

# CSV, Rinse, Repeat

## Javascript Data Exploration

Mathieu Jacomy  
Sciences Po Paris médialab  
Equipex DIME-SHS





Paris  
Sciences Po  
médiablab

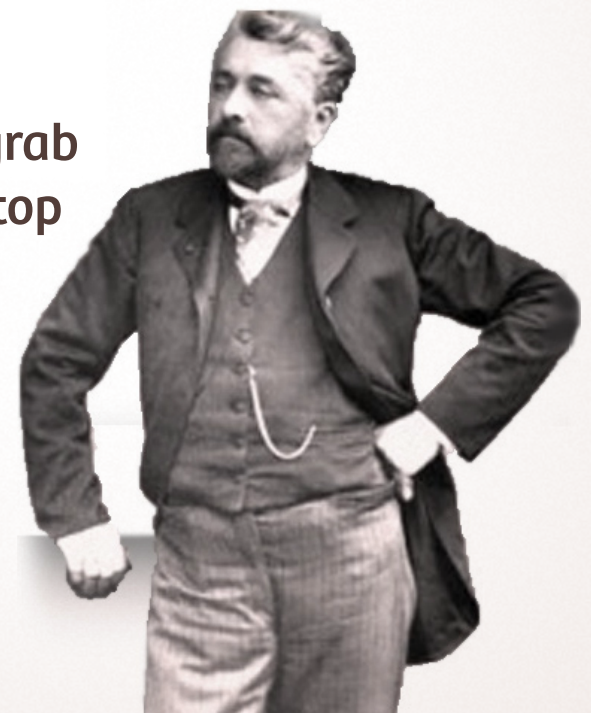
A hybrid laboratory  
for **social sciences**:

- Researchers
- Engineers
- Designers

I'm a sort of  
« social data scientist »



I just  
received  
a CSV



Let me grab  
my laptop



# Exploring data

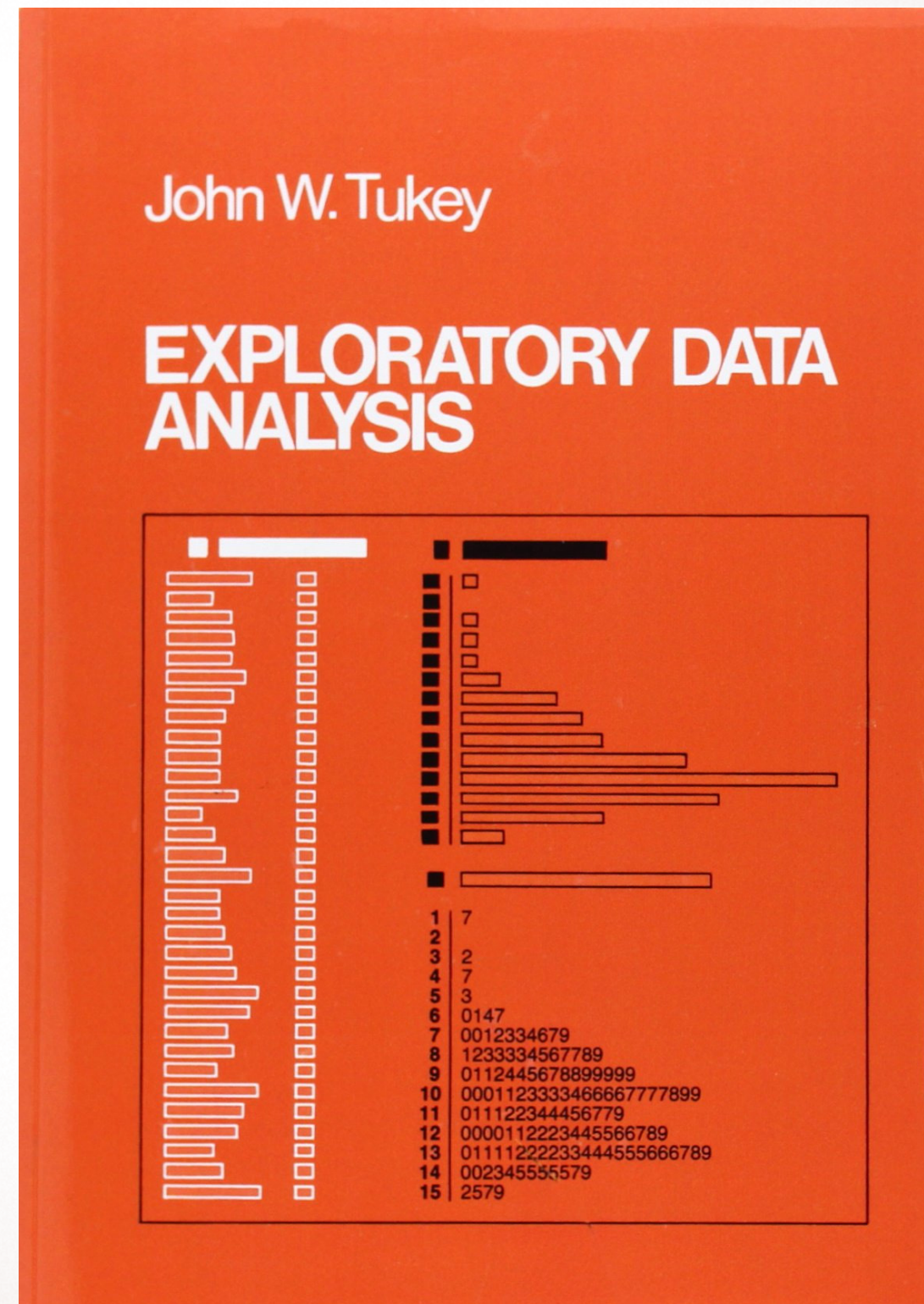
...is **not** about statistical metrics

The greatest value of a picture is  
when it forces us to notice what we  
**never expected to see.**


— John W. Tukey

Far better an approximate answer  
to **the right question**, which is often  
vague, than an exact answer to the  
wrong question, which can always  
be made precise.

—John W. Tukey







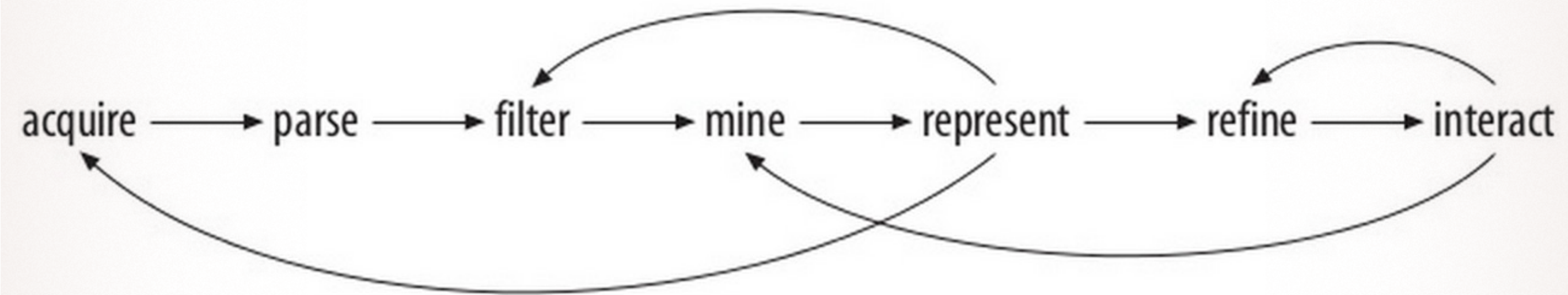
What is the  
right question?



# Exploring data

The chain of data mining  
by Ben Fry

...is about **iterating** facets of the data





# CSV Problem #1:

## Painful coding

In a spreadsheet environment,  
eg. Libre Office, Excel, Open Refine...  
coding features are designed for non-coders.

