



Gbif.es



Transferencia de capacidades técnicas asociadas a la publicación de datos de biodiversidad a través de GBIF



CESP: Regional Capacity Enhancement to Latinamerica by establishing Chile's node
23 al 25 abril 2018
Santiago - Chile



Leonardo Buitrago
Administrador Nodo GBIF Colombia,
Representante (S) Regional de Nodos GBIF (Latinoamérica y Caribe)
Líder Administración de contenidos SiB Colombia
albuitrago@humboldt.org.co

Conceptos básicos sobre calidad y curación de datos

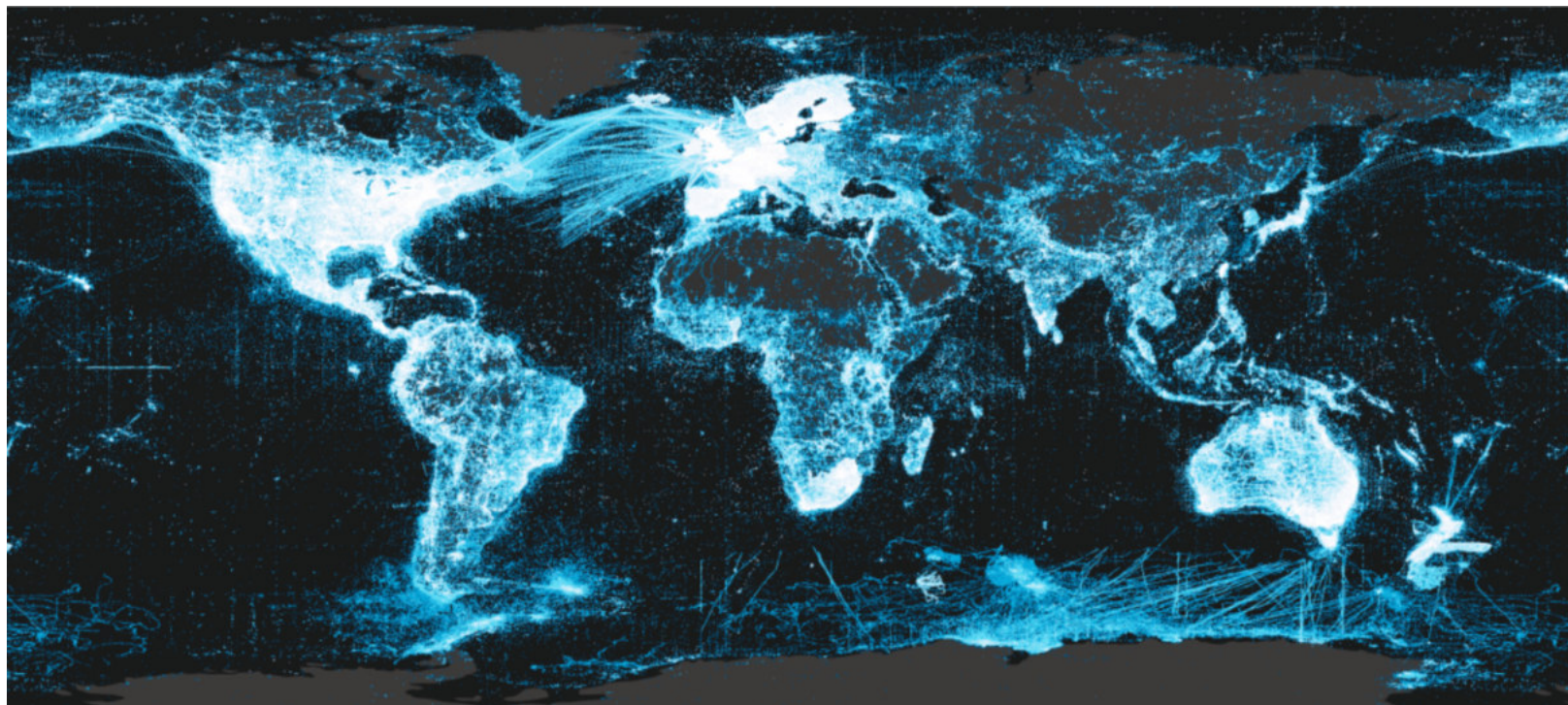


Gbif.es

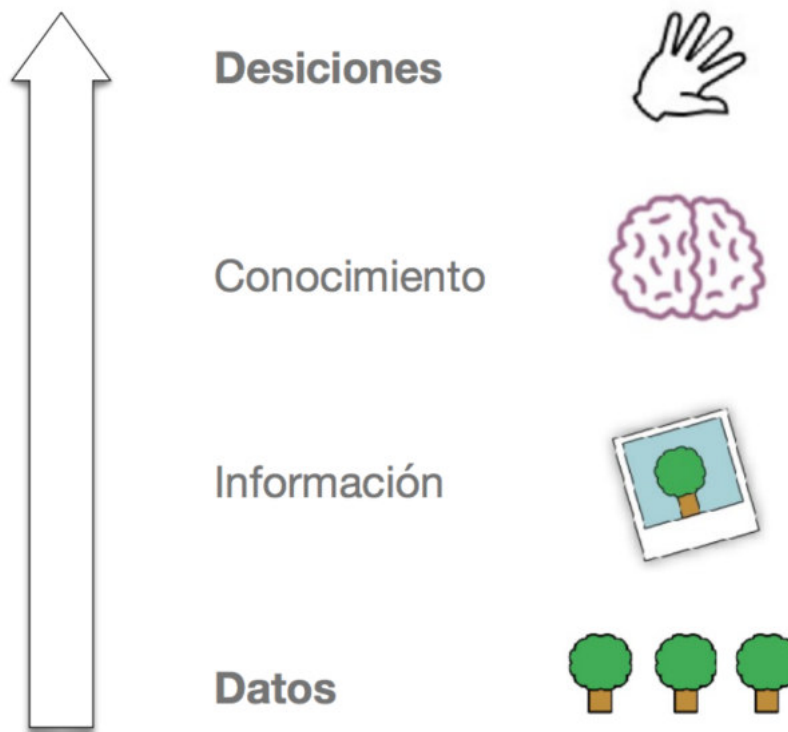


#GBIF1billion

Resources for the GBIF network to use in association with the upcoming big data milestone



DE LOS DATOS A LA TOMA DE DECISIONES



PERO...EXISTEN VACÍOS DE INFORMACIÓN

¿Cuántas y cuáles especies existen?

¿Cuál es el tamaño de sus poblaciones y dinámicas?

¿Cuál es su distribución temporal y espacial?

¿Cuántas están siendo afectadas por condiciones bióticas y abióticas?





**KEEP
CALM
AND
COLLECT
MORE DATA**



NECESITAMOS MÁS DATOS !

- Regiones pocas estudiadas o representadas
- Trabajo de campo y laboratorio
- Apoyo y financiación científica

NECESITAMOS USAR MEJOR LOS DATOS EXISTENTES !

