

Institute for Health and Society, Medical College of Wisconsin



Datum

NEWSLETTER
Division of Biostatistics



Data entry in Excel

Features and tips to reduce errors

Aniko Szabo, PhD

Introduction

Microsoft Excel remains a popular data entry tool, despite the availability of more reliable alternatives, such as RedCAP (MCW is a RedCAP partner). The main difference between dedicated data collection tools and Excel is that the latter does not help the user to enter the data consistently and in a form that is convenient for future analysis. For example, RedCAP always shows the name of the variable that is being entered right next to the entry box, and it would not let the user enter characters for a numeric outcome or an invalid date. It will also export data in a format that can be conveniently imported into statistical software. However, with self-restraint and use of some of less-known Excel features, the reliability and ease of use of Excel can be greatly improved. In this article, we will provide some tips on doing that. The screenshots use Excel 2010, but other versions work similarly.

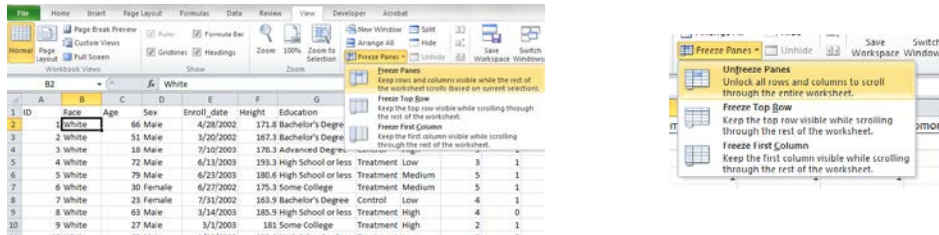
Data arrangement

In the *Database Basics* article (Datum Vol 20 No.1), Dan Eastwood described the general principles of setting up a data collection instrument in Excel. We always recommend that data be set up in a rectangular format, with rows representing separate observations/subjects and the columns representing variables. The first row should contain short, but descriptive variable names; more detailed explanations about each variable, including units, coding conventions, etc. can be placed on a separate "Data Dictionary" worksheet. Each row should be self-contained, so rearranging the order of the rows (by sorting, for example) should not result in any loss of information. As an example, information about the treatment group should be repeated in every row and not just once as a heading or a merged cell. Cell shading and text formatting (colors, bolding, etc.) should not be used to convey information as these attributes

are lost after importing the data into statistical software; instead use a separate column with a text or numeric indicator.

Seeing the variable names and subject IDs

If the data set has more than a few subjects and variables, it is often difficult to tell which subject and which variable does any given cell refer to, especially if the first row or the first column have “scrolled” out of view (as in Figure 1a). The “freeze panes” feature allows the user to specify that some rows and/or columns should remain visible when scrolling. In Figure 1b



this feature has been activated, so that the first column with patient IDs and the first row with variable names are always visible (separated by a thin black line), so it is easy to see that the selected cell refers to the Comorbidity_2 variable of subject #298.

	G	H	I	J	K	L	M	N
289	Some College	Treatment	Low	1	1	1	1	1
290	High School or less	Control	Low	1	1	1	1	0
291	High School or less	Treatment	High	1	1	1	1	1
292	Some College	Control	High	3	0	1	0	1
293	Advanced Degree	Control	Medium	2	1	1	1	0
294	Advanced Degree	Treatment	Medium	5	1	1	1	0
295	Advanced Degree	Control	Low	5	1	1	1	1
296	High School or less	Control	Medium	3	1	1	1	0
297	Advanced Degree	Control	Medium	2	0	1	1	1
298	Bachelor's Degree	Treatment	Medium	3	1	1	1	1
299	High School or less	Control	Low	2	1	1	1	1
300	High School or less	Control	Medium	4	1	0	0	1
301	High School or less	Treatment	Medium	5	1	1	1	1
302	High School or less	Treatment	Medium	2	0	1	1	1
303	Some College	Treatment	Low	2	1	1	1	1
304	Some College	Control	Low	4	1	1	0	0
305	Some College	Treatment	Medium	5	0	1	1	0