Consequence for non-coders:  
You invest your time in nonstandard, broken languages  
(you lose your time + it is still complicated)

Consequence for coders:  
Editing and filtering the data is painful  
(ranging from inefficient to WTF)



# Painful coding, simple example

Goal: retrieve years in movie titles

Rank	Rating	Movie
1	9.3/10	The Shawshank Redemption (1994)
2	9.2/10	The Godfather (1972)
3	9.0/10	The Godfather: Part II (1974)
4	9.0/10	The Dark Knight (2008)
5	8.9/10	12 Angry Men (1957)
6	8.9/10	Pulp Fiction (1994)
7	8.9/10	Schindler's List (1993)
8	8.9/10	The Lord of the Rings: The Return of the King (2003)
9	8.9/10	The Good, the Bad and the Ugly (1966)
10	8.9/10	Fight Club (1999)
11	8.8/10	The Lord of the Rings: The Fellowship of the Ring (2001)
12	8.8/10	Inception (2010)
13	8.8/10	Star Wars: Episode V - The Empire Strikes Back (1980)

Issue:  
the year is coded  
within the title



# Simple parsing with Libre Office

STXT					
=CHERCHE("\([0-9]{4}\)";C2;0)					
	A	B	C	E	F
1	Rank	Rating	Movie		
2	1	9.3/10	The Shawshank Redemption (1994)		
3	2	9.2/10	The Godfather (1972)	18	1972
4	3	9.0/10	The Godfather: Part II (1974)	27	1974

1. Find the year  
position

Look,  
it's coded  
in French!

2. Retrieve the  
string

STXT					
=STXT(C2;D2+1;4)					
	A	B	C	D	F
1	Rank	Rating	Movie		
2	1	9.3/10	The Shawshank Redemption (1994)	29	
3	2	9.2/10	The Godfather (1972)	18	1972
4	3	9.0/10	The Godfather: Part II (1974)	27	1974

In this situation the GUI is a problem, not a solution



# Simple parsing with Javascript

```
// Extract year data  
item.Year = item.Movie.match(/.*\(([0-9]{4})\)/)[1];
```

A real coding language is **more efficient**  
...if you can bring your CSV in the right **coding environment**



# CSV Problem #2:

## The filter/vis gap

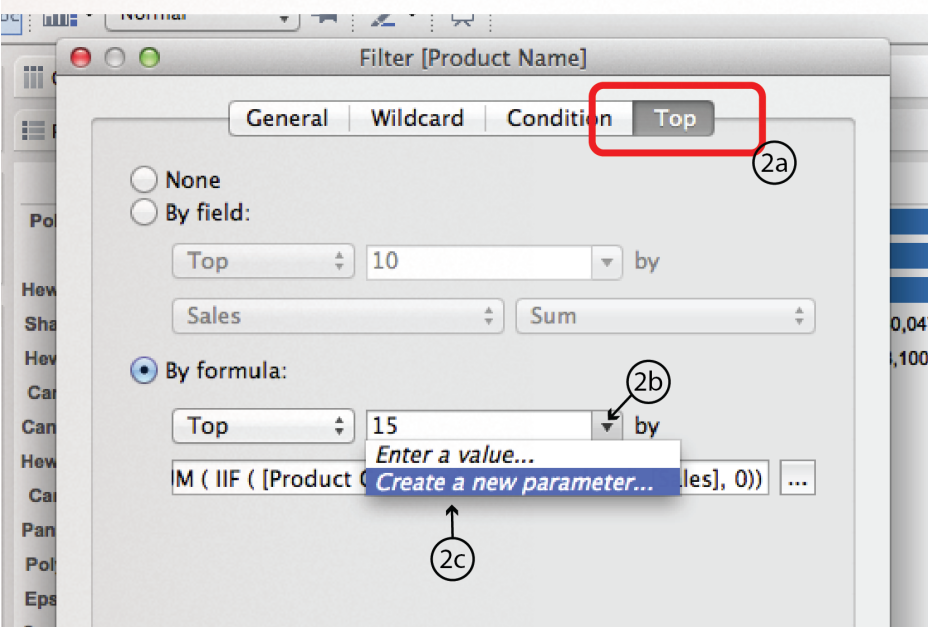
Code  
(Filter + Edit)

GAP

Visualize



# Simple filtering in Tableau Public



Formula + settings,  
in a form,  
inside a tab,  
of a modal  
that you open by a  
drag-and-drop  
(true story)

Edit the formula  
and/or the settings,  
apply and close

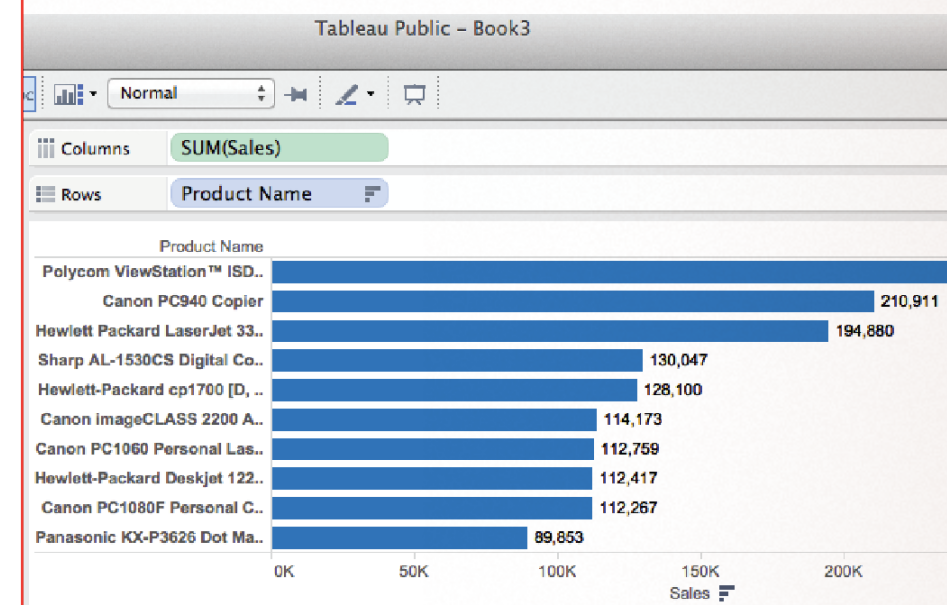


The modal hides  
the visualization



Reopen the modal,  
select the tab,  
select the field

Iterating  
is painful



## Visualization

Note:  
in Libre Office  
or Excel, it's even  
worse



Real  
Javascript  
coding

```
1 output = input.map(function(item, i){  
2  
3 // Extract year data  
4 item.Year = item.Movie.match(/.*\(([0-9]{4})\)/)[1];  
5  
6 // Clean title data (remove the year)  
7 item.Movie = item.Movie.replace(' ('+item.Year+')', '');  
8  
9 // Clean rating  
10 item.Rating = item.Rating.replace('/10', '');  
11  
12 return item;  
13  
14 }).filter(function(item, i){  
15  
16 // Only the movies with a number in their title  
17 return item.Movie.search(/[0-9]/gi) >= 0;  
18  
19 });  
20  
21 // Hit CTRL + ENTER to run the code
```

