

The first three steps in data entry, with examples in PASW/SPSS

Steve Simon
P.Mean Consulting
www.pmean.com

2. Abstract

- This **training class** will give you a general introduction to data management using PASW (formerly known as SPSS) software. This class is useful for anyone who needs to enter or analyze research data. **No statistical experience is necessary.**

3. Abstract

- There are three steps that will help you get started with data entry for a research project. First, arrange your data in a rectangular format. Second, create a name for each column of data and provide documentation on this column such as units of measurement. Third, create codes for categorical data and for missing values.

4. Abstract

- This class will show examples of data entry including the tricky issues associated with data entry of a two by two table and entry of dates.

5. Abstract

In this class, you will learn how to:

- document variables in a PASW data set;
- enter and manipulate dates in PASW; and
- import data into PASW from other programs.

6. Outline

1. Icebreaker
2. Pop quiz
3. Spreadsheet or database?
4. Documenting variables in PASW
5. Inputting two by two table
6. Inputting dates
7. Importing data
8. Repeat of pop quiz

7. Icebreaker: Your most complicated statistic

- Percentage
- Mean
- Standard deviation
- T-test
- Correlation coefficient
- Linear regression model
- Survival curve
- Logistic regression model
- Other

8. Pop quiz

1. PASW provides documentation about the individual levels of a categorical variable using
 - Format type
 - Missing value codes
 - Variable labels
 - Value labels

9. Pop quiz

2. In PASW, if you subtract one date from another to compute the number of days between two events you will get the following result.
 - An error message
 - A missing value
 - A result that is far too large to be correct
 - A warning message

10. Pop quiz

3. In PASW, you can simplify the data entry for a two by two table by using
 - Automatic recode
 - Crosstabs
 - Restructure wizard
 - Weight cases

11. Database or spreadsheet?

Fast answer

- Databases easily allow you to implement quality checks. They also allow you to easily integrate data from multiple sources. Finally, they are more effective in handling very large data sets.
- Spreadsheets are faster to set up and allow easier copying and duplication for data with repetitive patterns.

12. Database quality checks

- One of your variables might be gender. It might be coded 1=Male, 2=Female, 9=Unknown (though if gender is unknown, you might want to look at the credentials of the doctor doing the examination). With a database, you could set up data entry in that field so that it would beep anytime you tried to enter something other than a 1, 2, or a 9.

13. Database quality checks

- Another quality check found in databases is insuring that the same id code is not assigned to two different subjects. Database specialists refer to this as checking for unique primary keys.

14. Database quality checks

- It's also possible to program a database to check for consistencies in dates. If the birth date is in 1994, for example, and the examination date is in 1987, then either your data are in error or you have an extremely far-sighted pre-natal care program.

15. Database quality checks

- Yet another example is checking the gender or age of the subject before allowing certain data to be entered. Male subjects, for example, would not normally have a hysterectomy in their medical history. Five year olds are rarely married or widowed. The range of quality checks you can include in a database is limited only by your imagination.

16. Multiple data sources

- A database is effective at integrating data coming from a variety of sources. For example, you might have data coming from a laboratory, a questionnaire, and from the medical records. A database makes it easy to properly link the information from all three sources.

17. Multiple data sources

- Another example of where a database is extremely useful is in a multi-center clinical trial. The database offers a standard way for data entry that helps avoid the inconsistencies that can plague such studies.

18. Multiple data sources

- Of course, if you have a data set so complex as to take information from three different sources, then you should definitely consult an expert early in the design of your study. Databases are nice, but they are no substitute for careful planning.

19. Very large data sets

- A database is more effective at handling very large data sets. Unlike a spreadsheet, the entire data set does not have to fit into computer memory. Of course, this is a factor only when the data set on the order of tens of thousands of records or more. If your data set is smaller than this then fitting all of the data into computer memory is unlikely to be a problem.

20. Advantages of a spreadsheet

- Spreadsheets can be up and ready for data entry faster than databases. The extra time required by a database might be beneficial, but for a simple data entry situation, it might just as easily be overkill.

21. Advantages of a spreadsheet

- Spreadsheets also are more efficient at copying and duplicating blocks of information. This can be a time-saver for data sets with repetitive patterns, such as multi-factorial experiments.

