

# Data Science 1

STAT/CS 287

James Bagrow, UVM Dept of Math and Statistics

LECTURE 01

# Today's Schedule

Course logistics

Intro + Motivation

Intro/Review of programming

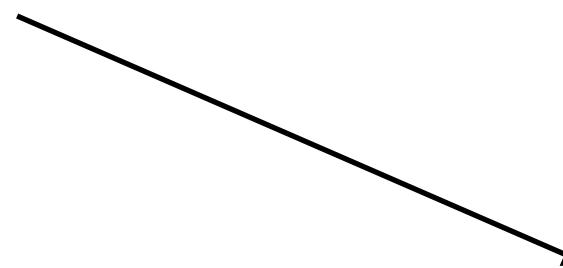
# Homework

1. Register and install Anaconda Python on your machine (see Python Setup slides)
  - Complete **HW01**
  - HW01 includes additional setup reading
2. Work through **reading**—**there's a lot** but the more prepared you are now, the faster you can work later

Quiz 01 will assess your preparation following this reading  
(will discuss schedule of assignments and quizzes shortly)

# Course logistics

Website: [bagrow.com/ds1/](http://bagrow.com/ds1/)



Syllabus  
Lecture notes/slides

## Assignments on Blackboard

- If you need to know when something is due, when is the quiz, etc., Blackboard has the truth!
- Check Blackboard regularly—*you're responsible for keeping track of assignment postings, updates to assignments, due dates, etc.*

### Data Science 1



#### Prof. James Bagrow

Email: [james.bagrow \[at\] uvm.edu](mailto:james.bagrow@uvm.edu)

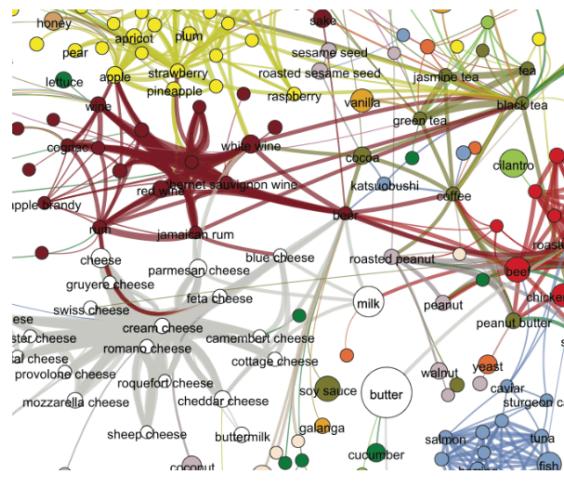
Lectures: Tu/Th 8:30–9:45 in [Perkins 107](#)

Office Hours: We/Th 10:00–11:00, or by appointment

Office: Innovation E426

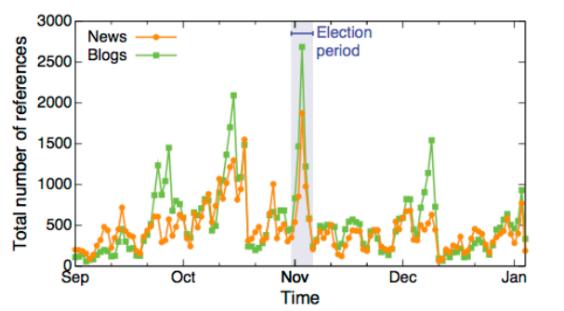
#### [Course syllabus](#)

Extracting meaning from data remains one of the biggest tasks of science and industry. The Internet and modern computers have given us vast amounts of data, so it is more important than ever to understand how to collect, process, and analyze these data. A picture is worth a thousand words, so visualizations, from scientific plots and infographics to interactive data explorers, are crucial to summarize and communicate new discoveries.



#### Resources

- The [course syllabus](#). Be sure to check this out.
- The course introductory reading, [A Whirlwind Tour of Python](#) (as a pdf) ([on GitHub](#)).
- [Anaconda](#). The scientific Python environment we will use.



#### Lecture materials and handouts

[Grid view](#) [List view](#)

<p>Preparing and submitting assignments</p> <p>Jim Bagrow STAT/CS 287 — Data Science 1</p>	<p>Python setup</p> <p>Jim Bagrow STAT/CS 287 — Data Science 1</p>
--	--

# Course logistics

## COVID-19



- Flexible and accommodating of attendance this semester → but keep on top of work!

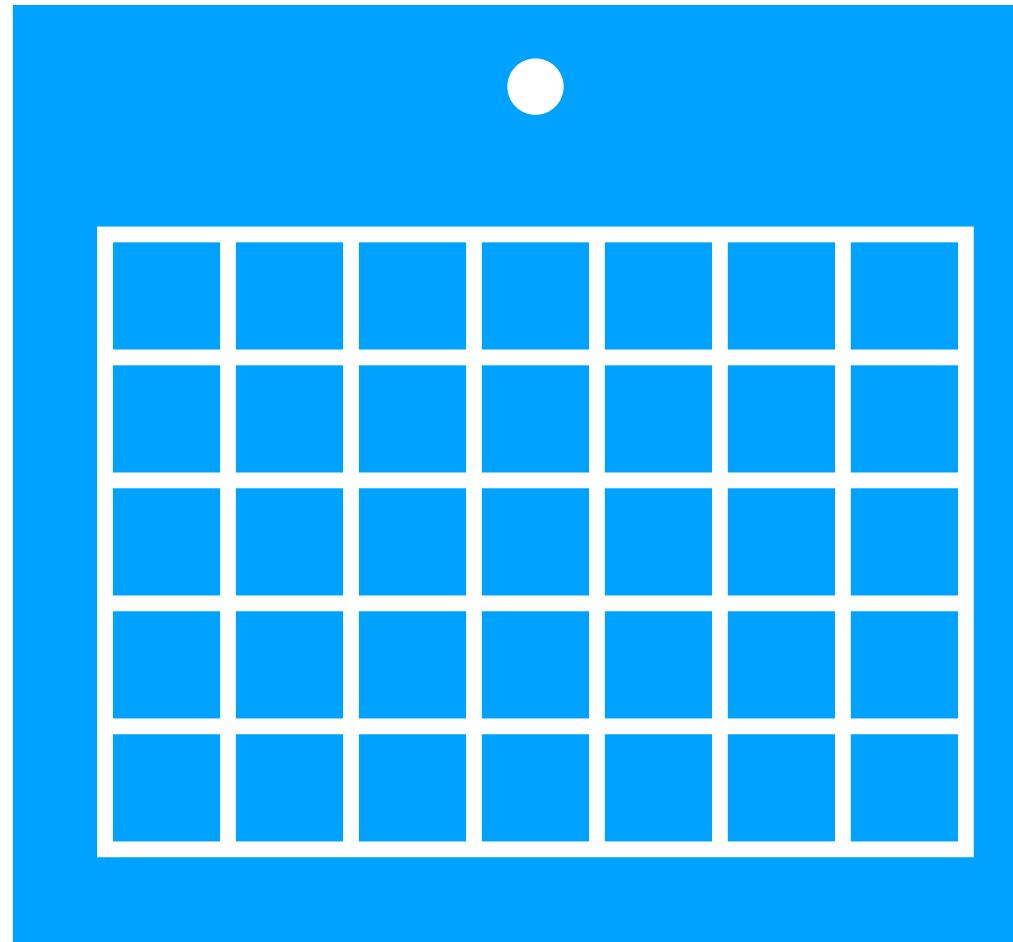


- Aim to **record lectures** (from laptop), post within 24 hours
- Future is uncertain, but right now nothing is so urgent that we need *real-time* remote

# Course logistics

Let's go over the [syllabus](#) right now!