CODE YOUR  
FILTER

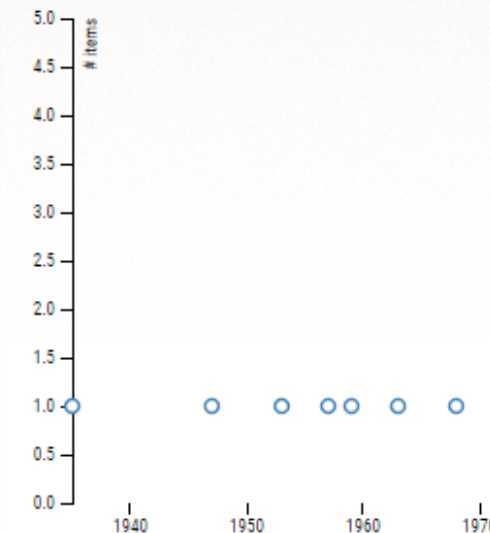
## CSV, Rinse Repeat

is about  
shortening the gap  
between code  
and visualization

to foster  
iterative  
exploration

Simple,  
ready-made  
visualizations

YEAR - YEARLY VOLUME



VISUALIZE  
RESULT

MOVIE - ITEMS TOP 50

1. Miracle on 34th Street (1)	the Deathly Hallows: Part 2 (1)	24. Se7en (1)	37. United 93 (1)
2. 12 Angry Men (1)	13. 3-Iron (1)	25. Big Hero 6 (1)	38. 28 Days Later... (1)
3. Terminator 2: Judgment Day (1)	14. Stalag 17 (1)	26. The 39 Steps (1)	39. Ghost in the Shell 2: Innocence (1)
4. 3 Idiots (1)	15. The Raid 2 (1)	27. (500) Days of Summer (1)	40. Apollo 13 (1)
5. Toy Story 3 (1)	16. Short Term 12 (1)	28. 3:10 to Yuma (1)	41. 13 Assassins (1)
6. 2001: A Space Odyssey (1)	17. Kill Bill: Vol. 2 (1)	29. 300 (1)	42. Ip Man 2 (1)
7. The 400 Blows (1)	18. Special 26 (1)	30. 21 Grams (1)	43. Mesrine Part 1: Killer Instinct (1)
8. The Legend of 1900 (1)	19. District 9 (1)	31. 50/50 (1)	44. 42(2013) (1)
9. Kill Bill: Vol. 1 (1)	20. 4 Months, 3 Weeks and 2 Days (1)	32. Harry Potter and the Deathly Hallows: Part 1 (1)	45. Despicable Me 2 (1)
10. 8½ (1)	21. Toy Story 2 (1)	33. 25th Hour (1)	46. X2 (1)
11. 12 Years a Slave (1)	22. How to Train Your Dragon 2: The Hidden World (1)	34. Cell 211 (1)	47. Buffalo '66 (1)

+ ADD VIZ



# CSV, Rinse, Repeat

...is a proposition to solve these problems **during exploration**

## A simple accessible tool

- Single web page

## A Javascript coding environment

- A standard coding panel
- CSV Import + Export
- Basic preview

## A layout designed to get rid of the filter/vis gap

- Input + code on the left
- Output + visualization on the right





Demo time!

<http://tools.medialab.sciences-po.fr/csu-rinse-repeat/>



# CSV, Rinse, Repeat

Sample datasets:

[Movies](#)

[NASA lab facilities](#)

[OECD BEPS consultation actors](#)

[+ more tools](#)

**SciencesPo**  
MÉDIALAB



LOAD CSV

CLICK or DRAG A FILE



INPUT 152610 r	created_at	from_user_name	text	filter
PREVIEW 3	937.0	2016-04-24T08:18:57	Jeremiahjones88	LIVE on #Periscope: Sonnet Sunday #shakespeare's bday <a href="https://www.periscope.tv/w/aesncDFIV0t5Tn pMTVZqQWd8MUJkR1I EVIZYem9LWG0EqIIRE S2vMPIQg5HxWYX6we _6DVazsiFQ70v06lqY">https://www.periscope.tv/w/aesncDFIV0t5Tn pMTVZqQWd8MUJkR1I EVIZYem9LWG0EqIIRE S2vMPIQg5HxWYX6we _6DVazsiFQ70v06lqY</a>
	382.0	2016-04-12T16:26:22	Lynchypoos	RT @britishlibrary: Here's a sneak peek at our exhibition, opening this Friday. What's your favourite #Shakespeare

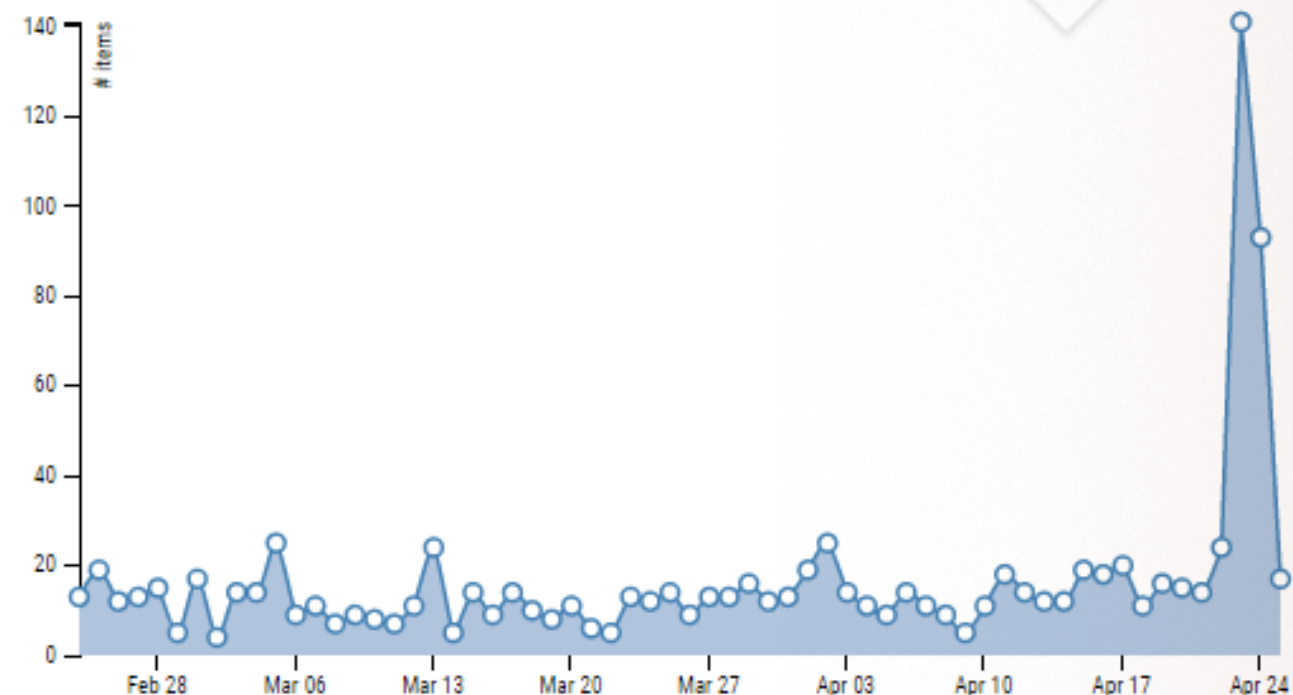
```

1 output = input.filter(function(item, i){
2
3     if (
4         // Check that latitude and longitude are valid
5         item.lat && item.lng
6         && !isNaN(+item.lat) && item.lat >= -90 && item.lat <
7         && !isNaN(+item.lng) && item.lng >= -180 && item.lng
8     ) {
9         // If they are, concatenate them for the vis module
10        item.coordinates = item.lat + '::' + item.lng
11        return true
12    } else {
13        // If not, filter them
14        return false
15    }
16
17 });
18
19

```

OUTPUT 1006 r	created_at	from_user_name	text
PREVIEW 3	321458.0	2016-04-22T10:37:38	jafrater
	595109.0	2016-04-25T14:38:29	MiguelWrites
			#Shakespeare's grave at Holy Trinity in #Stratford #shakespeare400 @ Holy Trinity Church,... <a href="https://www.instagram.com/p/BEf9uSrlUz/">https://www.instagram.com/p/BEf9uSrlUz/</a>
			We have indeed been celebrating

