# Data & Data Management

Jarrett Byrnes

9/6/2012

**DataONE**

---

http://dataone.org

---

# What is the Data Life Cycle?

Plan

Collect

Analyze

Assure

Integrate

Describe

Discover

Preserve

**DataONE**

---

# For Each Stage of the Data Lifecycle…

• …there are best practices…..and….tools to help!

• Your well-managed and accessible data can contribute to science in ways you may not even imagine today!

• `http://www.dataone.org/best-practices`

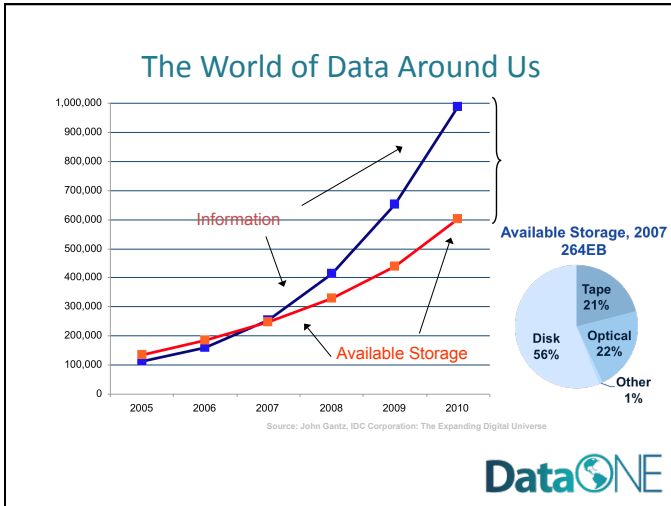• `http://www.dataone.org/education-modules`

**DataONE**

## Data deluge

Data is collected from sensors, sensor networks, remote sensing, observations, and more - this calls for increased attention to data management and stewardship
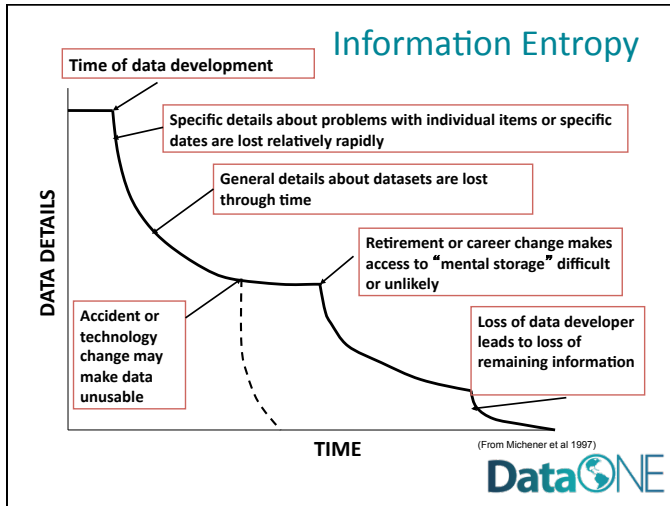


## The World of Data Around Us



Information

Available Storage

**Available Storage, 2007 264EB**

- Tape 21%
- Disk 56%
- Optical 22%
- Other 1%

Source: John Gantz, IDC Corporation: The Expanding Digital Universe

## The World of Data Around Us: Data Loss



- Natural disaster
- Facilities infrastructure failure
- Storage failure
- Server hardware/software failure
- Application software failure
- External dependencies (e.g. PKI failure)
- Format obsolescence
- Legal encumbrance
- Human error
- Malicious attack by human or automated agents
- Loss of staffing competencies
- Loss of institutional commitment
- Loss of financial stability
- Changes in user expectations and requirements

## Information Entropy

**Time of data development**

DATA DETAILS

**Specific details about problems with individual items or specific dates are lost relatively rapidly**

**General details about datasets are lost through time**

**Retirement or career change makes access to "mental storage" difficult or unlikely**

**Accident or technology change may make data unusable**

**Loss of data developer leads to loss of remaining information**

**TIME**

(From Michener et al 1997)

