions
- Representación finita
 - Muestreo de la geometría con una determinada frecuencia temporal
 - Interpolación para instantes intermedios





Guion

- Introducción
- Datos espacio-temporales
- Captura de datos
 - Procesos de observación (OGC SensorML)
- Análisis de datos espaciales
- Tecnologías
- Conclusiones



Procesos de Observación

OBSERVACIÓN

- Proceso (Normalmente Sensor)
- Entidad observada
- Propiedad observada
- Instante temporal
- Valor observado





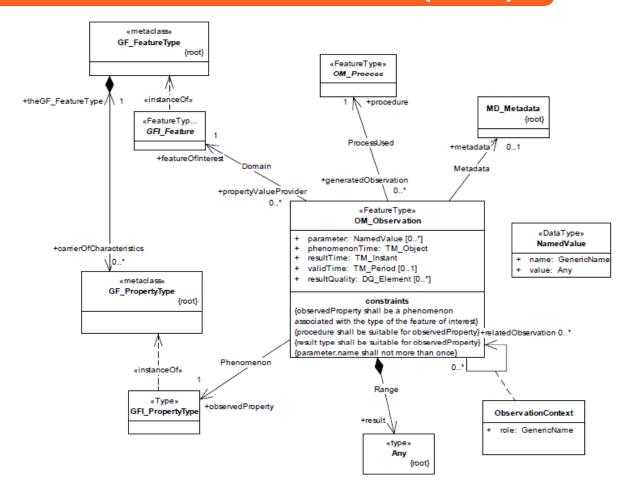






Procesos de Observación

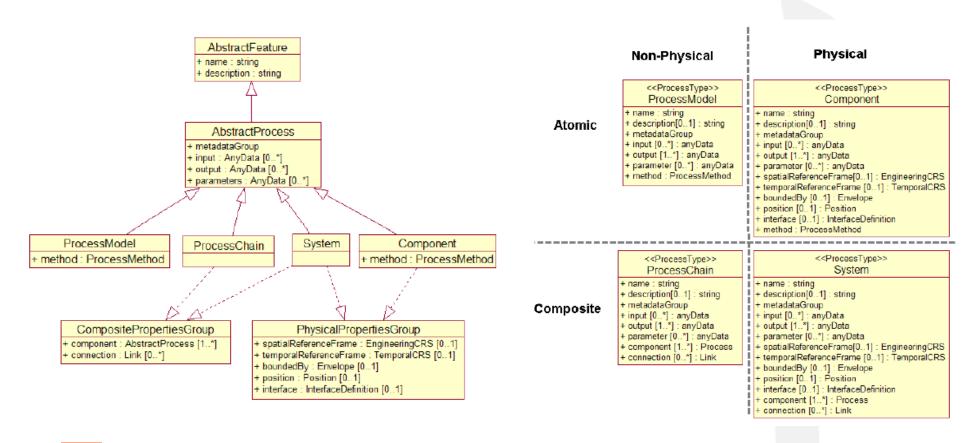
OGC - OBSERVATIONS & MEASUREMENTS (O&M)





Procesos de Observación

OGC - SENSOR MODEL LANGUAGE (SENSORML)





Procesos de Observación

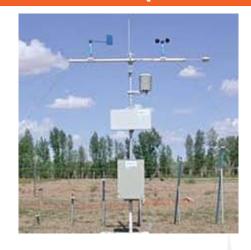
OGC - SENSOR MODEL LANGUAGE (SENSORML)