CREATED AT - DAILY VOLUME



+ ADD VIZ



INPUT 152610 r

PREVIEW 3

Input preview

	created_at	from_user_name	text	filter
937.0	2016-04-24T08:18:57	Jeremiahjones88	LIVE on #Periscope: Sonnet Sunday #shakespeare's bday <a href="https://www.periscope.tv/w/aesncDFIV0t5Tn pMTVZqQWd8MUJkR1I EVIZYem9LWG0EqIIRE S2vMPIQg5HxWYX6we _6DVazsiFQ70v06lqY">https://www.periscope.tv/w/aesncDFIV0t5Tn pMTVZqQWd8MUJkR1I EVIZYem9LWG0EqIIRE S2vMPIQg5HxWYX6we _6DVazsiFQ70v06lqY</a>	
382.0	2016-04-12T11:00:00		RT @britishlibrary: Here's a sneak peek at our exhibition, opening this Friday. What's your favourite #Shakespeare	

CODE YOUR FILTER

```

1 output = input.filter(function(item, i){
2
3   if (
4     // Check that latitude and longitude are valid
5     item.lat && item.lng
6     && !isNaN(+item.lat) && item.lat >= -90 && item.lat <
7     && !isNaN(+item.lng) && item.lng >= -180 && item.lng
8   ) {
9     // If they are, concatenate them for the vis module
10    item.coordinates = item.lat + '::' + item.lng
11    return true
12  } else {
13    // If not, filter them
14    return false
15  }
16
17 });
18
19

```

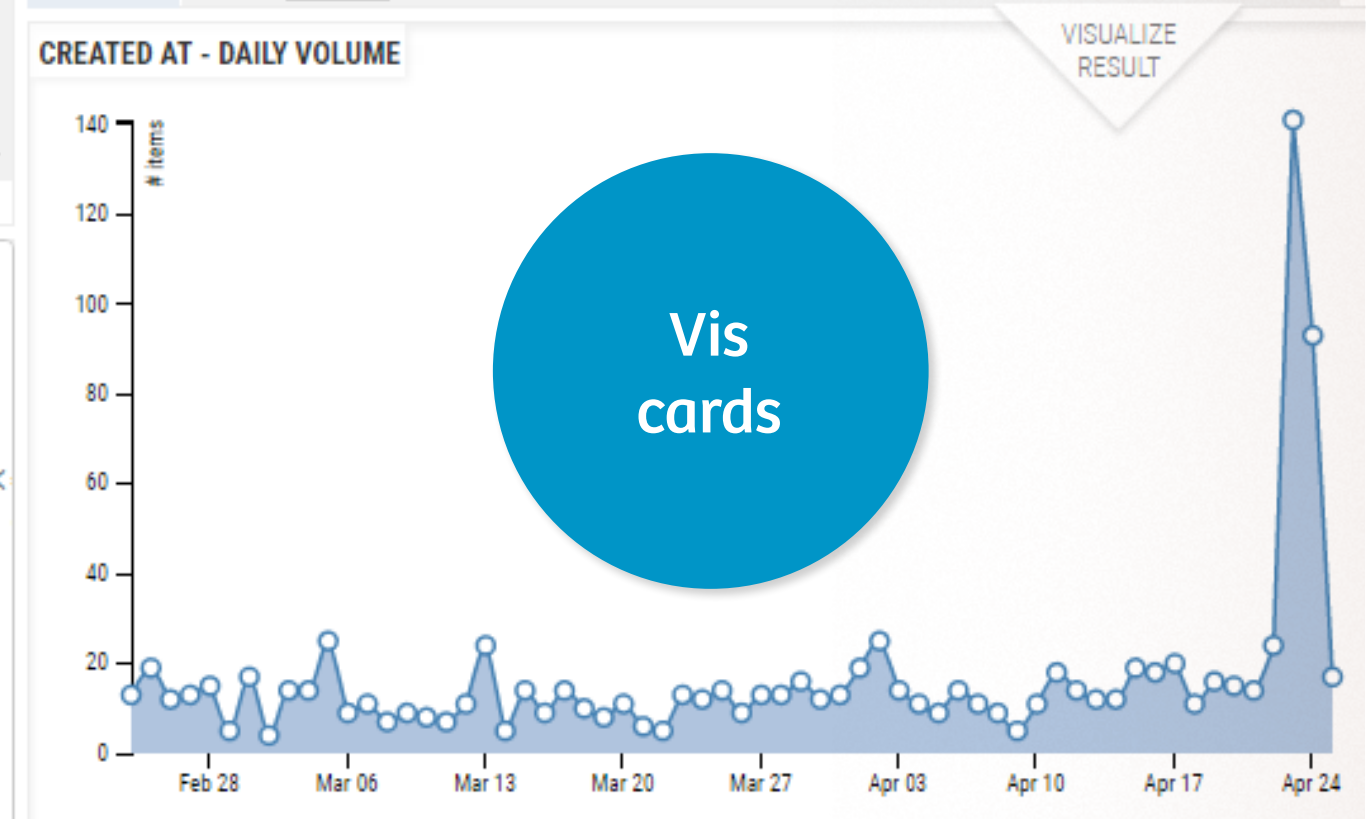
Coding panel

OUTPUT 1006 r

PREVIEW 3

Output preview

	created_at	from_user_name	text
321458.0	2016-04-22T14:38:29		#Shakespeare's grave at Holy Trinity in #Stratford #shakespeare400 @ Holy Trinity Church,... <a href="https://www.instagram.com/p/BEf9uSrlUz/">https://www.instagram.com/p/BEf9uSrlUz/</a>
595109.0	2016-04-25T14:38:29	MiguelWrites	We have indeed been celebrating



+ ADD VIZ



INPUT 152610 r	created_at	from_user_name	text	filter
PREVIEW 3	937.0	2016-04-24T08:18:57	Jeremiahjones88	LIVE on #Periscope: Sonnet Sunday #shakespeare's bday https://www.periscope .tv/w/aesncDFIV0t5Tn pMTVZqQWd8MUJkR1I EVIZYem9LWG0EqIIRE S2vMPIQg5HxWYX6we _6DVazsiFQ70v06lqY
	382.0	2016-04-12T16:26:22	Lynchypoos	RT @britishlibrary: Here's a sneak peek at our exhibition, opening this Friday. What's your favourite #Shakespeare

```

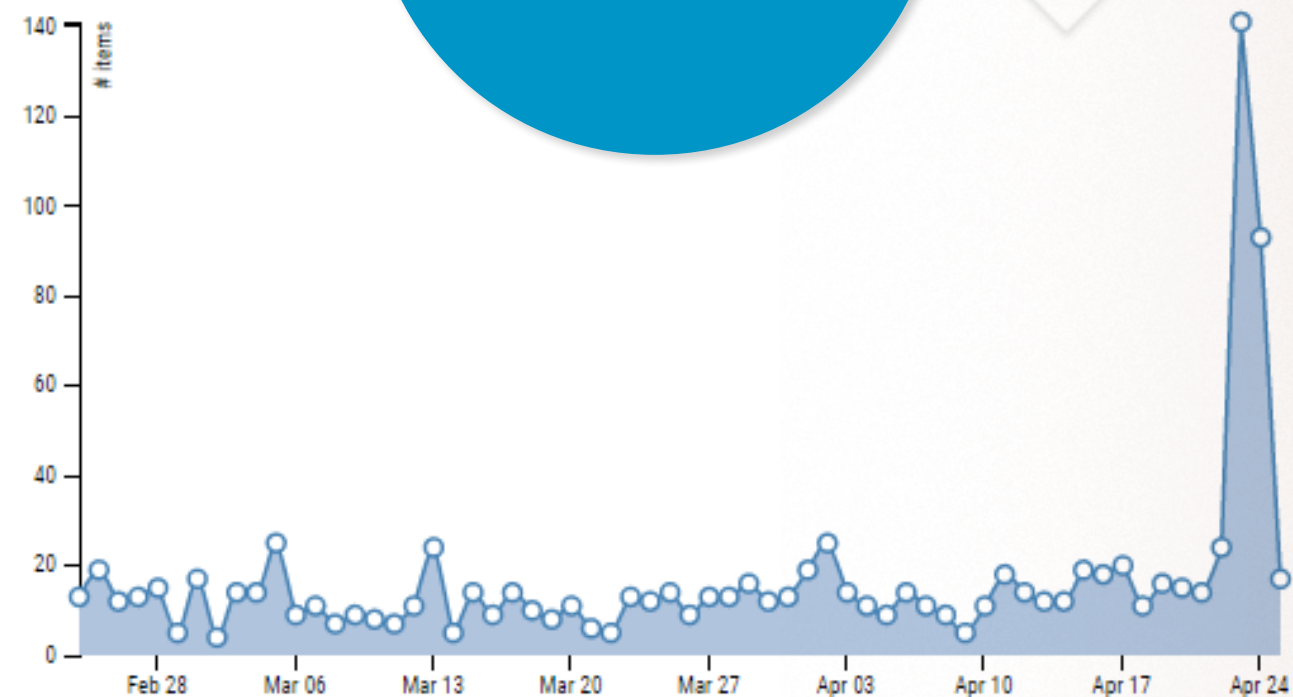
1 output = input.filter(function(item, i){
2
3     if (
4         // Check that latitude and longitude are valid
5         item.lat && item.lng
6         && !isNaN(+item.lat) && item.lat >= -90 && item.lat <
7         && !isNaN(+item.lng) && item.lng >= -180 && item.lng
8     ) {
9         // If they are, concatenate them for the vis module
10        item.coordinates = item.lat + '::' + item.lng
11        return true
12    } else {
13        // If not, filter them
14        return false
15    }
16
17 });
18
19

