

Data normalisation and formats

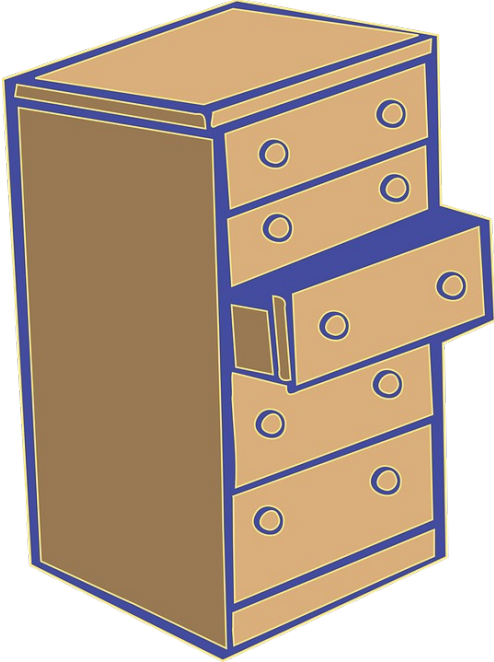
Sarah Faulwetter & Pieter Provoost

A Practical introduction to acquisition, validation,
quality control and access to (biodiversity) data

Vlissingen, Netherlands, 14.06.2015



Purpose



Storage & sharing



Analyses



Format depends on software,
type of analysis, etc.



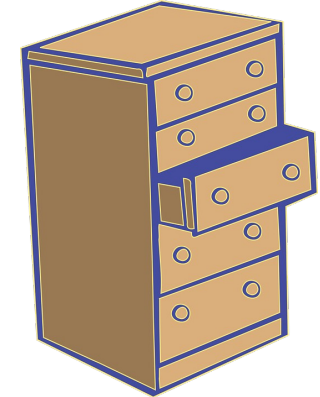
Format depends on software,
type of analysis, etc.

	A	B	C	D	E	F	
1	<u>scientificName</u>	Cardigan Bay	Celtic Deep	Central Channel	Humber	Irish Sea	Li
2	<u>Abra alba</u>	23	159	9	25		
3	<u>Abra nitida</u>						
4	<u>Abra prismatica</u>						
5	<u>Achelia echinata</u>						
6	<u>Acidostoma obesum</u>						
7	<u>Acidostoma sarsi</u>						
8	<u>Acteon tornatilis</u>						
9	<u>Aega crenulata</u>			7			
10	<u>Aequipecten opercularis</u>			4			
11	<u>Aglaophamus rubella</u>					2	
12	<u>Ammodytes marinus</u>						
13	<u>Ampelisca brevicornis</u>						
14	<u>Ampelisca diadema</u>						
15	<u>Ampelisca macrocephala</u>		3				
16	<u>Ampelisca spinipes</u>			10			
17	<u>Ampelisca tenuicornis</u>	1	6				

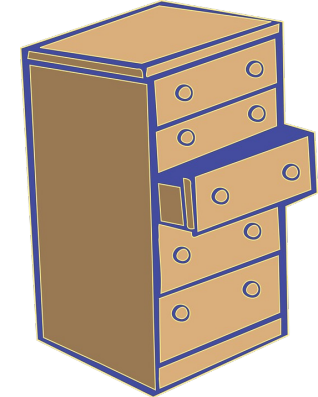


Format depends on software,
type of analysis, etc.

	A	B	C	D	E	F	
1	<u>scientificName</u>	Cardigan Bay	Celtic Deep	Central Channel	Humber	Irish Sea	Li
2	<u>Abra alba</u>	23	159		9	25	
3	<u>Abra nitida</u>						
4	<u>Abra prismatica</u>						
5	limited to three types of information						
6							
7							
8							
9							
10							
11	<u>Aglaophamus rubella</u>						2
12	<u>Ammodytes marinus</u>						
13	<u>Ampelisca brevicornis</u>						
14	<u>Ampelisca diadema</u>						
15	<u>Ampelisca macrocephala</u>			3			
16	<u>Ampelisca spinipes</u>				10		
17	<u>Ampelisca tenuicornis</u>	1	6				



