

Data Storage & Backup

LAS 6292: Data Collection & Management
Emilio M. Bruna, University of Florida

1

First some terminology...

Original Data



Operation Data

Backup Data



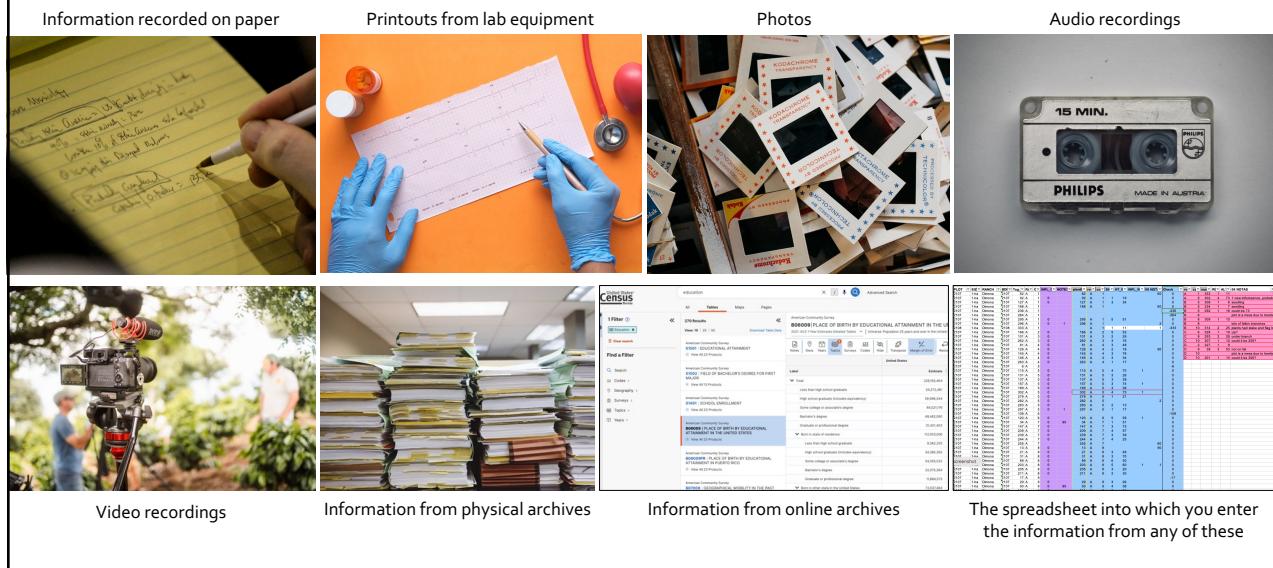
Data Archive



2

Original Data (i.e., 'Raw Data')

Information collected during research (plus the samples from which the info was extracted or the items on which you recorded it)



3

Operation Data (i.e., 'Clean Data')

Original data after it has been Processed, Edited, Corrected, or Transformed.

Spreadsheet after proofreading and correcting values

key	number	source	prop.rem.1	prop.rem.2	prop.rem.3	plot	destination	year
6	12591	High-N	0.489	0.489	0.439	2	Control	2009
7	12258	High-N	0.687	0.687	0.637	2	Control	2009
8	12295	High-N	0.76	0.76	0.71	2	Control	2009
9	9387	High-N	0.827	0.827	0.777	2	Control	2009
80	12352	Low-N	0.429	0.429	0.379	2	Control	2009
81	9403	Low-N	0.656	0.656	0.606	2	Control	2009
82	9466	Low-N	0.775	0.775	0.725	2	Control	2009
83	9481	Low-N	0.884	0.884	0.834	2	Control	2009
84	12609	Low-N	0.892	0.892	0.842	2	Control	2009
155	12771	Control	0.607	0.607	0.557	2	Control	2009
156	12563	Control	0.736	0.736	0.686	2	Control	2009
157	12114	Control	0.772	0.772	0.722	2	Control	2009
158	12601	Control	0.776	0.776	0.726	2	Control	2009
159	12598	Control	0.824	0.824	0.774	2	Control	2009
229	9371	High-N	0.84739379	0.84739379	0.797	2	Control	2010
230	9126	High-N	0.85311909	0.85311909	0.803	2	Control	2010
231	9280	High-N	0.87127609	0.87127609	0.831	2	Control	2010
232	12831	High-N	0.89589064	0.89589064	0.846	2	Control	2010
233	12375	High-N	0.91200673	0.91200673	0.863	2	Control	2010
305	12409	High-N	0.84452011	0.84452011	0.795	2	Control	2010
306	12048	Low-N	0.87297216	0.87297216	0.823	2	Control	2010
307	12807	Low-N	0.88313576	0.88313576	0.833	2	Control	2010
308	12564	Low-N	0.88713004	0.88713004	0.837	2	Control	2010
309	9088	Low-N	0.80986663	0.80986663	0.859	2	Control	2010
381	12630	Control	0.85197849	0.85197849	0.803	2	Control	2010
382	12388	Control	0.85995751	0.85995751	0.81	2	Control	2010
383	12760	Control	0.89649557	0.89649557	0.846	2	Control	2010

