



Instituto de Investigación Sanitaria
del Principado de Asturias

Los Diez Mandamientos de una buena base de datos

Patricio Suárez Gil
Coordinador Plataforma Bioestadística
y Epidemiología

Fuente:

<https://rtask.thinkr.fr/blog/the-ten-commandments-for-a-well-formatted-database/>

Commandement 1 : Tu feras tenir toutes les données dans un seul tableau

Commandment 1: all your data shall fit into one single dataframe

Spreading your data over several spreadsheets or several files is **not** an option. All your data should fit in one single dataframe.

Commandement 2 : Tu respecteras
une mise en page précise

Commandment 2: Thou shalt respect a precise formatting

Commandement 3 : Une ligne = un individu statistique

Commandment 3: A line = a statistical individual

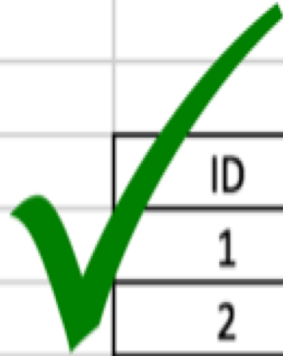
It is not always easy to define what a statistical individual is. In a nutshell, **an answer to a survey corresponds to one individual** in your experimental protocol. If a physical person answers two times, you can enter these on two lines.

Commandement 4 : Une colonne
= une variable

Commandment 4: A column = a variable



ID	Man	Woman
1	yes	no
2	no	yes
3	no	yes
4	yes	no
5	yes	no
6	yes	no
7	no	yes



ID	Gender
1	Man
2	Woman
3	Woman
4	Man
5	Man
6	Man
7	Woman



ID	Source of pain
1	foot and head
2	head and elbow
3	elbow
4	
5	foot and head
6	foot
7	elbow
8	no pain



ID	pain	source_foot	source_head	source_elbow
1	yes	yes	yes	no
2	yes	no	yes	yes
3	yes	no	no	yes
4	yes			
5	yes	yes	yes	no
6	yes	yes	yes	no
7	yes	no	no	yes
8	no			



ID	pain	source_foot	source_head	source_elbow
1	yes	x	x	
2	yes		x	x
3	yes			x
4	yes			
5	yes	x	x	
6	yes	x	x	x
7	yes			x
8	no			


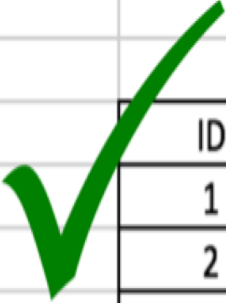


ID	pain	source_foot	source_head	source_elbow
1	yes	yes	yes	no
2	yes	no	yes	yes
3	yes	no	no	yes
4	yes			
5	yes	yes	yes	no
6	yes	yes	yes	no
7	yes	no	no	yes
8	no			

Commandement 5 : Tu ne coderas pas tes variables qualitatives

Commandment 5: Thou shalt not encode thy qualitative variables

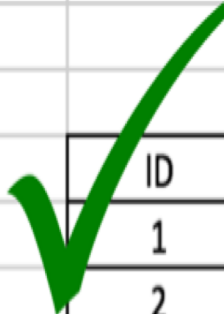

Avoid as much as possible encoding qualitative variables: a 1/2 coding for the gender is useless, prefer the use of "woman" and "man" for a better readability. If you really can't prevent yourself from using a code, then stick to 1/0 to mean presence/absence or yes/no. Not more !

	ID	disease			ID	disease	
	1	1			1	cold	
	2	1			2	cold	
	3	2			3	flu	
	4	2			4	flu	
	5	1			5	cold	
	6	2			6	flu	
	7	2			7	flu	

Commandement 6 : Ta base ne contiendra que les données

Commandment 6: Thy database shall only contain data

For Excel lovers: don't code your variables using colors or bold/italic/underlined. No "control in blue, test case in red": nothing worse than such a format to lose information when exporting those data.

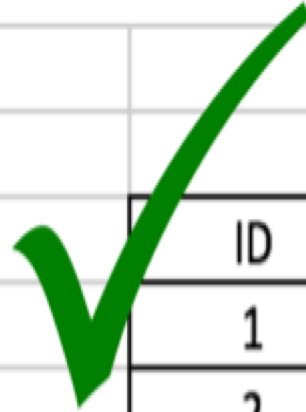


ID	Age
1	0,25
2	3
3	2
4	4
5	4,5
6	18
7	
8	

ID	Age	Status
1	0,25	healthy
2	3	healthy
3	2	sick
4	4	sick
5	4,5	healthy
6	18	sick
7		sick
8		healthy



