# Video 11 - Data management

Steve Simon

# Data management

- Data dictionary
  - Variable names, Variable labels, Value labels, Missing value codes
- Managing complex files
  - Multiple response, Longitudinal/repeated measures data
- Storage options
  - Spreadsheet, Text file, Database, REDCap

# Data dictionary

– Also called a code book

– Start before collecting data

– Revise as needed

# Variable names

– Brief, but descriptive explanation

– Roughly 4 to 16 characters

– No blanks and (almost) no symbols

– One to three words

# Good and bad variable names

– Names to avoid (www.writersexchange.com)
- systolic blood pressure
- systolic-blood-pressure

– Names that work
- systolic_blood_pressure
- systolic.blood.pressure
- SystolicBloodPressure

– NEVER USE ALL CAPS FOR VARIABLE NAMES
- Lower case ascenders (e.g., f and l)
- Lower case descenders (e.g., g and y)

# Variable labels

– Longer descriptions
- Can include spaces and punctuation
- Ideal length is 20-40 characters
- Mention units of measurement, special qualifiers

# Missing value codes

– Explain WHY the value is missing
– For a survey
  • Did not answer
  • Not applicable
– For a lab result
  • Below the limit of detection
  • Insufficient volume for testing
  • Dropped the test tube and it shattered making a huge mess

# Example of missing value codes

– Use extreme number code
  • 9, 99, 999
  • -1
– Use symbols
  • NA
  • (asterisk)
  • (dot)
– Never use blanks to designate missing
– Note missing value code on data dictionary

# Missing value example



Price tage from computer store

# First break

– What have you learned
  - variable names
  - Variable labels
  - Missing value codes
– What's coming up
  - Date formats
  - Categorical values

# Date formats



Cartoon showing variety of data formats

# Internal storage formats

- Excel - number of days since 1899-12-31 (1900-01-00)
- R - number of days since January 1, 1970
- SAS - number of days since January 1, 1960
- SPSS - number of seconds since October 14, 1582

# Gregorian calendar

| JULIAN 1582 | | October | | | Gregorian 1582 | |
|---|---|---|---|---|---|---|
| Sun | Mon | Tues | Wed | Thurs | Fri | Sat |
| | *1* | *2* | *3* | *4* | 15 | 16 |
| 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| 24 | 25 | 26 | 27 | 28 | 29 | 30 |
| 31 | | | | | | |

Transition to the Gregorian calendar

# Gregorian calendar



Painting of Pope Gergory XIII

# Categorical values

– Definition: small number of possible values
– Beware of ambiguities
  • YES, yes, and Yes are three distinct levels.
– Use number codes
  • 0, 1, 9 for binary variables
– Single letter codes
  • M, F, and U for gender
  • Potentially ambiguous
  • Consistent case is important.

# Example of ambiguous coding

**RaceID**

| | | Frequency | Percent | Valid Percent | Cumulative Percent |
|---|---|---|---|---|---|
| Valid | | 2 | 1.1 | 1.1 | 1.1 |
| | W | 1 | .6 | .6 | 1.7 |
| | A | 1 | .6 | .6 | 2.2 |
| | B | 11 | 6.1 | 6.1 | 8.3 |
| | C | 137 | 76.1 | 76.1 | 84.4 |
| | H | 9 | 5.0 | 5.0 | 89.4 |
| | O | 1 | .6 | .6 | 90.0 |
| | W | 18 | 10.0 | 10.0 | 100.0 |
| | Total | 180 | 100.0 | 100.0 | |

SPSS frequencies table for RaceID

# Reverse coding (1 of 2)

– Context specific
– Sequence of IF THEN ELSE statements
  • if (is.na(x)) then y=NA
  • else if (x=1) then y=4
  • else if (x=2) then y=3
  • else if (x=3) then y=2
  • else if (x=4) then y=1
  • else y=9

# Reverse coding (2 of 2)

– Functional transformations
  • 0,1 to 1,0 is $f(x)=1-x$
  • 1,2,3,4 to 4,3,2,1 is $f(x)=5-x$
  • 0,1,2,3,4 to 4,3,2,1,0 is $f(x)=4-x$
– Always check your results
– Watch out for missing value codes

