# Introduction to SPSS

**Mr. Kongmany Chaleunvong**

GFMER - WHO - UNFPA - LAO PDR
Training Course in Reproductive Health Research
Vientiane, 22 October 2009
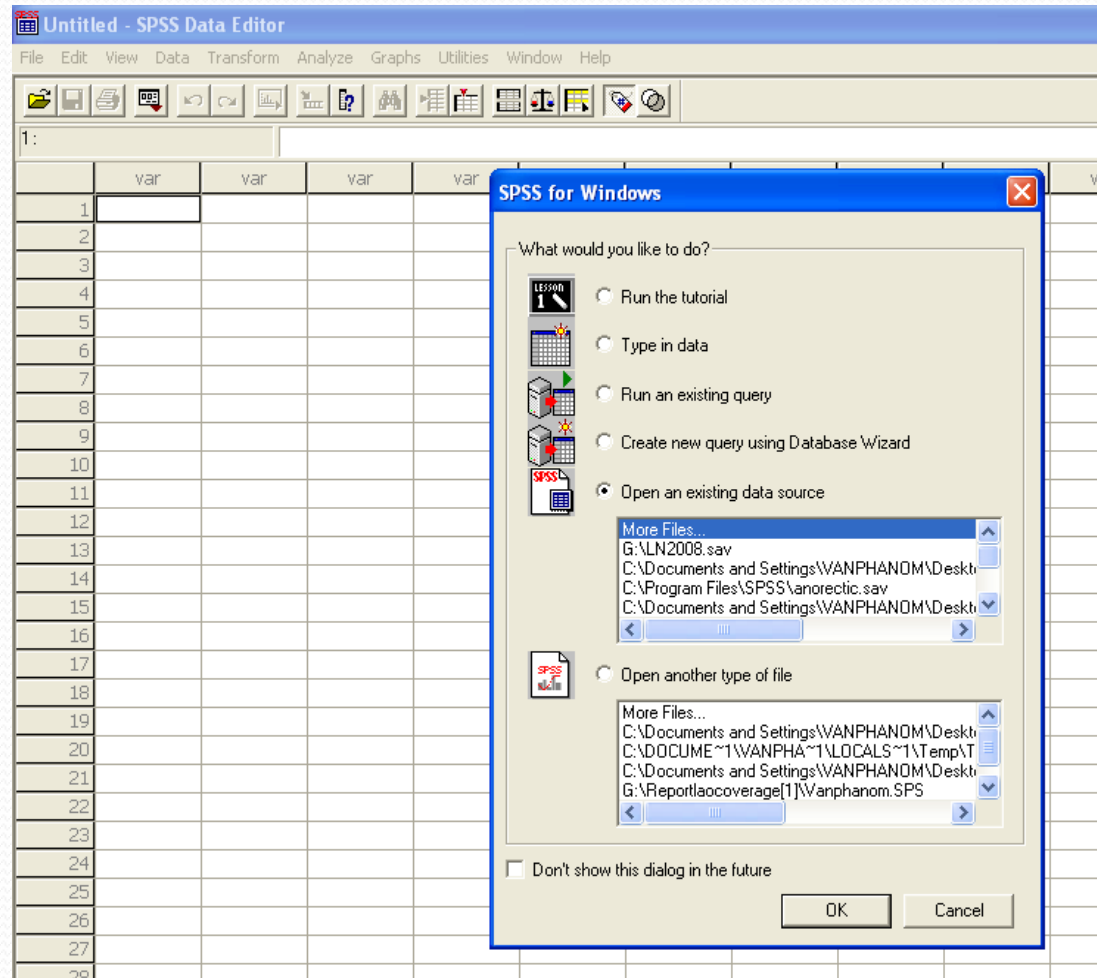
# Object of the Course

- Introduction to SPSS
- The basics of managing data files

# Introduction: What is SPSS?

- SPSS is a statistical package for beginning, intermediate, and advanced data analysis

- Originally it is an acronym of Statistical Package for the Social Science but now it stands for Statistical Product and Service Solutions

- One of the most popular statistical packages which can perform highly complex data manipulation and analysis with simple instructions
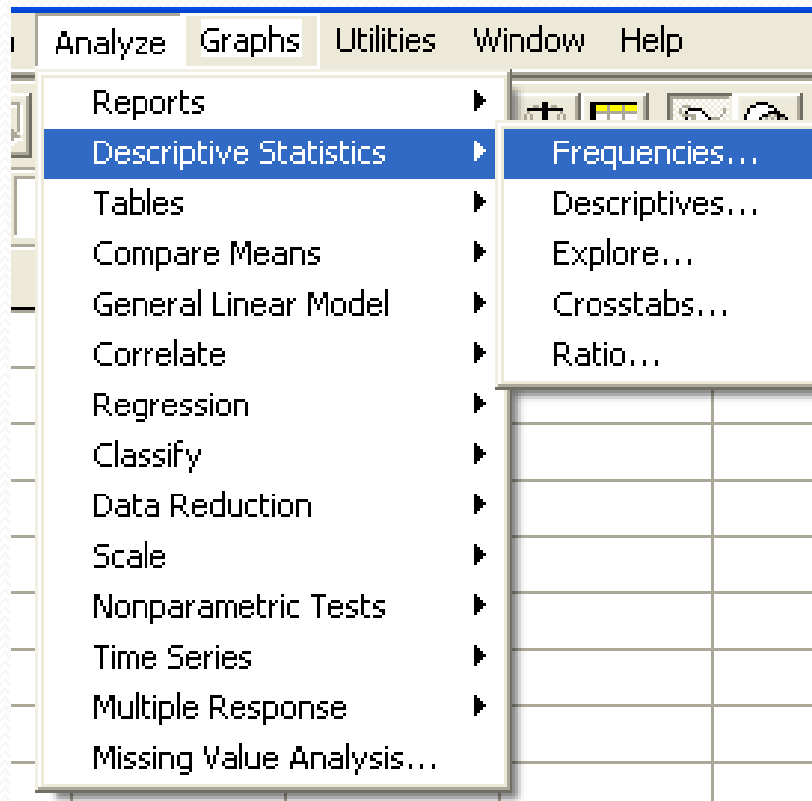
# Starting SPSS for Windows