22. What are your import options?

- Before you choose, check to make sure that the statistics software can import your version of the spreadsheet or database. PASW can import Excel and Lotus spreadsheets, dBase format databases like FoxPro, Microsoft Access, SQL, and many other databases.

23. Human factors

- Before you make your choice, be sure to factor in any human considerations. If the person doing data entry is much more comfortable with a spreadsheet than a database (or vice versa), that might outweigh some of the computer efficiencies.

24. Human factors

- On the other hand, keep in mind that software in general, and database software in particular, is getting easier to use. Don't let lack of experience keep you from trying a database. It's easier than you think.

25. Database or spreadsheet?

- In summary, databases allow for better error checking, for better integration of data from multiple sources, and for better handling of very large files.
- Spreadsheets are faster to get up and running, which can be an advantage for small tasks. Spreadsheets also have an advantage when there are repetitive patterns in the data.

26. First three steps in data entry

- Every data set is different, but most data entry procedures will start out the same. Here are the first three steps that you should follow to help insure a successful data entry.
 1. Arrange your data in rectangular format.
 2. Create variable names (8 characters or less).
 3. Assign number codes for categorical data and missing values.

27. Rectangular format

- Arrange your data in a rectangular format. **The intersection of each row and column should contain a single number.** Don't leave a cell empty if you can possibly avoid it. Don't try to squeeze two numbers into a single cell.

28. Problem with empty cells

- Empty cells are ambiguous (does it represent a missing value or is it a sign that data entry is not yet complete on this patient). Some computer programs (not PASW) will take an empty cell and convert it to zero, which can lead to disastrous results.

29. Don't squeeze two numbers into one cell

- You might be tempted, for example, to list blood pressure as 120/80 to represent the systolic and diastolic pressures. Don't do this! You will make it difficult for the computer to compute an average blood pressure of any type. Furthermore, some computer software programs might look at the entry 120/80, misinterpret the slash as a division sign, and replace the whole cell with 1.5.

30. Transforming to a rectangular format

- Here's an example of data which does not fit into a rectangular format. These data are loosely based on a study of breast feeding in pre-term infants.

```
REPORT: PARENTER STATUS AT 100 MONTHS
DO          DO          DO          DO          DO          DO
Mom's Gestational Birth Mom's Gestational Birth Mom's Gestational Birth
Age  Gestatic Weight Age  Gestatic Weight Age  Gestatic Weight
04  Normal 1.440  04  Normal 2.500  04  Normal 1.440
02  Single 1.100  04  Normal 1.100  04  Normal 2.430
01  Normal  20  Normal 2.500
01  Normal 1.640
```

31. Transforming to a rectangular format

- Notice that there is a 4 by 3 matrix (4 rows by 3 columns) for the “No” group, a 3 by 3 matrix for the “Yes” group, and a 2 by 3 matrix for the “Lost” group.

Screen Feeding status at 21A month

NO			YES			LOST TO FOLLOW-UP			
Mon's Marital	Birth	Mon's Marital	Birth	Mon's Marital	Birth	Mon's Marital	Birth	Mon's Marital	Birth
Age	Status	Weight	Age	Status	Weight	Age	Status	Weight	Age
33	Married	1.950	33	Single	1.980	33	Married	1.980	
32	Single	1.440	34	Single	1.180	34	Married	1.980	
34	Married		36	Married	2.040				
35	Married	1.640							

32. Transforming to a rectangular format

- If you stack these matrices one beneath the other rather than side by side, you will get closer to a rectangular format.
 - Notice that you have to add another column to denote which matrix is which.

Screen Feeding Status	Mon's Marital	Birth	Age	Status	Weight
No	33	Married	1.950		
No	32	Single	1.440		
No	34	Married			
No	36	Married	2.040		
Yes	33	Single	1.980		
Yes			1.180		
Yes	34	Married	1.980		
Lost	35	Married	1.640		
Lost			2.440		

33. Assigning codes

- If you have categorical data, assign a code to each category level. Use the code during data entry to save time and minimize errors.
- Here are some examples of codes:
 - Gender 1=Male, 2=Female, 9=Unknown;
 - Race 1=White, 2=Black, 3=Asian, 4=Hispanic, 5=Native American, 8=Multiracial, 9=Unknown;
 - Likert scale 1=Strongly Disagree, 2=Disagree, 3=Neutral, 4=Agree, 5=Strongly Agree, 9=No answer.