# Tentative schedule

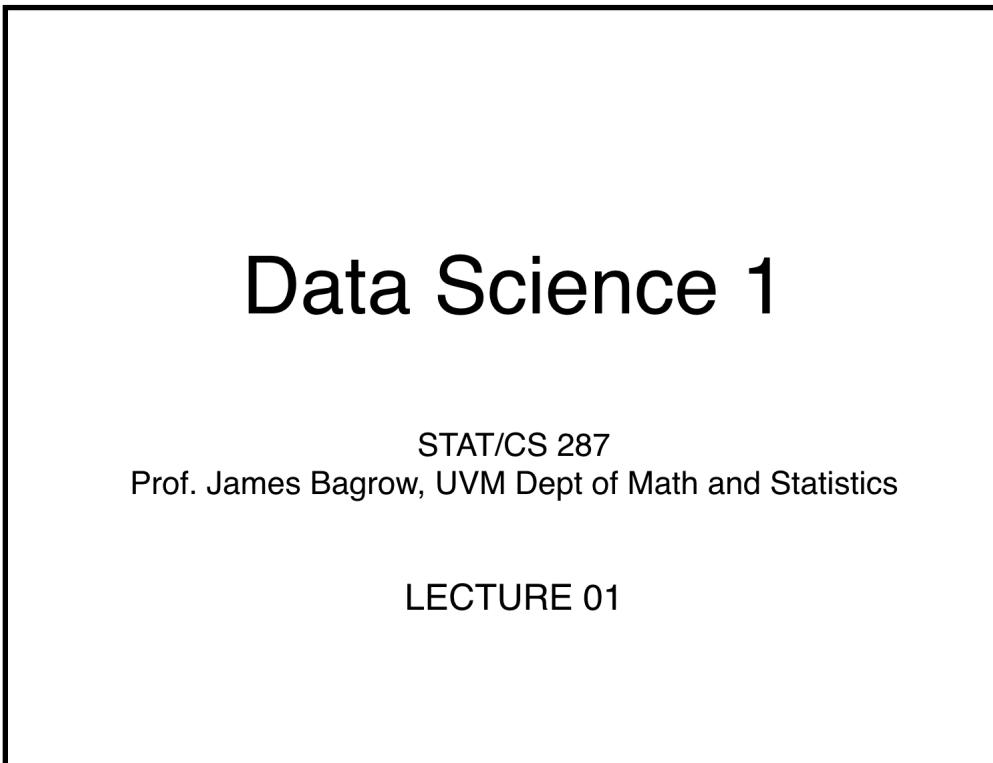


Project details coming soon!

Week	Assigned	Due	Quiz
1	HW 01		
2	HW 02	HW 01	Q 01
3	HW 03	HW 02	
4	HW 04	HW 03	Q 02
5	Project 1	HW 04	
6			Q 03
7	HW 05	Project 1	
8	HW 06	HW 05	Q 04
9	HW 07	HW 06	
10	Project 2	HW 07	
11			Q 05
12	HW 08	Project 2	
13	Thanksgiving Recess		
14	HW 09	HW 08	Q 06
15		HW 09	
Finals		Final project	

# Lectures consist of

Slides



Notebooks

**DS1 Lecture 01 supplement**

**Jim Bagrow**  
This [notebook](#) acts as an additional supplement on programming in Python. You are expected to learn this material during the first week of class, and the first quiz of the semester will assess this.

This document complements the course's "Whirlwind Tour" reading. You are also expected to complete this reading assignment. I recommend you review this supplement after reading the whirlwind tour.

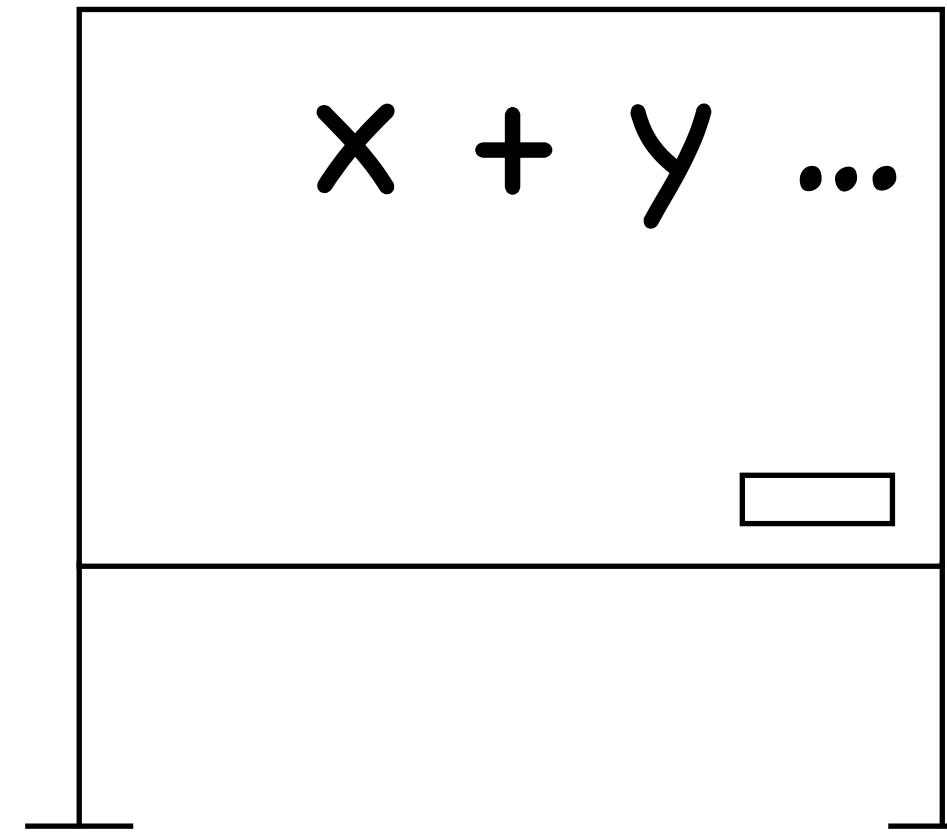
**About this document**

- This is a [Jupyter Notebook](#). I'll be using these for code-heavy lectures and lecture supplements. Notebooks contain IPython-style input and output cells in a document, interleaved with narrative (non-code) text.
- Your projects and assignments will be submitted as Python scripts (`.py` files) with associated writeups (`.docx` files).