	A	C	H	I	J	K	L	M	N
1	ID	Education	Group	Grade	Stage	Comorbidity_1	Comorbidity_2	Comorbidity_3	Comorbidity_4
289	289	Some College	Treatment	Medium	3	1	1	1	1
290	290	High School or less	Control	Low	1	1	1	1	1
291	291	High School or less	Treatment	High	1	1	1	1	1
292	292	Some College	Control	High	3	0	1	0	1
293	293	Advanced Degree	Control	Medium	2	1	1	1	0
294	294	Advanced Degree	Treatment	Medium	5	1	1	1	0
295	295	Advanced Degree	Control	Low	5	1	1	1	1
296	296	High School or less	Control	Medium	3	1	1	1	0
297	297	Advanced Degree	Control	Medium	2	0	1	1	1
298	298	Bachelor's Degree	Treatment	Medium	3	1	1	1	1
299	299	High School or less	Control	Low	2	1	1	1	1
300	300	High School or less	Control	Medium	4	1	1	0	1
301	301	High School or less	Treatment	Medium	5	1	1	1	1
302	302	High School or less	Treatment	Medium	2	0	1	1	1
303	303	Some College	Treatment	Low	2	1	1	1	1
304	304	Some College	Control	Low	4	1	1	1	0
305	305	Some College	Treatment	Medium	5	0	1	1	0

Figure 2 shows how to “freeze panes”. The first step is to select the anchor cell for the freezing action: the top left cell in the area of that should not be frozen. So to freeze one row and one column, the second cell in the second row (B2) is selected in Figure 2a. The “Freeze Panes” button is located on the “View” menu tab, and should be clicked only when the appropriate anchor cell has been selected. It is also possible to freeze just the top row or first column by using the appropriate options. Once the desired rows/columns have been locked, the “Freeze Panes” menu item becomes “Unfreeze Panes”, which can be used to deactivate this setup (Figure 2b).

Data types

A common occurrence in Excel spreadsheets is the presence of invalid or inconsistent data, whether due to a typo, column shift, or user misunderstanding. Common examples include text notes among numbers, impossible dates or codes, and inconsistent spelling. Excel has data

validation tools that allow checking of values that have already been entered, or stop the user from entering invalid values. Filtering is another feature to check quickly for inconsistent entries.

Consider the example in Figure 3. The highlighted column should contain dates, but there are two subtle typos. They can be easily detected by telling Excel that the values in this column should be dates. In the “Data Validation” menu of the Data tab one can specify the type of possible values for the highlighted area (here the entire column E has been highlighted). By selecting “Date” (Figure 3b) and then selecting the “Circle Invalid Data” menu item, invalid values are highlighted (Figure 3c). Now we can easily see the invalid entries: there are only 30 days in June, so 6/31/2003 is invalid, while 3/1//2003 had an extra slash. After setting this validation, attempt to enter non-date values in column E will result in an error message.

ID	Race	Age	Sex	Enroll_date	Height	Education	Group	Grade	Stage	Comorbidity_1	Comorbidity
1	White	66	Male	4/28/2002	171.8	Bachelor's Degree	Control	Low	3	1	
2	White	51	Male	3/20/2002	167.3	Bachelor's Degree	Treatment	Low	5	0	
3	White	18	Male	7/10/2003	178.3	Advanced Degree	Control	High	3	1	
4	White	72	Male	6/13/2003	193.3	High School or less	Treatment	Low	3	1	
5	White	79	Male	6/31/2003	180.6	High School or less	Treatment	Medium	5	1	
6	White	30	Female	6/27/2002	175.3	Some College	Treatment	Medium	5	1	
7	White	23	Female	7/31/2002	163.9	Bachelor's Degree	Control	Low	4	1	
8	White	63	Male	3/14/2003	185.9	High School or less	Treatment	High	4	0	
9	White	27	Male	3/1//2003	181	Some College	Treatment	High	2	1	
10	White	60	Male	1/10/2003	192.2	High School or less	Treatment	Low	2	0	
11	White	77	Male	12/7/2002	161.9	High School or less	Control	Medium	3	1	
12	White	39	Female	5/15/2003	168.1	Bachelor's Degree	Treatment	Low	5	1	
13	White	59	Male	8/18/2002	186.4	High School or less	Treatment	High	3	1	
14	White	51	Male	5/25/2003	170.3	Some College	Treatment	High	2	1	