- Una gran cantidad de datos no están disponible para su uso
 - No digitalizados, no compartidos
 - No fácilmente accesibles
 - Problemas de calidad



<http://goo.gl/pnjg82>

REFERENTES TEÓRICOS



Arthur D. Chapman¹

Although most data gathering disciplines treat error as an embarrassing issue to be expunged, the error inherent in [spatial] data deserves closer attention and public understanding ...because error provides a critical component in judging fitness for use. (Chrisman 1991).

Chapman, A. D. 2005. Principles of Data Quality, version 1.0. Report for the Global Biodiversity Information Facility, Copenhagen.



TALLER CALIDAD DE DATOS
MEJORANDO LOS DATOS PRIMARIOS
SOBRE BIODIVERSIDAD

Universidad de los Andes; Carrera 1 # 18A-12
Bogotá D.C., 18-21 Noviembre de 2014

Organizan:



Apoyan:



Saraiva A. 2014

BASURA ENTRA → BASURA SALE

- ❑ **Problemas de calidad:** conllevan a resultados de mala calidad: análisis, decisiones, etc.
- ❑ **Los problemas surgen de:** toma de datos, digitalización, falta de metadatos, ausencia de estándares.
- ❑ **Hay mucho por hacer:** limpieza de datos (corrección), prevención y políticas de calidad de datos