**DataONE**

## Information Entropy

DATA DETAILS

**Sound information management, including metadata development, can arrest the loss of dataset detail.**
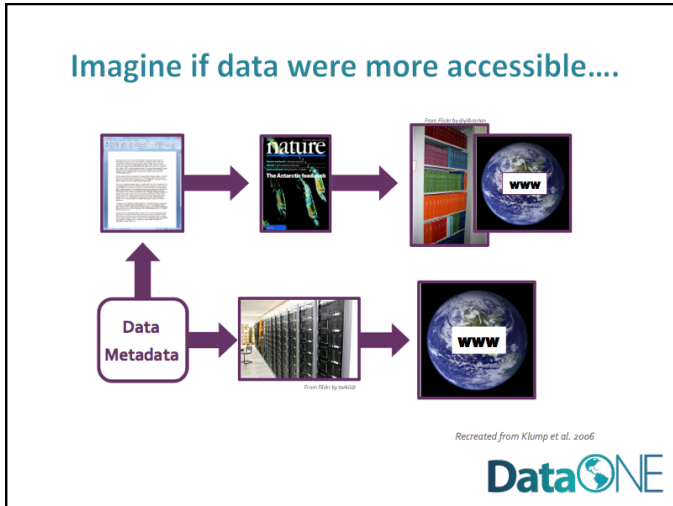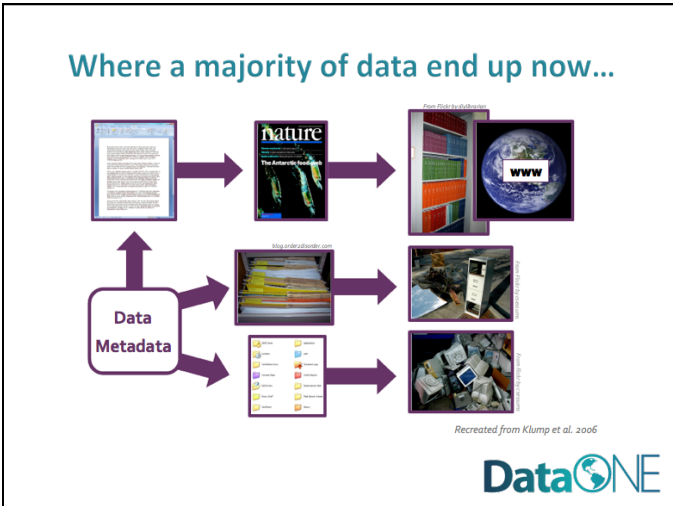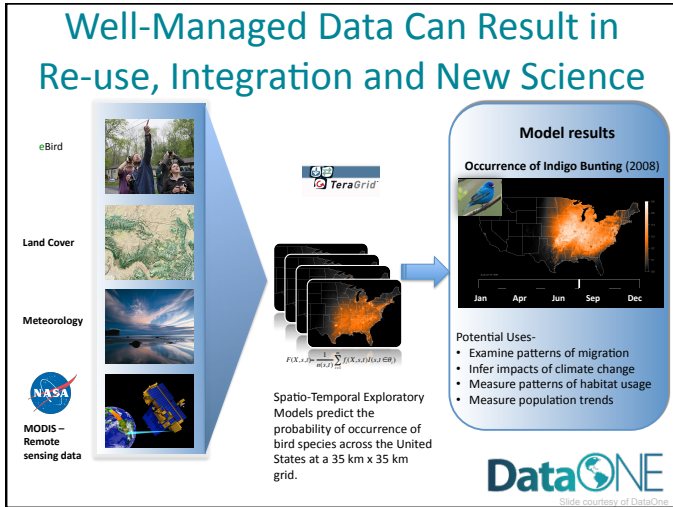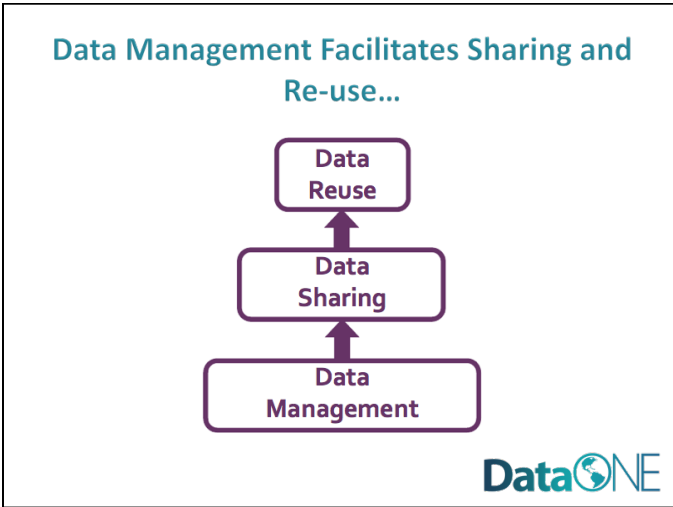
**TIME**