```

OUTPUT 1006 r	created_at	from_user_name	text
PREVIEW 3	321458.0	2016-04-22T10:37:38	jafrater
	595109.0	2	
			#Shakespeare's grave at Holy Trinity in #Stratford #shakespeare400 @ Holy Trinity Church,... https://www.instagram .com/p/BEf9uSrlUz/
			We have indeed been celebrating

You can  
download the  
modified CSV

CREATED AT - DAILY VOLUME

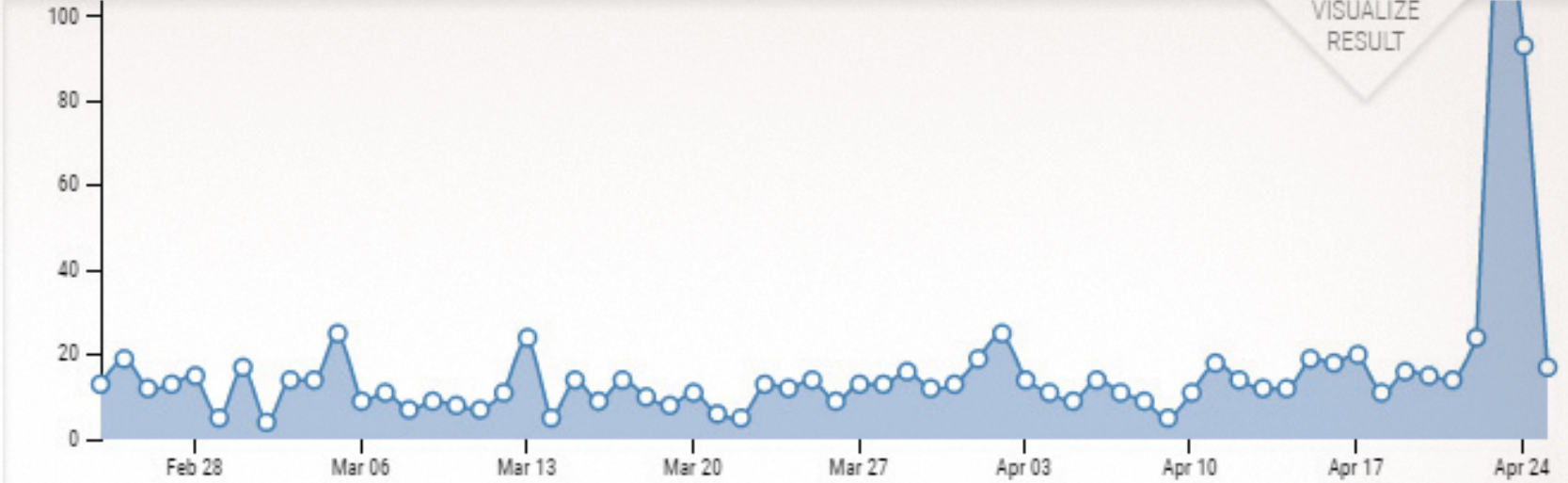


+ ADD VIZ



```
1 output = input.filter(func. CODE YOUR FILTER, i){
2
3   if (
4     // Check that latitude and longitude
5     item.lat && item.lng
6     && !isNaN(+item.lat) && item.lat >=
7     && !isNaN(+item.lng) && item.lng >=
8   ) {
9     // If they are, concatenate them for
10    item.coordinates = item.lat + '::' +
11    return true
12  } else {
13    // If not, filter them
14    return false
15  }
16
17 });
18
19
```

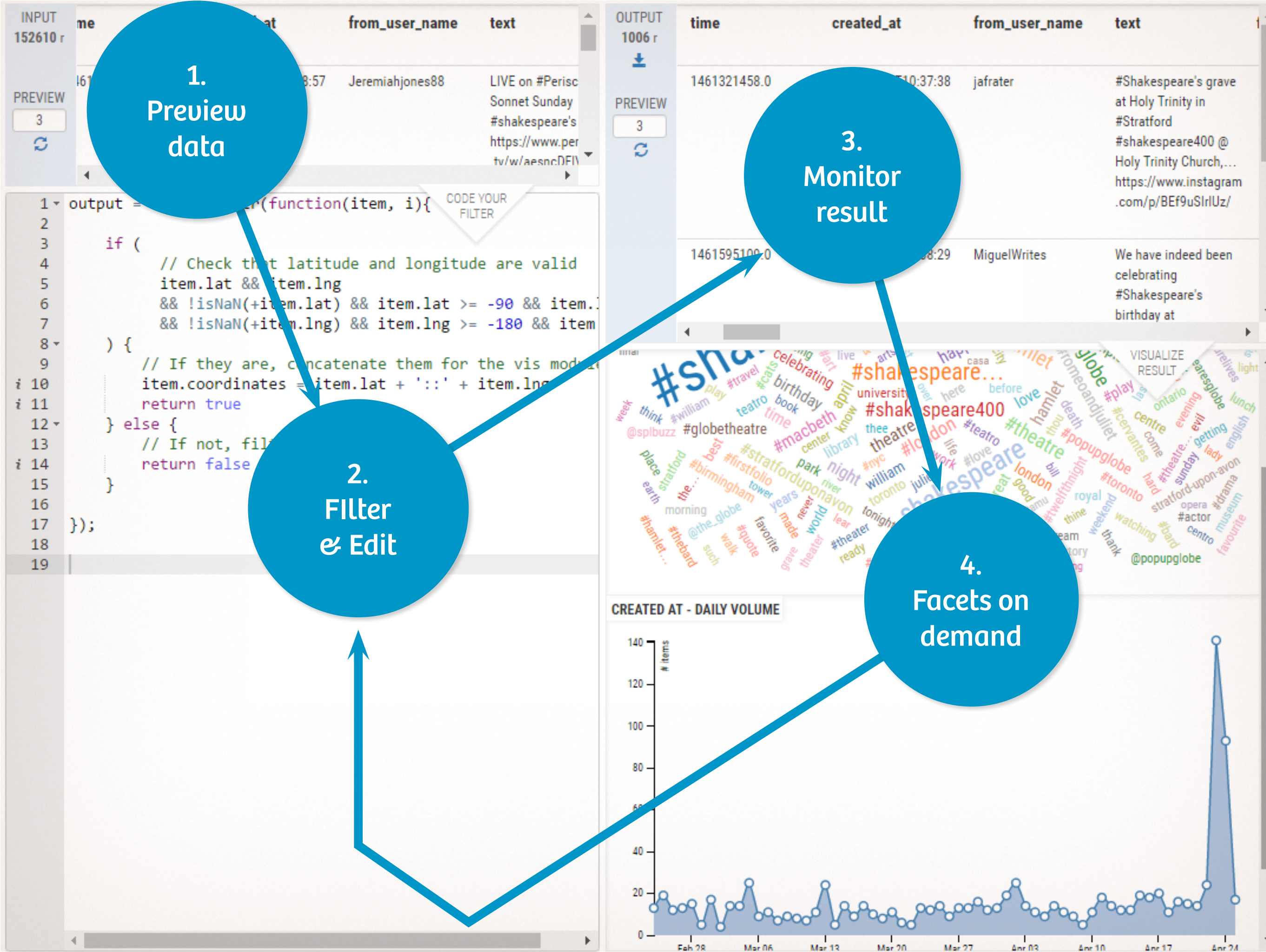
You can resize or hide the panels to fit your needs



