



Obtaining data, raw and tidy data

What does data really look like?

```
@HWI-EAS121:4:100:1783:550#0/1
CGTTACGAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGAGACGGATCTCGTATGCCGTCTGCTGCGTGACAAGACAGGGG
+HWI-EAS121:4:100:1783:550#0/1
aaaaa`b_aa`aa`YaX]aZ`aZM^Z]YRa]YSG[ [ZREQLHESDHNDHNMEEDDMPENITKFLFEEDDDHEJQMEDDD
@HWI-EAS121:4:100:1783:1611#0/1
GGGTGGGCATTTCCACTCGCAGTATGGGTGCGGCACGACAGGCAGCGGTCAGCCTGCGCTTTGGCCTGGCCTTCGGAAA
+HWI-EAS121:4:100:1783:1611#0/1
a``^\\_`_`_`^a``a`^a^_`_`]a_]`a`_____`_`^`^]X]_]XTV\\_]NX_XVX]]_TTTTG[VTHPN]VFDZ
@HWI-EAS121:4:100:1783:322#0/1
CGTTTATGTTTTTGAATATGTCTTATCTTAACGGTTATATTTTAGATGTTGGTCTTATTCTAACGGTCATATATTTTCTA
+HWI-EAS121:4:100:1783:322#0/1
abaa`^aaaaabbbbaababbbbbb`bbbb_bbbbbbbb`bbbaV^_a``a``]``aT]a__V\\]]_]`a`]a_abbaV__
@HWI-EAS121:4:100:1783:1394#0/1
GGGTCTTTATTGGTCTGGTGATCCCCATATTCTCCGGTTGTGTGGTTTAACCGATCATCGCGCATTA CTCCCGGCTGC
+HWI-EAS121:4:100:1783:1394#0/1
```[aa\\b^^[ ]aabbb][`a_abbb`a``bbbbbbabaabaaaab_VZa_^__bab_X`[a\\HV_[_]_[^_X\\T_VQQ
@HWI-EAS121:4:100:1783:207#0/1
CCCTGGGAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGAGACCGATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAACA
+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
GGGTAATTCAGGGACAATGTAATGGCTGCACAAAAAATACATCTTTCATGTTCCATTGCACCATTGACAAATACATATT
+HWI-EAS121:4:100:1783:455#0/1
abb_babbabaabbbbbbbbbbbbbbbba\\`b`\\abbbabbbbabbbbbbaabbbbb`bb`ab_O_bab_Q_bbabaa_a
```

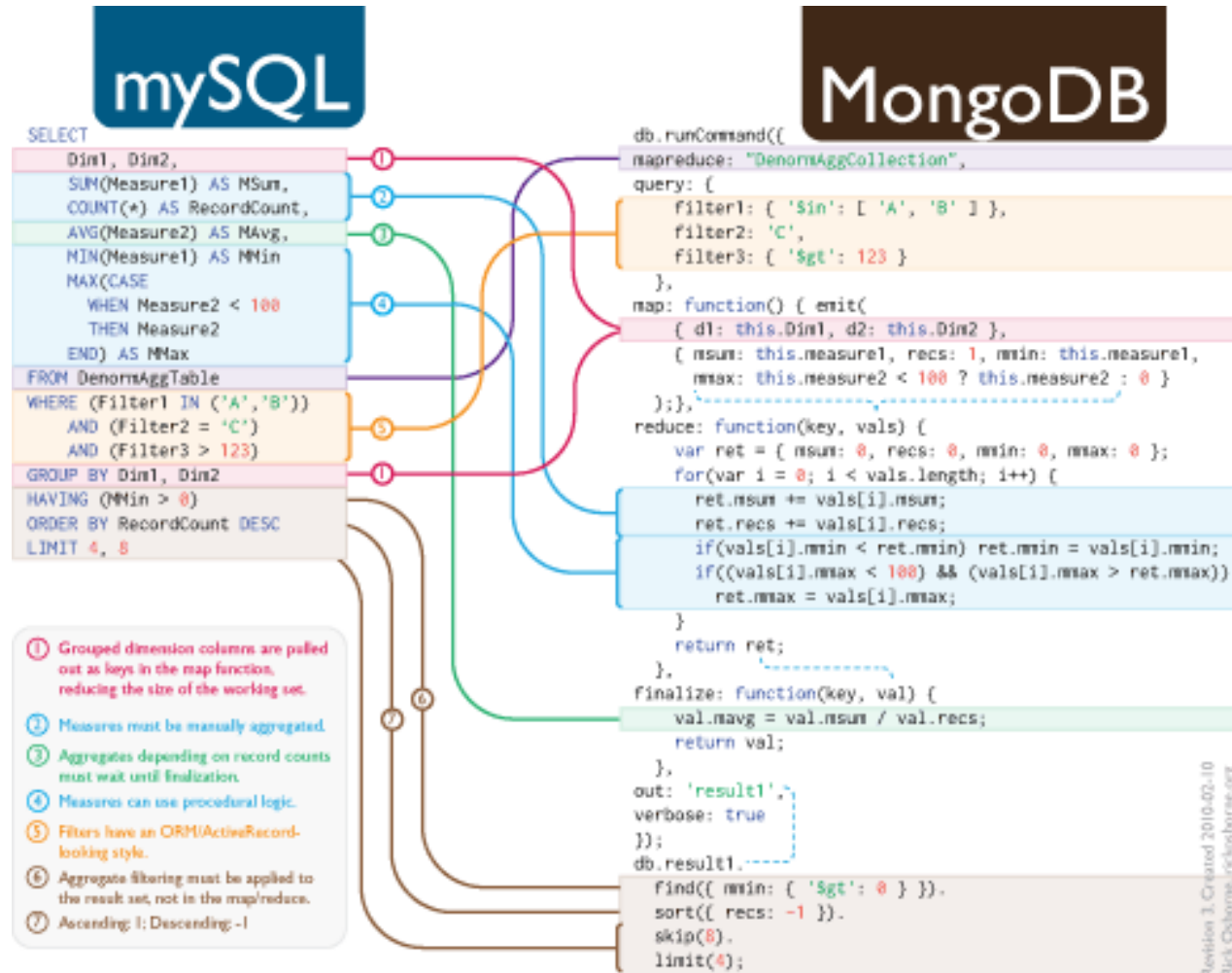
# What does data really look like?

```
-<breakfast_menu>
 -<food>
 <name>Belgian Waffles</name>
 <price>$5.95</price>
 -<description>
 Two of our famous Belgian Waffles with plenty of real maple syrup
 </description>
 <calories>650</calories>
 </food>
 -<food>
 <name>Strawberry Belgian Waffles</name>
 <price>$7.95</price>
 -<description>
 Light Belgian waffles covered with strawberries and whipped cream
 </description>
 <calories>900</calories>
 </food>
 -<food>
 <name>Berry-Berry Belgian Waffles</name>
 <price>$8.95</price>
 -<description>
 Light Belgian waffles covered with an assortment of fresh berries and whipped cream
 </description>
 <calories>900</calories>
 </food>
-</food>
```

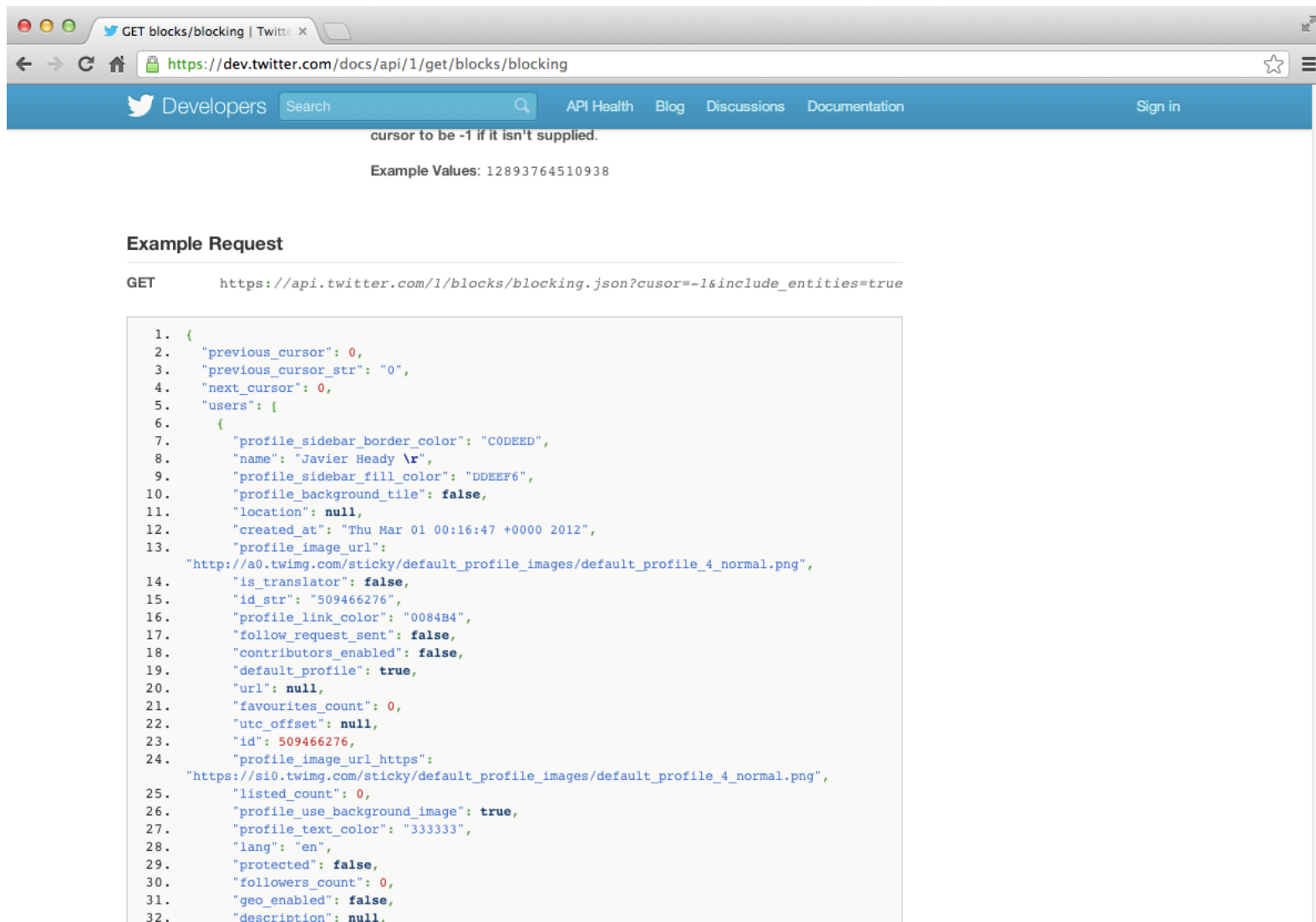
# What does data really look like?

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 1 HOUR BEFORE GRAPEFRUIT JUICE--
Date Entered:	09 Mar 2011	Status:	Active
Reaction:		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
		Pharmacy:	DAYTON
		Prescription Number:	2718953
Allergy Name:	TRAMADOL	Medication:	IBUPROFEN 600MG TAB
Location:	DAYT29	Instructions:	TAKE ONE TABLET BY MOUTH FOUR TIMES A DAY WITH FOOD
Date Entered:	09 Mar 2011	Status:	Active
Reaction:	URINARY RETENTION	Refills Remaining:	3
Allergy Type:	DRUG	Last Filled On:	20 Aug 2010
A Drug Class:	NON-OPIOID ANALGESICS	Initially Ordered On:	01 Jul 2010
Observed/Historical:	HISTORICAL	Quantity:	300
Comments:	gradually worsening difficulty emptying bladder		

# Where is data?



# Where is data?



The screenshot shows a web browser window with the URL `https://dev.twitter.com/docs/api/1/get/blocks/blocking`. The page is titled "GET blocks/blocking | Twitter" and includes a search bar and navigation links like "API Health", "Blog", "Discussions", and "Documentation". The main content area displays the endpoint details, including a note about the cursor and example values.