Launch SPSS either by double-clicking the SPSS icon on the desktop, or from the Start menu –SPSS will have a group under programs. The opening screen should appear as
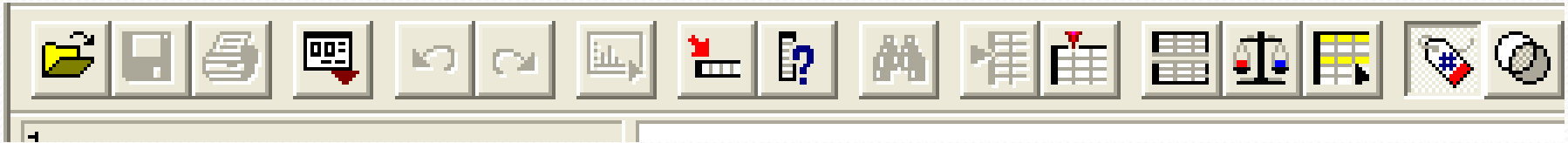
# The Menu bar

**Untitled - SPSS Data Editor**

File   Edit   View   Data   Transform   Analyze   Graphs   Utilities   Window   Help

The Menu bar lists 10 pull down menu, grouping the available SPSS commands. Some of these have sub-menus, the Analyze menu is like this
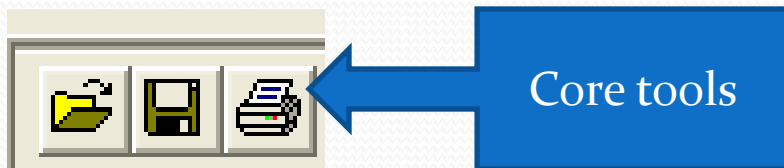
| Analyze | Graphs | Utilities | Window | Help |

Reports ▶

**Descriptive Statistics** ▶ — **Frequencies...**
Tables ▶ — Descriptives...
Compare Means ▶ — Explore...
General Linear Model ▶ — Crosstabs...
Correlate ▶ — Ratio...
Regression ▶
Classify ▶
Data Reduction ▶
Scale ▶
Nonparametric Tests ▶
Time Series ▶
Multiple Response ▶
Missing Value Analysis...

# The Toolbar

The toolbar, located just below the menu bar, provides quick and easy access to many frequently used facilities

Core tools

• Open File: Displays the Open File dialog box for the type of windows that is active.
• Save File: Saves the working file, if the file has no name, it displays the Save File dialog box for the type of document that is active.
• Print : Displays the print dialog box.
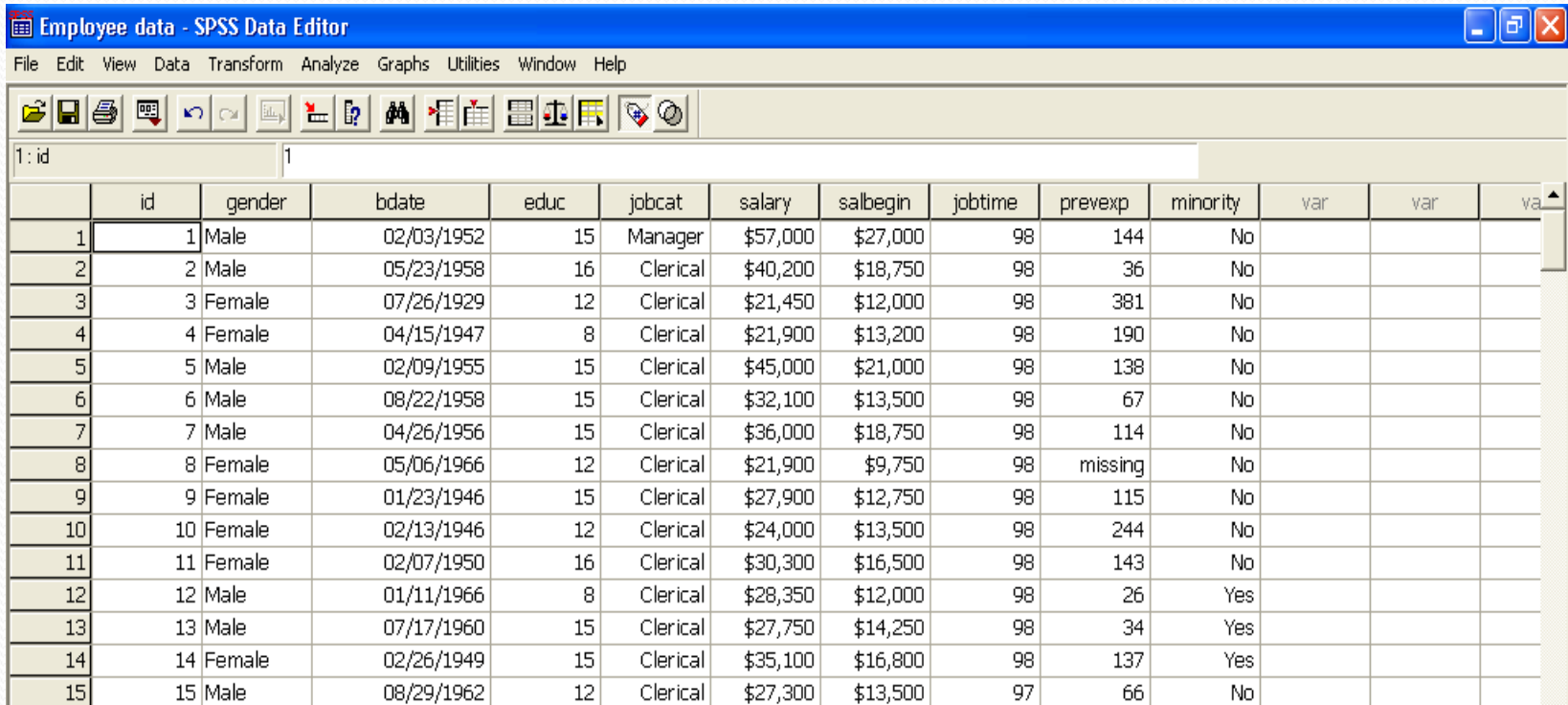
# About the four-windows in SPSS

