

Managing Research Data Effectively

Justin du Toit
Groofofontein Agricultural Development Institute

- From field data-sheets to electronic files
- Data entry and storage
- Managing data in Excel
- Managing visual data (photographs)


Entry format (base format)

Date	1 10 2008			
Site	Short grass			
		Number of grass plants	Number of flowering culms	Wattle seedlings
Quadrat 1	3.7 cm	6	0	12
Quadrat 2	5 cm	2	1	8 (11 dead)
Quadrat 3	5.3 cm	2	1	5
Quadrat 4	4.4 cm	2	2	8
Quadrat 5	6.9 cm	3	3	0
Date	1 10 2008			
Site	Tall grass			
		Number of grass plants	Number of flowering culms	Wattle seedlings
Quadrat 1	9 cm	2	2	0
Quadrat 2	11 cm	2	4	0
Quadrat 3	12.3 cm	Not recorded	0	3
Quadrat 4	8 cm	0	0	3
Quadrat 5	8 cm	4	4	1 (dead)

- Computers view data differently from the way we do
- Computers can analyse data *incredibly* fast, but it must be in the correct format
- Data sheets are laid out in a way that makes sense to us (see left)
- Data files must be laid out in a way that makes sense to a computer

Entry format (base format)

Date	1 10 2008			
Site	Short grass			
		Number of grass plants	Number of flowering culms	Wattle seedlings
Quadrat 1	3.7 cm	6	0	12
Quadrat 2	5 cm	2	1	8
Quadrat 3	5.3 cm	2	1	5
Quadrat 4	4.4 cm	2	2	8
Quadrat 5	6.9 cm	3	3	0
Date	1 10 2008			
Site	Tall grass			
		Number of grass plants	Number of flowering culms	Wattle seedlings
Quadrat 1	9 cm	2	2	0
Quadrat 2	11 cm	2	4	0
Quadrat 3	12.3 cm	Not recorded	0	3
Quadrat 4	8 cm	0	0	3
Quadrat 5	8 cm	4	4	1 (dead)



Dual entry



- To enter a list of 100 values takes between 0.01 and 0.1% of the time that it takes to run an experiment (person working-hours, not duration of the experiment)
- However, vested within these numbers is the *entire cost* of the experiment (usually hundreds of thousands of rands)
- Therefore, it is *vital* to ensure that the data are captured perfectly
- By far the most reliable way of doing this is for two different people to enter the data (having the same person enter it twice leads to unhealthy temptation with the Copy | Paste function!)

Spreadsheet or database?

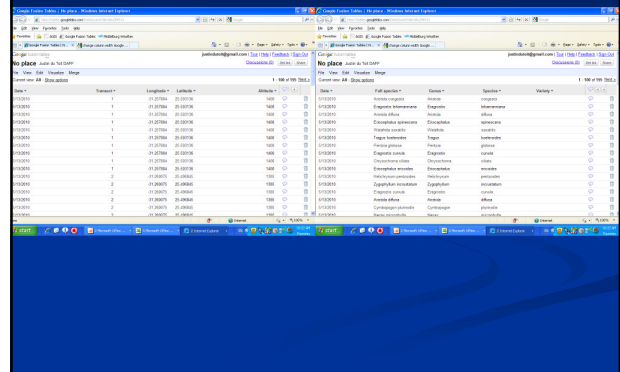
- Spreadsheets
 - Easy to use
 - Suitable for most applications
 - Lack of query capacity (hence less room for errors...)
 - Good analysis and illustration capabilities
 - Limited to 65536 rows x 256 columns
- Databases
 - Require expertise to run correctly
 - VERY EASY to ask for the wrong thing, and hence end up with junk data
 - Vital for large data sets, especially ongoing ones, and where >1 people enter data
 - Unlimited data storage

Backing up data

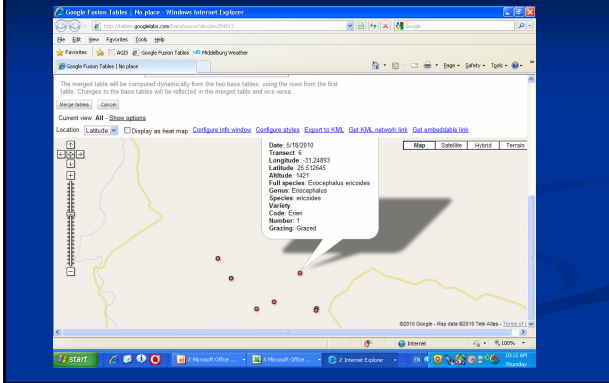
- Why backup? To have a copy of the data that is free from the hazards your original copy might face, e.g.:
 - Theft, a computer crash, losing the storage device, office burns down, etc.
- Therefore, the most commonly backups face exactly the same hazards as the original
 - Flashdrives, CDs, secondary hard-drives, laptops, etc
- The trick is to store data far, far away, and the easiest way to do that is electronically – set up a few internet-based email accounts, and Search & Send it to these regularly, and using e.g. Google Fusion
 - (Similarly, photograph or scan lab- or field-data sheets)



Example – Google Fusion



Example – Google Fusion



Make your life easier in Excel! Learn the basic functions!