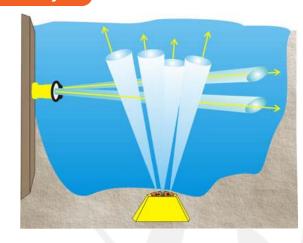
Tipos de Sensores

- Sensores In-Situ
- Sensores Remotos

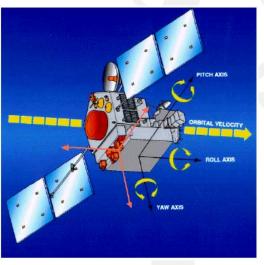
Plataformas

- Plataformas estáticas
- Plataformas móviles











Guion

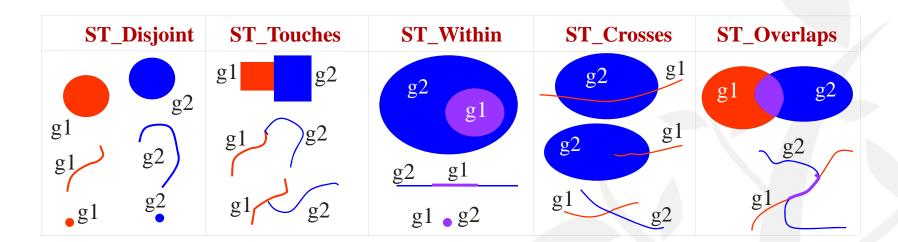
- Introducción
- Datos espacio-temporales
- Captura de datos
- Análisis de datos espaciales
 - **Colecciones de entidades espaciales**
 - **Coberturas espaciales**
- Tecnologías
- Conclusiones



Análisis de datos espaciales

Colecciones de entidades Espaciales

- Predicados topológicos
 - Calculus Based Method (CBM)



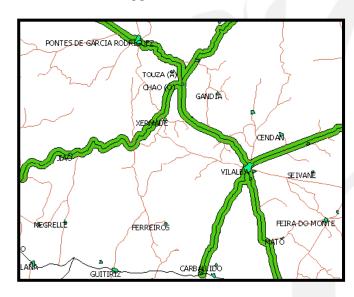


Análisis de datos espaciales

Colecciones de entidades Espaciales

- Propiedades de las geometrías
 - Tipo convencional: area, perímetro, longitud, etc.
 - Tipo espacial: mbr, borde, centroide, etc.
- Cálculo de distancias entre geometrías
- Operaciones de conjunto
 - Unión
 - Intersección
 - Diferencia
- Análisis de redes espaciales
 - Conectividad, caminos más cortos, etc.

- Rasterización
- Buffer
- Extrusión
- Convex Hull
- Voronoi
- Triangulación de Delaunay
- Etc.



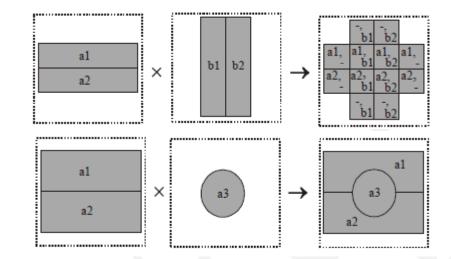


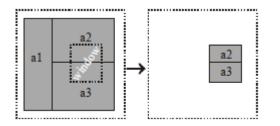
Análisis de datos espaciales

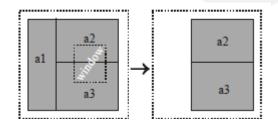
Coberturas Espaciales

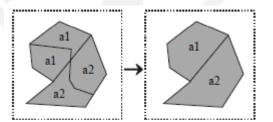
COBERTURAS DISCRETAS (VECTORIAL)

- Gestión de la topología
- Overlay
- Superimposición
- Clipping, Windowing
- Disolución de regiones
- Selección
- Cobertura











Análisis de datos espaciales

Coberturas Espaciales

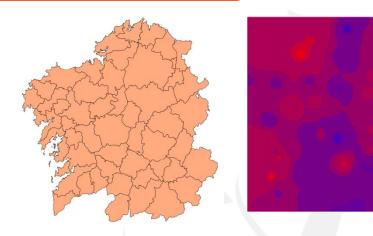
COBERTURAS CONTINUAS (RASTER)