**The Four Windows:**
Data editor
Output viewer
Chart editor window
Syntax editor

# The Four Windows: Data Editor

- Data Editor

  Spreadsheet-like system for defining, entering, editing, and displaying data. Extension of the saved file will be "sav."

# The Four Windows: Output Viewer

- Output Viewer

  Displays output and errors. Extension of the saved file will be "spo."

# The Four Windows: Chart editor window

- Output Viewer

  Displays output and errors. Extension of the saved file will be "spo."

# The Four Windows: Syntax editor

- Syntax Editor

  Text editor for syntax composition. Extension of the saved file will be "sps."

# Using the Syntax editor

- Click 'Analyze,' 'Descriptive statistics,' then click 'Frequencies.'

- Put 'Gender' in the Variable(s) box.

- Then click 'Charts,' 'Bar charts,' and click 'Continue.'

- Click 'Paste.'



Click

# The basics of managing data files

# Data Entry & Coding

- Before describing the process for defining variables, an important distinction should be made between two terms that are often confused: *variable* and *value*
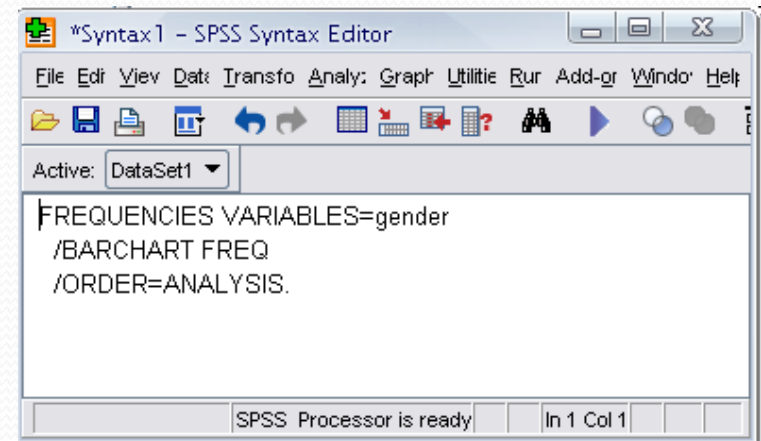
- A variable is a measure or classification scheme that can have several values

- Values are the numbers or categorical classification representing individual instances of the variable being measured

# Data Entry

- You may create a data file using one of your favorite text editors, or word processing packages (e.g., Word Perfect, MS-Word). Files created using word processing software should be saved in text format before trying to read them into an SPSS session.

- You may enter your data into a spreadsheet (e.g., Lotus 123, Excel, dBASE) and read it directly into SPSS for Windows.

- Finally, you may enter the data directly into the spreadsheet-like Data Editor of SPSS for Windows.

    - In this document we are going to examine one data entry methods: using the Data Editor of SPSS for Windows.

# The Data View



## The Variable View

# Define Information – The Variable View

- Name

    - Each variable name must be unique; duplication is not allowed.

    - Start with a letter.

    - May have up to 8 characters, including letters, numbers, and the symbols (@, #, _, or $).

    - Variable names cannot end with a period.

# The Variable View (con't)

- Name (con't)

  - Variable names that end with an underscore should be avoided.

  - The certain key words are reversed and may not be used as variable names, e.g. "compute", "sum" and so forth.

  - Ex. Subject_ID, but not "subject-ID", and not "Subject ID".

# The Variable View (con't)

- Type
  - Basic type – numeric and string
  - Maximum width for numeric variables is 40 characters, the maximum number of decimal positions is 16.
  - String variables may contain letters or numbers. For string values a blank is considered a valid value.
  - Numeric operations on the string variables will NOT be allowed, e.g. finding the mean, variance, standard deviation, etc…

– **If you select a string variable, you can tell SPSS how much "room" to leave in memory for each value, indicating the number of characters to b allowed for data entry in this string variable**

