



What is data?

Jeffrey Leek, Assistant Professor of Biostatistics
Johns Hopkins Bloomberg School of Public Health

Definition of data

“ Data are values of qualitative or quantitative variables, belonging to a set of items.”

<http://en.wikipedia.org/wiki/Data>

Definition of data

“Data are values of qualitative or quantitative variables, belonging to a **set of items**. ”

<http://en.wikipedia.org/wiki/Data>

Set of items: Sometimes called the population; the set of objects you are interested in

Definition of data

“ Data are values of qualitative or quantitative **variables**, belonging to a set of items.”

<http://en.wikipedia.org/wiki/Data>

Variables: A measurement or characteristic of an item.

Definition of data

“ Data are values of **qualitative** or **quantitative** variables, belonging to a set of items.”

<http://en.wikipedia.org/wiki/Data>

Qualitative: Country of origin, sex, treatment

Quantitative: Height, weight, blood pressure

Raw versus processed data

Raw data

- The original source of the data
- Often hard to use for data analyses
- Data analysis *includes* processing
- Raw data may only need to be processed once

http://en.wikipedia.org/wiki/Raw_data

Processed data

- Data that is ready for analysis
- Processing can include merging, subsetting, transforming, etc.
- There may be standards for processing
- All steps should be recorded

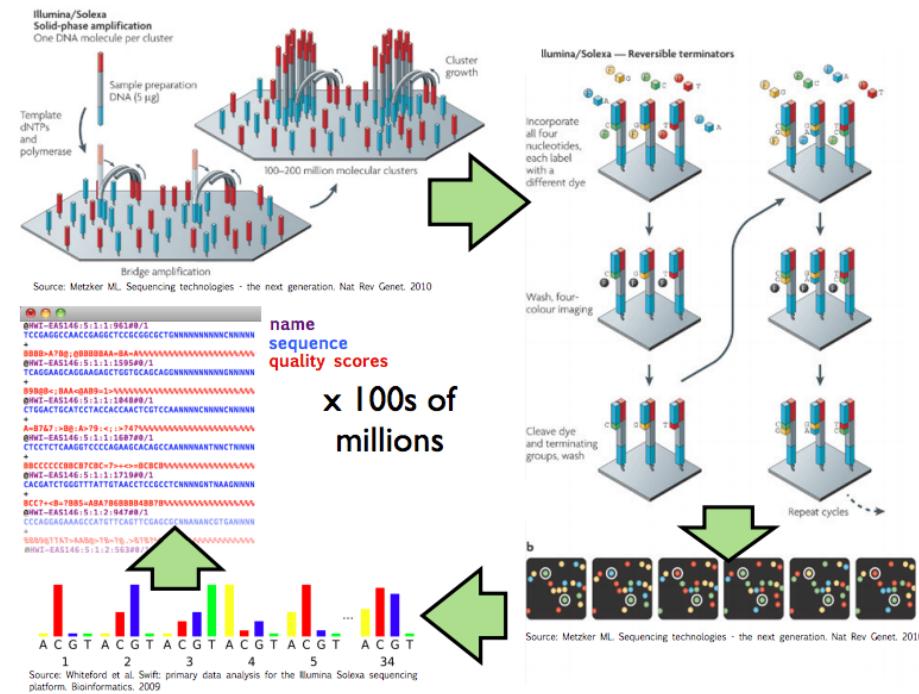
http://en.wikipedia.org/wiki/Computer_data_processing

An example of a processing pipeline



http://www.illumina.com.cn/support/sequencing/sequencing_instruments/hiseq_1000.asp