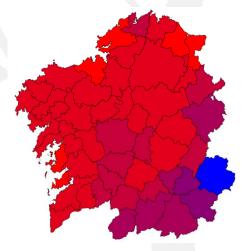
- Algebras de Mapas
 - Basadas en la aproximación de Dana Tomlin
 - Operaciones Locales
 - Operaciones Zonales
 - Operaciones Incrementales
 - Operaciones Focales













Análisis de datos espaciales

Coberturas Espaciales

COBERTURAS CONTINUAS (RASTER)

- Generalización
- Vectorización
- Análisis de terrenos
- Modelos Hidrológicos
- Generación de curvas de nivel
- Interpolaciones espaciales

- Transformadas de Fourier
- Álgebra lineal de matrices
- Funcionalidad compleja
 - Modelos meteorológicos
 - Modelos oceanográficos
- Etc.



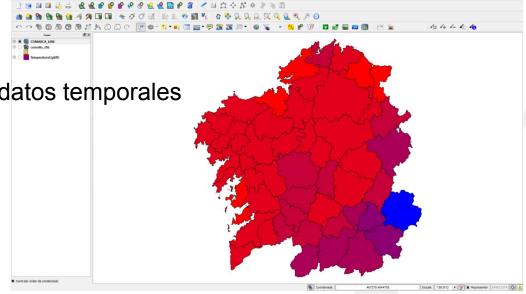
Guion

- Introducción
- Datos espacio-temporales
- Captura de datos
- Análisis de datos espaciales
- Tecnologías
 - Sistemas de Información Geográfica (SIG)
 - Bases de datos espaciales
 - Sistemas de gestión de arrays
- Conclusiones



Sistemas de Información Geográfica (SIG)

- Espacio geográfico
- Funcionalidad de análisis espacial de propósito general
- Interfaces de usuario final + entornos de programación
- No diseñadas para grandes volúmenes de datos
- Integran fuentes de datos
 - Archivos
 - Bases de datos espaciales
 - Servicios web
- Muy limitadas en la gestión de datos temporales
- Ejemplos
 - gvSIG
 - Quantum GIS
 - ArcGIS
 - Geomedia
 - Etc.