# The Variable View (con't)

- Width
  - The number of characters. SPSS will allow to be entered for the variable.
  - For a numerical value with decimals, this total width has to include a spot for each decimal, as well as one for the decimal point.
- Decimals
  - If more decimals have been entered or computed by SPSS, the additional information will be retained internally but not displayed on screen.
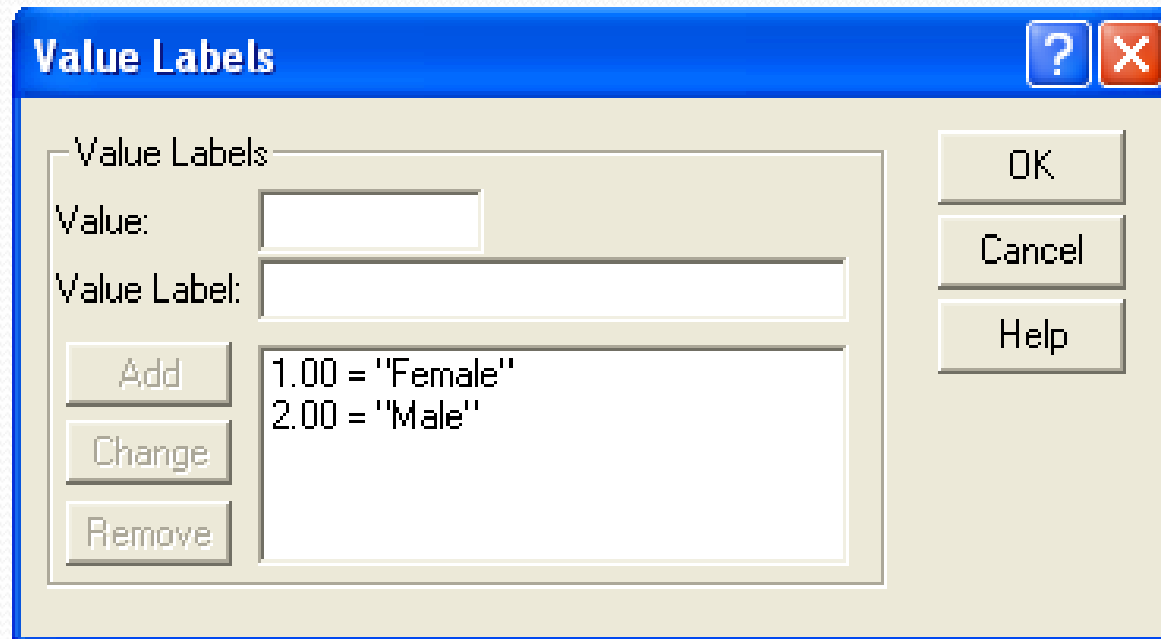
# The Variable View (con't)

- Label

  - A string to identify in detail what a variable represents.

  - Is limited to 255 characters

  - May contain spaces and punctuation.

# The Variable View (con't)

- Values
  - Indicate how the numbers are assigned for categorical data.
  - Instead of typing into the computer the full answer to each question, codes are typed in (e.g., 1 if the respondent is female, 2 if male).
  - Codes are usually numerical, because this is what most statistical software expects, and using only numerical codes makes data entry faster.
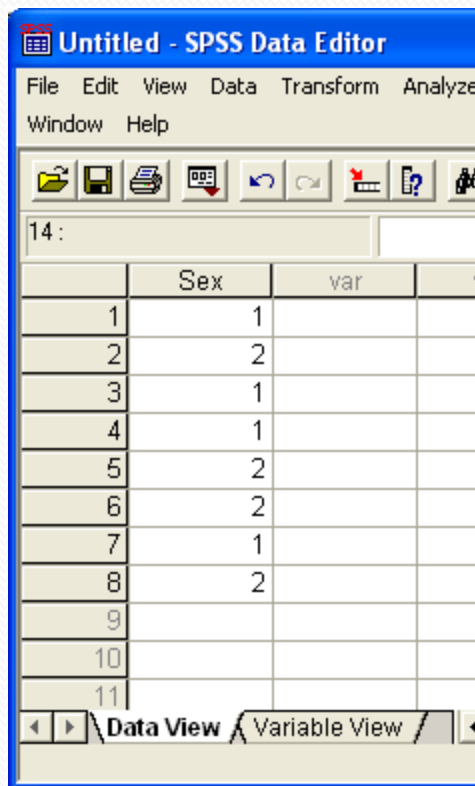  - These are easier to remember, and therefore tend to have lower error rates.

# The Variable View (con't)

- Values (con't)
  - To code categorical variables in numeric format.
  - The Value Labels will be used.

# The Variable View (con't)

The labels can be seen in the Data View by clicking on the "toe tag" icon in the tool bar  , which switches between the numeric values and their labels.
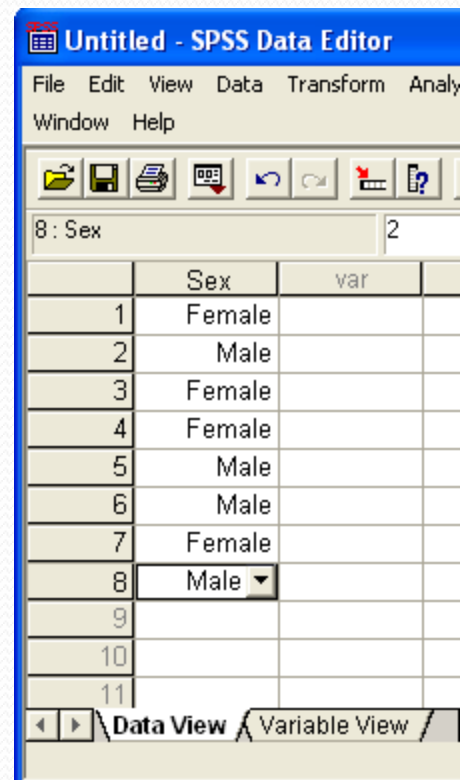
# The Variable View (con't)

- Missing
  - Signal to SPSS which data should be treated as missing.
  - System Missing data – SPSS display a single period.