34. Assigning codes

- **While I prefer to use number codes, there are some advantages to using short letter codes.**
- Here are some examples of letter codes:
 - Gender M, F, and U (Male, Female, and Unknown);
 - Race W, B, A, H, N, M, and U (White, Black, Asian-American, Hispanic, Native American, Mixed, and Unknown);
 - Likert scale SD, D, N, A, SA, NA (Strongly Disagree, Disagree, Neutral, Agree, Strongly Agree, No Answer).

35. Number versus letter codes

- Letter codes are easier to remember, and sometimes can be used effectively as plotting symbols.
- **I prefer number codes because they offer more flexibility during statistical analysis.** For example, SPSS will not allow you to draw a scatterplot when one of your variables uses letter codes.

36. Binary variables

- For binary variables, I prefer to use 0-1 coding rather than 1-2. It is similar to how computers work (0=off, 1=on). Here are some examples:
 - Treatment: 0=placebo, 1=active drug.
 - Exposure: 0=unexposed, 1=exposed.
 - Disease status: 0=healthy, 1=diseased.
 - Gender: 0=female, 1=male.

37. Example using number codes

- Let's assign number codes for the categorical variables in the breast feeding data example.

SPSS Variable Name	Mom's Age	Breast Feeding Status	Birth Weight
1	31	1	7.1123
1	32	1	7.1182
1	37	1	7.1411
1	34	1	7.1411
1	35	1	7.1767
1	36	1	7.1804
1	38	1	2.1322

38. Missing value codes

- Note also that there are still some blank spots in this data. These represent missing data. **Never let an empty field represent missing data.** Explicitly create a code for missing, and be sure to explain why the data are missing to anyone involved with analysis of your data.

39. Missing value codes

- Let -1 represent a missing value for Mom's Age and Birth Weight. Let 9 represent a missing value for Breast Feeding Status and Marital Status.

SPSS Variable Name	Mom's Age	Breast Feeding Status	Birth Weight
1	31	1	7.1123
1	45	1	7.1411
1	47	1	7.1411
1	46	1	7.1411
1	26	1	7.1767
1	3	1	7.1804
1	26	1	7.1804
1	2	1	2.1322

40. Create short names for each column of data

- If you are using a spreadsheet, place a descriptive variable name at the top of each column. If you are using a database, provide a descriptive name for each field. You will use this variable or field name in statistical software like SPSS to specify the variables that you want to analyze. **Try to be reasonably descriptive with your variable names; avoid generic names like VAR01, VAR02, etc.**

41. Guidelines for variable names

- Here are some general guidelines that will help avoid trouble with variable names.
 - Use a brief name (eight to sixteen characters long).
 - A mixture of numbers and letters is okay, but avoid special symbols such as \$, &, or %.
 - Don't rely on upper/lower case to distinguish among variable names
 - Avoid embedded blanks.

42. Use a brief name

- Use a brief name (eight to sixteen characters long). A long time ago, (version 11 and earlier of SPSS), you could not use a name longer than eight character long. Now you can use up to 255 characters, but you should show some restraint. It is convenient to have a short "handle" that you can refer to for any column of data in your data set.

43. Avoid special symbols

- A mixture of numbers and letters is okay, but avoid special symbols such as \$, &, or %. Most statistical software will reserve these special symbols for other purposes. The one major exception is the underscore (`_`), which is found usually paired on the same key with the minus sign.

44. Upper/lower case

- Don't rely on upper/lower case to distinguish among variable names. For example, don't name one variable `x` and the next one `X`. Some packages are case insensitive and even if they are not, having two variables with names that look almost identical is a formula for trouble.

45. Avoid embedded blanks

- In most statistical software, an embedded blank will cause problems. A variable with a name like "mom age" will possibly be interpreted as two variables "mom" and "age" in certain situations. This can lead to lots of problems.