An example of a processing pipeline



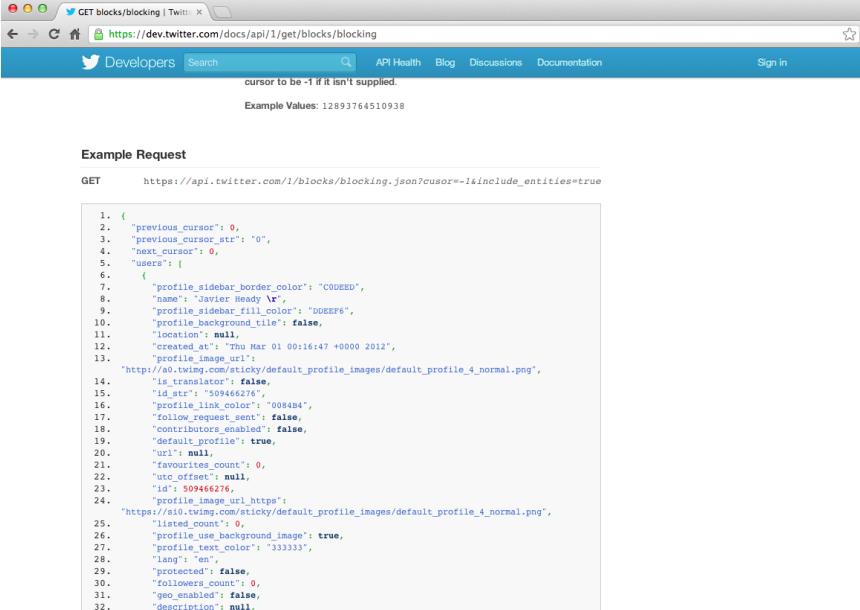
http://www.cbcn.umd.edu/~hcorrada/CMSC858B/lectures/lect22_seqIntro/seqIntro.pdf

What do raw data look like?

```
@HWI-EAS121:4:100:1783:550#0/1
CGTTACGAGATCGGAAGAGCGGGTTCAGCAGGAATGCCGAGACGGATCTGTATGCCGTCTGCTCCGTGACAAGACAGGG
+HWI-EAS121:4:100:1783:550#0/1
aaaaaa`b_aa`aa`YaX]aZ`aZM^Z]YRa]YSG[[ZREQLHESDHNDHNMEEDDM PENITKFLFEEDDDHEJQMEDDD
@HWI-EAS121:4:100:1783:1611#0/1
GGGTGGGCATTTCACACTCGCAGTATGGGTTGCCGCACGACAGGCAGCGGTCA GCGCTTGGCCTGGCCTTCGGAAA
+HWI-EAS121:4:100:1783:1611#0/1
a``^`_ ``^a``^a_`_ja_]`a____`_ ``^`]X_]XTV_\]NX_XVX]]_TTTG[VTHPN]VFDZ
@HWI-EAS121:4:100:1783:322#0/1
CGTTTATGTTTTGAATATGTCTTATCTAACGGTTATTTAGATGTTGGTCTTATTCTAACGGTCATATATTTCTA
+HWI-EAS121:4:100:1783:322#0/1
abaa``^aaaaabbbaaabbbbbbb`bbbb_bbbbbbbb`bbbaV^_a``a``]``aT]a__V\]]_`^a`_ja_abbaV_
@HWI-EAS121:4:100:1783:1394#0/1
GGGTCTTATTGGCTCTGGTGATCCCCCATATTCTCCGGTTGTGGTTAACCGATCATCGCGCATTACTCCGGCTGC
+HWI-EAS121:4:100:1783:1394#0/1
````[aa\b``^`]aabbb][`a_abbb`a``bbbbbabaabaaaab_Vza_``bab_X`[a\HV[_]_[^_X\T_VQQ
@HWI-EAS121:4:100:1783:207#0/1
CCCTGGGAGATCGGAAGAGCGGGTTCAGCAGGAATGCCGAGACCGATCTGTATGCCGTCTCTGTTGAAAAAAACAA
+HWI-EAS121:4:100:1783:207#0/1
abba`Xa``^`aa]ba_bba[a_O_a`aa`aa`a]^V]X_a^YS\R_\H[_]\ZTDUZZUSOPX]]POP\GS\WSHHD
@HWI-EAS121:4:100:1783:455#0/1
GGGTAATTCAAGGGACAATGTAATGGCTGCACAAAAAAATACATCTTCATGTTCCATTGCACCATTGACAAATACATATT
+HWI-EAS121:4:100:1783:455#0/1
abb_babbabaabbbbbbbbbbba``b`\abbabbabbabbbaabbbbb`bb`ab_O_bab_Q_bbabaa_a
```

[http://brianknaus.com/software/srtoolbox/s\\_4\\_1\\_sequence80.txt](http://brianknaus.com/software/srtoolbox/s_4_1_sequence80.txt)

# What do raw data look like?



The screenshot shows a web browser window with the URL [https://dev.twitter.com/docs/api/1/get\(blocks/blocking](https://dev.twitter.com/docs/api/1/get(blocks/blocking)). The page is titled "GET blocks/blocking | Twitter". The main content area displays an "Example Request" for the API endpoint. The request is a GET request to [https://api.twitter.com/1/blocks/blocking.json?cursor=-1&include\\_entities=true](https://api.twitter.com/1/blocks/blocking.json?cursor=-1&include_entities=true). The response body contains a large JSON object representing user blocking information. The JSON structure includes fields such as "previous\_cursor", "previous\_cursor\_str", "next\_cursor", "next\_cursor\_str", and a "users" array containing detailed user profiles.