**DataONE**

## Why Manage Data:
## Researcher Perspective

- Manage your data for yourself:
  - ○ Keep yourself organized – be able to find your files (data inputs, analytic scripts, outputs at various stages of the analytic process, etc)
  - ○ Track your science processes for reproducibility – be able to match up your outputs with exact inputs and transformations that produced them
  - ○ Better control versions of data – identify easily versions that can be periodically purged
  - ○ Quality control your data more efficiently

**DataONE**

## Why Data Management:
## Researcher Perspective

- Make backups to avoid data loss
- Format your data for re-use (by yourself or others)
- Be prepared: Document your data for your own recollection, accountability, and re-use (by yourself or others)
- Prepare it to share it – gain credibility and recognition for your science efforts!
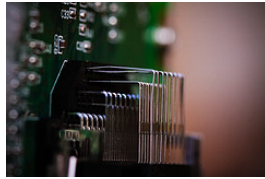
**DataONE**

## Data Management Facilitates Sharing and Re-use…



## Well-Managed Data Can Result in Re-use, Integration and New Science



## Where a majority of data end up now…



## Imagine if data were more accessible….

## Why Share Data?

Data sharing requires effort, resources, and faith in others. Why do it?

For the benefit of:
- o the public
- o the research sponsor
- o the research community
- o the researcher

## Value of Data Sharing:
## To the Scientist

Scientists that share data gain the benefit of:
- o Authority
- o Citation
- o Collaboration

## Concerns About Data Sharing

Even if the value of data sharing is recognized, concerns remain as to the impacts of increased data exposure.

## Concerns About Data Sharing

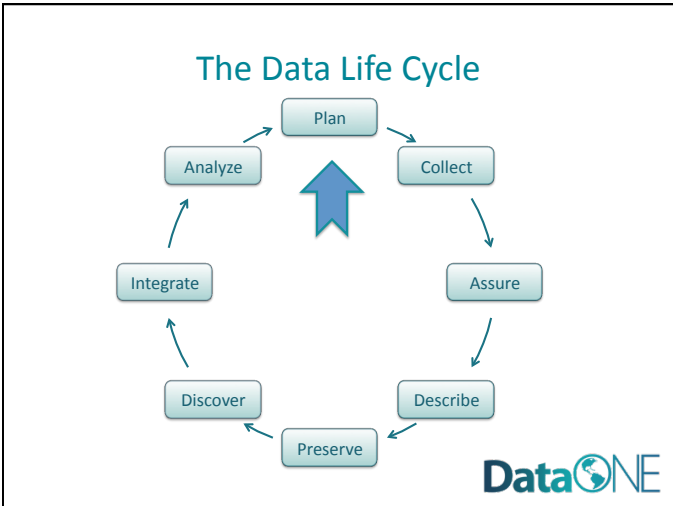| Concern | Solution |
|---|---|
| inappropriate use due to misunderstanding of research purpose or parameters | |
| security and confidentiality of sensitive data | |
| lack of acknowledgement / credit | |
| loss of advantage when competing for research dollars | |

## Concerns About Data Sharing

| Concern | Solution |
|---|---|
| inappropriate use due to misunderstanding of research purpose or parameters | ✓ metadata |
| security and confidentiality of sensitive data | ✓ metadata |
| lack of acknowledgement / credit | ✓ metadata |
| loss of advantage when competing for research dollars | ✓ metadata |

**DataONE**

---

# Data Management Planning & Meta-Data



**DataONE**

---

## The Data Life Cycle



Plan → Collect → Assure → Describe → Preserve → Discover → Integrate → Analyze

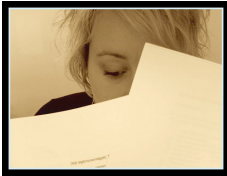**DataONE**

---

## What is a Data Management Plan?

- Formal document
- Outlines what you will do with your data **during** and **after** you complete your research
- Ensures your data is safe for the **present** and the **future**