cursor to be -1 if it isn't supplied.

Example Values: 12893764510938

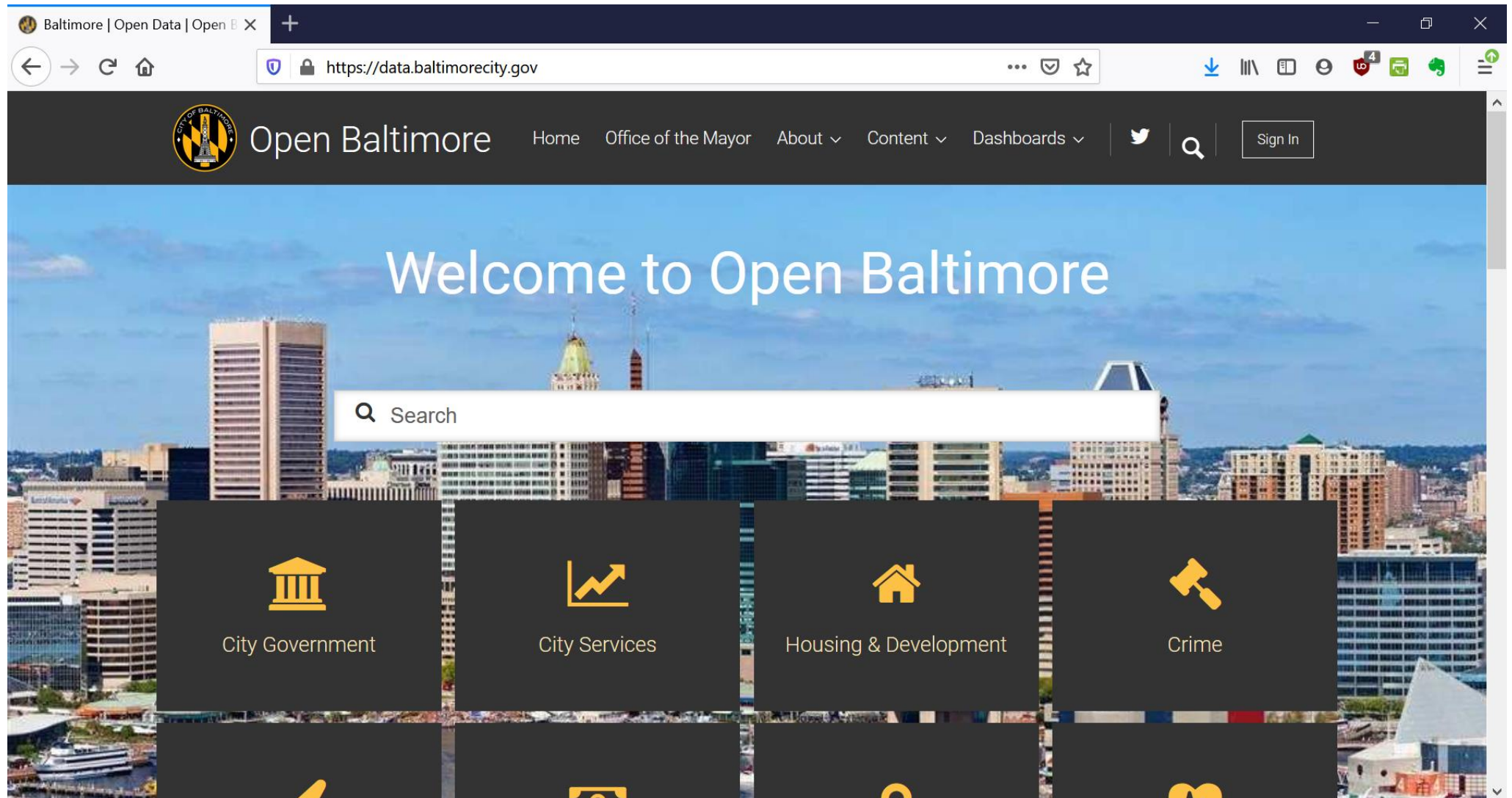
### Example Request

GET `https://api.twitter.com/1/blocks/blocking.json?cursor=-1&include_entities=true`

```
1. {
2. "previous_cursor": 0,
3. "previous_cursor_str": "0",
4. "next_cursor": 0,
5. "users": [
6. {
7. "profile_sidebar_border_color": "CODEED",
8. "name": "Javier Heady \r",
9. "profile_sidebar_fill_color": "DDEEF6",
10. "profile_background_tile": false,
11. "location": null,
12. "created_at": "Thu Mar 01 00:16:47 +0000 2012",
13. "profile_image_url":
14. "http://a0.twimg.com/sticky/default_profile_images/default_profile_4_normal.png",
15. "is_translator": false,
16. "id_str": "509466276",
17. "profile_link_color": "0084B4",
18. "follow_request_sent": false,
19. "contributors_enabled": false,
20. "default_profile": true,
21. "url": null,
22. "favourites_count": 0,
23. "utc_offset": null,
24. "id": 509466276,
25. "profile_image_url_https":
26. "https://s10.twimg.com/sticky/default_profile_images/default_profile_4_normal.png",
27. "listed_count": 0,
28. "profile_use_background_image": true,
29. "profile_text_color": "333333",
30. "lang": "en",
31. "protected": false,
32. "followers_count": 0,
33. "geo_enabled": false,
34. "description": null,
```



# Where is data?





# Data science pipeline

**Raw data -> Processing script -> tidy data -> data analysis -> data communication**

# Definition of Data

Data are a set of values of qualitative or quantitative variables about one or more persons or objects

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

**One or more persons or objects:** Sometimes called the population; the set of objects you are interested in.

**Variables:** A measurement or characteristic of an item.

- **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](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](http://en.wikipedia.org/wiki/Computer_data_processing)

# Components of tidy data

1. The raw data.
2. A tidy data set
3. A code book describing each variable and its values in the tidy data set.
4. An explicit and exact recipe you used to go from 1 -> 2,3.