```
1. {
2. "previous_cursor": 0,
3. "previous_cursor_str": "0",
4. "next_cursor": 0,
5. "next_cursor_str": "0",
6. "users": [
7. {
8. "profile_sidebar_border_color": "CODEED",
9. "name": "Twitter Meeny",
10. "profile_sidebar_fill_color": "DDEEF6",
11. "profile_background_tile": false,
12. "location": null,
13. "created_at": "Thu Mar 01 00:16:47 +0000 2012",
14. "profile_image_url": "http://si0.twimg.com/sticky/default_profile_images/default_profile_4_normal.png",
15. "id": 509466276,
16. "id_str": "509466276",
17. "profile_link_color": "#008484",
18. "follow_request_sent": false,
19. "contributors_enabled": false,
20. "default_profile": true,
21. "utc_offset": null,
22. "favourites_count": 0,
23. "id": 509466276,
24. "profile_image_url_https": "https://si0.twimg.com/sticky/default_profile_images/default_profile_4_normal.png",
25. "list_id": 0,
26. "profile_use_background_image": true,
27. "profile_text_color": "#333333",
28. "lang": "en",
29. "protected": false,
30. "followers_count": 0,
31. "geo_enabled": false,
32. "description": null,
33. }
34.]
35. }
```

[https://dev.twitter.com/docs/api/1/get\(blocks/blocking](https://dev.twitter.com/docs/api/1/get(blocks/blocking)

# What do raw data look like?

ALLERGIES		MEDICATION HISTORY	
Last Updated: 01 Dec 2011 @ 0851		Last Updated: 11 Apr 2011 @ 1737	
Allergy Name:	TRIMETHOPRIM	Medication:	AMLODIPINE BESYLATE 10MG TAB
Location:	DAYT29	Instructions:	TAKE ONE TABLET BY MOUTH TAKE ONE-HALF TABLET FOR GRAPEFRUIT JUICE--
Date Entered:	09 Mar 2011	Status:	Active
Action:		Refills Remaining:	3
Allergy Type:	DRUG	Last Filled On:	20 Aug 2010
A Drug Class:	ANTI-INFECTIVES, OTHER	Initially Ordered On:	13 Aug 2010
Observed/Historical:	HISTORICAL	Quantity:	45
Comments:	The reaction to this allergy was MILD (NO SQUELAE)	Days Supply:	90
Allergy Name:	TRAMADOL	Pharmacy:	DAYTON
Location:	DAYT29	Prescription Number:	2718953
Date Entered:	09 Mar 2011		
Action:	URINARY RETENTION		
Allergy Type:	DRUG		
A Drug Class:	NON-OPIOID ANALGESICS		
Observed/Historical:	HISTORICAL		
Comments:	gradually worsening difficulty emptying bladder		

<http://blue-button.github.com/challenge/>

# What do processed data look like?

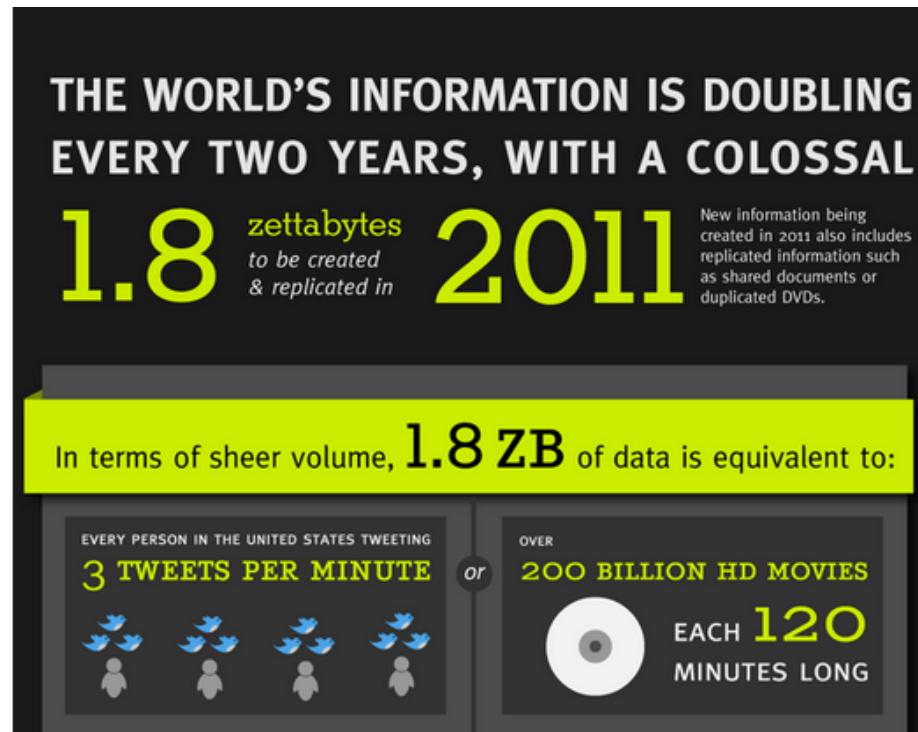
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	id	problem_id	subject_id	start	stop	time_left	answer									
2	1	498	17	1307119989	1307120016	2369	A									