# The Variable View (con't)

- Columns
  - How wide the column should be for each variable
  - Columns affect only the display of values in the Data Editor. Changing the column width does <u>not change</u> the defined width of a variable.
- Align

# The Variable View (con't)

- Measure
  - Indicates the level of measurement.
  - Since SPSS does NOT differentiate between interval and ratio levels of measurement, both of these quantitative variable types are lumped together as "Scale".
  - Nominal and ordinal levels of measurement ARE differentiated.

# Type of Measurement

The answers to the "numerical questions" are real numbers, not just arbitrary codes. There are four types of numerical scales that exist: nominal scales, ordinal scales, interval scales, and ratio scales.

- ◻ Scale
  - – A ratio scale is one in which the answers are real numbers, and an answer of zero means what it says. "What age are you?" - "How tall are you?" - "How many children do you have?"
  - – An interval scale (meaning equal-interval) - if there's a zero point, it's arbitrary, but the difference between two successive possible answers is the same. For example, the scale of temperature.

# Type of Measurement (con't)

- Ordinal
  - Frequently, categorical data responses represent more than two possible outcomes, and often these possible outcomes take on some inherent ordering.
  - No clue as to the relative distances between the levels.
  - For example, low – medium – high

    50% – 75% – 100% – 200%

    strong agree – agree – neutral – disagree – strongly disagree.

# Type of Measurement (con't)

- Nominal
  - A <u>nominal</u> scale isn't really a scale at all, but an arbitrary code value to distinguish the different groups.
  - No inherent ordering to the categories.
  - For example, "Do you prefer the beach, mountains, or lake for a vacation?"

    "Which color is your favorite?"

# Data Cleaning

- What most data entry programs will not do is warn the user when unlikely (but possible) codes occur. For example, if a respondent's age is shown as 99, this may be true, but it may also be a mistake.

- Therefore it's not only wild values that need to be checked. The first frequencies check from a program needs to be looked at very carefully to detect this kind of mistake.

# Data Cleaning (con't)

- Check missing values - If the question was "Which sex are you, male or female?" and the possible answers are 1 for male and 2 for female, these should be the only values for that variable - except perhaps for a few blanks for the missing values.

# Data Cleaning (con't)

- There are two types of missing values in SPSS: system-missing and user-defined.

- System-missing values are assigned by SPSS when, for example, you perform an illegal function, like dividing a number by zero. System-missing values can also be assigned in an input data set.

- User-defined missing values are numeric values that you can specify and SPSS will consider to be missing. For example, you may define -9999 to be a missing value.

# Data Cleaning (con't)

- You can assign many different missing values to a given variable, perhaps using the different values to indicate different reasons for the data point to be missing.
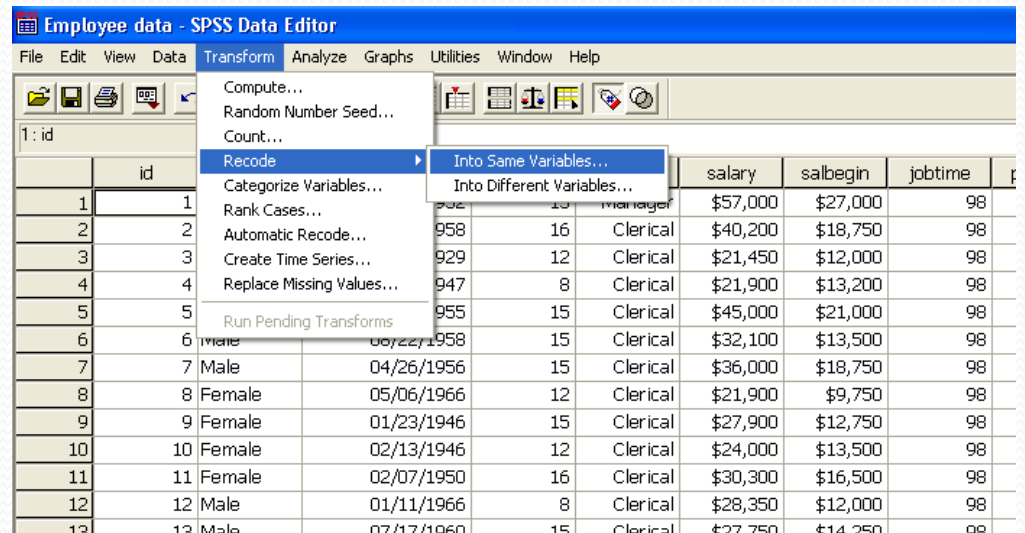
  For example, for an item on a survey, -9999 might indicate that the respondent skipped the item, -8888 might indicate that the item was not answered because it was part of skip pattern, and -7777 might indicate that a note was written in the margin instead of a standard response.

- You can specify up to three unique values for each variable. User-defined missing values can also be a range, such as 5 to 10. This is useful when you want to include only half of a scale, for example.