# Second break

– What have you learned
  • Dates
  • Value labels
  • Reverse coding
– What's coming next
  • Multiple response
  • Longitudinal/repeated measures data

# A multiple response example

Q1. What are a few of your favorite things?
☐ a. Raindrops on roses
☒ b. Whiskers on kittens
☒ c. Bright copper kettles
☒ d. Warm woolen mittens

Q1. What are a few of your favorite things?
☒ a. Raindrops on roses
☒ b. Whiskers on kittens
☒ c. Bright copper kettles
☐ d. Warm woolen mittens

Q1. What are a few of your favorite things?
☐ a. Raindrops on roses
☐ b. Whiskers on kittens
☒ c. Bright copper kettles
☒ d. Warm woolen mittens

Questionnaire with a multiple response question

# Coding multiple response with a single column



Multiple response coded into a single column

# A different way to code multiple response



A multiple response question coded into three columns

# The recommended way to code multiple response

| | q1.a | q1.b | q1.c | q1.d |
|---|---|---|---|---|
| Q1. What are a few of your favorite things?<br>☐ a. Raindrops on roses<br>☒ b. Whiskers on kittens<br>☒ c. Bright copper kettles<br>☒ d. Warm woolen mittens | 0 | 1 | 1 | 1 |
| Q1. What are a few of your favorite things?<br>☒ a. Raindrops on roses<br>☒ b. Whiskers on kittens<br>☒ c. Bright copper kettles<br>☐ d. Warm woolen mittens | 1 | 1 | 1 | 0 |
| Q1. What are a few of your favorite things?<br>☐ a. Raindrops on roses<br>☐ b. Whiskers on kittens<br>☒ c. Bright copper kettles<br>☒ d. Warm woolen mittens | 0 | 0 | 1 | 1 |

Multiple response coded with individual item indicators

# Longitudinal data, Repeated measures data

– Longitunal
  • Multiple time points per patient
– Repeated measurements
  • Measuring patient repeatedly under different conditions
– Tall and thin format
  • One line per visit/measurement
– Short and fat format
  • One line per patient

# Example of tall/thin, dictionary

Univariate format:

| Variable | Description |
|---|---|
| Subject | 1 to 40 |
| Sex | male or female |
| Age | Age of subject in years |
| Height | Height in cm |
| Weight | Weight in kg |
| Surface | normal or foam |
| Vision | eyes open, eyes closed, or closed dome |
| CTSIB | Qualitive measure of balance, 1 (stable) - 4 (unstable) |

Data dictionary for tall and thin format

# Example of short/fat, dictionary

Repeated measures format:

| Variable | Description |
|---|---|
| Subject | 1 to 40 |
| Sex | male or female |
| Age | Age of subject in years |
| Height | Height in cm |
| Weight | Weight in kg |
| NO1 | Balance measure on normal surface with eyes open, first replicate |
| NO2 | as above, second replicate |
| NC1 | Balance measure on normal surface with eyes closed, first replicate |
| NC2 | as above, second replicate |
| ND1 | Balance measure on normal surface with dome, first replicate |
| ND2 | as above, second replicate |
| FO1 | Balance measure on foam surface with eyes open, first replicate |
| FO2 | as above, second replicate |
| FC1 | Balance measure on foam surface with eyes closed, first replicate |
| FC2 | as above, second replicate |
| FD1 | Balance measure on foam surface with dome, first replicate |
| FD2 | as above, second replicate |

