

Motivation and pre-requisites

Jeffrey Leek

May 17, 2016

About this course

- ▶ This course covers the basic ideas behind getting data ready for analysis
- ▶ Finding and extracting raw data
- ▶ Tidy data principles and how to make data tidy
- ▶ Practical implementation through a range of R packages
- ▶ What this course depends on
- ▶ The Data Scientist's Toolbox
- ▶ R Programming
- ▶ What would be useful
- ▶ Exploratory analysis
- ▶ Reporting Data and Reproducible Research

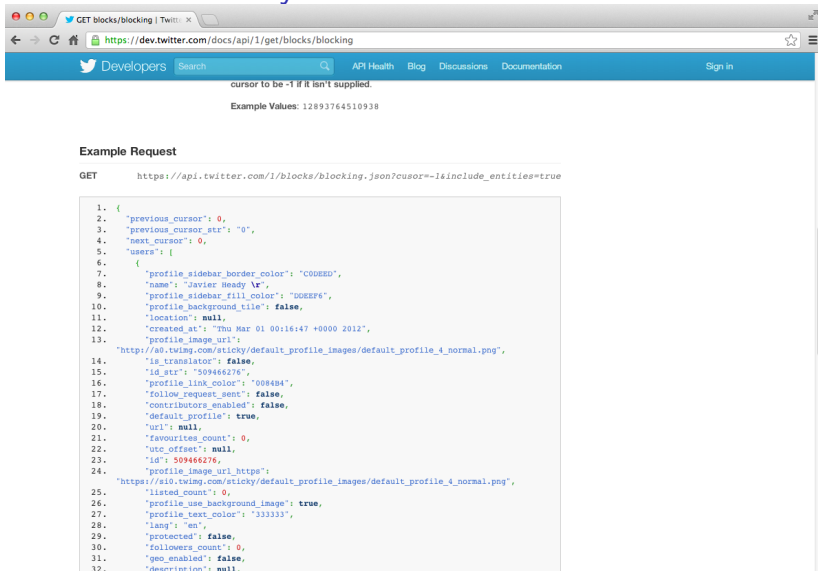
What you wish data looked like

What does data really look like?

```
@HWI-EAS121:4:100:1783:550#0/1
CGTTACGAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACGGATCTCGTATGCGGTCTGCTGCGTGACAAGACAGGGG
+HWI-EAS121:4:100:1783:550#0/1
aaaaa`b_aa`aa`YaX]az`aZM`Z]YRa]YSG[[ZREQLHESDHNDHDHNMEEDDMPENITKFLFEEDDDHEJQMEDDD
@HWI-EAS121:4:100:1783:1611#0/1
GGGTGGGCATTCCCATCGCAGTAGTGGGTTGCCGCACGACAGGCAGCGGTTCAGCCTGCGCTTTGGCGCTGGCCTTCGGAAA
+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`^aaaaabbbaababbbbbbb`bbbb_bbbbbbbb`bbbaV`a``a``]``aT[a__V\]]_]`a`]a_abbaV__
@HWI-EAS121:4:100:1783:1394#0/1
GGGTCTTTATTGGTCTGGTGATCCCCATATTCTCCGGTTGTGTGGTTTAACCGATCATCGCGCATTA TCTCCCGGCTGC
+HWI-EAS121:4:100:1783:1394#0/1
```[aa\b`^^[jaabbb][`a_abbb`a``bbbbbababaabaaaab_Vza_^____bab_X`[a\HV[_][_[^_X\T_VQQ
@HWI-EAS121:4:100:1783:207#0/1
CCCTGGGAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGATCTCGTATGCCGTCTTCTGCTTGAAAAAAAAAAAAACA
+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
GGGTAATTCAGGGACAATGTAATGGCTGCACAAAAAAATACATCTTTCATGTTCCATTGCACCATTGACAAATACATATT
+HWI-EAS121:4:100:1783:455#0/1
abb_babbababbbbbbbbbbbbbbbba\b`b\abbbabbbbabbbbbbaabbbbb`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 does data really look like?



The screenshot shows the Twitter API documentation page for the `GET blocks/blocking` endpoint. The page header includes the Twitter logo, a search bar, and navigation links for API Health, Blog, Discussions, and Documentation. The main content area displays the endpoint URL and an example JSON response. The JSON response is a list of one user object, containing details like name, location, created\_at, profile image URL, and various counts.