**Contents:**

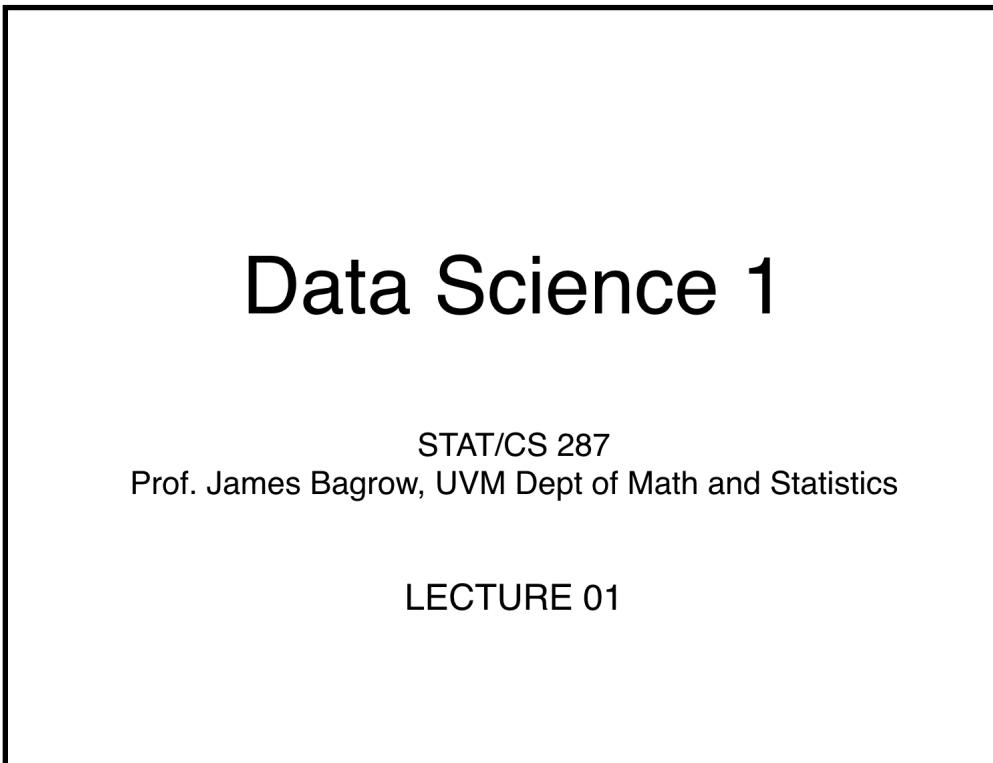
- A bit more on data structures and control flow (see also the Whirlwind Tour)
  - tuples vs. lists
    - multiassignment
    - sorting
    - zipping
- Working with [files](#)!
- Exploring the [Python standard library](#) (builtin modules and packages)
  - [numpy](#) and [scipy](#)
- A tour of useful third-party packages [time permitting]

Board



# Lectures consist of

Slides



Notebooks

**DS1 Lecture 01 supplement**

**Jim Bagrow**  
This [notebook](#) acts as an additional supplement on programming in Python. You are expected to learn this material during the first week of class, and the first quiz of the semester will assess this.

This document complements the course's "Whirlwind Tour" reading. You are also expected to complete this reading assignment. I recommend you review this supplement after reading the whirlwind tour.

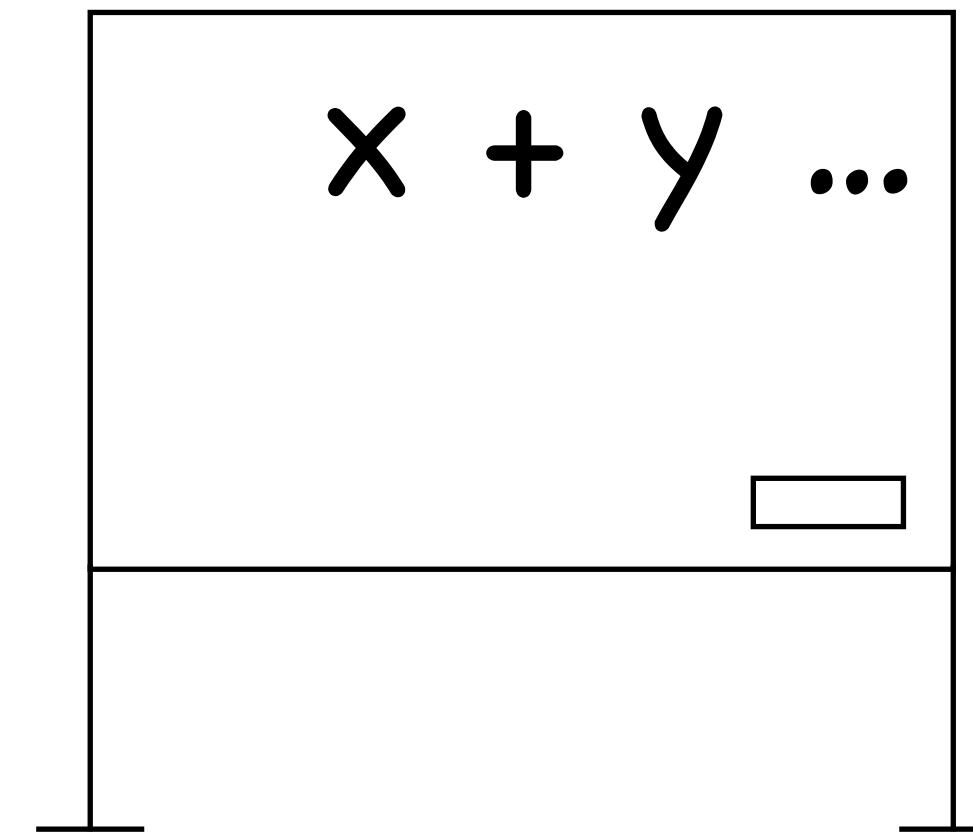
**About this document**

- This is a [Jupyter Notebook](#). I'll be using these for code-heavy lectures and lecture supplements. Notebooks contain IPython-style input and output cells in a document, interleaved with narrative (non-code) text.
- Your projects and assignments will be submitted as Python scripts (`.py` files) with associated writeups (`.docx` files).

**Contents:**

- A bit more on data structures and control flow (see also the Whirlwind Tour)
  - tuples vs. lists
    - multiassignment
    - sorting
    - zipping
- Working with [files](#)!
- Exploring the [Python standard library](#) (builtin modules and packages)
  - [numpy](#) and [scipy](#)
- A tour of useful third-party packages [time permitting]

Board



# Lectures consist of

## Notebooks

### DS1 Lecture 01 supplement

**Jim Bagrow**

This [notebook](#) acts as an additional supplement on programming in Python. You are expected to learn this material during the first week of class, and the first quiz of the semester will assess this.

This document complements the course's "**Whirlwind Tour**" reading. You are also expected to complete this reading assignment. I recommend you review this supplement after reading the whirlwind tour.

#### About this document

- This is a **Jupyter Notebook**. I'll be using these for code-heavy lectures and lecture supplements. Notebooks contain IPython-style input and output cells in a document, interleaved with narrative (non-code) text.
- Your **projects and assignments** will be submitted as Python scripts (`.py` files) with associated writeups (`.docx` files).

#### Contents:

- A bit more on data structures and control flow (see also the Whirlwind Tour)
  - tuples vs. lists
    - multiassignment
  - sorting
  - zipping
- Working with **files** !
- Exploring the [Python standard library](#) (builtin modules and packages)
- A tour of useful third-party packages [time permitting]
  - [numpy](#) and [scipy](#)

# Lectures consist of

## Notebooks

### DS1 Lecture 01 supplement

**Jim Bagrow**

This [notebook](#) acts as an additional supplement on programming in Python. You are expected to learn this material during the first week of class, and the first quiz of the semester will assess this.

