

Course Reminders

- **Due tonight (11:59 PM)**
 - Project Survey
 - Q1
 - D1
 - A1
- **No Class Monday (MLK, Jr. Day)**
- Projects
 - You will be assigned a GitHub repo this weekend - please accept the invitation (it will expire)
 - You will also be assigned a previous project to review (links on Canvas)

Data

Data Structures & Tidy Data

Shannon E. Ellis, Ph.D
UC San Diego

• • •

Department of Cognitive Science
sellis@ucsd.edu

Data Structures Review

Structured data

- can be stored in database SQL
- tables with rows and columns
- requires a relational key
- 5-10% of all data

Semi-structured data

- doesn't reside in a relational database
- has organizational properties (easier to analyze)
- CSV, XML, JSON

Unstructured

- non-tabular data
- 80% of the world's data
- images, text, audio, videos

(Semi-)Structured Data

Data that is stored in such a way that it is easy to search and work with. These data are stored in a particular format that adheres to organization principles imposed by the file format. These are the data structures data scientists work with most often.

CSVs

Each column separated by a comma

Has the extension ".csv"

Example CSV - Sheet1 — Notatnik				
Plik	Edycja	Format	Widok	Pomoc
Email	First Name	Last Name	Company	Snippet 1
example1@domain.com	John	Smith	Company 1	Snippet Sentence1
example2@gmail.com	Mary	Blake	Company 2	Snippet Sentence 2
example3@outlook.com	James	Joyce	Company 3	Snippet Sentence 3

Each row is separated by a new line



Example CSV



File Edit View Insert Format Data Tools Add-ons Help All changes saved in Drive

undo redo print filter | 100% | \$ % .0 .00 123 | Arial | 10 | B I S A | field

fx

	A	B	C	D	E	F
1	Email	First Name	Last Name	Company	Snippet 1	
2	example1@domain.com	John	Smith	Company 1	Snippet Sentence1	
3	example2@gmail.com	Mary	Blake	Company 2	Snippet Sentence 2	
4	example3@outlook.com	James	Joyce	Company 3	Snippet Sentence 3	

CSV file



Example CSV - Sheet1 — Notatnik
Plik Edycja Format Widok Pomoc
Email,First Name,Last Name,Company,Snippet 1
example1@domain.com,John,Smith,Company 1,Snippet Sentence1
example2@gmail.com,Mary,Blake,Company 2,Snippet Sentence 2
example3@outlook.com,James,Joyce,Company 3,Snippet Sentence 3