*From University of Virginia Library*

**DataONE**

## Why Prepare a DMP?

- Save time
  - o Less reorganization later
- Increase research efficiency
  - o Ensures you and others will be able to understand and use data in future

**Data**ONE

## Components of a General DMP

1. Information about data & data format
2. Metadata content and format
3. Policies for access, sharing and re-use
4. Long-term storage and data management
5. Budget

**Data**ONE

## 1. Information About Data & Data Format

1.1  Description of data to be produced

- Experimental
- Observational
- Raw or derived
- Physical collections
- Models and their outputs
- Simulation outputs
- Curriculum materials
- Software
- Images
- Etc…

**Data**ONE

## 1. Information About Data & Data Format

1.2  How data will be acquired

- When?
- Where?

1.3  How data will be processed
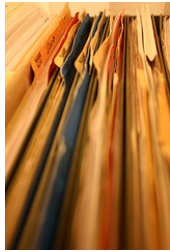
- Software used
- Algorithms
- Workflows

**Data**ONE

## 1. Information About Data & Data Format

1.4  File formats
- Justification
- Naming conventions

1.5  Quality assurance & control during sample collection, analysis, and processing

**Data**ONE

## 1. Information About Data & Data Format

1.6  Existing data
- If existing data are used, what are their origins?
- Will your data be combined with existing data?
- What is the relationship between your data and existing data?

1.7  How data will be managed in short-term
- Version control
- Backing up
- Security & protection
- Who will be responsible

**Data**ONE

## 2. Metadata Content & Format

**Metadata defined:**
- Documentation and reporting of data
- Contextual details: Critical information about the dataset
- Information important for using the data
- Descriptions of temporal and spatial details, instruments, parameters, units, files, etc.

**Data**ONE

## 2. Metadata Content & Format

2.1  What metadata are needed
- Any details that make data meaningful

2.2  How metadata will be created and/or captured
- Lab notebooks? GPS units?
- Auto-saved on instrument?

2.3  What format will be used for the metadata
- Standards for community
- Justification for format chosen

**Data**ONE

## 3. Policies for Access, Sharing, Reuse

3.4  Intellectual property & copyright issues
- Who owns the copyright?
- Institutional policies
- Funding agency policies
- Embargos for political/commercial reasons

3.5  Intended future uses/users for data

3.6  Citation
- How should data be cited when used?
- Persistent citation?

**Data**ONE

---

## 4. Long-term Storage & Data Management

4.1  What data will be preserved

4.2  Where will it be archived
- Most appropriate archive for data
- Community standards

3.6  Data transformations/formats needed
- Consider archive policies

4.4   Who will be responsible
- Contact person for archive

**Data**ONE

---

## 5. Budget

5.1  Anticipated costs
- Time for data preparation & documentation
- Hardware/software for data preparation & documentation
- Personnel
- Archive costs

5.2  How costs will be paid

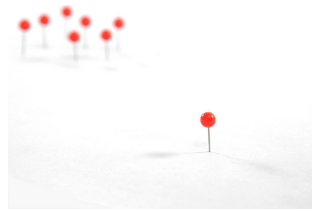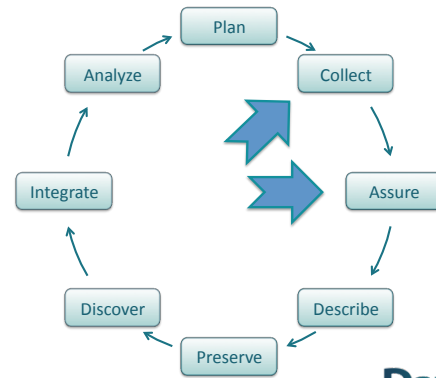**Data**ONE

---

## Tools for Creating Data Management Plans

dmp.cdlib.org

dmponline.dcc.ac.uk

JISC

**Data**ONE

## Data Entry and Quality Control



**Data**ONE

## The Data Life Cycle



- Plan
- Collect
- Assure
- Describe
- Preserve
- Discover
- Integrate
- Analyze

**Data**ONE

## Goals of Data Entry

- create data sets that are:
  - o Valid
  - o Organized to support ease of use



**Data**ONE

## Example: Poor Data Entry



- Inconsistency between data collection events
  - Location of Date information
  - Inconsistent Date format
  - Column names
  - Order of columns

**Data**ONE

## Example: Poor Data Entry



- Inconsistency between data collection events
  - Different site spellings, capitalization, spaces in site names—hard to filter
  - Codes used for site names for some data, but spelled out for others
  - Mean1 value is in Weight column
  - Text and numbers in same column – what is the mean of 12, "escaped < 15", and 91?