This document complements the course's "**Whirlwind Tour**" reading. You are also expected to complete this reading assignment. I recommend you review this supplement after reading the whirlwind tour.

#### About this document

- This is a **Jupyter Notebook**. I'll be using these for code-heavy lectures and lecture supplements. Notebooks contain IPython-style input and output cells in a document, interleaved with narrative (non-code) text.
- Your **projects and assignments** will be submitted as Python scripts (`.py` files) with associated writeups (`.docx` files).

#### Contents:

- A bit more on data structures and control flow (see also the Whirlwind Tour)
  - tuples vs. lists
    - multiassignment
  - sorting
  - zipping
- Working with **files** !
- Exploring the [Python standard library](#) (builtin modules and packages)
- A tour of useful third-party packages [time permitting]
  - [numpy](#) and [scipy](#)

# How to succeed in this class

- Be active
- Think critically
- Start early on assignments
- Focus on initial review materials
- Organize and summarize information yourself

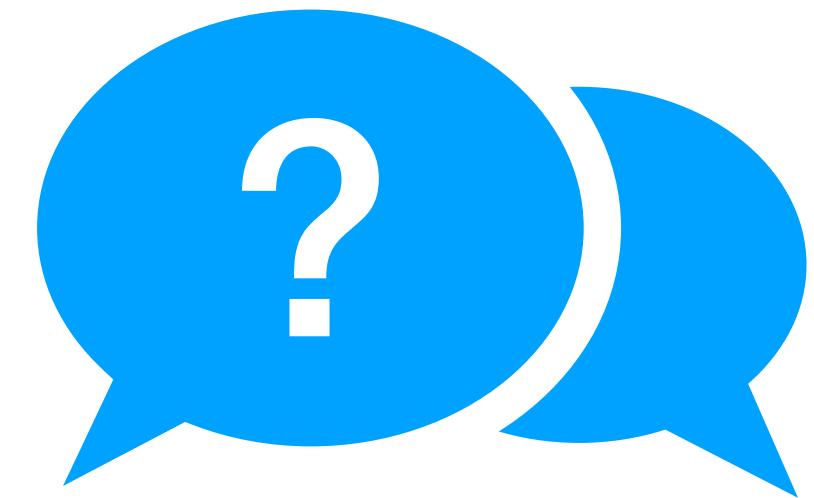
# How to succeed in this class

- Be active
- Think critically
- Start early on assignments
- Focus on initial review materials
- Organize and summarize information yourself



*Strongly advise* you to **avoid** electronic devices during class

# Course logistics



Questions?

# Today's Schedule

✓ Course logistics

Intro + Motivation

Intro/Review of programming

# Intro + Motivation

# What is data science?

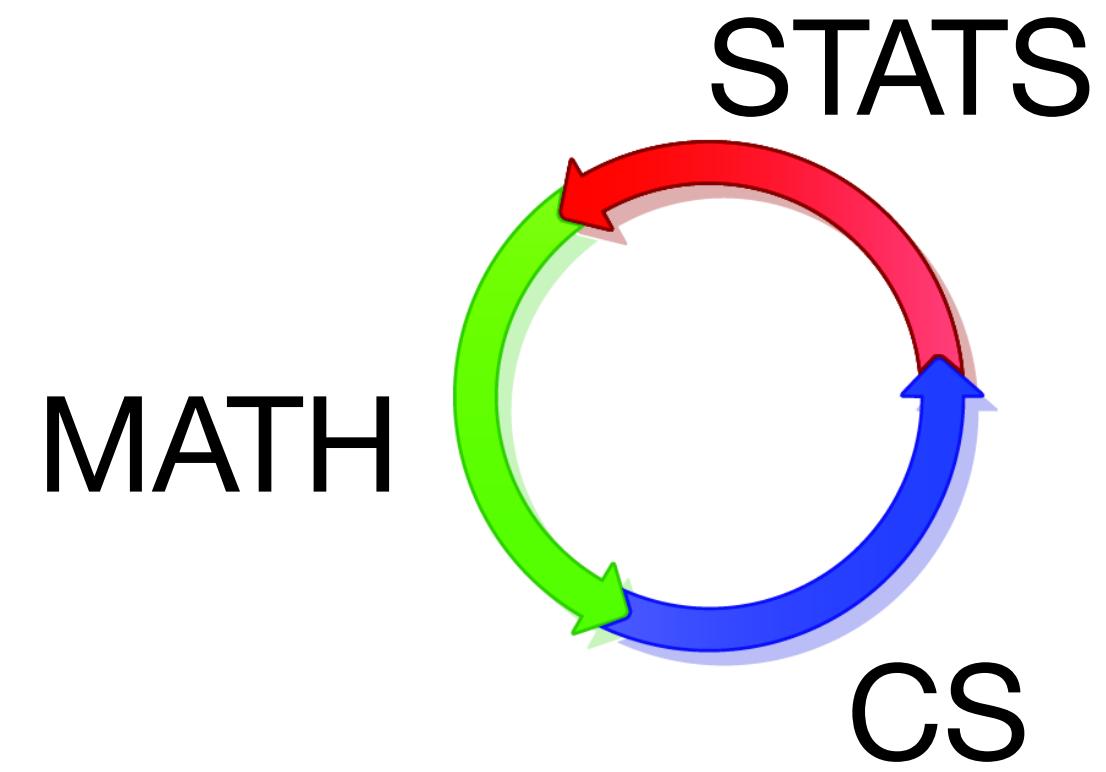
“Data science is the study of the generalizable extraction of knowledge from data”

[http://en.wikipedia.org/wiki/Data\\_science](http://en.wikipedia.org/wiki/Data_science)

In 2012, Harvard Business Review named **data scientist** the "**sexiest** job of the 21st century." More recently, Glassdoor named it the "best job of the year" for 2016. "It isn't a big surprise," Dr. Andrew Chamberlain, Glassdoor's chief economist, told Business Insider. Feb 25, 2016