Edited and Corrected Transcripts

The standard Lorem Ipsum passage, used since the 1500s
"Lorem ipsum dolor sit amet, consecetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum."
Section 1.10.32 of "de Finibus Bonorum et Malorum", written by Cicero in 45 BC
"Sed ut persicarpit deinde omnis iuste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritas et quasi architecto beatae vitae dicta sunt explicabo. Nemo enim ipsam voluptatem quia voluptas sit aspernatur aut odit aut fugit, sed quia consequuntur magni dolores eos qui ratione voluptatem sequi nesciunt. Neque porro quisquam est, qui dolorem ipsum quia dolor sit amet, consectetur, adipisci velit, sed quia non quamquam eius modi tempora incidunt ut labore et dolore magnam aliquam quaerat voluptatem. Ut enim ad minim veniam, quis nostrum exercitationem ullam corporis suscipit laboriosam, nisi ut aliquid vel ex commodi consequatur? Quis autem vel eum lures per se voluntariae vel esse quam nihil molestiae consequatur, vel illum qui dolorem eum fugit quo voluptas nulla paratur?"
1914 translation by H. Rackham
"But I must explain to you how all this mistaken idea of denouncing pleasure and praising pain was born and I will give you a complete account of the system, and expound the actual teachings of the great explorer of the truth, the master-builder of human happiness. No one rejects, dislikes, or avoids pleasure itself, because it is pleasure, but because those who do not know how to pursue pleasure rationally encounter consequences that are extremely painful. Nor again is there anyone who loves or pursues or desires to obtain pain of itself, because it is pain, but it occurs occasionally circumstances in which toll and pain can procure him some great pleasure. To take a trivial example, which of us ever undertakes laborious physical exercise, except to obtain some advantage from it? But who has any right to find fault with a man who chooses to enjoy a pleasure that has no annoying consequences, or one who avoids a pain that produces no resultant pleasure?"
Section 1.10.33 of "de Finibus Bonorum et Malorum", written by Cicero in 45 BC
"At vero eos et accusamus et iudicato dignissimos delictis ei blandius praesertim voluptatum dilectissimis corrupti quos dolores et quas molestias excepturi sint occaecati cupiditate non provident, similique sunt in cubo qui officia delectus molitla animi, id est laborum et dolorem fugit. Et harum quidem rerum facilis est et expedita distinctio. Nam libero tempore, cum soluta nobis est eligendi optio cumque nihil impedit quo minus id quod maxime placeat facere possimus, omnis voluptatis assumenda est, omnis dolor repellendus. Temporibus autem quibusdam et ut officiis debitis aut rerum necessitatibus saepe eveniet ut et voluptates repudiandae sint et molestiae non recusandae. Itaque earum rerum hic tenetur a sapiente delectus. Ut aut recidens voluptatus maiores alias consequatur aut perferendis doloribus asperiores repellat."

These are the ones on which you carry out 'operations', i.e., analyses.

4

Backup Data

Copies of 'Original' & 'Operation' Data allowing you to resume analyses if files are lost or damaged.

Electronic version stored on your computer, the cloud, or...



...a portable storage unit.



Paper Copies



5

Data Archive

Long-term storage of Original & Operational Data in a format & location where it is:

- (1) secure and (2) can be readily discovered and used by others.

Museum Collections



Library



Online Data Repository

The figure shows a screenshot of the Dryad Data Portal search results. The search term is 'Data from: Asymmetric dispersal and colonization success of Amazonian plant-ants queens'. The results page includes the title, authors (Bruno, Emilia M.; University of Florida; Smithsonian Tropical Research Institute), journal (Tropical Ecology and Systematics), volume (1), issue (1), and date (July 12, 2011). It also lists the DOI (10.5061/dryad.3t11), version (1), and URL (https://doi.org/10.5061/dryad.3t11). The page features sections for 'Data Files', 'Downloaded 1 times', 'Revised Works', 'Article', and 'Metrics'. Metrics include 198 views, 196 downloads, and 2 citations. A 'Keywords' section is also present.

6

Why bother backing up data?

Because bad things happen.
Especially to you.

7



8

Other common causes of data loss include...

9

DEMONIC INTERVENTION	natural disasters
pests	hackers
Cats	
document loss	ACCIDENTAL DELETION
VIRUSES & MALWARE	
software or media obsolescence	
THEFT	file corruption
ACCIDENTS	
human error	POWER FAILURE
HARD DRIVE FAILURE	coffee spilled on laptop

10

These things happen more often and on bigger scales than you think.