- String values can also be used as missing values, including a series of blanks (i.e., a null string).
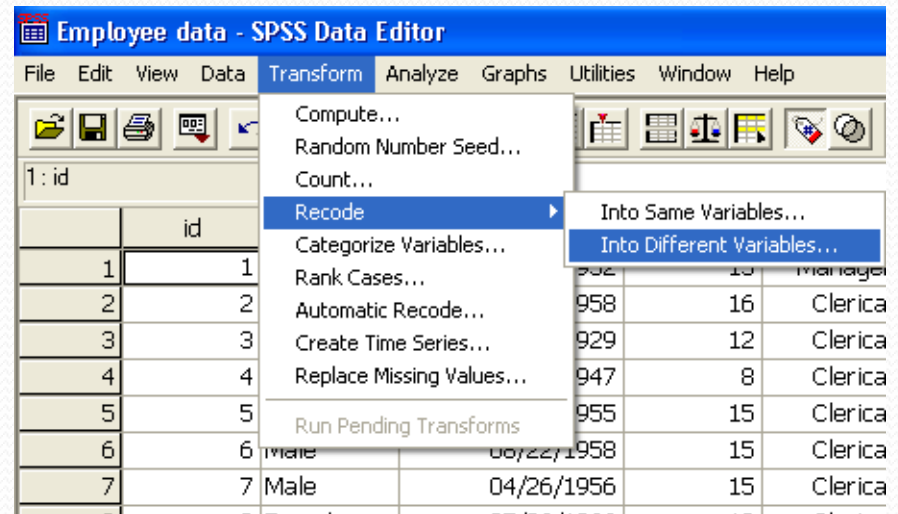
# Recode Procedure

Recode is used to

- to change the values of an existing variable

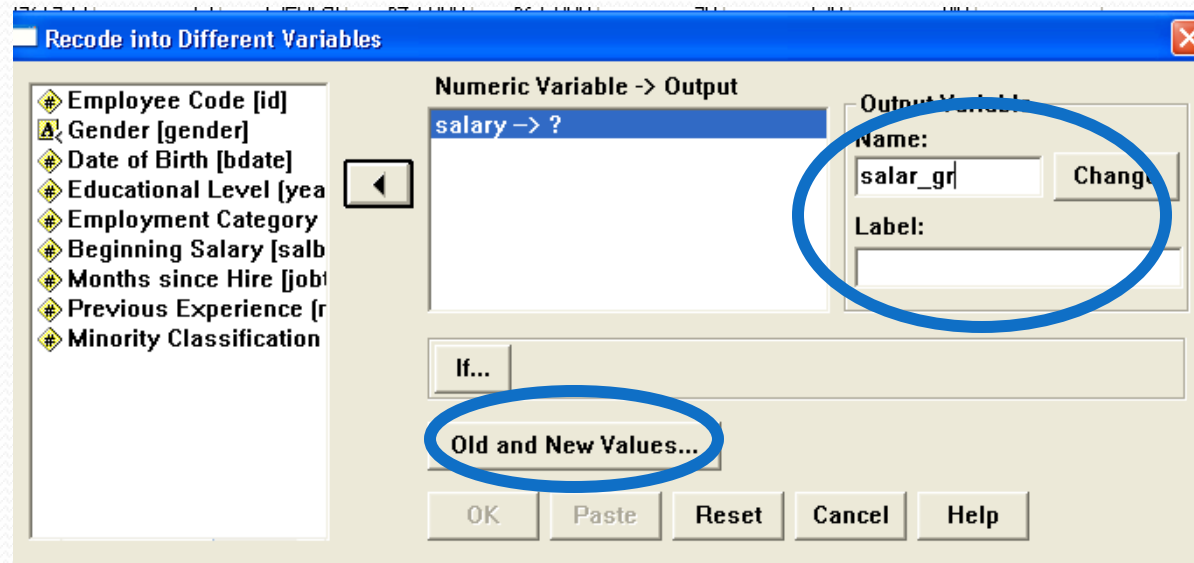- to create a new variable based on the values an existing variable

# Recode into New Variable

- In the menu, click Transform.
- Select Recode.
- Click Into Different Variable(s)

# Recode into New Variable

- Select and move variable(s) over.

- Name and label new variable.

- Click
Old and New Values

# Recode into New Variable

For each value of the existing variable

- Enter the new value
- Click Add
- Repeat for each value or range of values
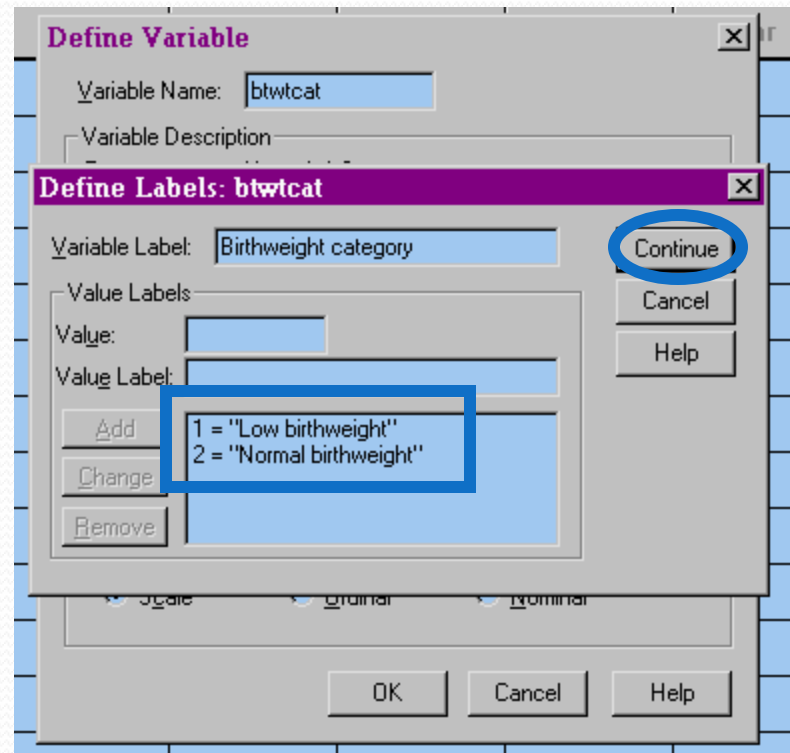- Click Continue