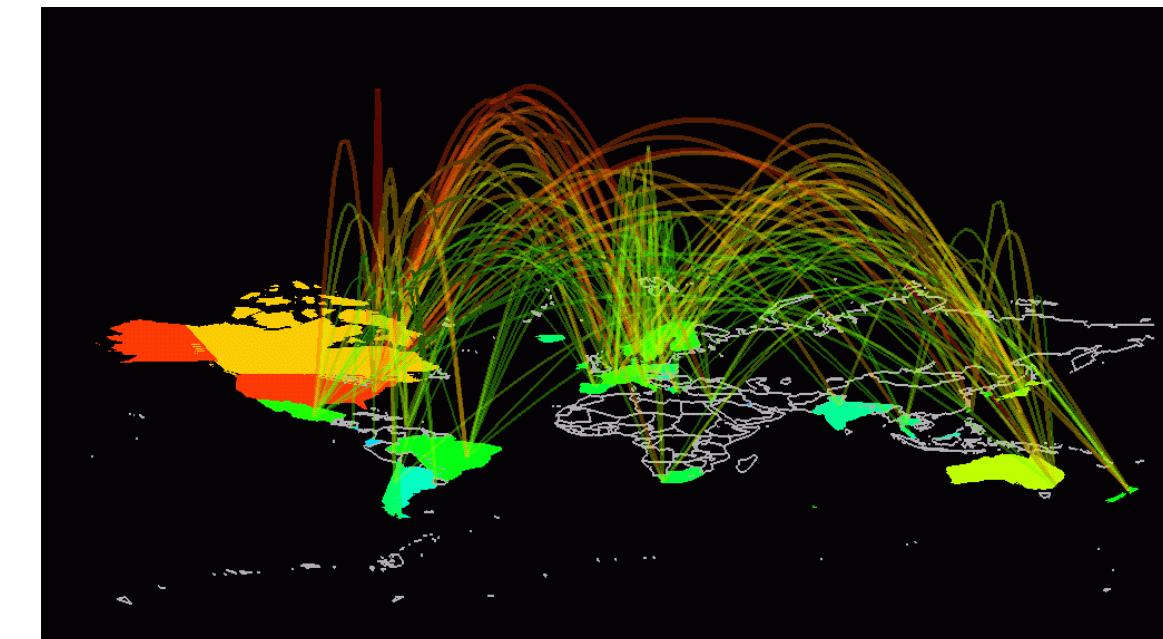
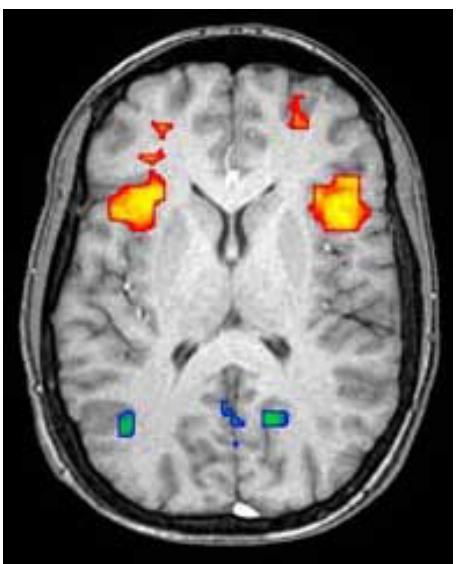
[How much money you earn in the 'sexiest job of the 21st century ...](http://www.businessinsider.com/how-much-money-you-earn-in-the-sexiest-job-of-the-21st-century-2016-2)  
[www.businessinsider.com/how-much-money-you-earn-in-the-sexiest-job... Business Insider ▾](http://www.businessinsider.com/how-much-money-you-earn-in-the-sexiest-job-of-the-21st-century-2016-2)

# What is data science?



mixed with **domain knowledge**:

- social networks/sociology
- biology or drug discovery
- neuroscience
- business
- etc...



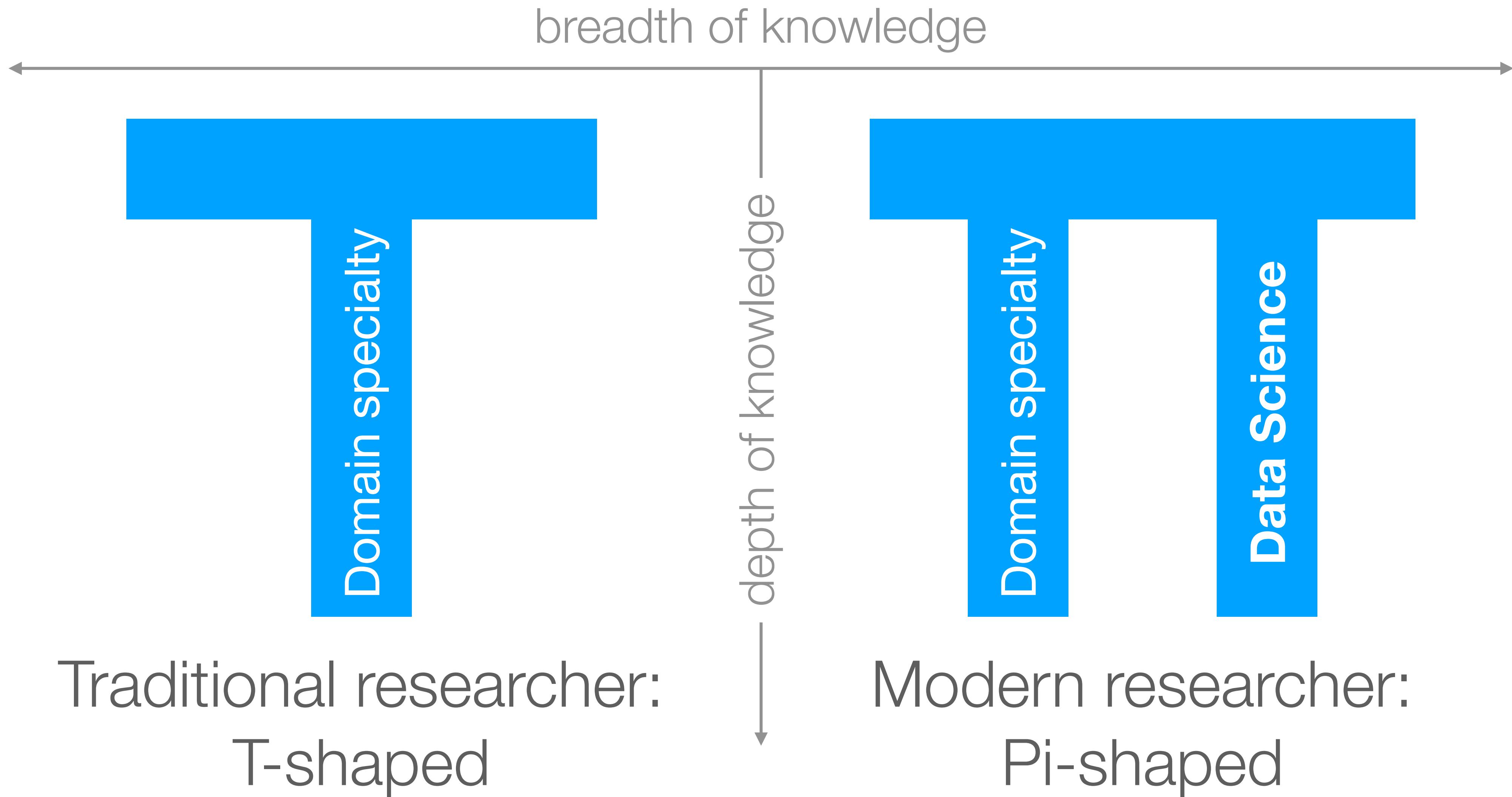
# What is data science?