46. Avoid embedded blanks

- But you shouldn't just strip out the blanks. There's a story about a website for a group known as Writer's Exchange. They used www.writersexchage.com for their website, but someone noticed that this could be read as writer sex change.

47. Avoid embedded blanks

- If your variable name consists of two or three short words, here are three strategies that work well.
 1. Use an upper case letter at the start of each new word: `MomAge`, `BreastFeedingStatus`.
 2. Separate words using a period: `mom.age`, `breast.feeding.status`
 3. Separate words using an underscore `mom_age`, `breast_feeding_status`.

48. Avoid embedded blanks

- Note that separating using a dash (minus sign) is not a good idea. A variable with the name `mom-age` looks too much to some computers like a subtraction calculation. Some software programs also use a dash to indicate a range of variables (`a1-a8`).

49. Example of variable names

- Here's what the data set looks like with variable names.

yr	total	new	old	avg	yr	birth	nl
2	28	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	
2	30	1	1	1	1	2.1250	

50. PASW Variable View tab

- Once you have names for each column of data, you should document what goes into each column, and control how it is displayed. In PASW, it starts with selecting the VARIABLE VIEW tab to add documentation to your data.

51. PASW Variable View tab



52. PASW Variable View tab

- From VARIABLE VIEW tab, you can tell SPSS **how to display the data in the SPSS data editor window** (how many decimal places shown, how dates are displayed, and how wide the columns are). You can also provide SPSS with **informational labels that will appear in your output window** (labels for the variable itself and, if needed, labels for category levels). You would also use the dialog box to **specify any codes that represent missing data**.

53. NAME

- When documenting your data, your first step should be to **provide a brief but descriptive variable name**. This goes into the NAME column of the VARIABLE VIEW tab. Please spend some time to provide descriptive variable names. As noted above, these names should be short (8 to 16 characters). You will get a chance to provide additional details in the LABEL field.

54. TYPE

- WhenClick in the TYPE column to add or change the format type. You will notice a gray button appear on the right hand side. Click on it to get VARIABLE TYPE dialog box.
- This dialog box has information about the type of data that you want to use.

55. TYPE (Numeric)

- The most common data type is **NUMERIC**, which is used for any data that can be represented solely by numbers. Unless you are dealing with unusually large numbers, the default width of 8 works well. For some situations, you might be tempted to use a smaller makes, but this can make it more difficult to view the variable name and the value labels.



56. TYPE (Decimal places)

- Be sure to set the number of decimal places. Please do not display decimal places that you don't really need.



57. TYPE (String)

- Select the **STRING** options for data that is all letters or a mixture of letters and numbers. When you select this option, SPSS provides a chance for you to tell how long the strings are.



58. TYPE (Date)

- If you click on the **DATE** option, you will be given choices between various display formats (month names versus month numbers, two digit versus four digit years, etc.).



59. LABEL

- Click on the label field to add a variable label. A variable label is a longer description of your data. **Variable labels appear in your output and make it easier to follow what is going on.** You can use a mixture of upper and lower case here, which I recommend for improving readability. **AVOID USING ALL UPPERCASE HERE BECAUSE IT IS FAR LESS READABLE THAN A MIXTURE OF CASES.**

60. LABEL

- You can put blanks and special symbols in your variable label.** If you are very excited about a variable, spice it up with a couple of exclamation points. Go ahead and type to your heart's content. Just a small warning though. A variable label that is too long can make your output look a bit unwieldy. Although you can type up to 255 characters here, it looks strange to have a six inch label underneath a two inch histogram. A variable label of around **20 to 40 characters in length** works well in practice.

61. VALUES

- **Value labels provide informative names for levels in any categorical variable.** Leave the value labels blank for continuous data like weight or height. They do make sense, though, for categorical data like gender. This will serve as a reminder that data values of 1 represents males and 2 females. The last thing you want is for people to think that you can't tell the difference between males and females.