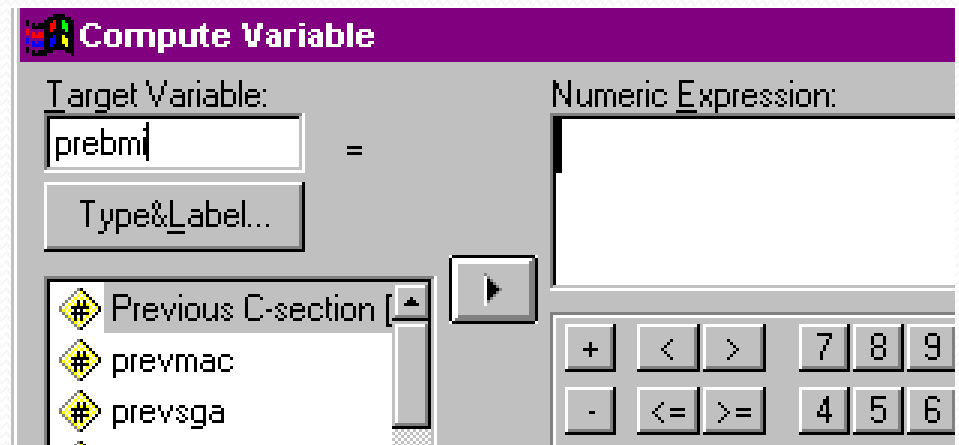# Recode into New Variable

- Click Change
- Click OK

# Define Labels for New Variable

- In the Data menu, click  Define Variable.

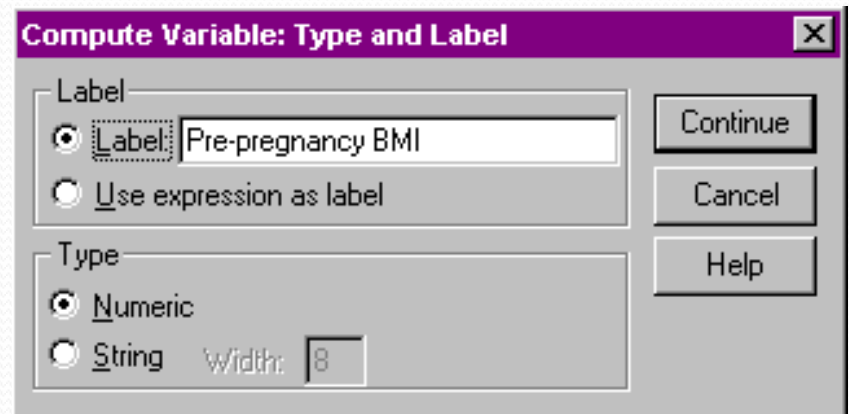- Click Labels.

- Enter value labels for the new variable.

# Compute Procedure

- Name the new variable.
- Click Type&Label to define the characteristics of the new variable.
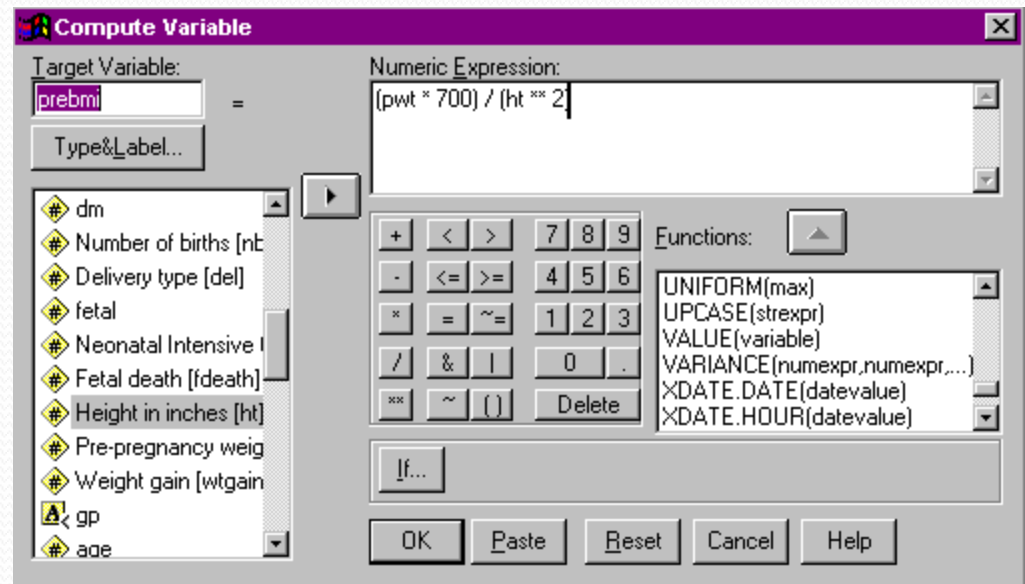
# Compute Procedure