ID	Age
1	0,25
2	3
3	2 (not sure)
4	4
5	4,5
6	18
7	
8	



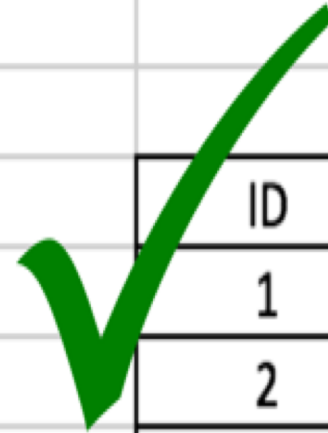
ID	Age	comment
1	0,25	
2	3	
3	2	not sure of the age
4	4	
5	4,5	
6	18	
7		
8		

Commandement 7 : Homogène, tu resteras

Commandment 7: Homogeneous thou shalt be



ID	Gender
1	Masculine
2	Woman
3	woman
4	M
5	Man
6	Man
7	Woman



ID	Gender
1	Man
2	Woman
3	Woman
4	Man
5	Man
6	Man
7	Woman

Commandement 8 : Tu respecteras tes variables numériques

Commandment 8: Thy numerical variables with respect thou shalt treat



ID	Age
1	3 months
2	3 years
3	between 1 and 3 years
4	4
5	4,5
6	18
7	??
8	-



ID	Age
1	0,25
2	3
3	2
4	4
5	4,5
6	18
7	
8	



ID	Date of birth
1	01-Nov-86
2	Thirteenth of April 1985
3	01/01/85
4	14/08/04
5	01/01/08



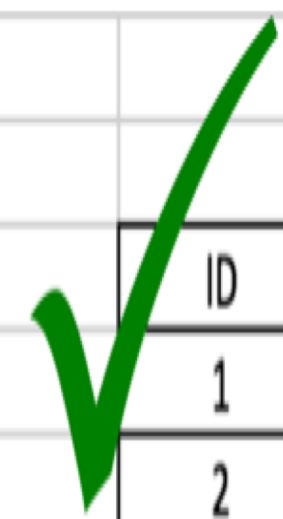
ID	Date of birth
1	1986/11/01
2	1985/04/13
3	1985/01/01
4	2004/08/14
5	2008/01/01

Commandement 9 : Tu anonymiseras ta base

Commandment 9: Anonymous thy database thou shalt keep



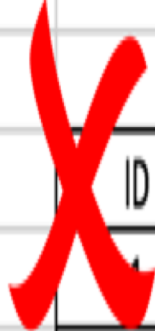
First name	Surname	Score	test
Alan	Smith	45	1
Kevin	Hasseloff	23	1
Helen	Marple	38	1
Alan	John	42	2
Judie	Roger	60	2



ID	Score	test
1	45	1
2	23	1
3	38	1
4	42	2
5	60	2

Commandement 10 : Tu garderas ta
base compréhensible par l'humain

Commandment 10: Human readable thy database shall be



ID	day session	score session	day session	score session
1	Monday	23	Thursday	56
2	Monday	54	Friday	43
3	Monday	12	Tuesday	56
4	Tuesday	23	Wednesday	89
5	Monday	56	Thursday	32
6	Monday	87	Thursday	34
7	Tuesday	45	Wednesday	5
8	Tuesday	3	Friday	17



ID	day Session 1	score Session 1	day Session 2	score Session 2
1	Monday	23	Thursday	56
2	Monday	54	Friday	43
3	Monday	12	Tuesday	56
4	Tuesday	23	Wednesday	89
5	Monday	56	Thursday	32
6	Monday	87	Thursday	34
7	Tuesday	45	Wednesday	5
8	Tuesday	3	Friday	17

Algo más...

- Hacer una copia en formato texto (.txt, .csv ...)
- Cuidado con los nombres de las variables: evitar símbolos excepto punto (.) y guión bajo (_)

Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA

ABSTRACT

Spreadsheets are widely used software tools for data entry, storage, analysis, and visualization. Focusing on the data entry and storage aspects, this article offers practical recommendations for organizing spreadsheet data to reduce errors and ease later analyses. The basic principles are: be consistent, write dates like YYYY-MM-DD, do not leave any cells empty, put just one thing in a cell, organize the data as a single rectangle (with subjects as rows and variables as columns, and with a single header row), create a data dictionary, do not include calculations in the raw data files, do not use font color or highlighting as data, choose good names for things, make backups, use data validation to avoid data entry errors, and save the data in plain text files.

ARTICLE HISTORY

Received June 2017
Revised August 2017

KEYWORDS

Data management; Data organization; Microsoft Excel; Spreadsheets



@patricsg



@patri_1871



www.ispasturias.es

Plataforma de Bioestadística y Epidemiología

Instituto de Investigación Sanitaria del Principado de Asturias (ISPA)

Fundación para la Investigación e Innovación Biosanitaria del Principado
de Asturias (FINBA)



Instituto de Investigación Sanitaria
del Principado de Asturias



FINBA
Fundación para la Investigación y la Innovación
Biosanitaria del Principado de Asturias