3	2	150	15	1307119991	1307120009	2376	D									
4	3	313	16	1307119994	1307120009	2376	E									
5	4	142	13	1307119995	1307120019	2360	B									
6	5	273	14	1307119996	1307120008	2357	A									
7	6	101	19	1307119996	1307120021	2364	B									
8	7	105	18	1307119998	1307120048	2337	B									
9	8	162	12	1307120004	1307120042	2343	C									
10	9	70	15	1307120011	1307120038	2347	C									
11	10	300	16	1307120012	1307120052	2293	B									
12	11	494	17	1307120013	1307120075	2310	D									
13	12	357	13	1307120021	1307120118	2367	A									
14	13	522	19	1307120025	1307120152	2233	D									
15	14	232	14	1307120030	1307120158	2227	C									
16	15	344	15	1307120041	1307120117	2268	B									
17	16	160	17	1307120079	1307120249	2136	D									
18	17	516	16	1307120080	1307120159	2226	B									
19	18	472	12	1307120119	1307120170	2215	A									
20	19	43	15	1307120122	1307120140	2245	C									
21	20	353	13	1307120144	1307120199	2186	C									
22	21	218	15	1307120152	1307120272	2113	E									
23	22	69	16	1307120163	1307120188	2197	D									
24	23	562	16	1307120164	1307120181	2090	D									
25	24	121	19	1307120252	1307120294	2091	E									
26	25	297	15	1307120277	1307120342	2043	B									
27	26	495	13	1307120281	1307120353	2032	E									
28	27	94	14	1307120288	1307120343	2042	E									
29	28	22	18	1307120310	1307120365	2020	C									
30	29	64	19	1307120311	1307120368	2000	B									
31	30	502	16	1307120322	1307120336	2049	B									
32	31	44	16	1307120339	1307120352	2033	A									
33	32	315	14	1307120348	1307120362	2023	B									
34	33	385	15	1307120352	1307120553	1832	E									
35	34	550	13	1307120356	1307120444	1941	B									
36	35	92	14	1307120371	1307120397	1984	B									
37	36	995	16	1307120377	1307120426	1959	D									
38	37	267	17	1307120382	1307120515	1870	E									
39	38	257	14	1307120401	1307120427	1958	C									
40	39	312	19	1307120407	1307120548	1837	D									
41	40	321	18	1307120431	1307120449	1936	A									
42	41	220	16	1307120437	1307120510	1875	A									

1. Each variable forms a column
2. Each observation forms a row
3. Each table/file stores data about one kind of observation (e.g. people/hospitals).

<http://vita.had.co.nz/papers/tidy-data.pdf>

Leek, Taub, and Pineda 2011 PLoS One

# How much is there?

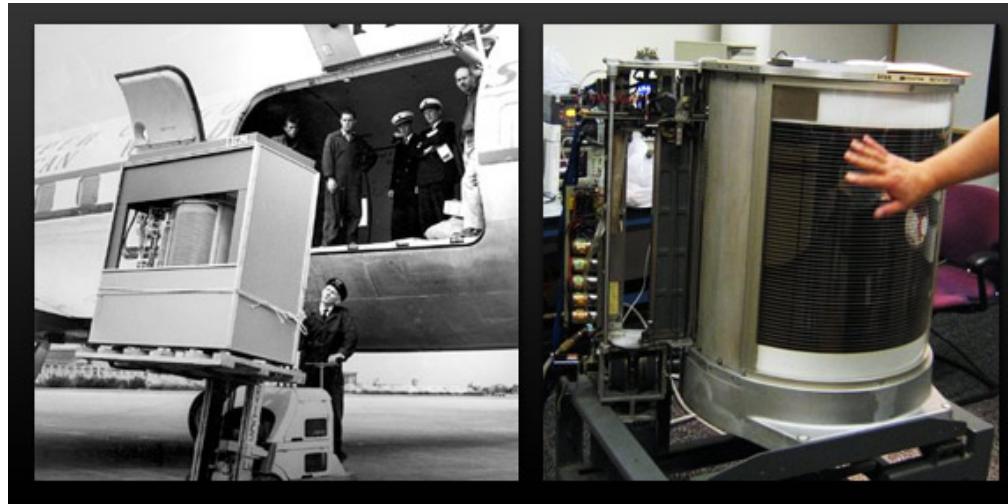


<http://mashable.com/2011/06/28/data-infographic/>