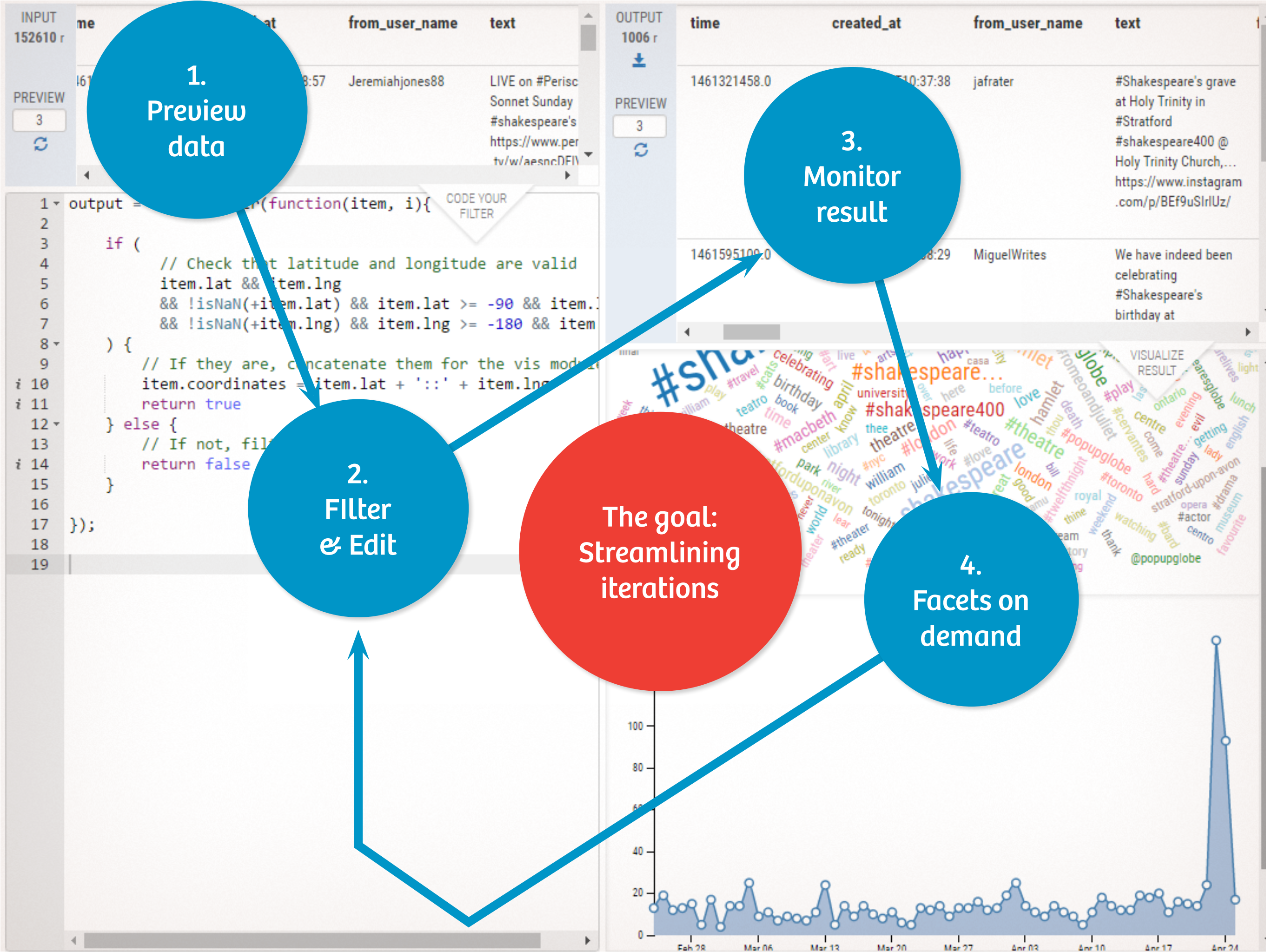
You can add & remove vis cards













# Example:

## Twitter data about Shakespeare

INPUT

152610 r

	time	created_at	from_user_name	text
8726097702912	1461452106.0	2016-04-23T22:55:06	ReDewhurst	Can someone please reassure me that I'm not a bad reader for not liking (as far as I know) #Shakespeare 😊
5302054563841	1461422679.0	2016-04-23T14:44:39	victortisp	RT @400Cervantes: En 2016 celebramos la universalidad y genialidad de los autores del #Quijote y #Hamlet #400Cervantes #Shakespeare400 http://twitter.com/400Cervantes/status/722454409788137473/o/1
6512314859520	1461406279.0	2016-04-23T10:11:19	PenguinIndia	RT @

PREVIEW

3

↺

```

1 // FILTER YOUR DATA HERE
2
3 // Just fill the "output" variable using "input"
4 output = input.filter(function(item, i){
5
6     // EXAMPLE: return the first 100 items
7     return i < 100;
8
9 });
10
11 // Hit CTRL + ENTER to run the code

```

Edit your code  
then press

**CTRL + ENTER**

*(works only from coding panel)*

Not much to see  
at upload



# Example:

## Twitter data about Shakespeare

```
1  
i 2 output = input  
3  
4 // Hit CTRL + ENTER to run the code
```



Let's filter nothing  
and take a look at  
different facets



# Example: Twitter data about Shakespeare

```

1
2 output = input
3
4 // Hit CTRL + ENTER to run the code

```

OUTPUT 152610 r

row #	id	time	created_at	from_user_na
128070	709048398062735361	1457885286.0	2016-03-13T16:08:06	JoeyFizzy


PREVIEW 3

Add  
visualization:  
which one?

Which one?

1. en (76332)  
2. fr (17222)  
3. es (10363)  
4. pt (8204)  
5. it (8116)

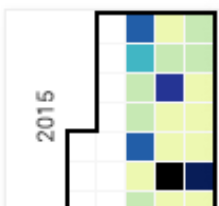
**ITEMS TOP 50**  
List of the 50 most occurrent items




**WORDS CLOUD**  
Word cloud obtained from basic text mining

16. excited (8369)  
17. info (8326)  
18. stay (8222)  
19. tuned (8136)  
20. #designer (77)


**WORDS TOP 50**  
The 50 most occurrent words, obtained from basic text mining




**CALENDAR**  
Daily count of items in a calendar view




**D DAILY VOLUME**  
Daily count of items as a curve



**M MONTHLY VOLUME**  
Monthly count of items as a curve



**Y YEARLY VOLUME**  
Yearly count of items as a curve



**MAP COORDINATES**  
Geographical distribution of items with dynamic clustering. Coordinates format: '0.0::0.0' (latitude - longitude)



# Example: Twitter data about Shakespeare

```

1
2 output = input
3
4 // Hit CTRL + ENTER to run the code

```

CODE YOUR FILTER

OUTPUT 152610 r

row #	id	time	created_at	from_user_na
128070	709048398062735361	1457885286.0	2016-03-13T16:08:06	JoeyFizzy

PREVIEW 3

We chose «Daily Volume». Data from which column?

DAILY VOLUME: Which column?

