

Plan:

1. Define tidy data & explain its benefits
2. Understand how data can be messy
3. Explain how tidy data are an intermediate step for many data science projects

Tidy Data



Shannon E. Ellis, Ph.D
UC San Diego

Department of Cognitive Science
sellis@ucsd.edu

"Good data scientists understand, in a deep way, that the heavy lifting of cleanup and preparation isn't something that gets in the way of solving the problem: it is the problem." - DJ Patil

untidy data

Australian Bureau of Statistics												
1800.0 Australian Marriage Law Postal Survey, 2017												
Released on 15 November 2017												
Table 5 Participation by Federal Electoral Division(a), Males and Age Gender apartheid												
Yeah NA												
Total participants 292 1,056 1,465 1,653 1,515 1,516 1,710 1,730 1,753 1,574 Eligible participants 572 2,910 3,789 3,996 3,607 3,506 3,645 3,331 2,960 2,456 Participation rate (%) 51.0 36.4 38.7 41.4 42.0 43.2 46.9 51.9 59.2 64.1												
Total participants 442 1,461 2,066 2,357 2,188 2,057 2,224 2,108 2,134 1,772 Eligible participants 750 2,991 3,994 4,155 3,634 3,358 3,427 3,066 2,931 2,355 Participation rate (%) 58.9 48.8 51.7 56.7 60.2 60.5 64.9 68.6 72.8 75.2												
Total participants 734 2,519 3,531 4,010 3,703 3,573 3,934 3,838 3,887 3,346 Eligible participants 1,322 5,901 7,783 8,151 7,241 6,904 7,072 6,397 5,891 4,811 Participation rate (%) 55.5 42.7 45.4 49.2 51.1 51.8 55.6 60.0 66.0 69.5												
Total participants 1,764 4,789 4,817 4,973 4,626 4,453 5,074 4,826 5,169 4,394 Eligible participants 2,260 6,471 6,448 6,509 5,953 5,805 6,302 5,902 6,044 5,057 Participation rate (%) 78.1 74.0 74.7 76.4 77.3 76.7 80.5 81.8 85.5 86.9												
Total participants 1,477 4,687 5,178 5,786 6,025 5,463 5,191 4,208 3,948 3,465 Eligible participants 1,904 6,354 7,121 7,802 7,960 7,155 6,480 4,692 4,692 3,945 Participation rate (%) 77.6 73.6 72.2 74.0 75.7 76.4 80.1 80.5 84.1 87.8												
Total participants 2,542 9,470 9,999 10,759 10,095 9,940 10,289 9,094 9,147 7,959 Eligible participants 4,164 12,825 13,569 14,331 13,943 12,960 12,782 11,108 10,736 9,000 Participation rate (%) 77.8 73.9 73.7 75.1 76.4 76.5 80.3 81.3 84.9 87.3												
Total participants 151,297 438,166 441,658 460,548 462,206 479,360 524,620 517,693 543,449 506,799 Eligible participants 201,439 635,909 646,916 665,250 656,446 660,841 693,850 659,150 664,720 597,386 Participation rate (%) 75.1 68.9 68.3 69.2 70.4 72.5 75.6 78.5 81.8 84.5												
(a) The Federal Electoral Divisions are current as at 24 August 2017 (b) Includes those whose age is unknown (c) Includes Christmas Island and the Cocos (Keeling) Islands (d) Includes Norfolk Island (e) Includes Jarvis Bay												

tidy data

area	gender	age	State	Area (sq km)	Eligible participants	Participation rate (%)	Total participants	Total Participants
Adelaide	Female	18-19 years	SA	76	1341	83.5	1120	1120
Adelaide	Female	20-24 years	SA	76	4620	81.2	3750	3750
Adelaide	Female	25-29 years	SA	76	4897	81.8	4004	4004
Adelaide	Female	30-34 years	SA	76	4784	79.8	3820	3820
Adelaide	Female	35-39 years	SA	76	4319	79	3411	3411
Adelaide	Female	40-44 years	SA	76	4310	80.6	3472	3472
Adelaide	Female	45-49 years	SA	76	4579	81.4	3728	3728
Adelaide	Female	50-54 years	SA	76	4475	84.7	3791	3791
Adelaide	Female	55-59 years	SA	76	4622	87.3	4033	4033
Adelaide	Female	60-64 years	SA	76	4342	89.3	3879	3879
Adelaide	Female	65-69 years	SA	76	3970	90.7	3602	3602
Adelaide	Female	70-74 years	SA	76	3009	90.3	2716	2716
Adelaide	Female	75-79 years	SA	76	2156	88.5	1908	1908
Adelaide	Female	80-84 years	SA	76	1673	85.1	1423	1423

data
→
wrangling

Tidy Data

1. Each **variable** you measure should be in a single column

	A	B	C	D	E	F	G
1	ID	LastName	FirstName	Sex	City	State	Occupation
2	1004	Smith	Jane	female	Frederick	MD	Welder
3	4587	Nayef	Mohammed	male	Upper Darby	PA	Nurse
4	1727	Doe	Janice	female	San Diego	CA	Doctor
5	6879	Jordan	Alex	male	Birmingham	AL	Teacher