A **meta-field** field encompassing:

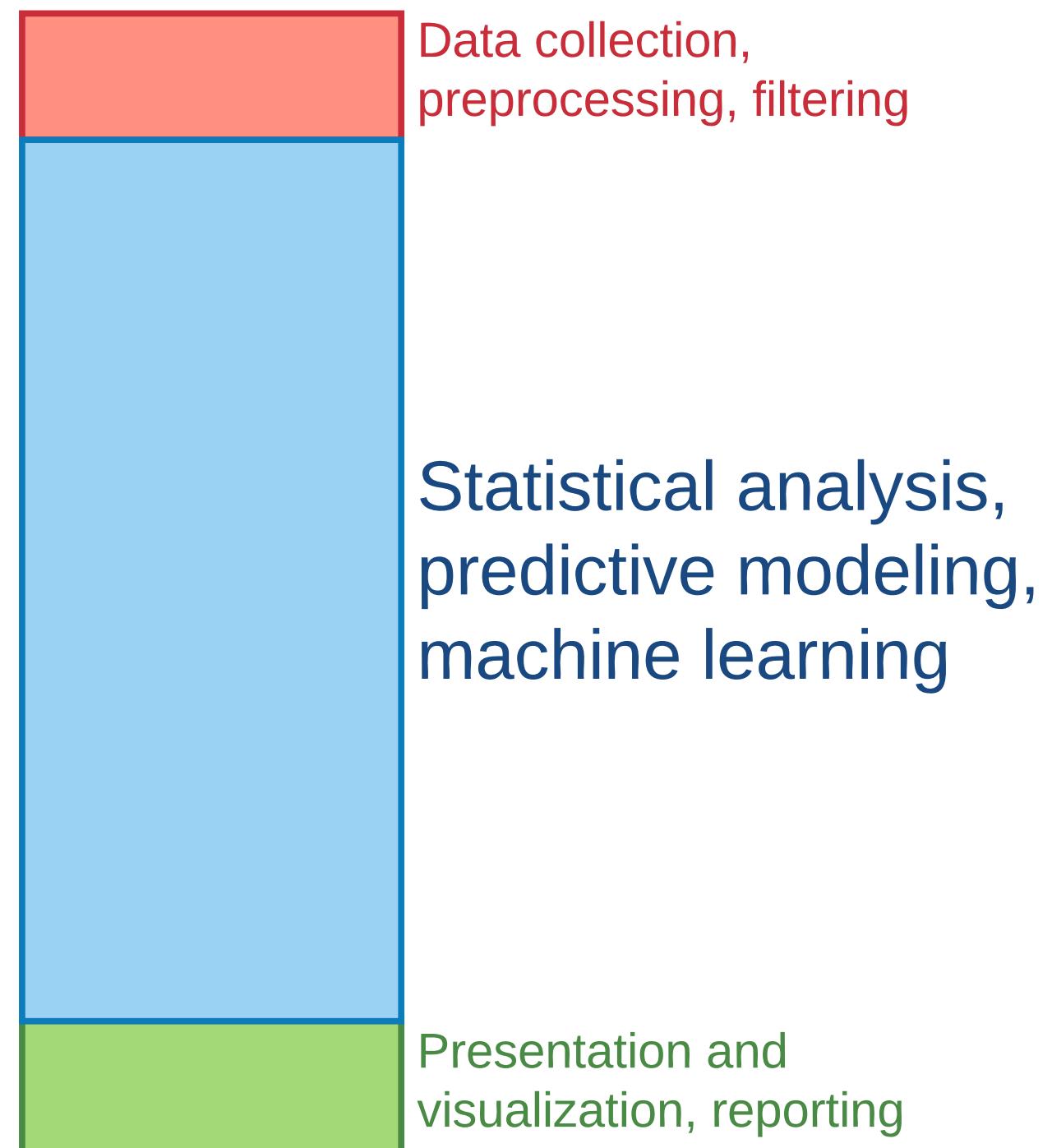
- statistics
- machine learning
- data mining
- knowledge discovery
- database engineering
- information retrieval
- visualization
- :

# What is data science?



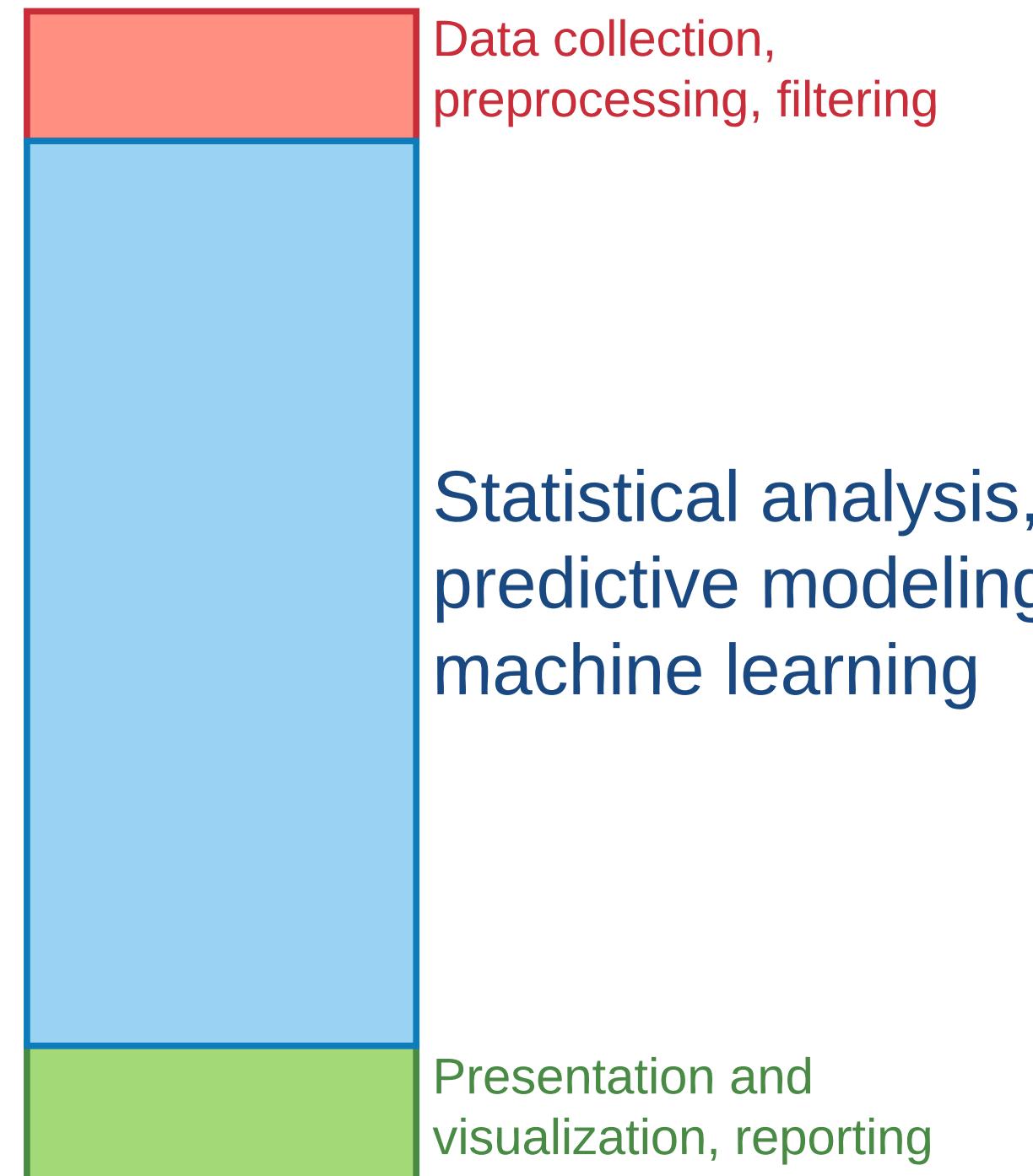
# What is data science?