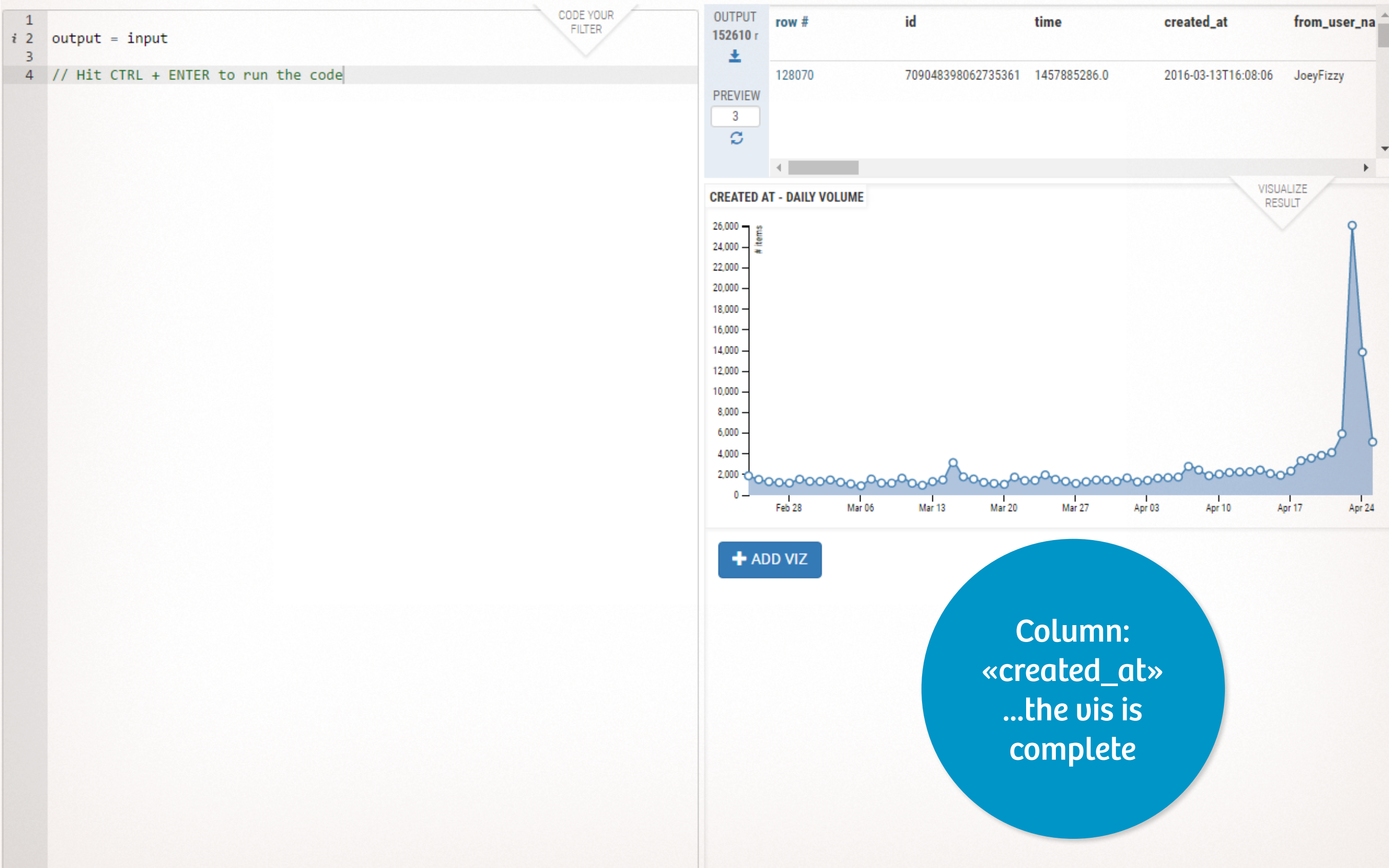
VISUALIZE RESULT

id
time
created\_at
from\_user\_name
text
filter\_level
possibly\_sensitive
withheld\_copyright
withheld\_scope
withheld\_countries
truncated
retweet\_count
favorite\_count
lang
to\_user\_name
in\_reply\_to\_status\_id
source
location
lat
lng
from\_user\_id
from\_user\_realname
from\_user\_verified
from\_user\_description
from\_user\_url
from\_user\_profile\_image\_url
from\_user\_utcoffset
from\_user\_timezone
from\_user\_lang
from\_user\_tweetcount
from\_user\_followercount
from\_user\_friendcount
from\_user\_favourites\_count
from\_user\_listed
from\_user\_withheld\_scope
from\_user\_withheld\_countries
from\_user\_created\_at



# Example:

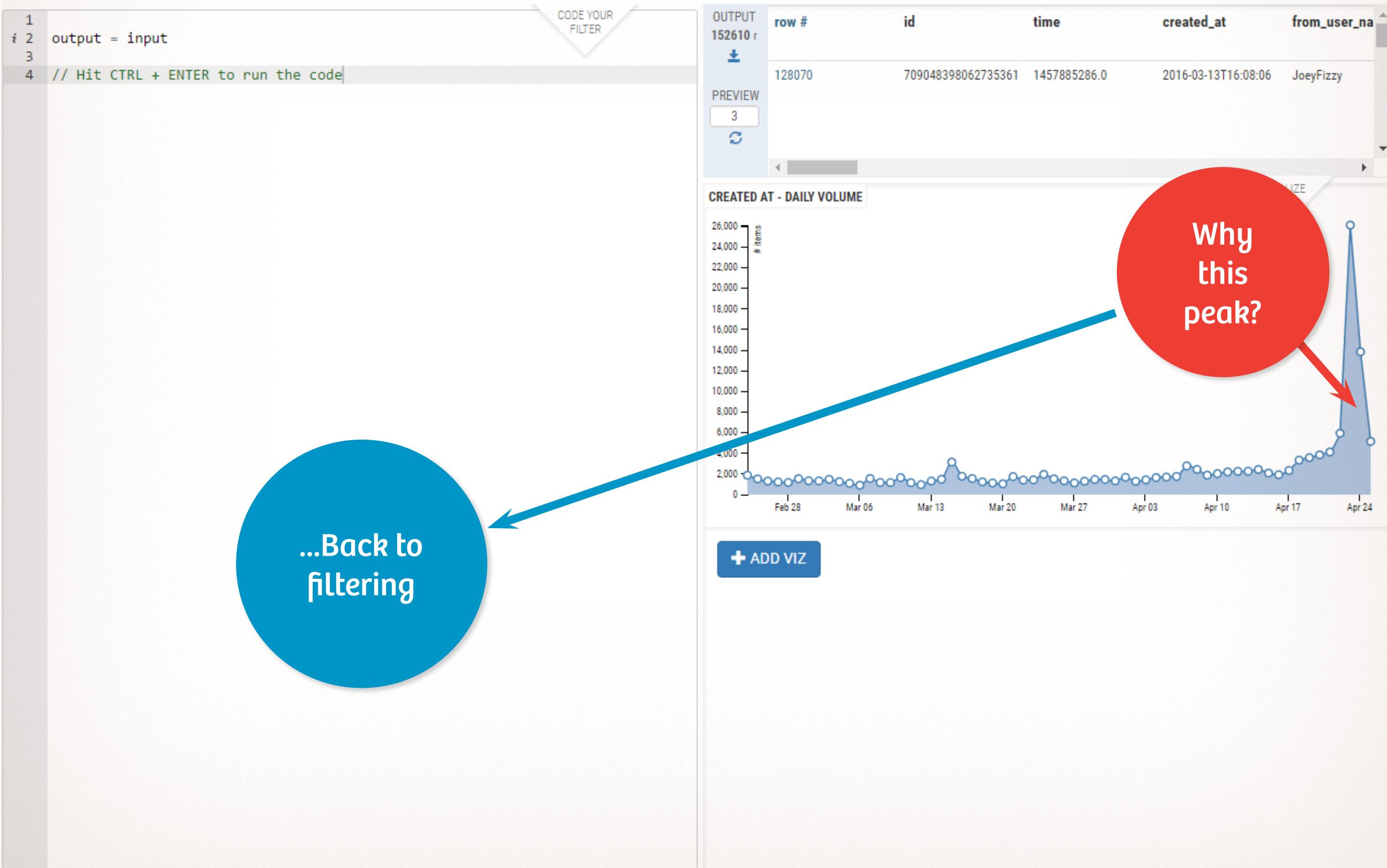
## Twitter data about Shakespeare





# Example:

## Twitter data about Shakespeare





# Example:

## Twitter data about Shakespeare

```
i 1  var limitDate = new Date('2016-04-20')
    2  output = input.filter(function(item, i){
    3
    4      var date = new Date(item.created_at)
    5      return date > limitDate
    6
    7  });
    8
```