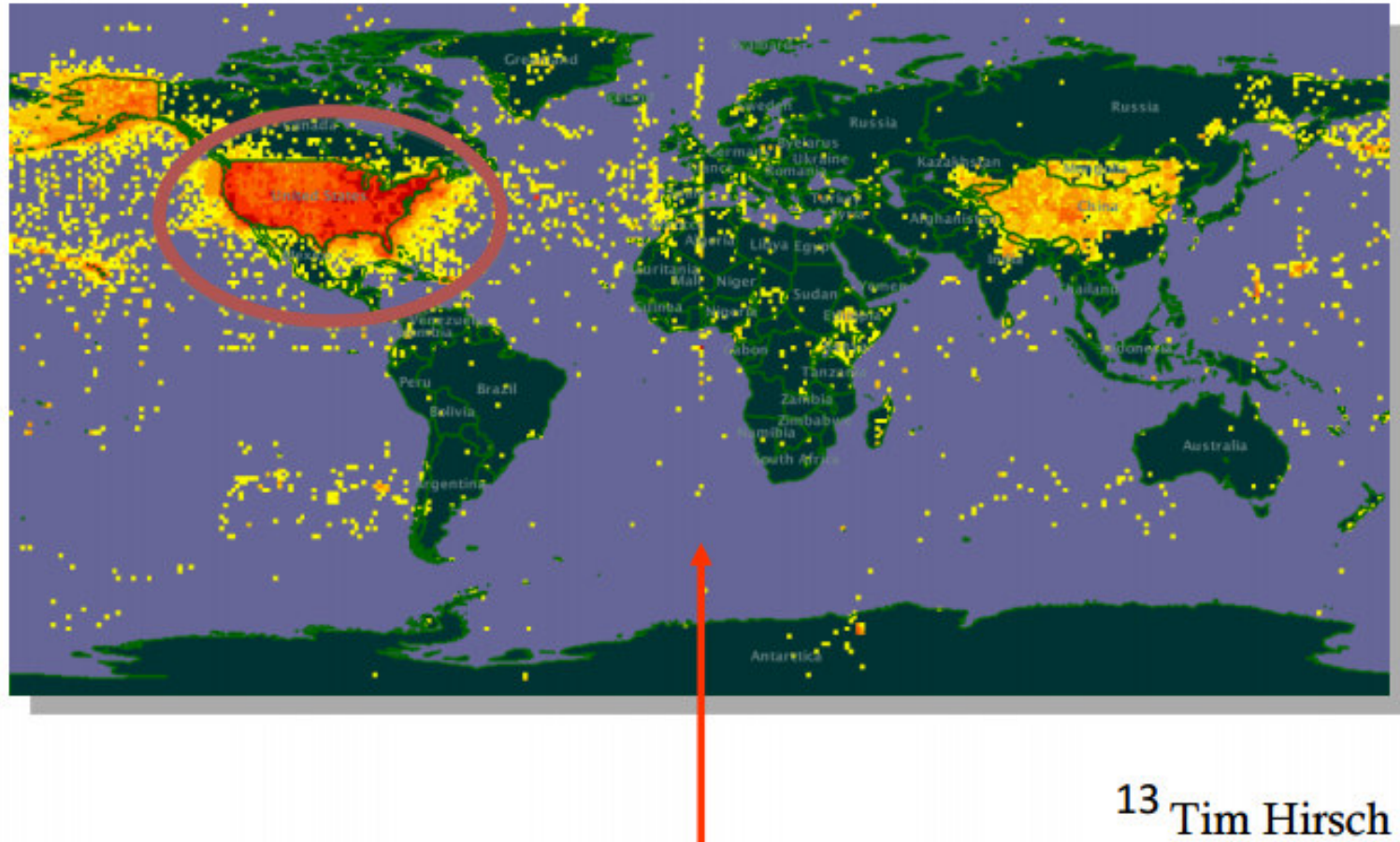


Artículos científicos

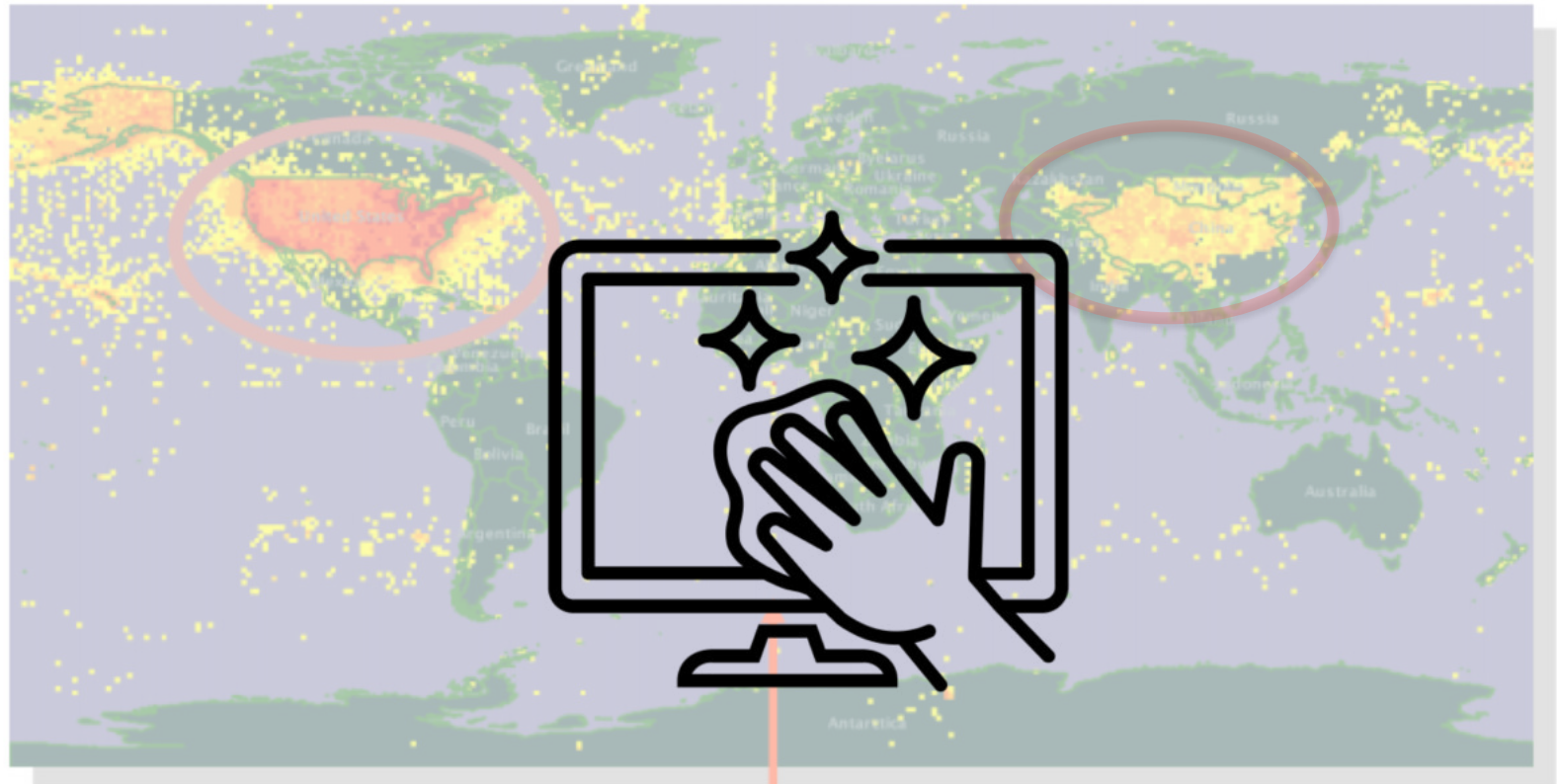
Modelamiento y análisis

Políticas de conservación

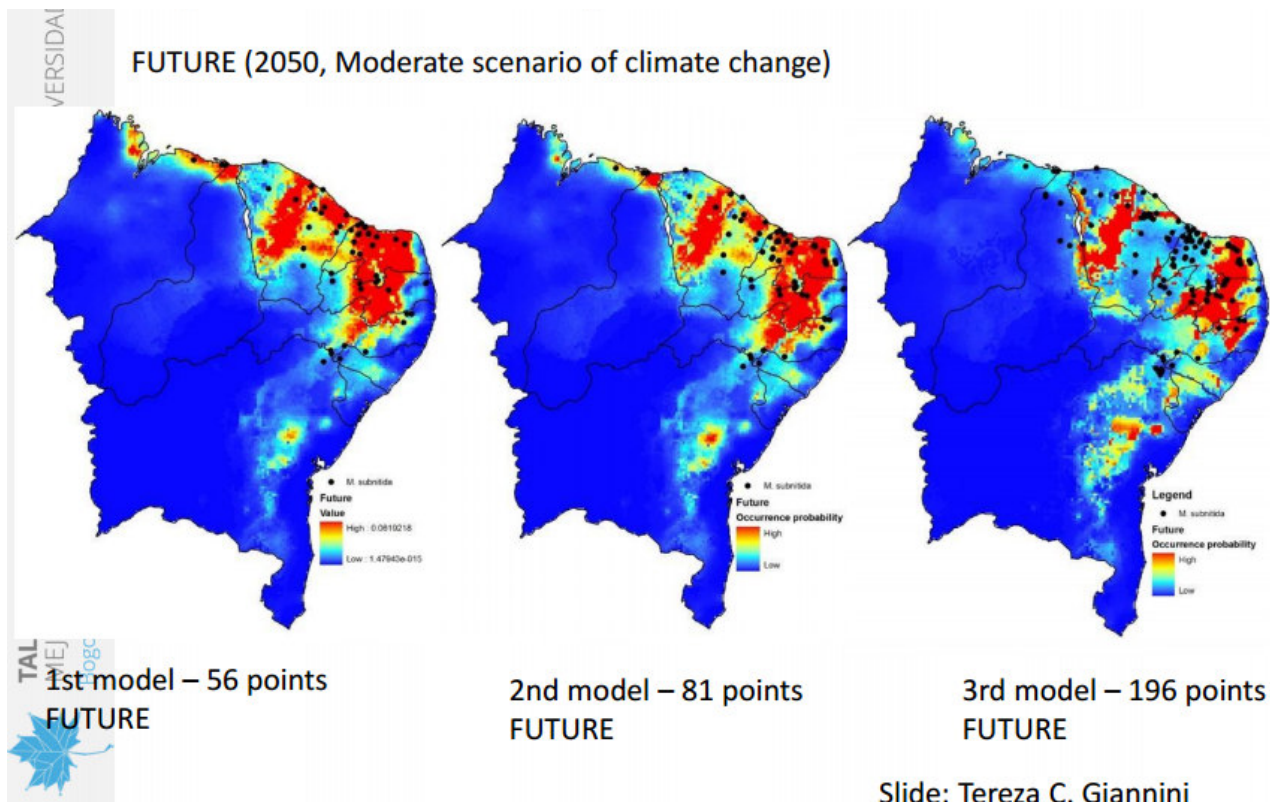
EJEMPLO



EJEMPLO



IMPACTO DE LOS DATOS EN LA GENERACIÓN DE MODELOS