taxonNameAsInFile	samplingLevel	value
Hydrobia	subtidal	1
Macoma balthica	subtidal	2
Mytilus	subtidal	1
Bathyporeia pilosa	subtidal	9
Polychaeta	subtidal	34
Mytilus	subtidal	1
Polychaeta	subtidal	10
Mytilus	subtidal	1
Bathyporeia pilosa	subtidal	36
Polychaeta	subtidal	23
Polychaeta	subtidal	14
Hydrobia	subtidal	3
Bathyporeia pilosa	subtidal	19
Polychaeta	subtidal	9
Hydrobia	subtidal	1
Bathyporeia pilosa	subtidal	20
Neomysis integer	subtidal	1
Polychaeta	subtidal	35
Hydrobia	subtidal	1
Bathyporeia pilosa	subtidal	16
Polychaeta	subtidal	20
Nephtys	intertidal	1
Oligochaeta	intertidal	1



taxonNameAsInFile	samplingLevel	value	sampleArea	samples.unit	replicateC	samplingMinDepth	sampledSedimentD
Hydrobia	subtidal	1	133	cm2	1	0.5	30
Macoma balthica	subtidal	2	133	cm2	1	0.5	30
Mytilus	subtidal	1	133	cm2	1	0.5	30
Bathyporeia pilosa	subtidal	9	133	cm2	1	0.5	30
Polychaeta	subtidal	34	133	cm2	1	0.5	30
Mytilus	subtidal	1	133	cm2	2	0.5	30
Polychaeta	subtidal	10	133	cm2	2	0.5	30
Mytilus	subtidal	1	133	cm2	3	0.5	30
Bathyporeia pilosa	subtidal	36	133	cm2	3	0.5	30
Polychaeta	subtidal	23	133	cm2	3	0.5	30
Polychaeta	subtidal	14	133	cm2	1	0.5	30
Hydrobia	subtidal	3	133	cm2	1	0.5	30
Bathyporeia pilosa	subtidal	19	133	cm2	1	0.5	30
Polychaeta	subtidal	9	133	cm2	1	0.5	30
Hydrobia	subtidal	1	133	cm2	2	0.5	30
Bathyporeia pilosa	subtidal	20	133	cm2	2	0.5	30
Neomysis integer	subtidal	1	133	cm2	2	0.5	30
Polychaeta	subtidal	35	133	cm2	2	0.5	30
Hydrobia	subtidal	1	133	cm2	2	0.5	30
Bathyporeia pilosa	subtidal	16	133	cm2	2	0.5	30
Polychaeta	subtidal	20	133	cm2	2	0.5	30
Nephtys	intertidal	1	133	cm2	2		30
Oligochaeta	intertidal	1	133	cm2	2		30



“

What's the most used tool for data sharing in biology?

Excel.

What's the most hated tool for data sharing?

Excel.

— Bruce Wilson, DataONE

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	
1	Mixed	Data based on the 1981-2018 FYE categories for 12 data. Certain excluded from measured individuals in MIPF category (both MIPF and only not)														
2	Non categories	After from Matrix in office, use of copy from Regis to Brown														
3		Total Douglas														
4	Spring (M-F)				Non-Road	Other	Transferable	Transferable	Transferable	Transferable	Transferable	Transferable	Transferable	Transferable	Transferable	Transferable
5	Summer (J-A)				0.00	0.04									1000	10.00
6	Autumn (S-N)				0.00	0.04										
7	Winter (D-F)				0.00	0.04										
8																
9																
10																
11																
12																
13																
14																
15																
16																
17																
18																
19																
20																
21																
22																
23																
24																
25																
26																
27																
28																
29																
30																
31																
32																
33																
34																
35																
36																
37																
38																
39																
40																
41																
42																

source: @tomjwebb 

1

do	merge cells
not	

1

do	merge cells
not	

2

do not	use
colour	coding

1

do	merge cells
not	

2

do not	use
colour	coding

3

use "talking" column headers

st.	ab_t	d	s
B1	12	0.33	4
B2	5	0.135	6



station	abundanceTotal	density	nrOfSpecies
B1	12	0.33	4
B2	5	0.135	6

1

do	merge cells
not	

2

do not	use
colour	coding

3

use "talking" column headers

st.	ab_t	d	s
B1	12	0.33	4
B2	5	0.135	6



station	abundanceTotal	density	nrOfSpecies
B1	12	0.33	4
B2	5	0.135	6

4

include units:

B	C
Temp_C	Oxygen_%
12	86.9

1

do	merge cells
not	

2

do not	use
colour	coding

3

use "talking" column headers

st.	ab_t	d	s
B1	12	0.33	4
B2	5	0.135	6



station	abundanceTotal	density	nrOfSpecies
B1	12	0.33	4
B2	5	0.135	6

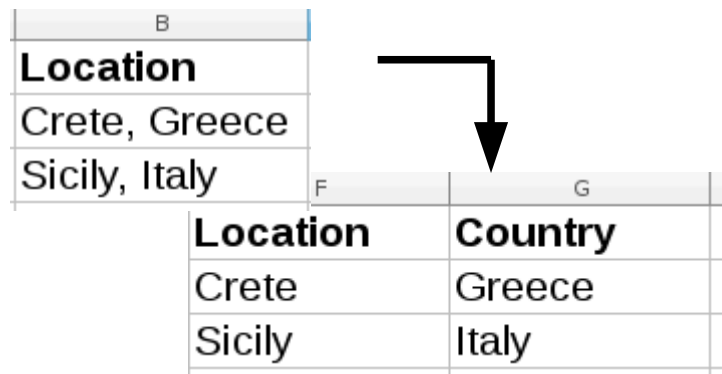
4

include units:

B	C
Temp_C	Oxygen_%
12	86.9

5

Only one piece of information per cell



1

do	merge cells
not	

2

do not	use
colour	coding

3 use "talking" column headers

st.	ab_t	d	s
B1	12	0.33	4
B2	5	0.135	6



station	abundanceTotal	density	nrOfSpecies
B1	12	0.33	4
B2	5	0.135	6

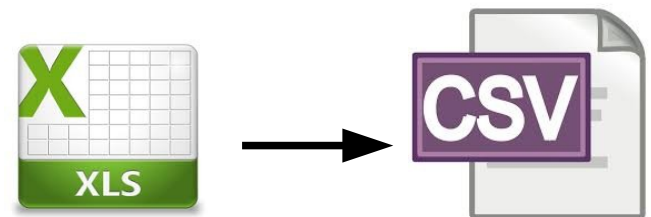
4 include units:

B	C
Temp_C	Oxygen_%
12	86.9

5 Only one piece of information per cell

B	F	G
Location	Location	Country
Crete, Greece	Crete	Greece
Sicily, Italy	Sicily	Italy

6 final version not in proprietary format



atomization

Name	Contact details
Robert Ingram	robert@ingram.com , New York, Tel. 555-861-2025
Jane Wright	jwright@gmail.com , New York, Tel. 555-861-2033
Maria Fernandez	mfernandez@gmail.com , Chicago, Tel. 393-3345-2235

Name	Contact details
Robert Ingram	robert@ingram.com , New York, Tel. 555-861-2025
Jane Wright	jwright@gmail.com , New York, Tel. 555-861-2033
Maria Fernandez	mfernandez@gmail.com , Chicago, Tel. 393-3345-2235

Sort by last name

All employees from New York

→ Problem

Name	Contact details
Robert Ingram	robert@ingram.com , New York, Tel. 555-861-2025
Jane Wright	jwright@gmail.com , New York, Tel. 555-861-2033
Maria Fernandez	mfernandez@gmail.com , Chicago, Tel. 393-3345-2235