JSON: key-value pairs

nested/hierarchical data

```
{"Name": "Isabela"}
```

key

value

JSON

These are all
nested within
attributes

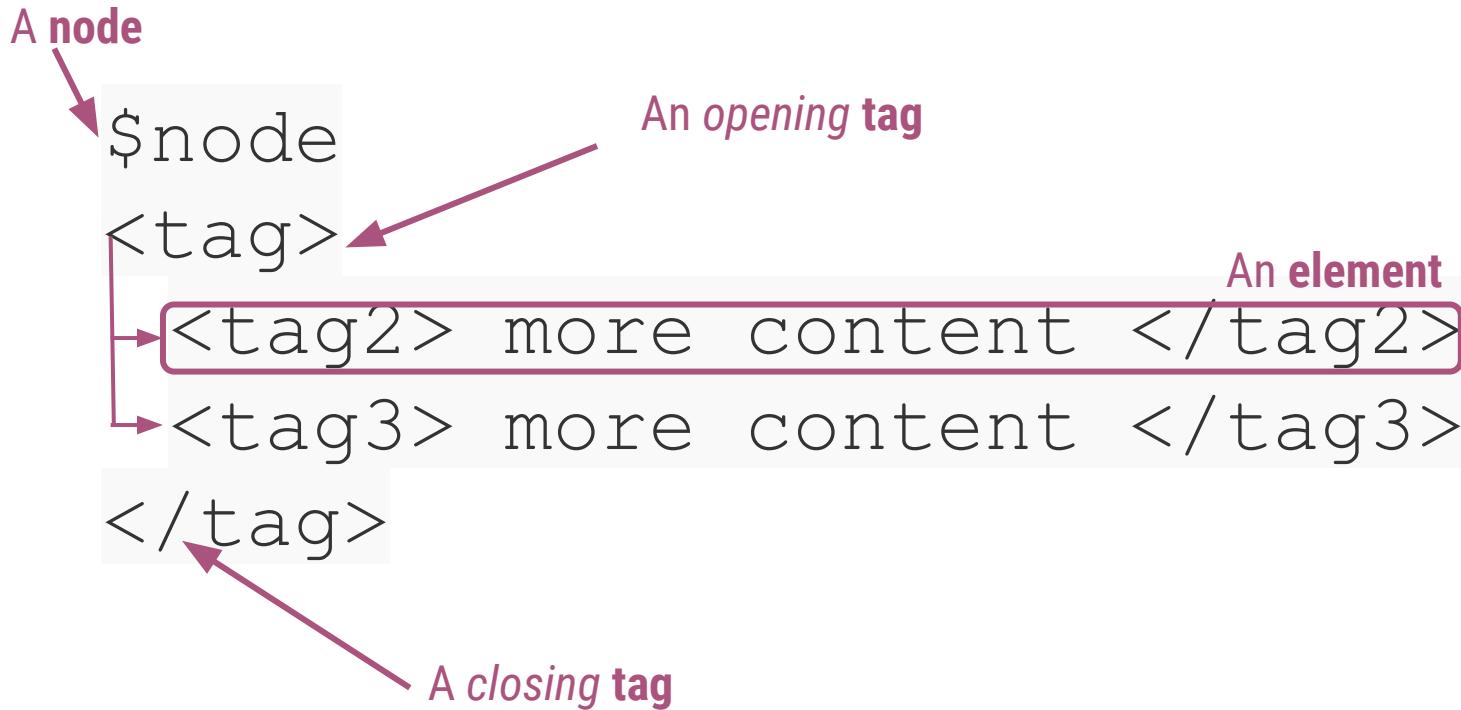
```
"attributes": {  
    "Take-out": true,  
    "Wi-Fi": "free",  
    "Drive-Thru": true,  
    "Good For": {  
        "dessert": false,  
        "latenight": false,  
        "lunch": false,  
        "dinner": false,  
        "breakfast": false,  
        "brunch": false  
    },
```

These are all
nested within
"Good For"

JSON

Extensible Markup Language (XML): nodes, tags, and elements

nested/hierarchical data



XML

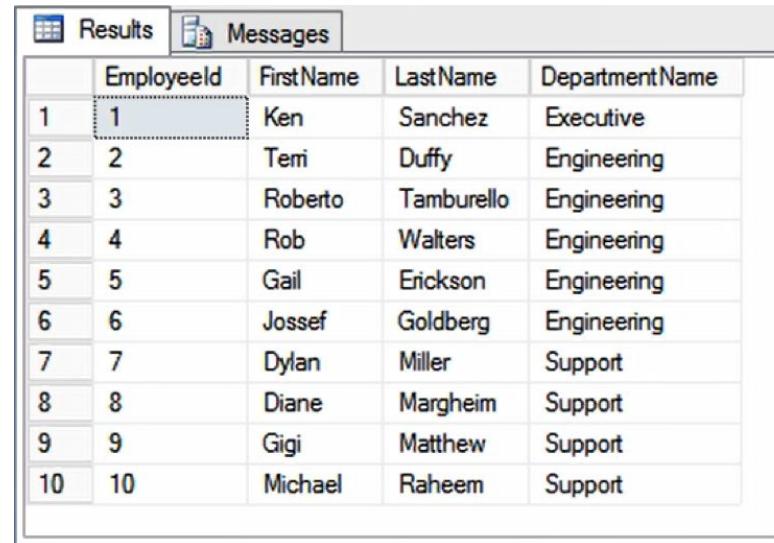
```
<?xml version="1.0" encoding="UTF-8"?>
<customers>
    <customer>
        <customer_id>1</customer_id>
        <first_name>John</first_name>
        <last_name>Doe</last_name>
        <email>john.doe@example.com</email>
    </customer>
    <customer>
        <customer_id>2</customer_id>
        <first_name>Sam</first_name>
        <last_name>Smith</last_name>
        <email>sam.smith@example.com</email>
    </customer>
    <customer>
        <customer_id>3</customer_id>
        <first_name>Jane</first_name>
        <last_name>Doe</last_name>
        <email>jane.doe@example.com</email>
    </customer>
</customers>
```