## Best Practices



- Columns of data are consistent: only numbers, dates, or text
- Consistent Names, Codes, Formats (date) used in each column
- Data are all in one table, which is much easier for a statistical program to work with than multiple small tables which each require human intervention

## Best Practices

- Create descriptive column names without spaces or special characters
  - Soil T30 → Soil_Temp_30cm
  - Species-Code → Species_Code (avoid using -,+,*,^ in column names. Some software may interpret these symbols as an operator)
- Use a descriptive file name. For instance, a file named SEV_SmallMammalData_v.5.25.2010.csv indicates the project the data is associated with (SEV), the theme of the data (SmallMammalData) and also when this version of the data was created (v.5.25.2010). This name is much more helpful than a file named mydata.xls.

## Best Practices

- Missing data
  - Preferably leave field empty (NULL = no value)
  - In numeric fields, use a distinct value such as 9999 to indicate a missing value – but only if this is in your meta-data!
  - In text fields, use NA ("Not Applicable" or "Not Available")
  - Use Data flags in a separate column to qualify missing value

| Date | Time | NO3_N_Conc | NO3_N_Conc_Flag |
|------|------|------------|-----------------|
| 20081011 | 1300 | 0.013 | |
| 20081011 | 1330 | 0.016 | |
| 20081011 | 1400 | | M1 |
| 20081011 | 1430 | 0.018 | |
| 20081011 | 1500 | 0.001 | E1 |

M1 = missing; no sample collected

E1 = estimated from grab sample

## Best Practices

- Enter complete lines of data



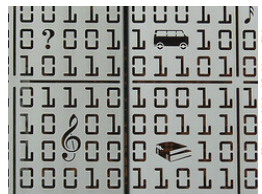Sorting an Excel file with empty cells is not a good idea!

## Best Practices

- For the long term, store data in a consistent format that can be read well in to the future and that can be used by any application now or in the future
- Appropriate file types include:
  - Non-proprietary: Open, documented standard
  - Common usage by research community: Standard representation (ASCII, Unicode)
  - Unencrypted
  - Uncompressed
- ASCII formatted files will be readable into the future
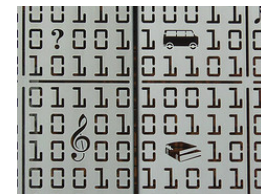  - Use ASCII (comma-separated) for tabular data

## Quality Control



## Definitions

**Data Contamination**

- Process or phenomenon, other than the one of interest, that affects the variable value
- Erroneous values

## Definitions: Types of Errors

- Errors of **Commission**
  - o Incorrect or inaccurate data entered
  - o Examples: malfunctioning instrument, mistyped data
- Errors of **Omission**
  - o Data or metadata not recorded
  - o Examples: inadequate documentation, human error, anomalies in the field

**Data**ONE

## Defining QA/QC

- Strategies for preventing errors from entering a dataset
- Activities to ensure quality of data before collection
- Activities that involve monitoring and maintaining the quality of data during the study

**Data**ONE

## QA/QC Before Collection

- Define & enforce standards
  - – Formats
  - – Codes
  - – Measurement units
  - – Metadata
- Assign responsibility for data quality
  - – Be sure assigned person is educated in QA/QC

**Data**ONE

## QA/QC  During Data Entry

- Double entry
  - – Data keyed in by two independent people
  - – Check for agreement with computer verification
- Record a reading of the data and transcribe from the recording
- Use text-to-speech program to read data back

**Data**ONE

## QA/QC During Data Entry

- Design data storage well
  - Minimize number of times items that must be entered repeatedly
  - Use consistent terminology
  - Atomize data: one cell per piece of information
- Document changes to data
  - Avoids duplicate error checking
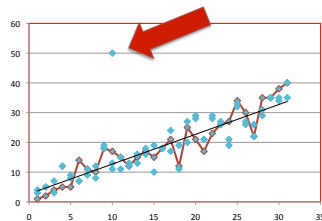  - Allows undo if necessary