# The raw data

- The strange binary file your measurement machine spits out
- The unformatted Excel file with 10 worksheets the company you contracted with sent you
- The complicated JSON data you got from scraping the Twitter API
- The hand-entered numbers you collected looking through a microscope

# The tidy data

1. Each variable you measure should be in one column
2. Each different observation of that variable should be in a different row
3. There should be one table for each "kind" of variable
4. If you have multiple tables, they should include a column in the table that allows them to be linked

## *Some other important tips*

- Include a row at the top of each file with variable names.
- Make variable names human readable AgeAtDiagnosis instead of AgeDx
- In general data should be saved in one file per table.



# The code book

1. Information about the variables (including units!) in the data set not contained in the tidy data
2. Information about the summary choices you made
3. Information about the experimental study design you used

## *Some other important tips*

- A common format for this document is a Word/text file.
- There should be a section called "Study design" that has a thorough description of how you collected the data.
- There must be a section called "Code book" that describes each variable and its units.

# The instruction list

- Ideally a computer script (in R :-), but I suppose Python is ok too...)
- The input for the script is the raw data
- The output is the processed, tidy data
- There are no parameters to the script

In some cases it will not be possible to script every step. In that case you should provide instructions like:

- 1.Step 1 - take the raw file, run version 3.1.2 of summarize software with parameters a=1, b=2, c=3
- 2.Step 2 - run the software separately for each sample
- 3.Step 3 - take column three of outputfile.out for each sample and that is the corresponding row in the output data set

This presentation is based on materials from Coursera Getting and Cleaning Data course

<https://www.coursera.org/learn/data-cleaning/home/welcome>