XML

adapted from Chris Keown

Relational Databases: A set of interdependent tables

1. Efficient Data Storage
2. Avoid Ambiguity
3. Increase Data Privacy



The screenshot shows a Microsoft SQL Server Management Studio (SSMS) interface with the 'Results' tab selected. The results grid displays a table of employee data with columns: EmployeeId, FirstName, LastName, and DepartmentName. The data shows 10 employees, each assigned to the 'Engineering' or 'Support' department. The first row (EmployeeId 1) is highlighted with a dashed border.

	EmployeeId	FirstName	LastName	DepartmentName
1	1	Ken	Sanchez	Executive
2	2	Teni	Duffy	Engineering
3	3	Roberto	Tamburello	Engineering
4	4	Rob	Walters	Engineering
5	5	Gail	Erickson	Engineering
6	6	Jossef	Goldberg	Engineering
7	7	Dylan	Miller	Support
8	8	Diane	Margheim	Support
9	9	Gigi	Matthew	Support
10	10	Michael	Raheem	Support

relational database

Information is stored across tables

	unique_identifier	
	AH13JK	
	JJ29JJ	
	CI21AA	

	unique_identifier	
	AH13JK	
	JJ29JJ	
	JJ29JJ	
	XJ11AS	
	CI21AA	

	unique_identifier	
	AH13JK	
	SE92FE	
	CI21AA	

entries are *related* to one another by their unique identifier

relational database

restaurant

name	id	address	type
Taco Stand	AH13JK	1 Main St.	Mexican
Pho Place	JJ29JJ	192 Street Rd.	Vietnamese
Taco Stand	XJ11AS	18 W. East St.	Fusion
Pizza Heaven	CI21AA	711 K Ave.	Italian

health inspections

id	inspection_date	inspector	score
AH13JK	2018-08-21	Sheila	97
JJ29JJ	2018-03-12	D'eonte	98
JJ29JJ	2018-01-02	Monica	66
XJ11AS	2018-12-16	Mark	43
CI21AA	2018-08-21	Anh	99

rating

id	stars
AH13JK	4.9
JJ29JJ	4.8
XJ11AS	4.2
CI21AA	4.7

relational database

restaurant

name	id	address	type
Taco Stand	AH13JK	1 Main St.	Mexican
Pho Place	JJ29JJ	192 Street Rd.	Vietnamese
Taco Stand	XJ11AS	18 W. East St.	Fusion
Pizza Heaven	CI21AA	711 K Ave.	Italian

Two different restaurants with
the same name will have
different unique identifiers

health inspections

id	inspection_date	inspector	score
AH13JK	2018-08-21	Sheila	97
JJ29JJ	2018-03-12	D'eonte	98
JJ29JJ	2018-01-02	Monica	66
XJ11AS	2018-12-16	Mark	43
CI21AA	2018-08-21	Anh	99

rating

id	stars
AH13JK	4.9
JJ29JJ	4.8
XJ11AS	4.2
CI21AA	4.7

relational database

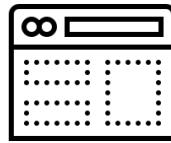
Unstructured Data

Some datasets record information about the state of the world, but in a more heterogeneous way. Perhaps it is a large text corpus with images and links like Wikipedia, or the complicated mix of notes and test results appearing in personal medical records.