**DataONE**

## QA/QC After Data Entry

- Make sure data line up in proper columns
- No missing, impossible, or anomalous values
- Perform statistical summaries



**DataONE**

## QA/QC After Data Entry

- Look for outliers
  - Outliers are extreme values for a variable given the statistical model being used
  - The goal is not to eliminate outliers but to identify potential data contamination



**DataONE**

## R and Data Screening

**DataONE**

## R and Data: Use Data in a Long Table Format



## R doesn't Play Nice with Excel – save data in a comma or tab delimited file



Note .csv in file name

Select comma Separated values (.csv) in the "Format" dropdown.

## Step 1) Set your Working Directory



Click on the … to select a directory

## Step 1) Set your Working Directory



The "More" button

## Loading Data

```
westNile <- read.csv("./data/SwaddleWestNile2002NCEAS-BAD.csv")
```

Note:

- ▶ File path (./ is this directory, ../ is back one directory)
- ▶ Quotes
- ▶ Our data is now an object in R

## Look at Your Data

```
head(westNile)
```

```
##    State                Infected.County WNV.incidence
## 1    AL        Autauga, AL                     2.290
## 2    AL      Calhoun , AL                      0.891
## 3    AL   Chambers, AL                         2.734
## 4    AL   Dallas , AL                          2.157
## 5    AL   Marengo , AL                         8.874
## 6    AL   Marion, AL                           3.204
##    Species.Richness Corvid.Abundance
## 1              66                 8
## 2              67                64
## 3              41                69
## 4              60                66
## 5              69                64
## 6              NA     NOT AVAILABLE
```

## Look at Columns 3 through 4

```
head(westNile[, 3:4])
```

```
##    WNV.incidence Species.Richness
## 1         2.290              66
## 2         0.891              67
## 3         2.734              41
## 4         2.157              60
## 5         8.874              69
## 6         3.204              NA
```

- ▶ Data Frame is treates as a Matrix.
- ▶ [$rows, columns$]

## Look at Your Individual Columns

```
names(westNile)
```

```
## [1] "State"           "Infected.County"  "WNV.incidence"
## [4] "Species.Richness" "Corvid.Abundance"
```

(Note that spaces are now .s)

```
westNile$Species.Richness
```

```
##    [1] 66 67 41 60 69 NA 56 65 54 52 81 51 47 59 49 51 72 53
##   [19] 54 49 61 81 62 70 71 57 87 64 50 62 71 70 59 63 58 51
##   [37] 46 66 53 59 58 56 58 43 65 51 51 63 54 60 53 39 62 67
##   [55] 68 82 70 76 58 60 72 59 72 62 82 63 68 39 67 66 63 47
##   [73] 59 61 65 79 54 56 30 48 56 68 58 42 51 64 73 55 61 65
##   [91] 61 74 65 61 51 93 42 63 68 58 68 61 56 60 81 66 53 49
##  [109] 68 72 76 57 76 55 76 56 73 59 73 57 90 50 73 64 78 75
##  [127] 61 80 59 69
```

## Missing Data is NA

```
westNile$Species.Richness
```

```
##    [1] 66 67 41 60 69 NA 56 65 54 52 81 51 47 59 49 51 72 53
##   [19] 54 49 61 81 62 70 71 57 87 64 50 62 71 70 59 63 58 51
##   [37] 46 66 53 59 58 56 58 43 65 51 51 63 54 60 53 39 62 67
##   [55] 68 82 70 76 58 60 72 59 72 62 82 63 68 39 67 66 63 47
##   [73] 59 61 65 79 54 56 30 48 56 68 58 42 51 64 73 55 61 65
##   [91] 61 74 65 61 51 93 42 63 68 58 68 61 56 60 81 66 53 49
##  [109] 68 72 76 57 76 55 76 56 73 59 73 57 90 50 73 64 78 75
##  [127] 61 80 59 69
```

Note the NA. This is missing data.

```
westNile$Species.Richness[6]
```

```
## [1] NA
```