Repeated measures example in short and fat format

# Example of tall/thin, data



Repeated measures data in tall and thin format

# Example of short/fat, data



Repeated measures data in short and fat format

# A simple alternative to both tall/thin and short/fat

- Disadvantages of tall/thin
  - Too much repetition
- Disadvantages of short/fat
- Database format
  - Time constant table
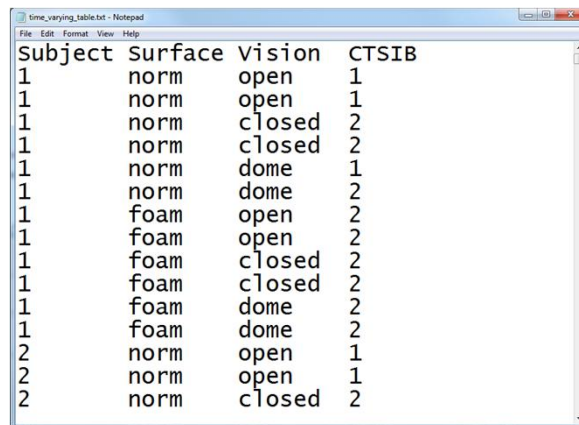  - Time varying table

---

# Time constant data



| Subject | Sex | Age | Height | Weight |
|---------|--------|-----|--------|--------|
| 1 | male | 22 | 176 | 68.2 |
| 2 | male | 22 | 181 | 67.6 |
| 3 | male | 22 | 175.5 | 72 |
| 4 | male | 21 | 180 | 73.2 |
| 5 | female | 20 | 166 | 63.8 |
| 6 | male | 18 | 177 | 78.8 |
| 7 | male | 29 | 183 | 86.4 |
| 8 | female | 22 | 150 | 44.6 |
| 9 | female | 29 | 154 | 57.8 |
| 10 | male | 31 | 176.5 | 80.8 |
| 11 | male | 24 | 176 | 91 |
| 12 | male | 33 | 184 | 89.8 |
| 13 | male | 18 | 187 | 85 |
| 14 | female | 34 | 168 | 54.4 |
| 15 | female | 27 | 173 | 60.8 |

Table listing time constant data only

# Time varying data



Table listing time varying data only

# Contents of a data dictionary

– Variable names

– Variable labels

– Units of measurement

– Permissible/impermissible values

– Value labels

– Missing value codes

– Source

– License

# Third break

- – What have you learned
  - Multiple response variables
  - Longitudinal/repeated measures data
- – What's next
  - Double entry coding
  - Excel files

# Double entry coding

- – Great quality check
  - If you can afford it
- – Prepare a code book first
  - Count the proportion of discrepancies
- – If too many discrepancies
  - Revise the code book and re-do the data entry.
- – If discrepancies small enough
  - Report this number in your publication

# If you enter data into Excel

- Do not use colors
- Do not include summary statistics
- Rectangular grid
- Don't squeeze two data values into one cell
  - Systolic/diastolic blood pressures
  - 44M for a 44 year old male
- Variable names in first row
- No blank cells
  - Contradicts your book

# A poorly structured spreadsheet



A spreadsheet with data

# Revisions to this spreadsheet



A revised and better organized spreadsheet

# The codebook from this spreadsheet



A codebook associated with revised spreadsheet

# Fourth break

– What have you learned
  • Double entry
  • Excel files
– What's coming next
  • Text files
  • Database files

# Text files

– Fixed width
– Delimited
  • Commas
  • Spaces
  • Tabs
  • "Quotes around text"

# Data dictionary for aboriginal prison death study

**StatSci.org** / Home

OzDASL

**Aboriginal Deaths in Custody**

Keywords: binomial regression.

**Description**

The data give the number of deaths in prison custody in Australia in each of the six years 1990 to 1995, given separately for Aboriginal and Torres Strait Islanders (indigenous) and others (non-indigenous).

| Variable | Description |
|---|---|
| Year | 1990 through 1995 |
| Indigenous | Yes = Aboriginal or Torres Strait Islander, No = Non-indigenous |
| Prisoners | Total number in prison custody |
| Deaths | Number of deaths in prison custody |
| Population | Adult population (15+ years) |

The data were collected in response to the Royal Commission into Aboriginal Deaths in Custody, the final report of which was tabled in the Federal Parliament on the 9 May 1991.