Unstructured Data Types



Text files and documents



Websites and applications



Sensor data



Image files



Audio files



Video files



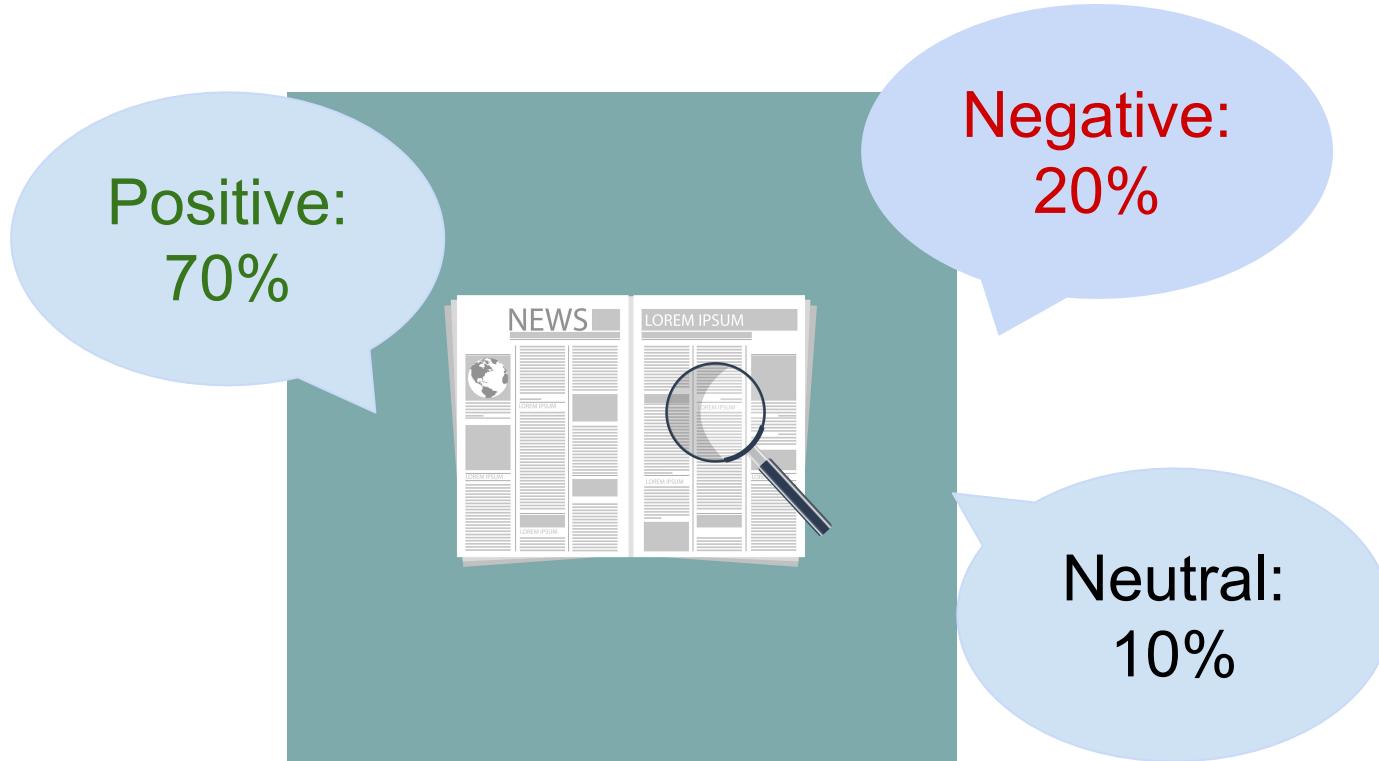
Email data

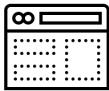


Social media data



Text: Sentiment Analysis