First name	Last name	e-mail	City	Phone Number
Robert	Ingram	robert@ingram.com	New York	555-861-2025
Jane	Wright	jwright@gmail.com	New York	555-861-2033
Maria	Fernandez	mfernandez@gmail.com	Chicago	393-3345-2235

First name	Last name	e-mail	City	Phone Number
Robert	Ingram	robert@ingram.com	New York	555-861-2025
Jane	Wright	jwright@gmail.com	New York	555-861-2033
Maria	Fernandez	mfernandez@gmail.com , mfernandez@company.com	Chicago	393-3345-2235

First name	Last name	e-mail	City	Phone Number
Robert	Ingram	robert@ingram.com	New York	555-861-2025
Jane	Wright	jwright@gmail.com	New York	555-861-2033
Maria	Fernandez	mfernandez@gmail.com, mfernandez@company.com	Chicago	393-3345-2235



Bad format!

First name	Last name	e-mail	City	Phone Number
Robert	Ingram	robert@ingram.com	New York	555-861-2025
Jane	Wright	jwright@gmail.com	New York	555-861-2033
Maria	Fernandez	mfernandez@gmail.com	Chicago	393-3345-2235
Maria	Fernandez	mfernandez@company.com	Chicago	393-3345-2235

First name	Last name	e-mail	City	Phone Number
Robert	Ingram	robert@ingram.com	New York	555-861-2025
Jane	Wright	jwright@gmail.com	New York	555-861-2033
Maria	Fernandez	mfernandez@gmail.com	Chicago	393-3345-2235
Maria	Fernandez	mfernandez@company.com	Chicago	393-3345-2235

Bad format!
(Duplicate information)

First name	Last name	e-mail	City	Phone Number
Robert	Ingram	robert@ingram.com	New York	555-861-2025
Jane	Wright	jwright@gmail.com	New York	555-861-2033
Maria	Fernandez	mfernandez@gmail.com	Chicago	393-3345-2235
Maria	Fernandez	mfernandez@company.com	Chicago	393-3345-2235

Bad format!
(Duplicate information)

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Species/Station	I- 50	I-60	I-70	I-72	I-75	I-77	I-79	I-81	I-83	I-87	I-89	I-91
2													
3	PORIFERA												
4	<u>Axinella domicornis</u>	0	0	1	0	0	0	0	0	0	0	0	0
5	<u>Suberites carnosus</u>	0	0	0	0	0	0	0	0	0	0	0	0
6	'-----												
7	HYDROZOA												
8	<u>Laomedea angulata</u>	0	0	0	0	0	0	0	0	0	0	0	0
9	<u>Obelia geniculata</u>	0	0	0	0	0	0	0	0	0	0	0	0
10	<u>Podocoryna carnea</u>	0	0	0	0	0	0	0	0	0	0	0	0
11	'-----												
12	ANTHOZOA												
13	<u>Amphianthus dohrnis</u>	1	0	0	0	0	0	0	1	0	1	0	0
14	<u>Edwardsia calimorpha</u>	0	0	0	0	0	0	0	0	0	0	0	0
15	<u>Sarcodictyon coralloides</u>	0	0	0	0	0	0	0	0	0	0	0	0
16	'-----												
17	TUBELLARIA												
18	<u>Stylochus pilidium</u>	0	0	0	0	0	0	0	0	0	0	0	0
19	'-----												
20	NEMERTINEA												
21	<u>Sp. 1</u>	0	0	2	0	0	0	0	0	0	0	2	0
22	<u>Sp. 2</u>	0	2	3	0	0	0	0	0	0	0	0	0
23	<u>Sp. 3</u>	0	0	0	0	0	0	0	0	0	0	0	0
24	'-----												
25	POLYCHAETA												
26	<u>Arabella iricolor</u>	0	0	0	1	0	0	0	0	0	0	0	0
27	<u>Drilonereis filum</u>	0	0	0	0	0	0	0	0	0	0	0	0
28	<u>Eteone siphonodonta</u>	0	0	1	0	0	0	2	0	0	0	0	0