11

Insects destroy Idaho county's historic documents

by Associated Press | Wednesday, October 3rd 2012



AP
Audit finds Baltimore workers using outdated data storage

September 28, 2019

U.S. News **World News** **Politics** **Sports** **Entertainment** **Business** **Technology** **Health** **Science** **Additio...**

INDEPENDENT

Pixar's billion-dollar delete button nearly lost Toy Story 2 animation

Gillian Orr • Thursday 17 May 2012 11:05 + [Comments](#)



REUTERS
[World](#) [Business](#) [Markets](#) [Breakfast](#) [Video](#) [More](#)

Millions of websites offline after fire at French cloud services firm

PARIS (Reuters) - A fire at a French cloud services firm has disrupted millions of websites, knocking out government agencies' portals, banks, shops, news websites and taking out a chunk of the .FR web space, according to internet monitors.



12

Why else should we bother backing up data?

Because you are required to.

13

Compliance with: Public Records Requests, lawsuits, funder mandates...



National Institutes
of Health

And finally....

14

Because it makes research and collaboration easier, more efficient, and more productive.



15

OK...so what should we be doing to **back up** our data?

16

3-2-1

3 copies of your data stored on
2 different types of media, with
1 copy stored off-site

17

What does 'off-site' mean?



If you **WORK** in your office

Then your office could be **OFF-SITE**



► home could be **OFF-SITE**

↔ If you **WORK** from home

18

“Different” Storage Media could be:

- Desktop Computer
- Laptop Computer
- External hard drive
- Campus server
- Commercial Cloud Backup
- Optical Storage (CD, DVD)
- Paper

*original notebooks & data sheets,
printouts of spreadsheets,
xerox copies of originals, etc.*



19

Keep in mind that some media can become obsolete.....



20

While others do not.



21

WARNING!!

22

Synchronized cloud storage is *not*
“off-site” or “different”.



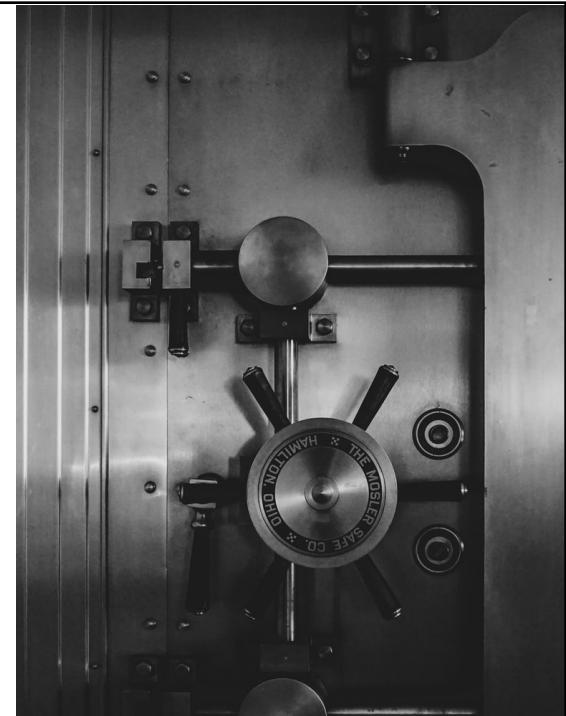
Why?

Deletion from Location 1 **automatically** triggers deletion from Location 2.

23

IMPORTANT: There are special considerations for “**Restricted Data**”

- **Find out** if your data require extra protection, and
- **Get guidance** from your institution on the additional steps and resources needed to protect original and backups of Restricted Data



24

And finally....

To ensure that your data can be read by as many computer programs as possible for many years into the future :

USE OPEN (a.k.a. 'non-proprietary') DATA FORMATS

25

And be sure to...

- Test Your Storage Media Regularly
- Beware of Early Hardware Failures
- Determine the Life of Your Hard Drives
- Routinely Inspect and Replace Data Discs
- Handle and Store Your Media With Care



26

Eerlijk gezegd is
gaangvloebement?



27

Make a plan. Write it down.

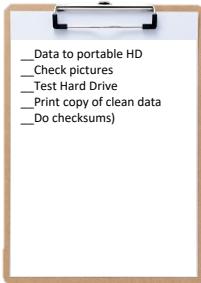


- what needs to be prepared and for how long
- where backups are located
- who can access backups and how contacted
- how often data should be backed up
- what kind of backups are performed
- who is responsible for performing the backups
- hardware and software used for performing backups
- how / how often to check if backups are successful
- the media are used to backup data
- a list of any data that are *not* archived or backed up

Then be sure to...

28

1. Make it easy to follow your plan
2. Automate as much of the process as you can.
3. Revisit & Revise your plan regularly



*And remember: synced computer and cloud drive are **not** independent copies*

29

The
End

(off-site)

Photos: Unsplash or EM Bruna
Music: www.bensound.com

30

What about the The Originals?

Consider the volume and security needs – both short-term & long-term – of specimens, samples, documents and data sheets, photographs, or other physical items. Depending on the situation you might use:

- Campus Office
- Lab or Museum/Herbarium
- Departmental Office
- Commercial storage facility
- Home
- Library



31

Photo, video, and audio backup:

More challenging, in part because transcriptions & captioning can be important for interpretation, discovery, and accessibility.



Be sure to back-up both the original files *and* the transcriptions.

Think ahead: storage can also be challenging because the files are large and there are often large numbers of them.

32