Tidy and Messy Data

"Tidy" Data

- Idea formulated by Hadley Wickham
- "Tidy" does not mean "good" or "reliable" or "useful."
- It now has a technical meaning when applied to a data set.

Which one is tidy? (AMA)

		treatmenta	treatmentb
	John Smith		2
D1:	Jane Doe	16	11
	Mary Johnson	3	1

person	treatment	result
John Smith	a	
Jane Doe	a	16
Mary Johnson	a	3
John Smith	b	2
Jane Doe	b	11
Mary Johnson	b	1

D2:

Tidy data:

- 1. Each variable forms a column
- 2. Each observation forms a row
- 3. Each type of observational unit forms a table

Closely related to Codd's 3rd Normal Form (database literature)

Variable? Observation? Type of OU?

		treatmenta	treatmentb
	John Smith		2
D1:	Jane Doe	16	11
	Mary Johnson	3	1

person	treatment	result
John Smith	a	
Jane Doe	a	16
Mary Johnson	a	3
John Smith	b	2
Jane Doe	b	11
Mary Johnson	b	1

D2:

religion	<\$10k	\$10-20k	\$20–30k	\$30–40k	\$40–50k	\$50-75k
Agnostic	27	34	60	81	76	137
Atheist	12	27	37	52	35	70
Buddhist	27	21	30	34	33	58
Catholic	418	617	732	670	638	1116
Don't know/refused	15	14	15	11	10	35
Evangelical Prot	575	869	1064	982	881	1486
Hindu	1	9	7	9	11	34
Historically Black Prot	228	244	236	238	197	223
Jehovah's Witness	20	27	24	24	21	30
Jewish	19	19	25	25	30	95

Table 4: The first ten rows of data on income and religion from the Pew Forum. Three columns, \$75–100k, \$100–150k and >150k, have been omitted.

What would a tidy version look like?

"Melting" data

				ľ	OW	column	value
					4	a	1
				I	3	a	2
row	a	b	С	(С	a	3
А	1	4	7	I	4	b	4
В	2	5	8	I	3	b	5
С	3	6	9	(С	b	6
(a)]	Raw	data	ł	ľ	4	С	7
				I	3	С	8
				(2	С	9

(b) Molten data

First rows of tidy version of PEW data

religion	income	freq
Agnostic	<\$10k	27
Agnostic	\$10–20k	34
Agnostic	\$20–30k	60
Agnostic	\$30–40k	81
Agnostic	\$40–50k	76
Agnostic	50-75k	137
Agnostic	\$75–100k	122
Agnostic	\$100–150k	109
Agnostic	>150k	84
Agnostic	Don't know/refused	96

country	year	m014	m1524	m2534	m3544	m4554	m5564	m65	mu	f014
AD	2000	0	0	1	0	0	0	0		
AE	2000	2	4	4	6	5	12	10		3
AF	2000	52	228	183	149	129	94	80		93
AG	2000	0	0	0	0	0	0	1		1
AL	2000	2	19	21	14	24	19	16		3
AM	2000	2	152	130	131	63	26	21		1
AN	2000	0	0	1	2	0	0	0		0
AO	2000	186	999	1003	912	482	312	194		247
AR	2000	97	278	594	402	419	368	330		121
AS	2000					1	1			

country	year	column	cases
AD	2000	m014	0
AD	2000	m1524	0
AD	2000	m2534	1
AD	2000	m3544	0
AD	2000	m4554	0
AD	2000	m5564	0
AD	2000	m65	0
AE	2000	m014	2
AE	2000	m1524	4
AE	2000	m2534	4
AE	2000	m3544	6
AE	2000	m4554	5
AE	2000	m5564	12
AE	2000	m65	10
AE	2000	f014	3

country	year	sex	age	cases
AD	2000	m	0–14	0
AD	2000	m	15 - 24	0
AD	2000	m	25 - 34	1
AD	2000	m	35 - 44	0
AD	2000	m	45 - 54	0
AD	2000	m	55 - 64	0
AD	2000	m	65 +	0
AE	2000	m	0 - 14	2
AE	2000	m	15 - 24	4
AE	2000	m	25 - 34	4
AE	2000	m	35 - 44	6
AE	2000	m	45 - 54	5
AE	2000	m	55 - 64	12
AE	2000	m	65 +	10
AE	2000	f	0-14	3

id	year	month	element	d1	d2	d3	d4	d5	d6	d7	d8
MX17004	2010	1	tmax								
MX17004	2010	1	tmin								
MX17004	2010	2	tmax		27.3	24.1					
MX17004	2010	2	tmin		14.4	14.4					
MX17004	2010	3	tmax					32.1			
MX17004	2010	3	tmin					14.2			
MX17004	2010	4	tmax								
MX17004	2010	4	tmin								
MX17004	2010	5	tmax								
MX17004	2010	5	tmin								

Table 11: Original weather dataset. There is a column for each possible day in the month. Columns d9 to d31 have been omitted to conserve space.

id	date	element	value
MX17004	2010-01-30	tmax	$\overline{27.8}$
MX17004	2010-01-30	tmin	14.5
MX17004	2010-02-02	tmax	27.3
MX17004	2010-02-02	tmin	14.4
MX17004	2010-02-03	tmax	24.1
MX17004	2010-02-03	tmin	14.4
MX17004	2010-02-11	tmax	29.7
MX17004	2010-02-11	tmin	13.4
MX17004	2010-02-23	tmax	29.9
MX17004	2010-02-23	tmin	10.7

id	date	tmax	tmin
MX17004	2010-01-30	27.8	14.5
MX17004	2010-02-02	27.3	14.4
MX17004	2010-02-03	24.1	14.4
MX17004	2010-02-11	29.7	13.4
MX17004	2010-02-23	29.9	10.7
MX17004	2010-03-05	32.1	14.2
MX17004	2010-03-10	34.5	16.8
MX17004	2010-03-16	31.1	17.6
MX17004	2010-04-27	36.3	16.7
MX17004	2010-05-27	33.2	18.2

Why tidy?

- Analysis tools typically want tidy data.
- Visualization tools typically want tidy data.

Why not tidy?

References

- Idea of Tidy:
- <u>http://vita.had.co.nz/papers/tidy-data.html</u>