## Let's look at another

```
westNile$Corvid.Abundance
```

```
##    [1] 8           64          69
##    [4] 66          64          NOT AVAILABLE
##    [7] 59          129         54
##   [10] 100         62          82
##   [13] 102         35          31
##   [16] 13          51          60
##   [19] 10          87          53
##   [22] 9999        34          86
##   [25] 75          102         216
##   [28] 71          43          57
##   [31] 98          84          44
##   [34] 109         165         44
##   [37] 68          48          34
##   [40] 63          9999        52
##   [43] 24          39          41
##   [46] 32          47          23
```

## Cleaner Data

```
westNile <- read.csv("./data/SwaddleWestNile2002NCEAS-BAD.csv",
    na.strings = "NOT AVAILABLE")
```
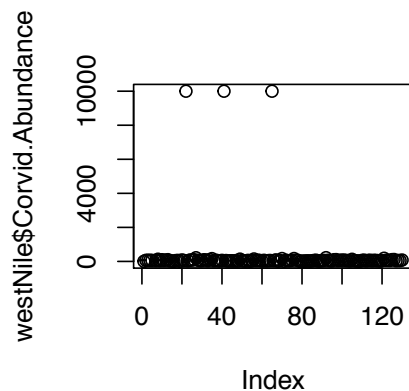
## And, Fixed!

```
westNile$Corvid.Abundance
```

```
##    [1]    8.00   64.00   69.00   66.00   64.00      NA
##    [7]   59.00  129.00   54.00  100.00   62.00   82.00
##   [13]  102.00   35.00   31.00   13.00   51.00   60.00
##   [19]   10.00   87.00   53.00 9999.00   34.00   86.00
##   [25]   75.00  102.00  216.00   71.00   43.00   57.00
##   [31]   98.00   84.00   44.00  109.00  165.00   44.00
##   [37]   68.00   48.00   34.00   63.00 9999.00   52.00
##   [43]   24.00   39.00   41.00   32.00   47.00   23.00
##   [49]  135.00   49.00   32.00   27.00   63.00   15.00
##   [55]   45.00  144.00   61.00   71.00   57.00   29.00
##   [61]   66.00   36.00   46.00   57.00 9999.00   54.00
##   [67]   91.00   19.00   56.00  168.00   14.00   71.00
##   [73]   43.00   48.00   70.00  170.00   75.00   63.00
##   [79]    6.00   18.00   21.00   29.00   34.00   18.00
##   [85]   39.00   57.00   71.00   26.00   31.00   47.00
##   [91]   63.00  220.00   70.00   42.00   36.00  101.00
```

## What about fixing many bad values?

```
plot(westNile$Corvid.Abundance)
```



## What about fixing many bad values?

```
which(westNile$Corvid.Abundance == 9999)
```

```
## [1] 22 41 65
```

$==$ is makes a COMPARISON and returns a logical value
Can also use $<$, $>$, and more.

```
westNile$Corvid.Abundance == 9999
```

```
##   [1] FALSE FALSE FALSE FALSE FALSE    NA FALSE FALSE FALSE
##  [10] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
##  [19] FALSE FALSE FALSE  TRUE FALSE FALSE FALSE FALSE FALSE
##  [28] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
##  [37] FALSE FALSE FALSE FALSE  TRUE FALSE FALSE FALSE FALSE
##  [46] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
##  [55] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
##  [64] FALSE  TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
##  [73] FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE
```

## Replace the 9999s

```
westNile$Corvid.Abundance[which(westNile$Corvid.Abundance ==
    9999)] <- NA
```
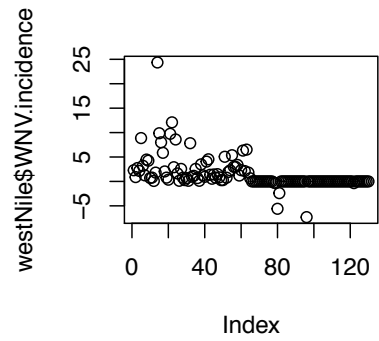
The which approach is often good, as once you spot a single problem observation, there may be others like it.