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": "C0DEED",
8. "name": "Javier Heady \r",
9. "profile_sidebar_fill_color": "D0E0F6",
10. "profile_background_tile": false,
11. "location": null,
12. "created_at": "Thu Mar 01 00:16:47 +0000 2012",
13. "profile_image_url":
14. "https://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,
```

https:

`//dev.twitter.com/docs/api/1/get/blocks/blocking`

# What does data really look like?

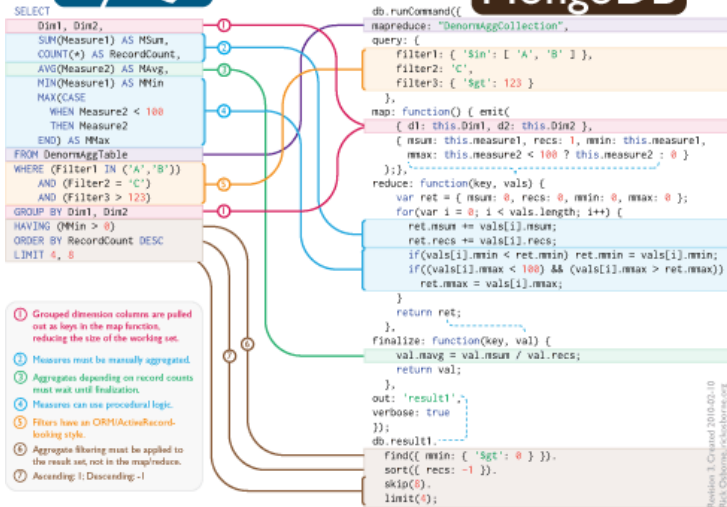
ALLERGIES		MEDICATION HISTORY	
Last Updated: 01 Dec 2011 @ 0851		Last Updated: 11 Apr 2011 @ 1737	
Allergy Name:	TRIMETHOPRIM	Medication:	AMLODIPINE BESYLATE 10MG TAB
Location:	DAYT29	Instructions:	TAKE ONE TABLET BY MOUTH TAKE ONE-HALF TABLET FOR GRAPEFRUIT JUICE--
Date Entered:	09 Mar 2011	Status:	Active
Reaction:		Refills Remaining:	3
Allergy Type:	DRUG	Last Filled On:	20 Aug 2010
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	Medication:	IBUPROFEN 600MG TAB
Reaction:	URINARY RETENTION	Instructions:	TAKE ONE TABLET BY MOUTH FOUR TIMES A DAY WITH FOOD
Allergy Type:	DRUG	Status:	Active
Drug Class:	NON-OPIOD ANALGESICS	Refills Remaining:	3
Observed/Historical:	HISTORICAL	Last Filled On:	20 Aug 2010
Comments:	gradually worsening difficulty emptying bladder	Initially Ordered On:	01 Jul 2010
		Quantity:	300

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

# Where is data?

mySQL

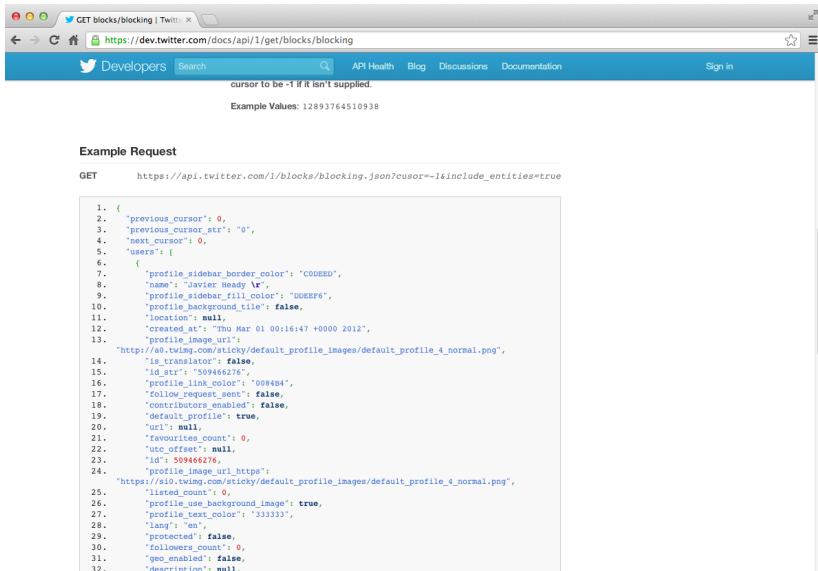
MongoDB



<http://rickosborne.org/blog/2010/02/>

[infographic-migrating-from-sql-to-mapreduce-with-mongodb/](http://infographic-migrating-from-sql-to-mapreduce-with-mongodb/)

# Where is data?



GET blocks/blocking | Twitter X

https://dev.twitter.com/docs/api/1/get/blocks/blocking

Developers Search API Health Blog Discussions Documentation Sign in

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": "C0DEED",
8. "name": "Javier Heady \r",
9. "profile_sidebar_fill_color": "D0E0F6",
10. "profile_background_tile": false,
11. "location": null,
12. "created_at": "Thu Mar 01 00:16:47 +0000 2012",
13. "profile_image_url":
14. "https://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,
```

https:

//dev.twitter.com/docs/api/1/get/blocks/blocking



# Where is data?



<https://data.baltimorecity.gov/>

# The goal of this course

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