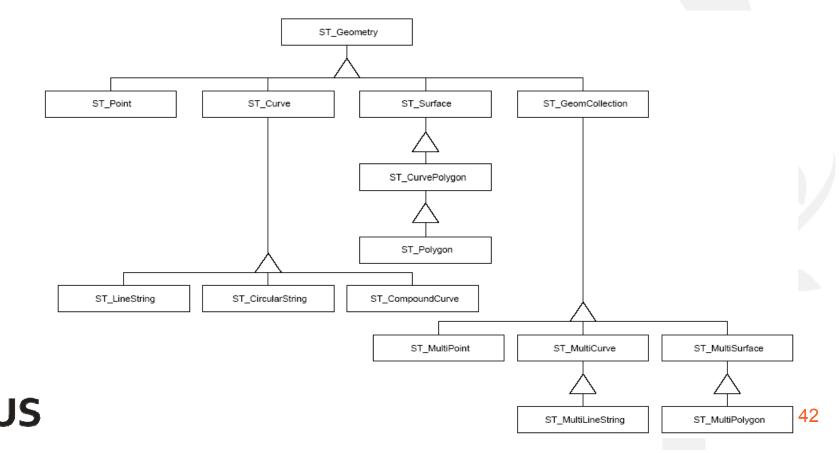




Bases de datos espaciales

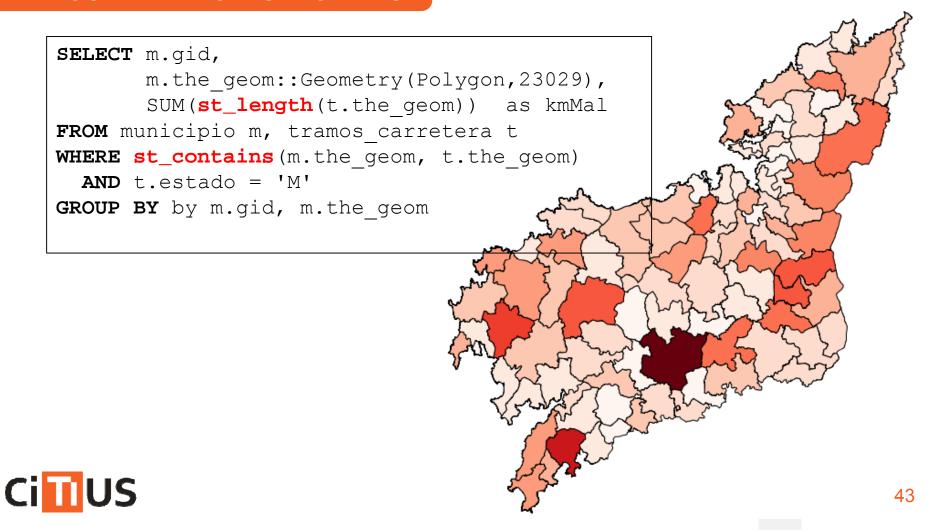
TIPOS DE DATO ESPACIALES

Tipos de datos geométricos para colecciones de entidades espaciales SQL Multimedia and Application Packages (SQL/MM) — Part 3: Spatial



Bases de datos espaciales

TIPOS DE DATO ESPACIALES



Bases de datos espaciales

TIPOS DE DATO TEMPORALES

- SQL:2003, SQL:2008
 - TimeStamp, con y sin zona horaria
 - Time, con y sin zona horaria
 - Date
 - Intervalos
 - Año, mes. Ejemplo: 3 años y 10 meses
 - Día, hora. Ejemplo: 123 días, 12 horas, 15 minutos y 5.45 segundos
- SQL:2011
 - Incluye períodos de tiempo de validez y tiempo de transacción a nivel de tupla, que forman parte de la clave (Bitemporal)
 - Predicados temporales para períodos
 - Contains, Overlaps, Precedes, etc.
 - Consulta en un instante específico o en un período de tiempo de transacción



Bases de datos espaciales

IMPLEMENTACIONES

- SGBD de almacenamiento por filas
 - Implementan estándares de ISO
 - Parte espacial del SQL/MM
 - Pronto parte temporal del SQL:2011
 - Oracle Spatial (http://www.oracle.com/es/products/database/options/spatial/index.html)
 - ▶ IBM DB2 Spatial Extender (<u>http://www-03.ibm.com/software/products/us/en/db2spaext/</u>)
 - PostgreSQL + PostGIS (http://www.postgis.org/)
 - Algunos incluyen tipos de datos para almacenar Raster
 - Oracle, PostGIS
- SGBD de almacenamiento por columnas
 - Extensión geospacial de MonetDB (http://www.monetdb.org/Documentation/Extensions/GIS)
 - Implementa la especificación OGC Simple Features for SQL



Bases de datos espaciales

IMPLEMENTACIONES

- Bases de datos NoSQL
 - Implementaciones todavía muy limitadas en funcionalidad
 - Gestión de documentos GeoJSON (http://geojson.org/)
 - Formato para codificar en JSON colecciones de entidades espaciales con propiedades de tipos geométricos
 - MongoDB (http://www.mongodb.org/)
 - Geocouch (https://github.com/couchbase/geocouch/)
 - Seguramente pronto veremos más
 - Cassandra
 - voltDB (new SQL)



Sistemas de gestión de arrays

ARCHIVOS

- Formatos para imagen geográfica
 - GeoTIFF (formato TIFF con cabecera geográfica)
 - ecw (Enhanced Compressed Wavelet)
- Arrays e imágenes en entornos científicos (Ejemplos)
 - Flexible Image Transport System (FITS)
 - Hierarchical Data Format (HDF4, HDF5)
 - GRIdded Binary or General Regularly-distributed Information in Binary form (GRIB)
 - Network Common Data Form (NetCDF)
- BIG Data?
 - - Implementaciones ad-hoc