Data dictionary

# Comma separated values (csv)

aboriginal_data_comma - Notepad

File  Edit  Format  View  Help

```
Year,Indigenous,Prisoners,Deaths,Population
1990,Yes,2041,6,168317
1991,Yes,2166,8,172462
1992,Yes,2223,2,176827
1993,Yes,2416,7,181341
1994,Yes,2742,11,185836
1995,Yes,2907,17,190438
1990,No,12264,27,13141817
1991,No,12855,31,13326044
1992,No,13336,34,13501987
1993,No,13450,42,13649262
1994,No,14302,42,13810095
1995,No,14501,42,13995940
```

Data set using a comma separated value format

# Comma separated values with quotes



Data set using a quoted format

# Fixed width format



Data set using a fixed width format

# Spaced format

```
aboriginal_data_space - Notepad                                    —  □  ×
File  Edit  Format  View  Help
Year Indigenous Prisoners Deaths Population
1990        Yes      2041      6    168317
1991        Yes      2166      8    172462
1992        Yes      2223      2    176827
1993        Yes      2416      7    181341
1994        Yes      2742     11    185836
1995        Yes      2907     17    190438
1990         No     12264     27  13141817
1991         No     12855     31  13326044
1992         No     13336     34  13501987
1993         No     13450     42  13649262
1994         No     14302     42  13810095
1995         No     14501     42  13995940
```
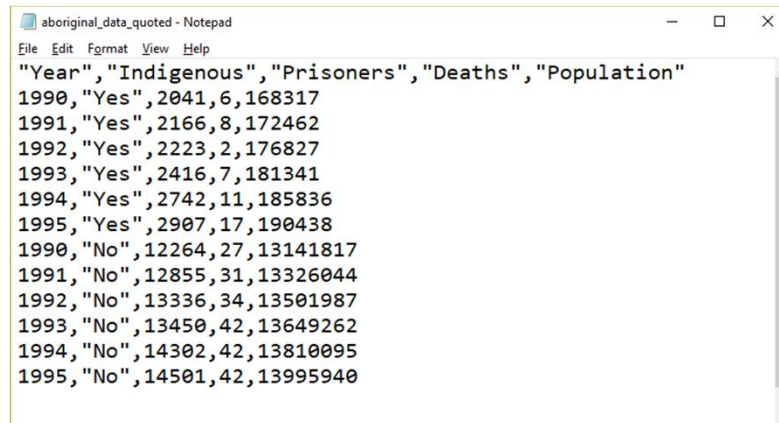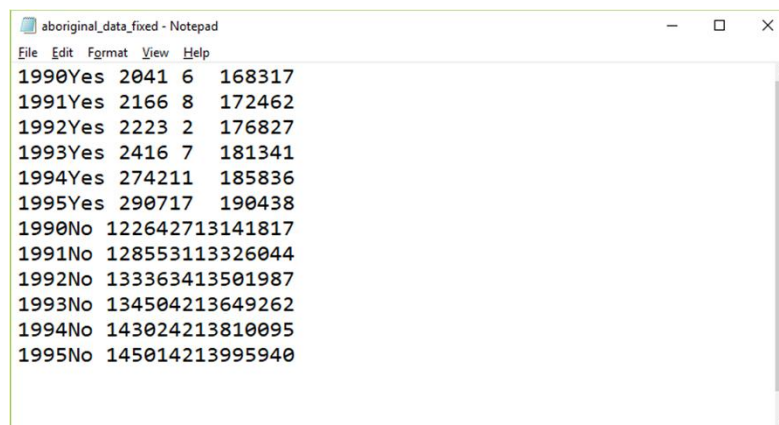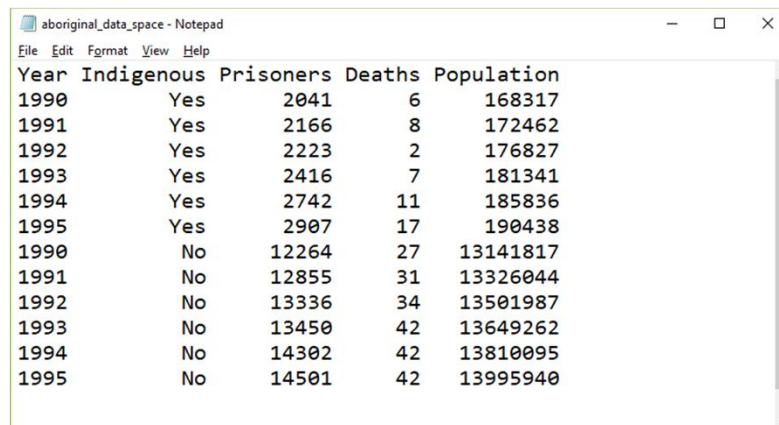
Data set using a spaced format