62. VALUES



63. MISSING

- If needed, click on the MISSING VALUES button to designate missing value codes. **Missing value codes are useful for designating data in SPSS where the value is unknown, not applicable or otherwise not provided.**

64. MISSING

- Be careful about missing values. Make sure you understand why your data is missing and discuss this issue with anyone you are consulting with. **The statistical handling of missing values can vary greatly depending on how the value came to be missing.**

65. MISSING

- When you are planning your project, it is a good idea to **select a very clearly impossible code for your missing value.** For example, use -1 for a birth weight because any infant with a negative birth weight would float up to the ceiling after delivery. Use a value of 9 to code missing for gender, since it is obvious to most of us that the number of possible genders is much smaller than 9.

66. VALUES



67. Summary

- When you are planning your project, it is a good idea to **select a very clearly impossible code for your missing value**. For example, use -1 for a birth weight because any infant with a negative birth weight would float up to the ceiling after delivery. Use a value of 9 to code missing for gender, since it is obvious to most of us that the number of possible genders is much smaller than 9.

68. Summary

- To get started with data entry, follow these three steps.
 1. Arrange your data in rectangular format.
 2. Create codes for category levels and missing values.
 3. Create variable names (8 characters or less).

69. Two by two table

- Data in two by two tables occur commonly in research, but they are a bit tricky to handle. Here is an example of a two by two table.

	D+	D-	Total
F+	34	23	57
F-	139	119	258
Total	173	142	315

70. Two by two table

- For data like this, you have to re-arrange things and then apply weights. To re-arrange the data, you need to specify three variables: F, D, and COUNT.
 - F takes the value of 1 for F+ and 0 for F-.
 - D takes the value of 1 for D+ and 0 for D-.
 - COUNT represents the total number of subjects for each combination of F and D.

71. Two by two table

- Here's what your re-arranged data would look like.

	f	d	count
1	1	1	34
2	1	0	23
3	0	1	139
4	0	0	119

72. Two by two table

- Enter the data, and tell SPSS that **W** represents a **weighting variable**, and you're ready to rock and roll. You do this by selecting **Data | Weight Cases** from the SPSS menu.



73. Two by two table

- Here's what a typical output from PSAW would look like.

f * d Crosstabulation

		d		Total
		0	1	
f	0	Count 118	Count 57	Count 175
		% within f 66.8%	% within f 32.2%	100.0%
1	0	Count 128	Count 24	Count 152
		% within f 84.2%	% within f 15.8%	100.0%
Total		Count 246	Count 81	Count 327
		% within Total 75.2%	% within Total 24.8%	100.0%

74. Two by two table

- If you forgot to use WEIGHT CASES, you would get the following.

f * d Crosstabulation

		d		Total
		0	1	
f	0	Count 1	Count 1	Count 2
		% within f 50.0%	% within f 50.0%	100.0%
1	0	Count 1	Count 1	Count 2
		% within f 50.0%	% within f 50.0%	100.0%
Total		Count 2	Count 2	Count 4
		% within Total 50.0%	% within Total 50.0%	100.0%

75. Dates in PASW

- A while back I got the following email inquiry:
 - Dear Professor Mean, I am trying to use dates in SPSS for certain calculations. For example, I want to use a compute statement in SPSS to create a new variable called duration of injury (during). I know that I must subtract the date of injury from the date of interview. However, when I do this, I get a number in the millions. What am I doing wrong? -- Stumped Sharon
 - Dear Stumped, Maybe your patients were waiting for their HMO to approve a visit to a specialist.

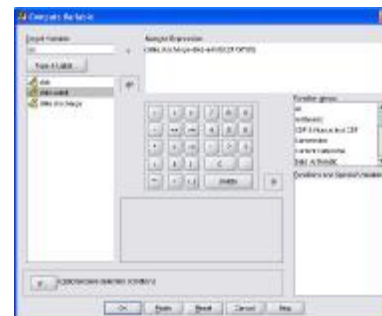
76. Dates in PASW

- Dates in SPSS are actually a bit tricky. **SPSS stores date/time values as the number of seconds since October 14, 1582** (the start of the Gregorian calendar). If you specify only a date and not a time, then SPSS sets the time to midnight. **When you subtract two dates, you get the duration of injury in seconds.** Divide by 86,400 (=24*60*60) to get the duration of injury in days. Divide again by 7, 30, or 365.25 to get duration in weeks, months, or years.