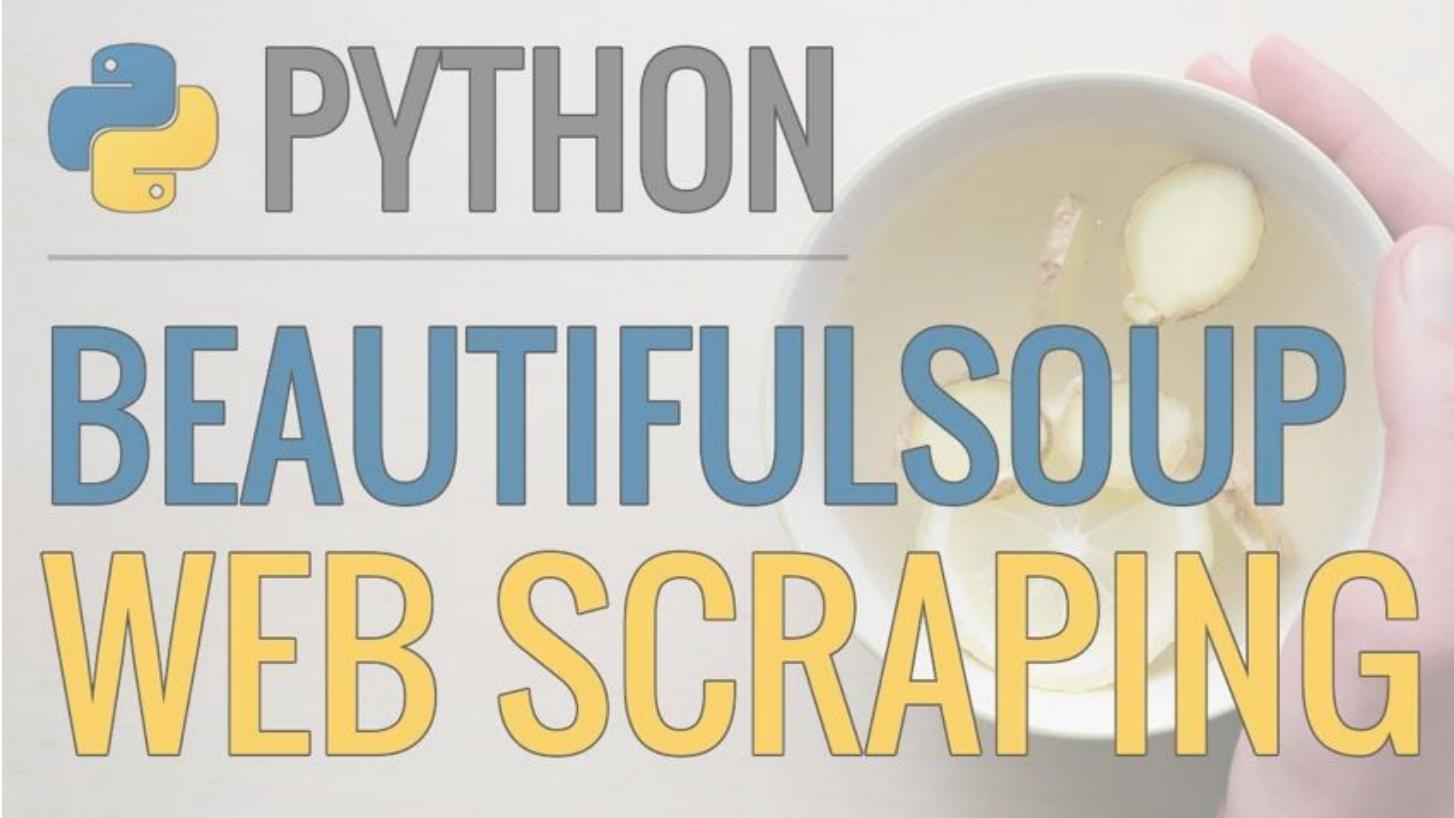




PYTHON

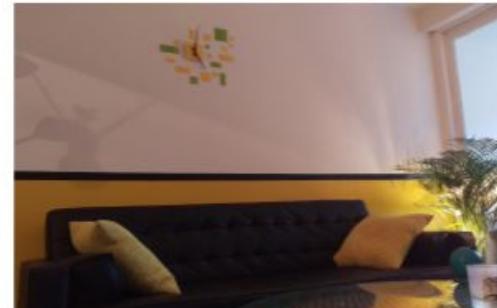
BEAUTIFULSOUP

WEB SCRAPING





Bedroom Or Not?



"The left two photos were correctly predicted as bedrooms; The right two photos were correctly predicted NOT as bedrooms."