Slide: Tereza C. Giannini

La calidad de datos puede afectar los indicadores, análisis, toma de decisiones y políticas



IUCN Red List Index

Guidance for national and regional use



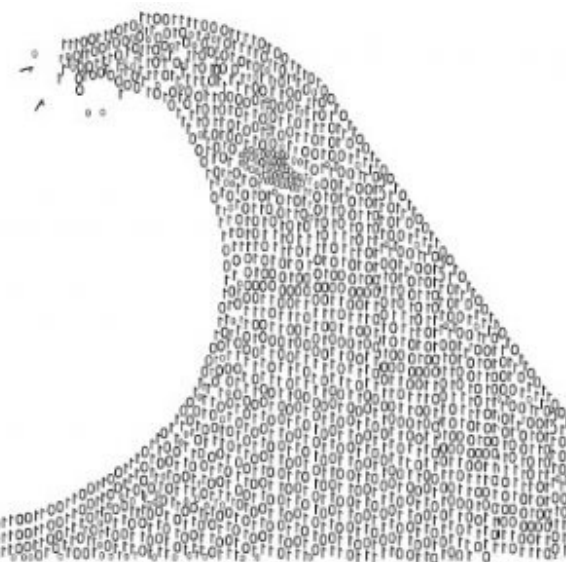
...y ahora quién podrá ayudarnos?

CALIDAD DE DATOS



ALGUNOS CONCEPTOS