77. Date example

deb	date.admit	date.discharge
1	01/22/1995	02/01/1995
2	01/26/1995	02/07/1995
3	05/10/1995	05/14/1995
4	05/07/1995	05/15/1995

78. Date example



79. Date example

	date	date.admit	date.discharge	los
1	01/22/1995	02/01/1995	02/01/1995	5.00
2	01/26/1995	02/07/1995	02/05/1995	2.00
3	05/10/1995	05/14/1995	06/17/1995	3.00
4	05/07/1995	05/15/1995	06/15/1995	4.00

80. Date and time wizard



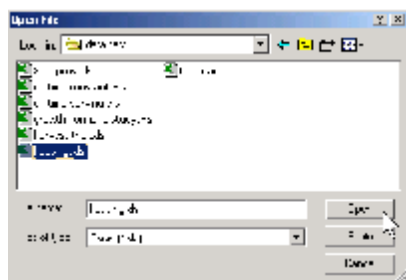
81. Date and time wizard



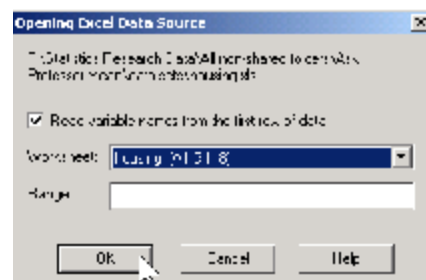
82. Importing data from Excel to PASW

- You should also do some prep work before you import Excel data. Excel is an extremely flexible program that allows you to put your data in just about any way you like.
 1. Arrange the data in a **rectangular grid**
 2. **Don't mix** strings and numbers in a single column.
 3. Put **descriptive names** in your first row.

83. Importing data from Excel to PASW



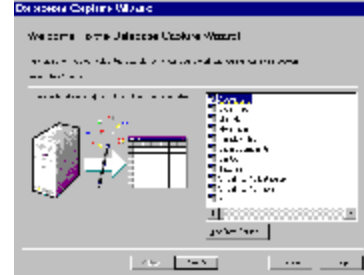
84. Importing data from Excel to PASW



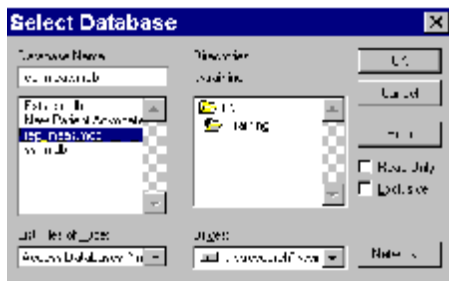
85. Importing data from Access to PASW

- SPSS can import data from a variety of sources using a system known as ODBC (Object Data Base Connectivity). **ODBC has links to just about every database that you would ever need to use.**
- I'll show you an example using Microsoft Access, but this would work just as well on other database systems, such as Oracle and Informix. **To import data from Access, select FILE | DATABASE CAPTURE | NEW QUERY from the SPSS menu.**

86. Importing data from Access to PASW



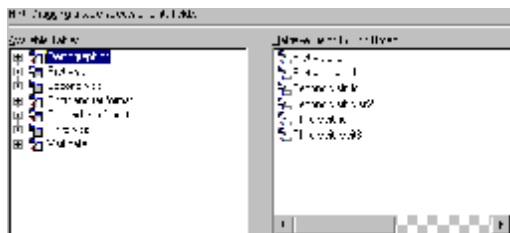
87. Importing data from Access to PASW



88. Importing data from Access to PASW



89. Importing data from Access to PASW



90. Pop quiz

1. PASW provides documentation about the individual levels of a categorical variable using
 - Format type
 - Missing value codes
 - Variable labels
 - Value labels

91. Pop quiz

2. In PASW, if you subtract one date from another to compute the number of days between two events you will get the following result.
- An error message
 - A missing value
 - A result that is far too large to be correct
 - A warning message

92. Pop quiz

3. In PASW, you can simplify the data entry for a two by two table by using
- Automatic recode
 - Crosstabs
 - Restructure wizard
 - Weight cases