We parse dates  
and filter after the  
2016-04-20



# Example: Twitter data about Shakespeare

```

1 var limitDate = new Date('2016-04-20')
2 output = input.filter(function(item, i){
3
4     var date = new Date(item.created_at)
5     return date > limitDate
6
7 });
8

```

CODE YOUR  
FILTER

OUTPUT  
58848 r



PREVIEW

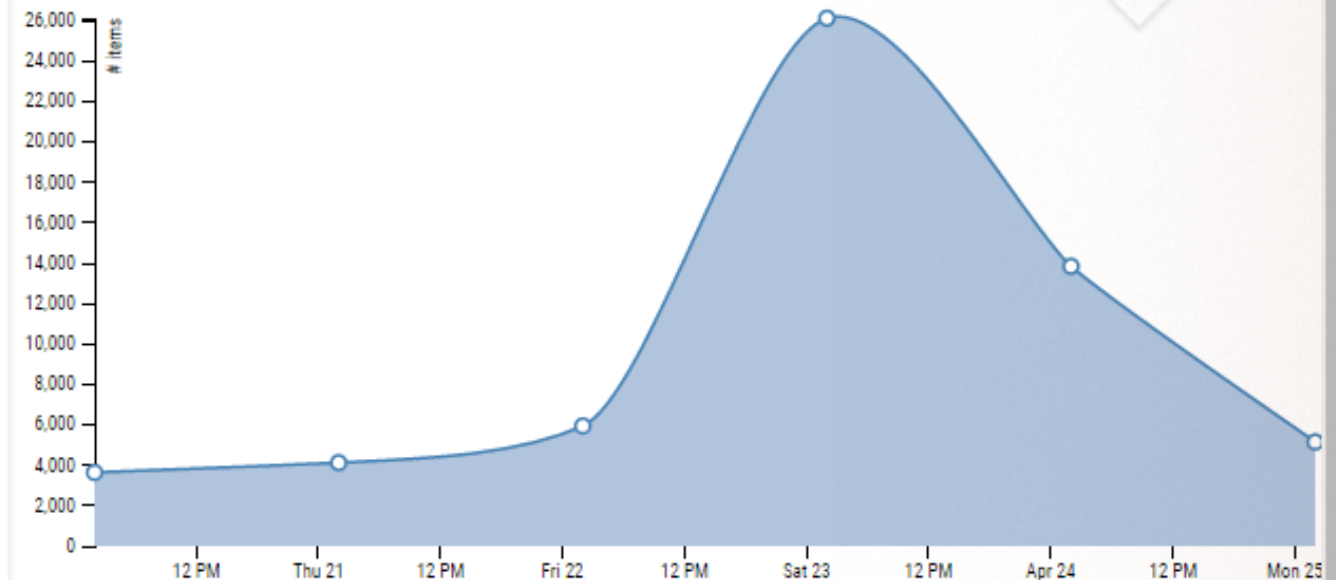
3



time	created_at	from_user_name	text	filter
797359617	1461536744.0	2016-04-24T22:25:44	peiweotck	RT @AkashaGarnier: "Boldness be my friend." ~ @Wwm_Shakespeare #Shakespeare400

VISUALIZE  
RESULT

CREATED AT - DAILY VOLUME



TEXT - WORDS TOP 50

1. #shakespeare (51389)	13. anniversary (2547)	26. here (1615)	39. play (1079)
2. #shakespeare400 (14623)	14. #cervantes (2400)	27. merchant (1596)	40. #onthisday (1054)
3. today (4496)	15. 400th (2220)	28. venice (1591)	41. prince (1041)
4. shakespeare (4124)	16. think (2010)	29. love (1535)	42. still (1025)
5. #hamlet (3824)	17. celebrate (1934)	30. gold (1534)	43. phrases (969)
6. #shakespearelives (3548)	18. #macbeth (1912)	31. glisters (1526)	44. @thersc (946)
7. wise (3516)	19. birthday (1898)	32. @chakerkhazaal (1518)	45. great (942)
8. fool (3515)	20. happy (1872)	33. more (1491)	46. course (907)
9. years (3428)	21. himself (1821)	34. plays (1417)	47. friend (860)
10. death (3206)	22. #rsclive (1801)	35. bard (1221)	48. quotes (846)
11. william (3086)	23. knows (1782)	36. @bbc1e (1194)	49. learn (843)
12. died (2928)	24. doth (1781)	37. celebrating (1147)	50. read (842)
	25. "the (1710)	38. @dailyshakes (1143)	

+ ADD VIZ

And we add a  
vis card to look at  
the content of the  
tweets



# Example:

## Twitter data about Shakespeare

### Before the peak

1. #shakespeare (84829)
2. shakespeare (10099)
3. #shakespeare400 (4721)
4. here (3376)
5. #hamlet (3141)
6. today (3128)
7. more (2902)
8. love (2679)
9. william (2588)
10. #cervantes (2489)
11. april (2469)
12. #macbeth (2190)
13. first (2058)

By iterating,  
we can  
compare

### During the peak

1. #shakespeare (51389)
2. #shakespeare400 (14623)
3. today (4496)
4. shakespeare (4124)
5. #hamlet (3824)
6. #shakespearelives (3548)
7. wise (3516)
8. fool (3515)
9. years (3428)
10. death (3206)
11. william (3086)
12. died (2928)



# Example:

## Twitter data about Shakespeare

### Before the peak

1. #shakespeare (84829)
2. shakespeare (10099)
3. #shakespeare400 (4721)
4. here (3376)
5. #hamlet (3141)
6. today (3128)
7. more (2902)
8. love (2679)
9. william (2588)
10. #cervantes (2489)
11. april (2469)
12. #macbeth (2190)
13. first (2058)

Its the anniversary  
of Shakespeare's  
death

### During the peak

1. #shakespeare (51389)
2. #shakespeare400 (14623)
3. today (4496)
4. shakespeare (4124)
5. #hamlet (3824)
6. #shakespearelives (3548)
7. wise (3516)
8. fool (3515)
9. years (3428)
10. death (3206)
11. william (3086)
12. died (2928)



# Example:

## Twitter data about Shakespeare

created\_at

from\_user\_name

text

2016-04-23T18:45:45

ru\_arena

RT @Libroantiguo:  
"The fool doth think he  
is wise, but the wise  
man knows himself to  
be a fool."  
#Shakespeare died  
#OnThisDay in 1616.  
<http://twitter.com/Libroantiguo/status/723785774420865025/photo/1>

The output monitoring  
helps validating the  
hypothesis



# Wrap up

Exploring data requires **iterating**

That is why **CSV, Rinse, Repeat**  
is about **constantly rewriting filters**

The visualizations are too basic, the preview is not comfortable...  
That's fine! You don't really need more **during exploration**.  
Our design aims at being «KISS»: Keep It Simple Stupid

Exploration is over when you have **hypotheses**.  
At this point, just switch to a more analytical environment:  
Libre Office, Tableau, R, Stata...



# Thank you for your attention

Mathieu.Jacomy@sciencespo.fr

**SciencesPo**  
MÉDIALAB

<http://medialab.sciences-po.fr>