- **Información:** morfè (forma) / éidos (concepto)
- Es la **representación** de la **realidad**
- La realidad es diferente de la “representación de la realidad”



EJEMPLO



Dato 1: *Saguinus*

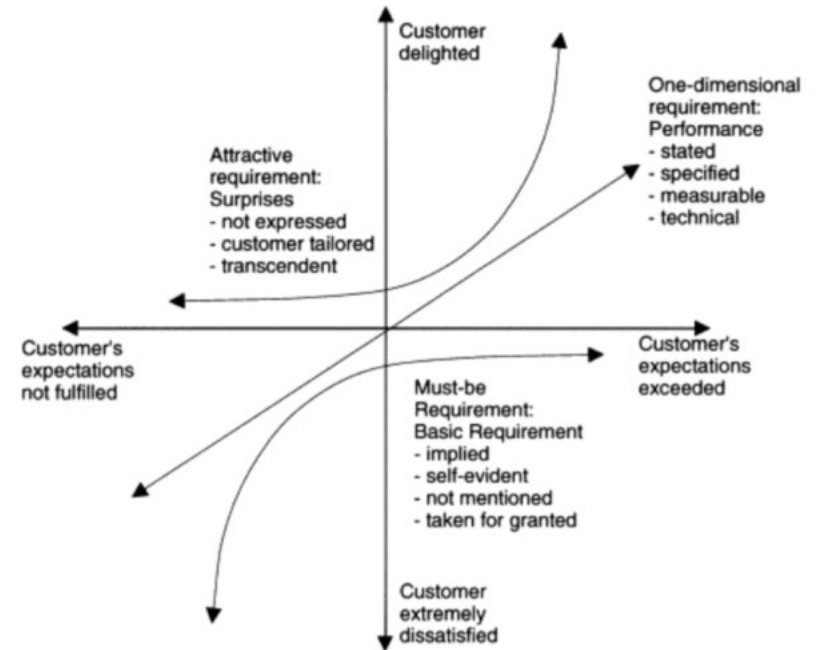
Dato 2: Mico tití

¿tienen la misma calidad?