- Label the new variable.
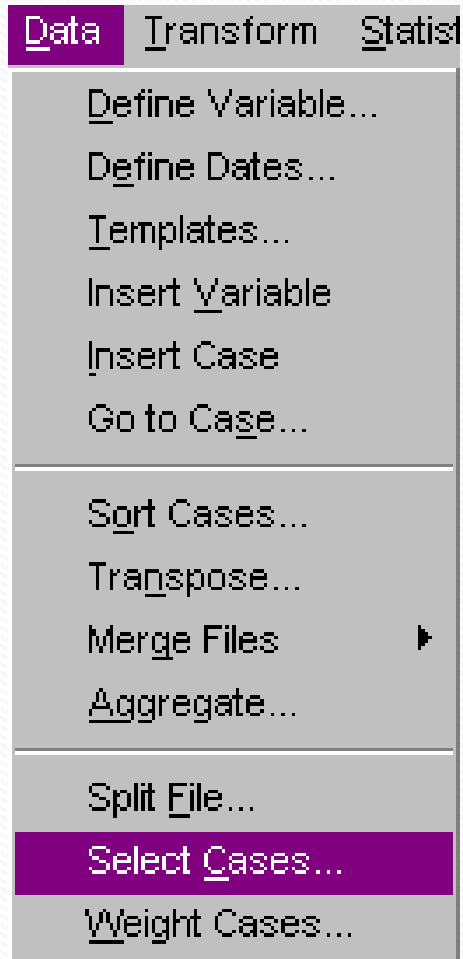- Enter the variable type.

# Compute Procedure

- Enter the numeric expression that will determine the values of the new variable.
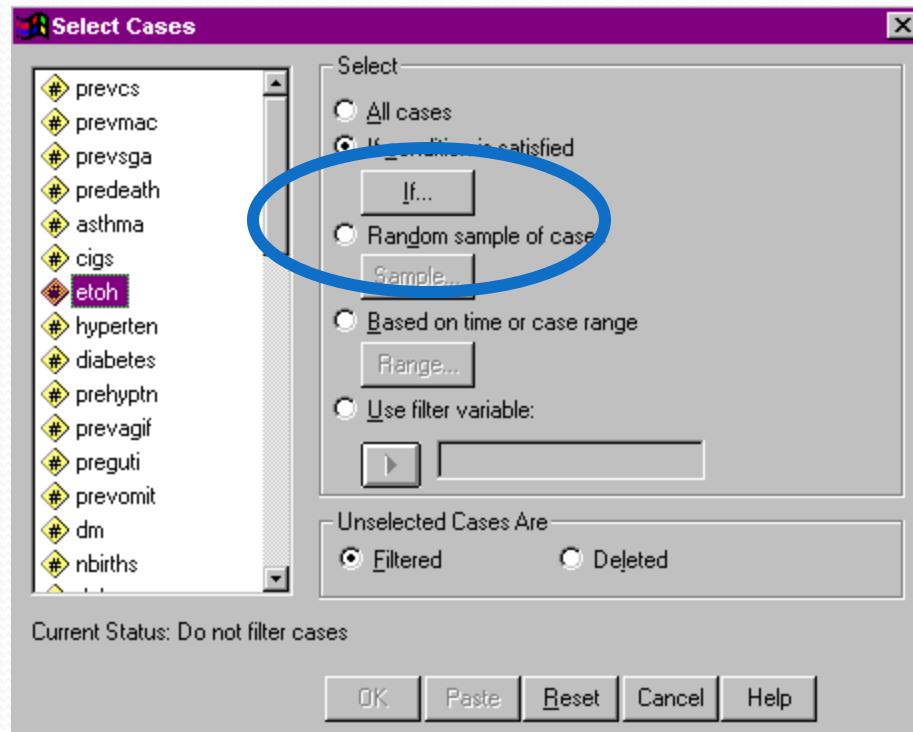- Click OK.

# Select Cases



For a subset of the datafile, use Select Cases.

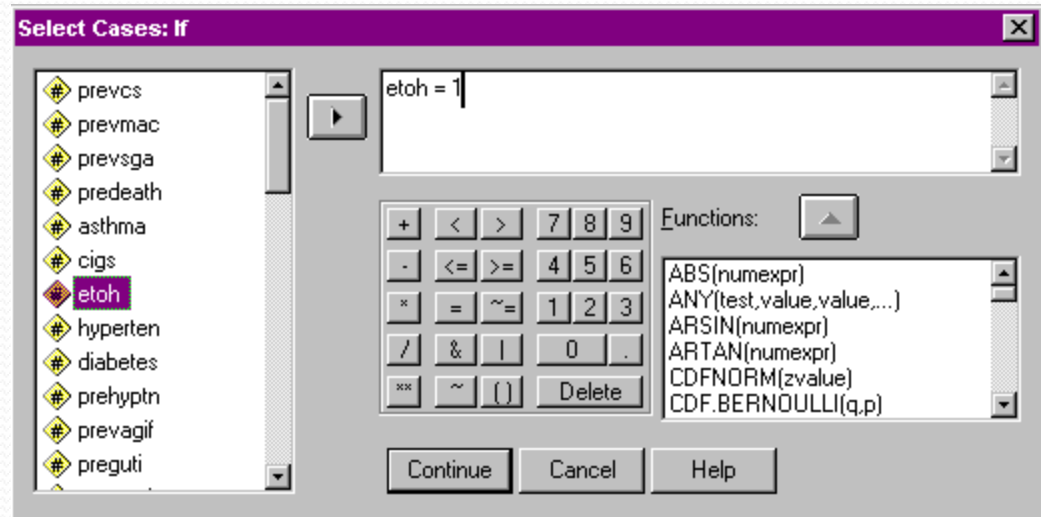- In the menu, click Data.
- Click Select Cases...

# Select Cases - Alcohol drinkers only

To select only those cases which meet certain criteria, choose the If option.

# Select Cases - Alcohol drinkers only

- Enter the expression that will determine which variables will be selected.
- Click Continue.

# Select Cases - Alcohol drinkers only

When you've finished specifying selection criteria, click OK.