ID	Race	Age	Sex	Enroll_date	Height	Education	Group	Grade	Stage	Comorbidity_1	Comorbidity_2	Comorbidity_3	Comorbidity_4	Comorbidity_5	LOS
1	White	66	Male	4/28/2002	171.8	Bachelor's Degree	Control	Low	3	1					13
2	White	51	Male	3/20/2002	167.3	Bachelor's Degree	Treatment	Low	5	0					5
3	White	18	Male	7/10/2003	178.3	Advanced Degree	Control	High	3	1					13
4	White	72	Male	6/13/2003	193.3	High School or less	Treatment	Low	3	1					13
5	White	79	Male	6/31/2003	180.6	High School or less	Treatment	Medium	5	1					11
6	White	30	Female	6/27/2002	175.3	Some College	Treatment	Medium	5	1					7
7	White	23	Female	7/31/2002	163.9	Bachelor's Degree	Control	Low	4	1					10
8	White	63	Male	3/14/2003	185.9	High School or less	Treatment	High	4	0					8
9	White	27	Male	3/1//2003	181	Some College	Treatment	High	2	1					8
10	White	60	Male	1/10/2003	192.2	High School or less	Treatment	Low	2	0					12
11	White	77	Male	12/7/2002	161.9	High School or less	Control	Medium	3	1					17
12	White	39	Female	5/15/2003	168.1	Bachelor's Degree	Treatment	Low	5	1					8
13	White	59	Male	8/18/2002	186.4	High School or less	Treatment	High	3	1					8
14	White	51	Male	5/25/2003	170.3	Some College	Treatment	High	2	1					16
15	White														9
16	White														11
17	White														3
18	White														8
19	White														12
20	White														17
21	White														8
22	White														16
23	White														4

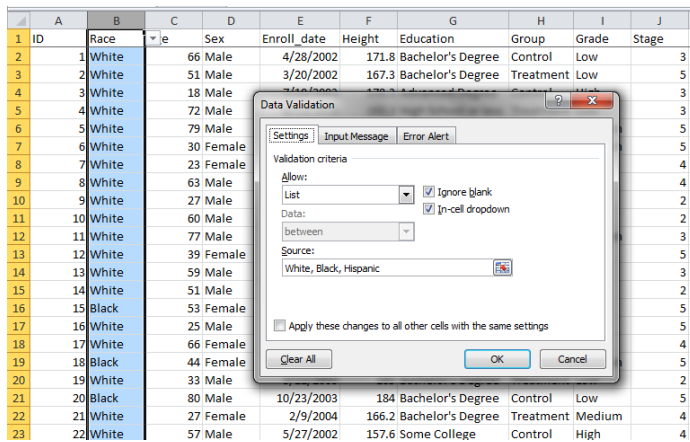
Age	Sex	Enroll_date	Height	Education
66	Male	4/28/2002	171.8	Bachelor's Degree
51	Male	3/20/2002	167.3	Bachelor's Degree
18	Male	7/10/2003	178.3	Advanced Degree
72	Male	6/13/2003	193.3	High School or less
79	Male	6/31/2003	180.6	High School or less
30	Female	6/27/2002	175.3	Some College
23	Female	7/31/2002	163.9	Bachelor's Degree
63	Male	3/14/2003	185.9	High School or less
27	Male	3/1//2003	181	Some College
60	Male	1/10/2003	192.2	High School or less
77	Male	12/7/2002	161.9	High School or less
39	Female	5/15/2003	168.1	Bachelor's Degree
59	Male	8/18/2002	186.4	High School or less
51	Male	5/25/2003	170.3	Some College

Using the same “Data Validation” menu we can specify that the values in column J, which represent disease stage, should be whole numbers between 1 and 5, avoiding illegal entries such as “2b” or “0” (assuming that is the intention).

Multiple choice dropdown

The data validation feature can also be used to ensure consistent entry of categorical variables such as race or sex. A common approach is to use integer coding for such variables (eg 1=male, 2=female), however such conventions are error-prone as they are difficult to remember during data entry. In fact, modern statistical software has no difficulty handling character variables as long as they are entered consistently with the same spelling and capitalization.

Figure 4 shows how to set up a dropdown box of possible options for data entry for the race/ethnicity variable. Such a list makes it easy to remember the possible options and avoid typos. In the same “Data Validation” menu that we saw before, we would select the “List” validation option, and either enter the possible values separated by comma in the “Source” box, or use the red arrow to highlight a range of cells (on another worksheet) in which all the possible values are listed.



7	White	23	Femal
8	White	63	Male
9	White	27	Male
10	White	60	Male
11	Hispanic	77	Male
12	White	39	Femal
13	White	59	Male
14	White	51	Male

Summary

As this article shows, Excel provides a variety of tools to simplify data entry. We recommend showing the planned data collection sheet to a statistician before starting data entry, perhaps with a few filled rows, to ensure that the data can be easily utilized for statistical analyses at the end of the project.