Calidad de datos

Definición #2

Satisfacción del consumidor. Si un consumidor está satisfecho con un servicio producto, este servicio o producto tiene calidad para este consumidor.



EJEMPLO



Requerimiento: el dato debe tener nombre científico y debe ser suministrado a nivel de especie

Nombre: *Saguinus*

Categoría: *Genero*

¿este dato tiene calidad?

¿puede ser usado en el estudio de la distribución de primates en Suramérica?

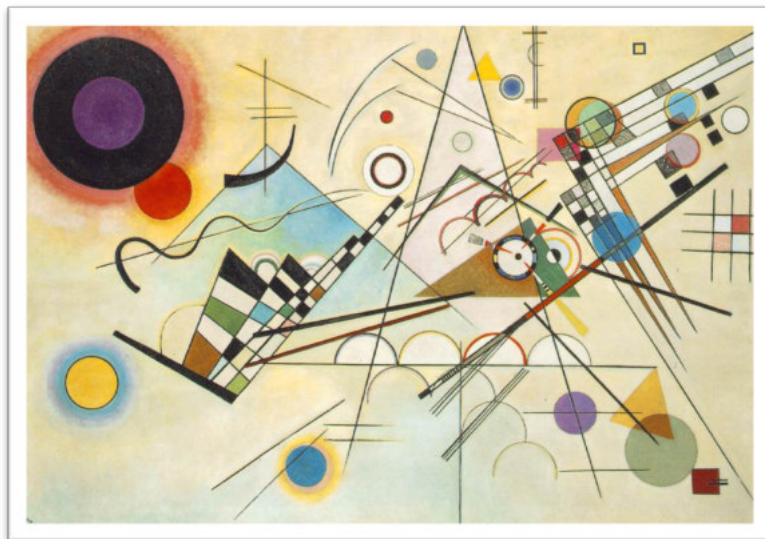
Calidad de datos

Definición #3

Usabilidad. Un dato tiene calidad si es adecuado para ser usado. Si el dato no sirve para el propósito del que lo usa, **puede ser útil para otros.**



ALGUNOS CONCEPTOS



La calidad de datos es un concepto idiosincrásico

“La idiosincracia es algo distintivo y propio de un individuo”

Definir calidad de datos es similar a definir qué es bonito, bueno divertido o valioso.

PALABRAS CLAVE

La palabra clave y la definición mas aceptada:

Usabilidad de los datos



USO: Calidad en relación a un propósito.