- Basic arithmetic functions (=, -, x, /, sum, average, count)
- Conditionally manipulate data (=if(CONDITION, TRUE, FALSE))
- Simplify data – e.g. convert continuous data to binomial (=if(VALUE>x, 1, 0))
- Split data into separate records e.g. “Eragrostis curvula” into two columns titled “Genus” and “Species”
- Transform data – e.g. (=log(VALUE+1))
- Truncate text e.g. (=left(“Eragrostis”, 3)) gives “Era” – combined with the same for “curvula” (cur) can be joined (=“Era”&“cur”) to give a code **Eracur**
- Look up information about a common value or term (=vlookup(VALUE, TABLE, COLUMN))
- You NEVER have to re-enter values, and NEVER have to use a calculator in Excel
- Remove unwanted spaces (=trim(TEXT))

	A	B	C	D
1	SpeciesFull	Genus	Species	Code
2	Gnidia polycephala	Gnidia	polycephala	Gniala
3	Helichrysum cerastoides	Helichrysum	cerastoides	Heldes
4	Helichrysum zeyheri	Helichrysum	zeyheri	Heleri
5	Hermannia coccocarpa	Hermannia	coccocarpa	Herrpa
6	Hertia pallens	Hertia	pallens	Herens
7	Heteropogon contortus	Heteropogon	contortus	Hettus

Original names

Split using Data | Text to Columns

6-letter code (1st 3 letters of genus, 1st 3 letters of species) using: =left(A2,3)&right(B2,3)

	A	B	C	D
1	SpeciesFull	RelAb	>10%?	log(RelAb)
2	Helichrysum zeyheri	16.2	1	1.21
3	Hermannia coccocarpa	7.7	0	0.89
4	Heteropogon contortus	31.2	1	1.49
5	Hyparrhenia hirta	22.6	1	1.35
6	Lepidium africanum	2.3	0	0.36
7	Limeum aethiopicum	13.9	1	1.14
8	Oropetium capense	7.6	0	0.88

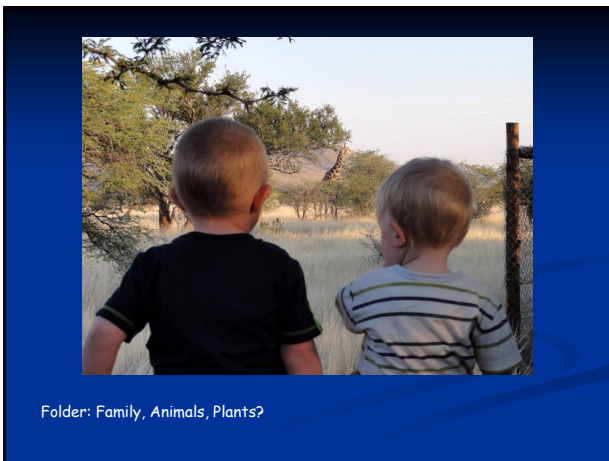
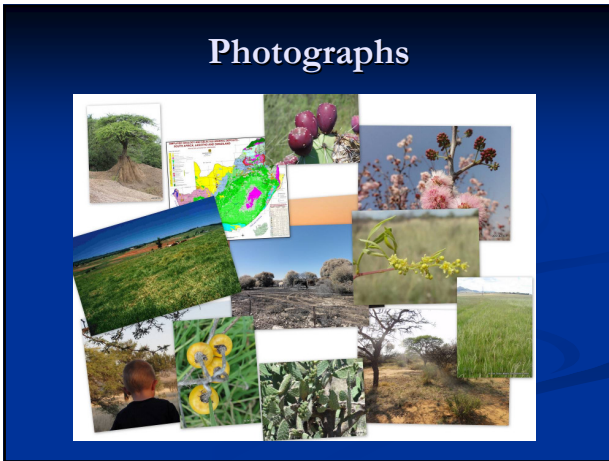
Data

Conditional manipulation: 1 if greater than 10, otherwise 0
=if(B2>10,1,0)

Arithmetic conversion
Log of original data
=log(B2)

	A	B	C	D	E	F	G	H
1	Date	Rain (mm)	Cumulative rainfall	Rain?	Cumulative sequential wet days	Cumulative sequential dry days	Cumulative wet days	Cumulative dry days
2	01-Jan	0	0	0	0	1		1
3	02-Jan	4	4	1	1	0		
4	03-Jan	33	37	1	2	0		
5	04-Jan	14	51	1	3	0		
6	05-Jan	4	55	1	4	0	4	
7	06-Jan	0	55	0	0	1		
8	07-Jan	0	55	0	0	2		
9	08-Jan	0	55	0	0	3		3
10	09-Jan	2	57	1	1	0		
11	10-Jan	8	65	1	2	0	2	
12	11-Jan	0	65	0	0	1		
13	12-Jan	0	65	0	0	2		
14	13-Jan	0	65	0	0	3		
15	14-Jan	0	65	0	0	4		4

=B3+C2	=IF(B3>0,1,0)	=IF(D3=1,E2+1,0)	=IF(D3=0,E2+1,0)	=IF(E3>E4,E3,"")	=IF(F3>E4,E3,"")
--------	---------------	------------------	------------------	------------------	------------------



Rather use tagging