## Exercise

1. Is everything OK with West Nile Virus Incidence?

2. Let's say a database overwrote some 0 values - fix these values!

## The Fix

```
plot(westNile$WNV.incidence)
```



```
westNile$WNV.incidence[which(westNile$WNV.incidence < 0)] <- 0
```
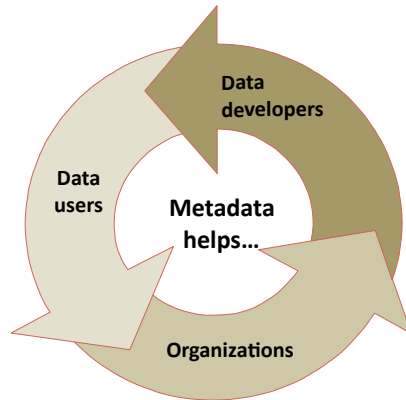
# Metadata

**Data**O**NE**

---

## What is Metadata?

**Metadata is: Data 'reporting'**

- **WHO** created the data?
- **WHAT** is the content of the data?
- **WHEN** were the data created?
- **WHERE** is it geographically?
- **HOW** were the data developed?
- **WHY** were the data developed?

**Data**O**NE**

---

## The Value of Metadata

Data developers

Data users

**Metadata helps…**

Organizations

**Data**O**NE**

---

## Multiple Metadata Standards Exist: Examples

- **Ecological Metadata Language (EML)**
  - Focus on ecological data
  - http://knb.ecoinformatics.org/eml_metadata_guide.html
- **Darwin Core**
  - Emphasis on museum specimens
  - http://rs.tdwg.org/dwc/index.htm
- **Geography Markup Language (GML)**
  - Emphasis on geographic features (roads, highways, bridges)
  - http://www.opengeospatial.org/standards/gml

**Data**O**NE**

## Tips for Writing Quality Metadata

- Do not use jargon
- Define technical terms and acronyms:
  - CA, LA, GPS, GIS : what do these mean?
- Clearly state data limitations
  - E.g., data set omissions, completeness of data
  - Express considerations for appropriate re-use of the data
- Use "none" or "unknown" meaningfully
  - None usually means that you knew about data and nothing existed (e.g., a "0" cubic feet per second discharge value)
  - Unknown means that you don't know whether that data existed or not (e.g., a null value)



**Data**ONE

## Tips for Writing Quality Metadata

Titles, Titles, Titles…

- Titles are critical in helping readers find your data
  - While individuals are searching for the most appropriate data sets, they are most likely going to use the title as the first criteria to determine if a dataset meets their needs.
  - Treat the title as the opportunity to sell your dataset.
- A complete title includes: What, Where, When, Who, and Scale
- An informative title includes: topic, timeliness of the data, specific information about place and geography

**Data**ONE

## Tips for Writing Quality Metadata

- A Clear Choice: Which title is better?

- *Rivers*

  OR

- *Greater Yellowstone Rivers from 1:126,700 U.S. Forest Service Visitor Maps (1961-1983)*

Greater Yellowstone (where) Rivers (what) from 1:126,700 (scale) U.S. Forest Service (who) Visitor Maps (1961-1983) (when)



**Data**ONE

## Tips for Writing Quality Metadata

- Be specific and quantify when you can! The goal of a metadata record is to give the user enough information to know if they can use the data without contacting the dataset owner.

  Vague:  We checked our work and it looks complete.

  Specific:   We checked our work using a random sample of 5 monitoring sites reviewed by 2 different people. We determined our work to be 95% complete based on these visual inspections.



**Data**ONE

## Tips for Writing Quality Metadata

- Select keywords wisely
- Use descriptive and clear writing
- Fully qualify geographic locations
- Use thesauri for keywords whenever possible
- Example: USGS Biocomplexity Thesaurus (over 9,500 terms)

**DataONE**

## Tips for Writing Quality Metadata

- Remember: a computer will read your metadata

- Do not use symbols that could be misinterpreted:
  Examples: ! @ # % { } | / \ < > ~

- Don't use tabs, indents, or line feeds/carriage returns

- When copying and pasting from other sources, use a text editor (e.g., Notepad) to eliminate hidden characters
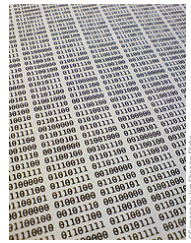
**DataONE**

## 3. Policies for Access, Sharing, Reuse

3.1  Obligations for sharing
- Funding agency
- Institution
- Other organization
- Legal

3.2  Details of data sharing
- How long?
- When?
- How access can be gained?
- Data collector rights

3.2  Ethical/privacy issues with data sharing

**DataONE**