- *Modelos de distribución, lista nacional de especies, etc.*

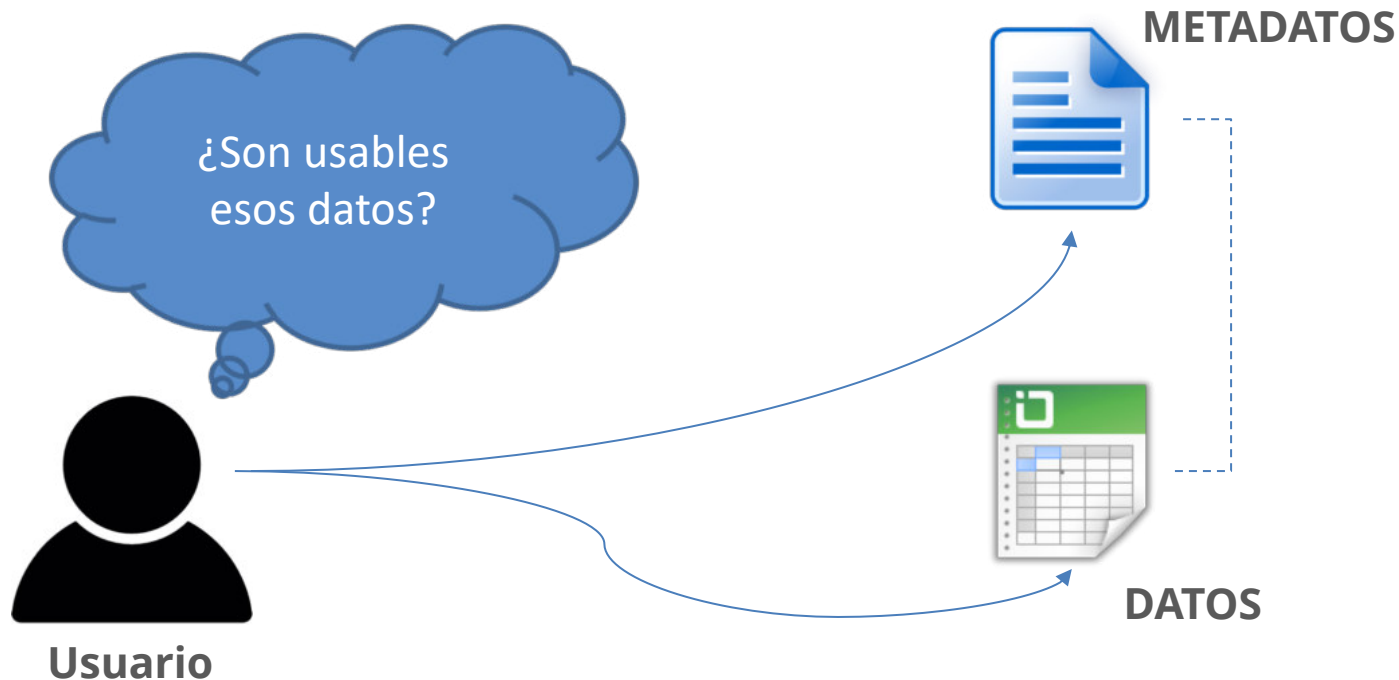
DATOS: Para cada propósito existe un tipo de datos.

- *Modelamiento: coordenadas y nombres de especies.*

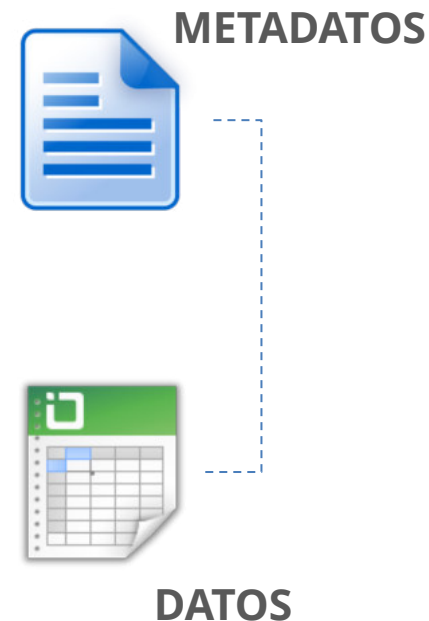
USABILIDAD: Para cada tipo de dato existen atributos a cumplir.