Sistemas de gestión de arrays

BASES DE DATOS OBJETO-RELACIONALES (SQL:2003)

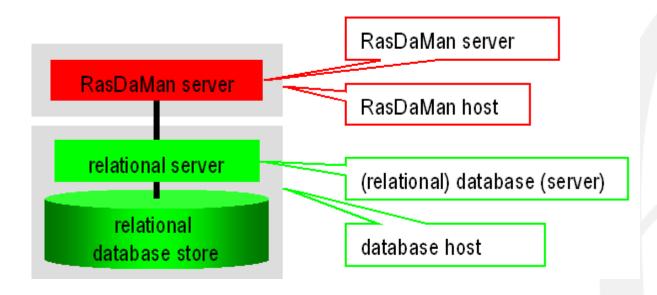
- Tipos de dato ARRAY
- Permiten anidar arrays dentro de tuplas del modelo relacional
- No diseñado para gestionar los arrays de gran tamaño que aparecen en aplicaciones de tipo científico

```
CREATE TABLE Empleado (
  id_emp INTEGER PRIMARY KEY,
  nombre VARCHAR(50),
  salarios DECIMAL(8,2) ARRAY[12]
)
```



Sistemas de gestión de arrays

- Gestor de arrays multidimensionales
- Implementado sobre un gestor relacional
 - Típicamente PostgreSQL
 - Arrays almacenados en tipos BLOB





Sistemas de gestión de arrays

- Lenguaje de consulta declarativo sobre arrays (rasql)
 - Lenguaje de definición de arrays (rasdl)
 - Definición del tipo de las celdas del array
 - Tipos simples y Tipos complejos (struct)
 - Definición de arrays sobre un tipo de celda
 - Definición de colecciones de arrays del mismo tipo

```
struct RGBPixel {char red, green, blue; };

typedef marray <RGBPixel,[0:799, 0:599]> RGBImage;

typedef set <RGBImage> RGBSet;
```



Sistemas de gestión de arrays

- Lenguaje de consulta declarativo sobre arrays (rasql)
 - Lenguaje de manipulación de arrays (rasml)
 - Similar a SQL.

```
select RGBSet[ 120:160, 55:75 ]
from RGBSet
```

```
select RGBSet
from RGBSet
where all_cells(RGBSet.green > 20)
```

```
select RGBSet.red * 2
from RGBSet
```

```
select png(RGBSet)
from RGBSet
```



Sistemas de gestión de arrays

- Lenguaje de consulta declarativo sobre arrays (rasql)
 - Lenguaje de manipulación de arrays (rasml)
 - Similar a SQL.

```
select RGBSet[ 120:160, 55:75 ]
from RGBSet
```

```
select RGBSet.red * 2
from RGBSet
```

```
select RGBSet
from RGBSet
where all_cells(RGBSet.green > 20)
```

```
select png(RGBSet)
from RGBSet
```



Sistemas de gestión de arrays

SciDB(http://www.scidb.org/)

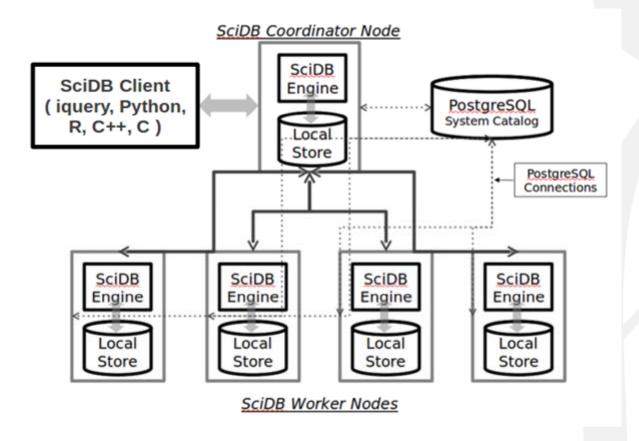
- Optimizado para Big Data y Big Analytics
- Arrays y Vectores son objetos de primera clase con operadores específicamente optimizadas para ellos.
- Arquitectura paralela de tipo share-nothing
- Extensible: Tipos y funciones definidos por el usuario
- Interacción con el paquete estadístico R
- Multiversión (tiempo de transacción)
- Comparativa entre Hadoop + HDFS y SciDB realizada por la NASA
 - Respective Strengths and Weaknesses of SciDB, MapReduce-HDFS, and a Custom Technique for a Data-Intensive Analysis System
 - https://ams.confex.com/ams/93Annual/flvgateway.cgi/id/23219?recordingid=23219
 - Cálculo de agregados sobre 6GB de datos
 - Hadoop: 180 segundos. 89 segundos después de optimizar
 - SciDB: 39 segundos out-of-box