# So what about big data?



# Depends on your perspective



# Why big data now?

## An Experimental Study of the Small World Problem\*

JEFFREY TRAVERS

Harvard University

AND

STANLEY MILGRAM

The City University of New York

*Arbitrarily selected individuals (N=296) in Nebraska and Boston are asked to generate acquaintance chains to a target person in Massachusetts, employing "the small world method" (Milgram, 1967). Sixty-four chains reach the target person. Within this group, the mean number of intermediaries between starters and targets is 5.2. Boston starting chains reach the target*

Travers and Milgram (1969) Sociometry

# Why big data now?

arXiv.org > physics > arXiv:0803.0939

Search or A

Physics > Physics and Society

## Planetary-Scale Views on an Instant-Messaging Network

Jure Leskovec, Eric Horvitz

(Submitted on 6 Mar 2008)

We present a study of anonymized data capturing a month of high-level communication activities within the whole of the Microsoft Messenger instant-messaging system. We examine characteristics and patterns that emerge from the collective dynamics of large numbers of people, rather than the actions and characteristics of individuals. The dataset contains summary properties of 30 billion conversations among 240 million people. From the data, we construct a communication graph with 180 million nodes and 1.3 billion undirected edges, creating the largest social network constructed and analyzed to date. We report on multiple aspects of the dataset and synthesized graph. We find that the graph is well-connected and robust to node removal. We investigate on a planetary-scale the oft-cited claim that people are separated by ``six degrees of separation'' and find that the average path length among Messenger users is 6.6. We also find that people tend to communicate more with each other when they have similar age, language, and location, and that cross-gender conversations are both more frequent and of longer duration than conversations with the same gender.

Leskovec and Horvitz WWW '08

# Big or small - you need the right data

“The data may not contain the answer. The combination of some data and an aching desire for an answer does not ensure that a reasonable answer can be extracted from a given body of data...”

Tukey

# Big or small - you need the right data

“ ...no matter how big the data are. ”

Leek