- *Compleitud, consistencia, precisión, exactitud, etc.*

USABILIDAD



USABILIDAD



EVALUACIÓN DE LA CALIDAD



DATOS



EVALUACIÓN



NO USABLES



USABLES

VALIDACIÓN DE CALIDAD



DATOS



VALIDACIÓN

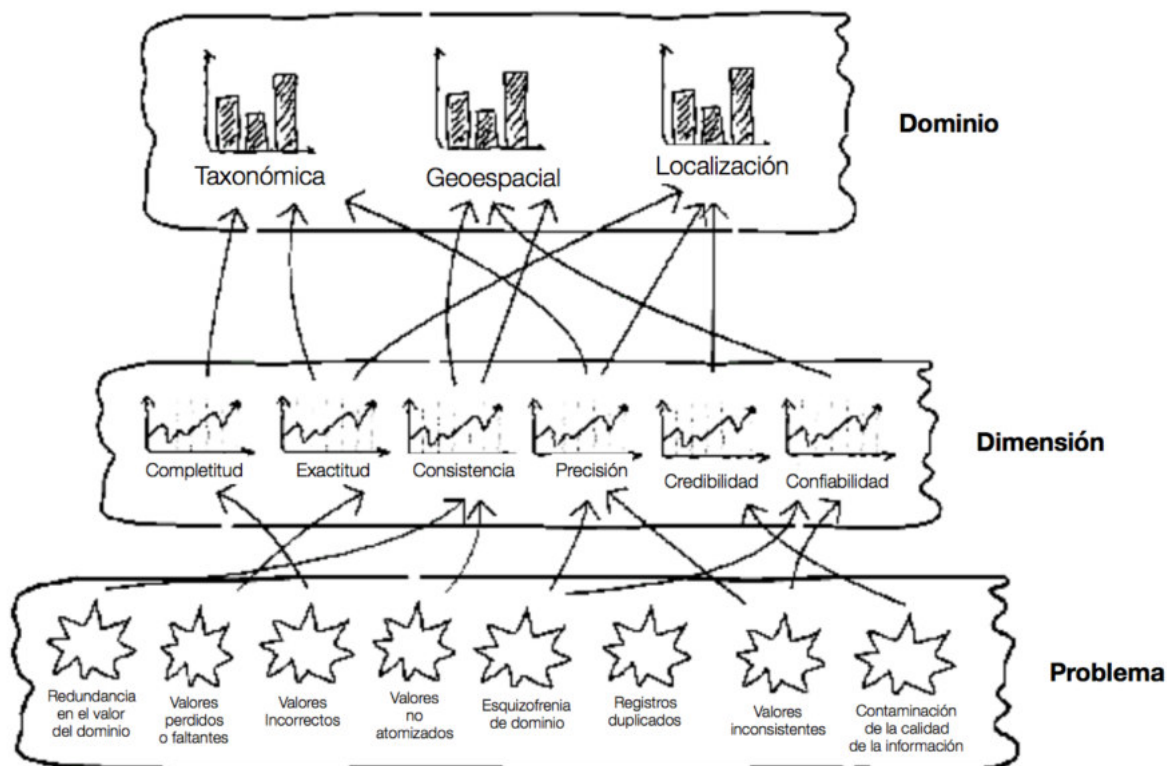


NO USABLES



USABLES

EVALUACIÓN DE LA CALIDAD



DOMINIOS



Geo-espacial: Datos con coordenadas o georreferenciados y regiones político-administrativas documentadas.



Taxonomía (nomenclatura): Nombres científicos, niveles taxonómicos



Formato (Estandarización): Errores de tipeo, formatos de fecha, formatos de coordenadas, caracteres especiales y demás.

DIMENSIÓN

Es el aspecto **medible** de la calidad del dato.

*Ejemplo: **Precisión***

* La calidad de datos es un concepto multidimensional

Dominio Geoespacial
-23.98 es menos preciso que -23.98**74**

Dominio Taxonómico
Taxón A: reino= X; filo= Y; clase=Z
Taxón B: reino= X; filo= Y; clase=?

PROBLEMAS

Todo lo que pueda **degradar** la calidad para una o mas **dimensiones**.

- **Domain value redundancy:** Ex.: Brazil x Brasil, São Paulo x Sao Paulo.
- **Missing data value:** absence of latitude or longitude.
- **Incorrect data values:** misspelling errors.
- **Nonatomic data values:** Ex.: São Paulo, SP.
- **Domain schizophrenia:** latitude and longitude equal "0" (zero) when they are not known.
- **Duplicate occurrences:** More than one record representing the same fact.
- **Inconsistent data values:** Ex.: lat=43; long=44; country=Brazil.
- **Information quality contamination:** create a new data based on an already existent wrong data.

CADENA DE LA INFORMACIÓN

CADENA DE LA INFORMACIÓN



CADENA DE LA INFORMACIÓN



1



2

CADENA DE LA INFORMACIÓN



1



2



3

CADENA DE LA INFORMACIÓN



Planificación

1



Captura de datos y Documentación

2



Digitalización

3



Control de calidad

4

CADENA DE LA INFORMACIÓN



Planificación

1



Captura de datos y Documentación

2




Digitalización

3



Control de calidad

4



Publicación

5

CADENA DE LA INFORMACIÓN

COSTO DE LA CORRECCIÓN DE ERRORES



Planificación

1



Captura de datos y Documentación

2



Digitalización

3



Control de calidad

4



Publicación

5

MECANISMOS DE MEJORA

 **PREVENCIÓN:** Evitar que se presenten errores previo a la creación de los datos

 **DETECCIÓN Y LIMPIEZA:** Detectar errores en el conjunto de datos y corregirlos

 **DETECCIÓN Y RECOMENDACIONES:** Detectar errores en el conjunto de datos y generar recomendaciones de limpieza



Gbif.es



CEESP: Regional Capacity Enhancement to Latinamerica by establishing Chile's node
23 al 25 abril 2018
Santiago - Chile

Muchas gracias por su atención

Preguntas?

Leonardo Buitrago
Administrador Nodo GBIF Colombia,
Representante (S) Regional de Nodos GBIF (Latinoamérica y Caribe)
Líder Administración de contenidos SiB Colombia
albuitrago@humboldt.org.co