Sistemas de gestión de arrays

SciDB(http://www.scidb.org/)

Arquitectura

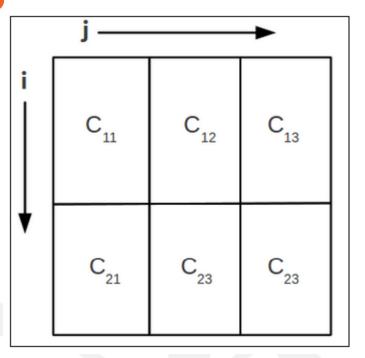




Sistemas de gestión de arrays

SciDB(http://www.scidb.org/)

- Modelo de datos de arrays
 - Array
 - Nombre
 - Dimensiones
 - Atributos
- Chunking
 - Cada dimensión se divide en Chunks
 - Dimensión i con tamaño 10 y tamaño de chunk 5
 - Dimensión j con tamaño 30 y tamaño de chunk 10
 - Los chunks se distribuyen por los nodos usando un esquema basado en hashing
 - Solapamiento de los chuncks
 - $_$ A <a: int32> [i=1:10,5,1, j=1:30,10,5]
 - Mejora rendimiento en operaciones de vecindario

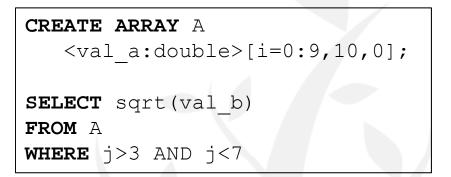




Sistemas de gestión de arrays

SciDB(http://www.scidb.org/)

- Almacenamiento
 - Dentro de cada chunk se almacena por atributos (particionamiento vertical)
 - Compresion con Run LengthEncoding (RLE) en cada atributo
 - Los datos nunca se sobrescriben, se generan nuevas versiones
- Catalogo almacenado en PostgreSQL
- Lenguajes
 - Array Query Language (AQL)
 - DDL + DML
 - Array Functional Language (AFL)
 - Operadores sobre arrays que incluyen Algebra Lineal de vectores y matrices



3.4

3.1

4.3



Guion

- Introducción
- Datos espacio-temporales
- Captura de datos
- Análisis de datos espaciales
- Tecnologías
- Conclusiones



Conclusiones

- Gestión de datos espaciales y espacio-temporales presente en un gran número de dominios de aplicación
 - Muchas aplicaciones de gestión de datos científicos
 - Se necesitan extensiones espaciales
- Gestión de entidades espaciales resuelto con bases de datos espaciales
 - Generalmente no necesitan tecnologías de Big Data
- Big Data?
 - Almacenamiento y análisis de datos de sensores
- Herramientas SIG no soportan grandes cantidades de datos
- Almacenamiento de datos crudos de sensores
 - Muchas inserciones por unidad de tiempo
 - Aproximaciones NoSQL: Cassandra, etc.
 - Nuevas soluciones SQL: voltDB, etc.
- Análisis sobre muestreos temporales y espaciales (coberturas)



Gracias por su atención

José Ramón Ríos Viqueira: jrr.viqueira@usc.es

citius.usc.es