Tidy Data

"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 15 November 2017													
Table junk													
1	Australian Bureau of Statistics	Yeah NA	18-19 years	20-24 years	25-29 years	30-34 years	35-39 years	40-44 years	45-49 years	50-54 years	55-59 years	60-64 years	
2	Lingua(c)	Total participants	292	1,058	1,460	1,653	1,515	1,516	1,710	1,730	1,753	1,574	
3	Eligible participants	572	2,910	3,040	3,096	3,607	3,506	3,645	3,331	2,960	2,456		
4	Participation rate (%)	51.0	36.4	38.7	41.4	42.0	43.2	46.8	51.9	58.2	64.1		
5	Primary keynotes	Comma on											
6	Merged cells	Total participants	442	1,461	2,068	2,357	2,188	2,057	2,224	2,108	2,134	1,772	
7	Solomon	Eligible participants	750	2,991	3,994	4,155	3,634	3,398	3,427	3,066	2,931	2,355	
8	Participation rate (%)	56.9	48.8	51.7	56.7	60.2	60.5	64.9	68.8	72.8	75.2		
9	Northern Territory	Total participants	734	2,519	3,531	4,010	3,703	3,573	3,934	3,838	3,887	3,346	
10	(Total)	Eligible participants	1,322	5,901	7,783	8,151	7,241	6,904	7,072	6,397	5,891	4,811	
11	Participation rate (%)	55.5	42.7	45.4	49.2	51.1	51.8	55.6	60.0	66.0	69.5		
12	Australian Capital Territory Divisions	Subheading											
13	Covariate as Subheading	Summary of data inside data											
14	Canberra(d)	Total participants	1,764	4,789	4,817	4,973	4,626	4,453	5,074	4,826	5,169	4,394	
15	Eligible participants	2,260	6,471	6,446	6,509	5,983	5,805	6,302	5,902	6,044	5,057		
16	Participation rate (%)	78.1	74.0	74.7	76.4	77.3	76.7	80.5	81.8	85.5	86.9		
17	Fisher(e)	Total participants	1,477	4,687	5,178	5,786	6,025	5,463	5,191	4,208	3,948	3,465	
18	Eligible participants	1,904	6,354	7,121	7,822	7,960	7,155	6,480	5,206	4,692	3,945		
19	Participation rate (%)	77.6	73.8	72.7	74.0	75.7	76.4	80.1	80.8	84.1	87.8		
20	NA Yeah												
21	Australian Capital Territory (Total)	Total participants	4,242	9,476	9,895	10,155	10,054	9,219	10,205	9,854	9,411	7,659	
22	Eligible participants	4,164	12,825	15,569	14,331	13,943	12,960	12,782	11,108	10,736	9,002		
23	Participation rate (%)	77.8	73.9	73.7	75.1	76.4	76.5	80.3	81.3	84.9	87.3		
24	Australia	Total participants	151,297	438,166	441,658	460,548	462,206	479,360	524,620	517,693	543,449	506,799	
25	Eligible participants	201,439	635,909	646,916	665,250	656,446	660,841	693,850	659,150	664,720	597,386		
26	Participation rate (%)	75.1	68.9	68.3	69.2	70.4	72.5	75.6	78.5	81.8	84.8		
27	a) The Federal Electoral Divisions are current as at 24 August 2017												
28	b) Includes those whose age is unknown												
29	c) Includes Christmas Island and the Cocos (Keeling) Islands												
30	d) Includes Norfolk Island												
31	e) Includes Jervis Bay												
32	Return of the table junk												
33	MS Excel or Die												

tidy data

data
wrangling

area	gender	age	State	Area (sq km)	Eligible participants	Participation rate (%)	Total participants	Total Participants
1	Adelaide	Female	SA	76	1341	83.5	1120	1120
2	Adelaide	Female	SA	76	4620	81.2	3750	3750
3	Adelaide	Female	SA	76	4897	81.8	4004	4004
4	Adelaide	Female	SA	76	4784	79.8	3820	3820
5	Adelaide	Female	SA	76	4319	79	3411	3411
6	Adelaide	Female	SA	76	4310	80.6	3472	3472
7	Adelaide	Female	SA	76	4579	81.4	3728	3728
8	Adelaide	Female	SA	76	4475	84.7	3791	3791
9	Adelaide	Female	SA	76	4622	87.3	4033	4033
10	Adelaide	Female	SA	76	4342	89.3	3879	3879
11	Adelaide	Female	SA	76	3970	90.7	3602	3602
12	Adelaide	Female	SA	76	3009	90.3	2716	2716
13	Adelaide	Female	SA	76	2156	88.5	1908	1908
14	Adelaide	Female	SA	76	1673	85.1	1423	1423

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
3		4587	Nayef	Mohammed	male	Upper Darby	PA
4		1727	Doe	Janice	female	San Diego	CA
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
3		4587	75	215	1.46
4		1727	62	124	0.72
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.