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Species/Station	I-50	I-60	I-70	I-72	I-75	I-77	I-79	I-81	I-83	I-87	I-89	I-91
2													
3	PORIFERA												
4	<u>Axinella domicornis</u>	0	0	1	0	0	0	0	0	0	0	0	0
5	<u>Suberites carnosus</u>	0	0	0	0	0	0	0	0	0	0	0	0
6	-----												
7	HYDROZOA												
8	<u>Laomedea angulata</u>	0	0	0	0	0	0	0	0	0	0	0	0
9	<u>Obelia geniculata</u>	0	0	0	0	0	0	0	0	0	0	0	0
10	<u>Podocoryna carnea</u>	0	0	0	0	0	0	0	0	0	0	0	0
11	-----												
12	ANTHOZOA												
13	<u>Amphianthus dohrnis</u>	1	0	0	0	0	0	0	1	0	1	0	0
14	<u>Edwardsia calimorpha</u>	0	0	0	0	0	0	0	0	0	0	0	0
15	<u>Sarcodictyon coralloides</u>	0	0	0	0	0	0	0	0	0	0	0	0
16	-----												
17	TUBELLARIA												
18	<u>Stylochus pilidium</u>	0	0	0	0	0	0	0	0	0	0	0	0
19	-----												
20	NEMERTINEA												
21	<u>Sp. 1</u>	0	0	2	0	0	0	0	0	0	0	2	0
22	<u>Sp. 2</u>	0	2	3	0	0	0	0	0	0	0	0	0
23	<u>Sp. 3</u>	0	0	0	0	0	0	0	0	0	0	0	0
24	-----												
25	POLYCHAETA												
26	<u>Arabella iricolor</u>	0	0	0	1	0	0	0	0	0	0	0	0
27	<u>Drilonereis filum</u>	0	0	0	0	0	0	0	0	0	0	0	0
28	<u>Eteone siphonodonta</u>	0	0	1	0	0	0	2	0	0	0	0	0

1: "Headers" in file 2: graphical elements, not data 3: Data only understandable in combination with the "header" 4: White spaces in or before taxon names

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Species/Station	I-50	I-60	I-70	I-72	I-75	I-77	I-79	I-81	I-83	I-87	I-89	I-91
2												
3												
4												
5												
6												
7	<u>Amphianthus dohrnis</u>	1	0	0	0	0	0	0	1	0	1	0	0
8	ANTHOZOA												
9	<u>Arabella iricolor</u>	0	0	0	1	0	0	0	0	0	0	0	0
10	<u>Axinella domicornis</u>	0	0	1	0	0	0	0	0	0	0	0	0
11	<u>Drilonereis filum</u>	0	0	0	0	0	0	0	0	0	0	0	0
12	<u>Edwardsia calimorpha</u>	0	0	0	0	0	0	0	0	0	0	0	0
13	<u>Eteone siphonodonta</u>	0	0	1	0	0	0	2	0	0	0	0	0
14	HYDROZOA												
15	<u>Laomedea angulata</u>	0	0	0	0	0	0	0	0	0	0	0	0
16	NEMERTINEA												
17	<u>Obelia geniculata</u>	0	0	0	0	0	0	0	0	0	0	0	0
18	<u>Podocoryna carnea</u>	0	0	0	0	0	0	0	0	0	0	0	0
19	POLYCHAETA												
20	PORIFERA												
21	<u>Sarcodictyon coralloides</u>	0	0	0	0	0	0	0	0	0	0	0	0
22	<u>Sp. 1</u>	0	0	2	0	0	0	0	0	0	0	2	0
23	<u>Sp. 2</u>	0	2	3	0	0	0	0	0	0	0	0	0
24	<u>Sp. 3</u>	0	0	0	0	0	0	0	0	0	0	0	0
25	<u>Stylochus pilidium</u>	0	0	0	0	0	0	0	0	0	0	0	0
26	<u>Suberites carnosus</u>	0	0	0	0	0	0	0	0	0	0	0	0
27	TUBELLARIA												
28													

Same file, sorted alphabetically