What research and coursework make data science **look like**

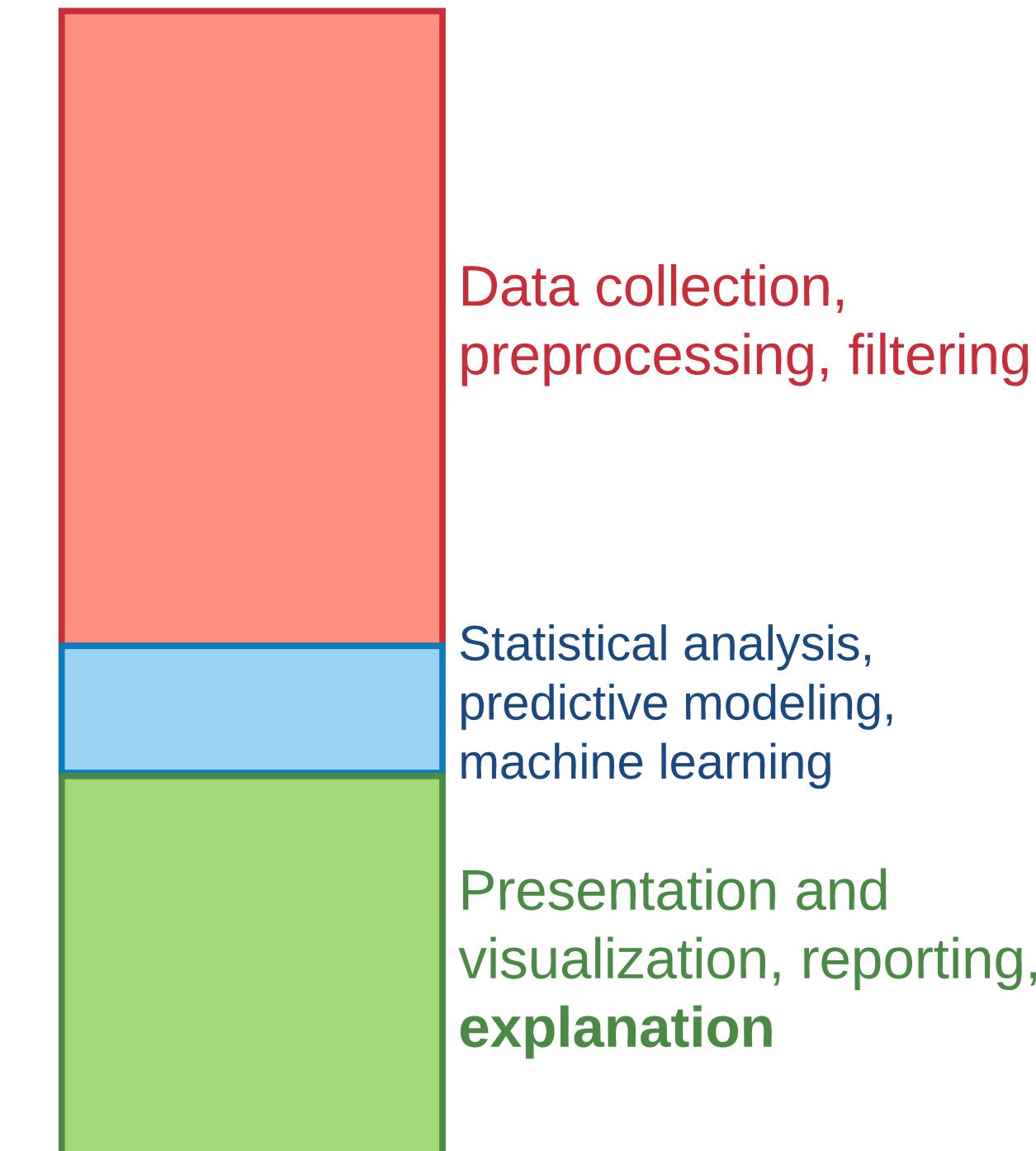


# What is data science?

What research and coursework make data science **look like**



What data science **actually is**

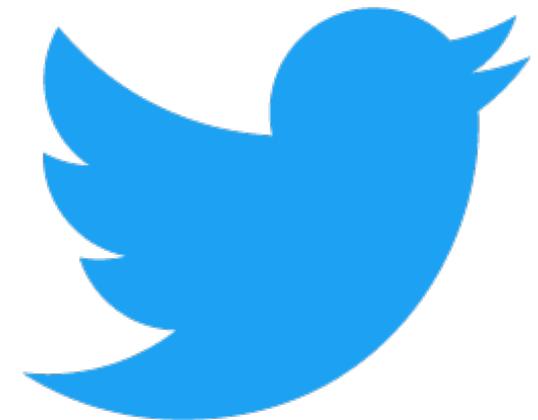


# Course topics include

- "good" practices for scientific computing
- practical data exploration and analysis
- clean and process data (data dexterity)
- web access to data sets, data "APIs"
- Natural language processing (NLP) and machine learning
- Classic regression and inference methods
- working with missing and malformed data
- Bayesian statistics and computation
- visualization principles and design
- Interpreting and communicating results!

A data science question

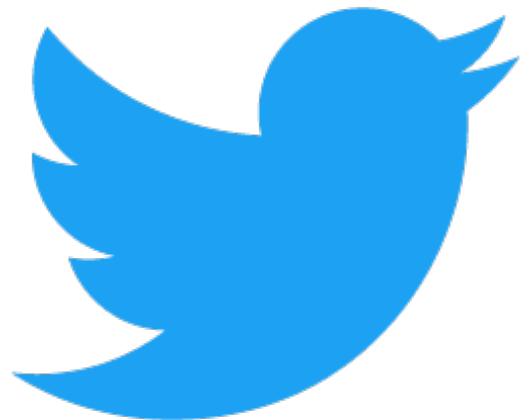
# A data science question



Twitter is famous for using very short messages called tweets. Tweets were originally limited to 140 characters

Short messages are easy to read, but not always easy to write

# A data science question

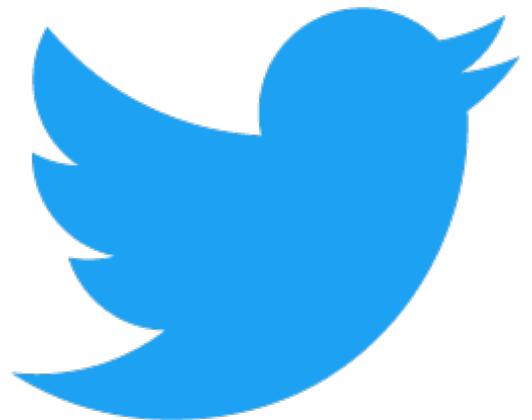


Twitter is famous for using very short messages called tweets. Tweets were originally limited to 140 characters

Short messages are easy to read, but not always easy to write

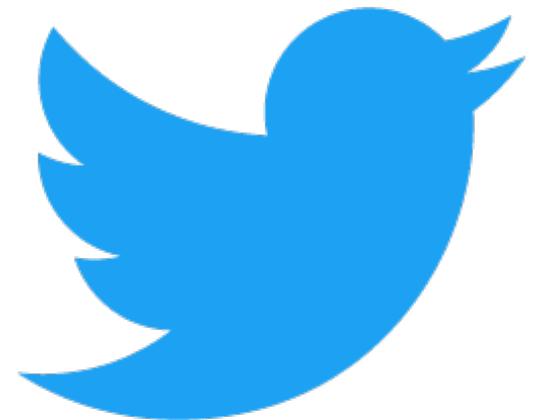
Question: Is 140 characters too short?