# Tab separated values

```
aboriginal_data_tab - Notepad                                      —  □  ×
File  Edit  Format  View  Help
Year    Indigenous      Prisoners       Deaths  Population
1990    Yes     2041    6       168317
1991    Yes     2166    8       172462
1992    Yes     2223    2       176827
1993    Yes     2416    7       181341
1994    Yes     2742    11      185836
1995    Yes     2907    17      190438
1990    No      12264   27      13141817
1991    No      12855   31      13326044
1992    No      13336   34      13501987
1993    No      13450   42      13649262
1994    No      14302   42      13810095
1995    No      14501   42      13995940
```
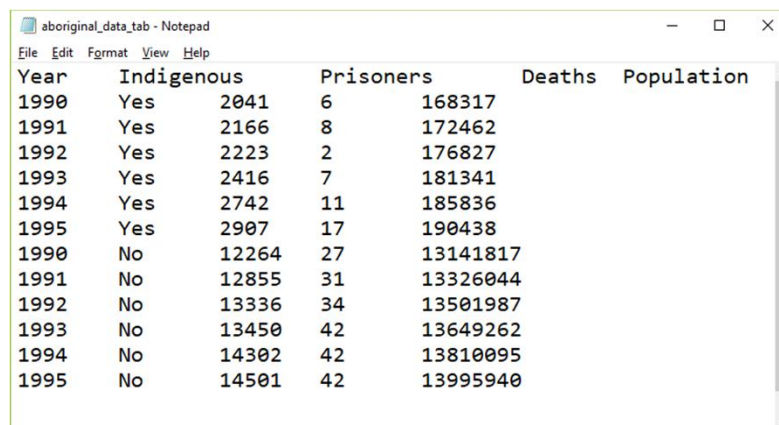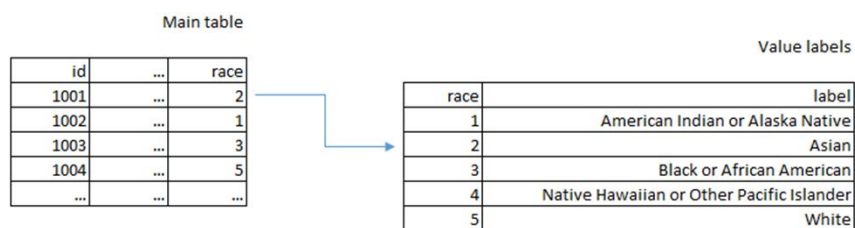
Data set using a tab separated value format

# Database systems

– Terminology
  - Tables
  - Fields
  - Records
  - Primary key
  - Foreign key

# Database table for value labels

| Main table | | |
| --- | --- | --- |
| id | ... | race |
| 1001 | ... | 2 |
| 1002 | ... | 1 |
| 1003 | ... | 3 |
| 1004 | ... | 5 |
| ... | ... | ... |

Value labels

| race | label |
| --- | --- |
| 1 | American Indian or Alaska Native |
| 2 | Asian |
| 3 | Black or African American |
| 4 | Native Hawaiian or Other Pacific Islander |
| 5 | White |

Database linkage between race code and race labels

# REDCap

– Research Electronic Data Capture
– Not open source, but freely distributed by Vanderbilt
– Software components
  • PHP
  • JavaScript
  • MySQL
– Case report forms
– Strongly recommended

# Conclusion

– Data dictionary
  • Variable names, Variable labels, Value labels, Missing value codes
– Managing complex files
  • Multiple response, Longitudinal/repeated measures data
– Storage options
  • Spreadsheet, Text file, Database, REDCap