2. Every **observation** of a variable should be in a different row

	A	B	C	D	E	F	G
1	ID	LastName	FirstName	Sex	City	State	Occupation
2	1004	Smith	Jane	female	Frederick	MD	Welder
3	4587	Nayef	Mohammed	male	Upper Darby	PA	Nurse
4	1727	Doe	Janice	female	San Diego	CA	Doctor
5	6879	Jordan	Alex	male	Birmingham	AL	Teacher

3. There should be one table for each type of data

Demographic Survey Data

	A	B	C	D	E	F	G
1	ID	LastName	FirstName	Sex	City	State	Occupation
2	1004	Smith	Jane	female	Frederick	MD	Welder
3	4587	Nayef	Mohammed	male	Upper Darby	PA	Nurse
4	1727	Doe	Janice	female	San Diego	CA	Doctor
5	6879	Jordan	Alex	male	Birmingham	AL	Teacher

Doctor's Office Measurements Data

	A	D	E	F	G
1	ID	Height_inches	Weight_lbs	Insulin	Glucose
2	1004	65	180	0.60	163
3	4587	75	215	1.46	150
4	1727	62	124	0.72	177
5	6879	77	160	1.23	205

4. If you have multiple tables, they should include a column in each *with the same column label* that allows them to be joined or merged

	A	B	C	D	E	F	G
1	ID	LastName	FirstName	Sex	City	State	Occupation
2	1004	Smith	Jane	female	Frederick	MD	Welder
3	4587	Nayef	Mohammed	male	Upper Darby	PA	Nurse
4	1727	Doe	Janice	female	San Diego	CA	Doctor
5	6879	Jordan	Alex	male	Birmingham	AL	Teacher

	A	D	E	F	G
1	ID	Height_inches	Weight_lbs	Insulin	Glucose
2	1004	65	180	0.60	163
3	4587	75	215	1.46	150
4	1727	62	124	0.72	177
5	6879	77	160	1.23	205

Tidy data == rectangular data

A

	A	B	C	D	E
1	id	sex	glucose	insulin	triglyc
2	101	Male	134.1	0.60	273.4
3	102	Female	120.0	1.18	243.6
4	103	Male	124.8	1.23	297.6
5	104	Male	83.1	1.16	142.4
6	105	Male	105.2	0.73	215.7

Tidy Data Benefits

1. consistent data structure
2. foster tool development
3. require only a small set of tools to be learned
4. allow for datasets to be combined

Common Problems with Messy Data Sets

1. Column headers are values but should be variable names.
2. A single column has multiple variables.
3. Variables have been entered in both rows and columns.
4. Multiple "types" of data are in the same spreadsheet.
5. A single observation is stored across multiple spreadsheets.

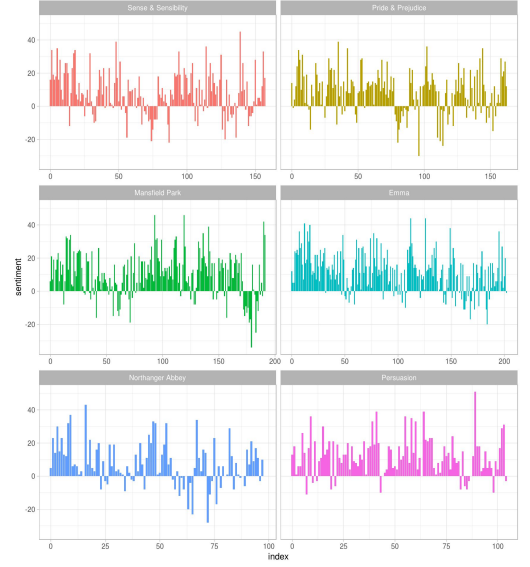


 text

tidy dataset

Word	Novel	Frequency
good	Emma	359
young	Emma	192
friend	Emma	166

results

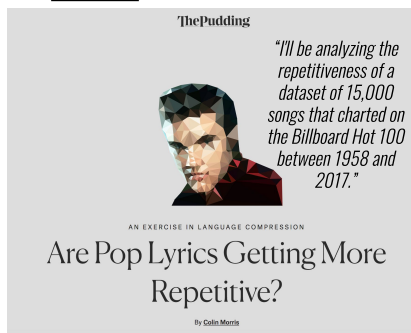




■ Told a public lie
 ■ Told a public falsehood
 ■ Didn't tell a public lie or falsehood



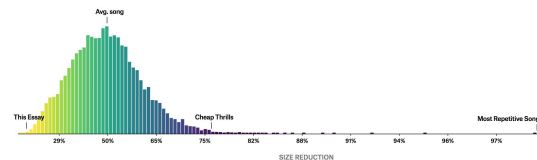
text (lyrics)



tidy dataset

song	Artist	Released	Reduction
Cheap Thrills	Sia	2016	76
Around The World	Daft Punk	1997	98
Everybody Dies	J. Cole	2018	27

results



What are these uber-repetitive outliers? *Around The World* by Daft Punk gets reduced a whopping 98%. It goes from 2,610 characters to 61. Small enough to fit in a tweet - twice!