Tabular Data Time

A

ID	Last	First	height_m	height_f
1004	Smith	Jane	NA	65
4587	Nayef	Mohammed	72	NA
1727	Doe	Janice	NA	60
6879	Jordan	Alex	55	NA

B

ID	Last	First	height_m	height_f
1004	Smith	Jane		65
4587	Nayef	Mohammed	72	
1727	Doe	Janice		60
6879	Jordan	Alex	55	

C

ID	Last	First	sex	height
1004	Smith	Jane	female	65
4587	Nayef	Mohammed	male	72
1727	Doe	Janice	fem	60
6879	Jordan	Alex	male	55

D

ID	Last	First	sex	height
1004	Smith	Jane	F	65
4587	Nayef	Mohammed	M	72
1727	Doe	Janice	F	60
6879	Jordan	Alex	M	55

Which of these tables stores data best?



A



B



C



D



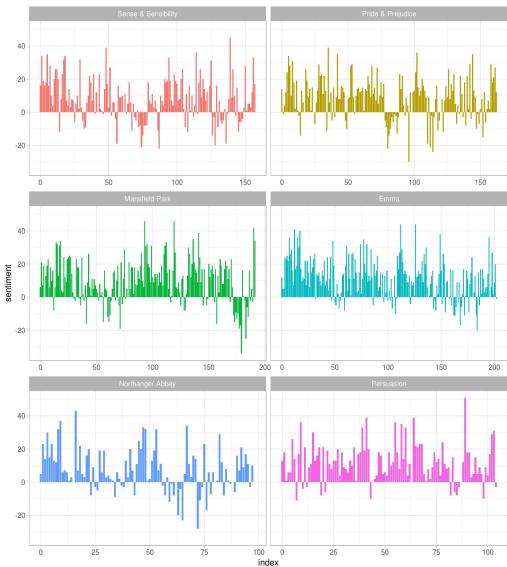
text

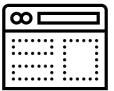
tidy dataset

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



results





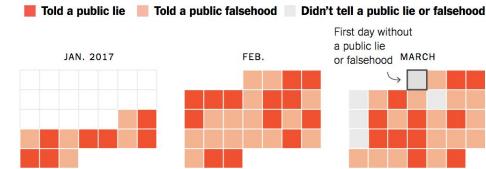
website

A black arrow pointing to the right, indicating the direction of the next page or section.

tidy dataset

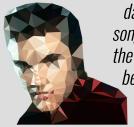
date	lie	explanation	url
0 Jan 21, 2017	I wasn't a fan of Iraq. I didn't want to go in...	He was for an invasion before he was against it.	https://wwwbuzzfeedcomandrewkaczynskiin2020
1 Jan 21, 2017	A reporter for Time magazine — and I have been...	Trump was on the cover 11 times and Nixon appeared...	http://nationtimecom2013/11/06/10-things-yo
2 Jan 23, 2017	Between 3 million and 5 million illegal votes ...	There's no evidence of illegal voting.	https://wwwnytimescom2017/01/23/us/politics...
3 Jan 25, 2017	Now, the audience was the biggest ever. But th...	Official aerial photos show Obama's 2009 inauguration...	https://wwwnytimescom2017/01/21/us/politics...
4 Jan 25, 2017	Take a look at the Pew reports (which show vot...	The report never mentioned voter fraud.	https://wwwnytimescom2017/01/24/us/politics...

results



text (lyrics)

The Pudding



"I'll be analyzing the repetitiveness of a dataset of 15,000 songs that charted on the Billboard Hot 100 between 1958 and 2017."

AN EXERCISE IN LANGUAGE COMPRESSION
Are Pop Lyrics Getting More Repetitive?

By Colin Morris

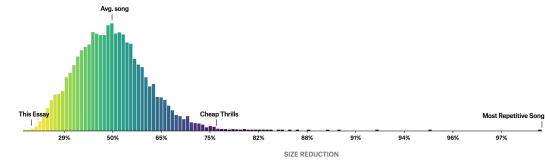


tidy dataset

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!