# A data science question

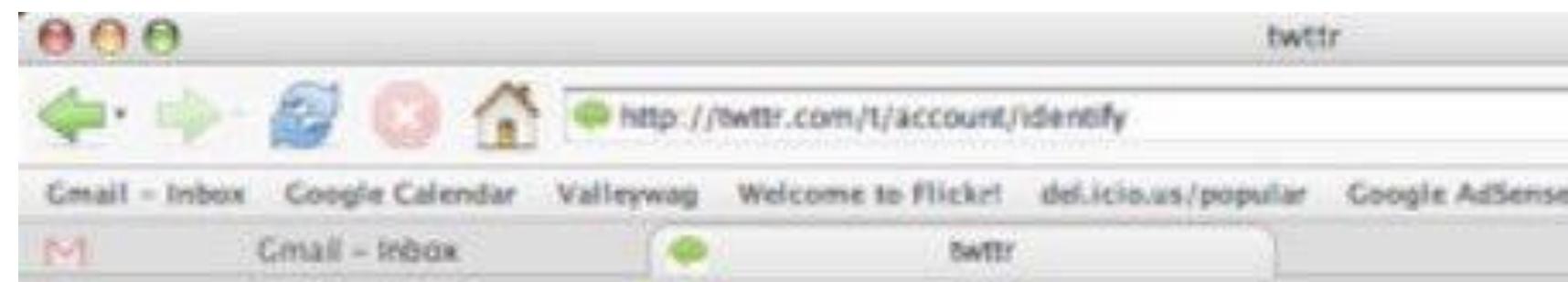


Twitter went many years with 140 characters. Why 140?

# A data science question



Twitter went many years with 140 characters. Why 140?



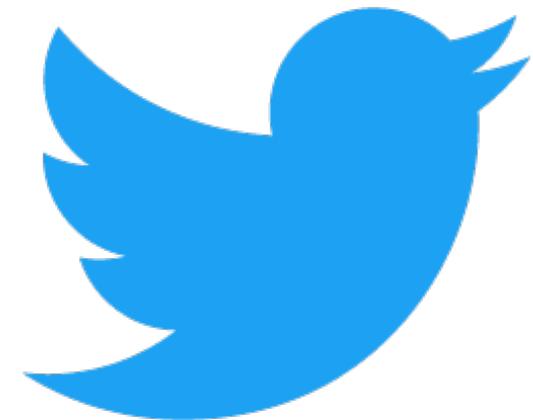
"Twittr" in 2006:



Use twittr to stay in touch with your friends all the time. If you have a cell and can txt, you'll never be bored again...E V E R !

A screenshot of the Twittr website. It features two main sections: "Timeline" on the left showing updates from friends like "Karen-Brennan", "Lenny", and "Katie", and "txt" (or) "What are you doing?" on the right where users can enter text and send updates. Below these are buttons for "what's up?" and "send updates from your cell".

# A data science question



Twitter went many years with 140 characters. Why 140?



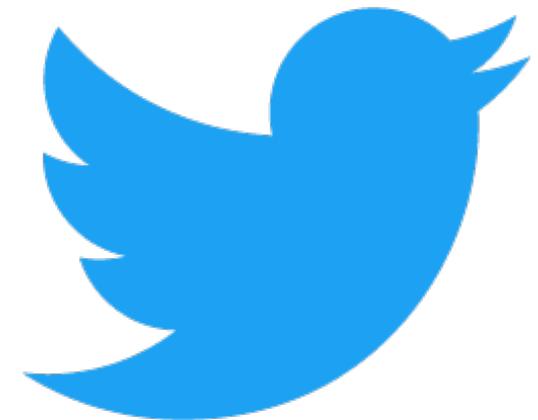
"Twittr" in 2006:



The screenshot shows the Twittr homepage from 2006. The title "twittr" is displayed in a large, green, stylized font. Below it is a tagline: "Use twittr to stay in touch with your friends all the time. If you have a cell and can txt, you'll never be bored again...E V E R !". Two main features are shown: "Timeline" (with a preview of friend posts) and "what up?" (a text input field for sending updates). The URL "http://twittr.com" is visible at the bottom.

"If you have a cell and  
can txt..."

# A data science question

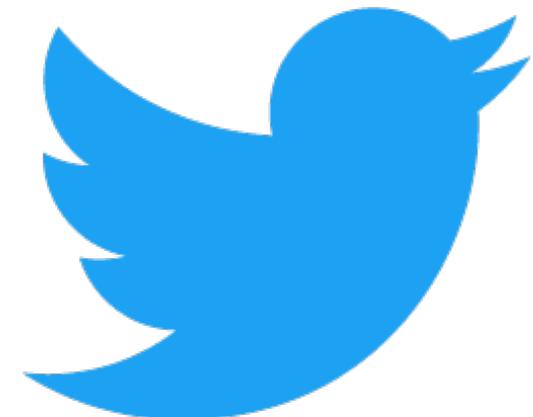


Twitter went many years with 140 characters. Why 140?

Twitter was originally focused on text messaging (SMS)

The SMS communication protocol (over 30 years old!) specifies the **"payload" size** for each message:

# A data science question



Twitter went many years with 140 characters. Why 140?

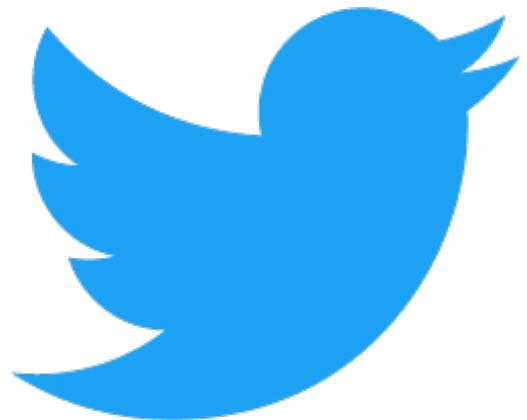
Twitter was originally focused on text messaging (SMS)

The SMS communication protocol (over 30 years old!) specifies the "payload" size for each message:

Annex 5 - Breakdown of SMS PDU lengths	
Mobile Originated	Mobile Terminated
User Data 140	User Data 140
MTI, VPF, SRR 1	MTI, MMS, SRI 1
Message Ref 1	Orig Addr 12
Dest Addr 12	Protocol Id 1
Protocol Id 1	Data Code Schm 1
Data Code Schm 1	SC Timestamp 7
Validity Per 7	User Data Len 1
User Data Len 1	
TPDU TOTAL 164	TPDU TOTAL 163
TPDU 164	
TPDU Length 1	
RP-Message type 1	
RP-Message ref 1	
RP-OA, RP-DA 14	
RPDU TOTAL 181	

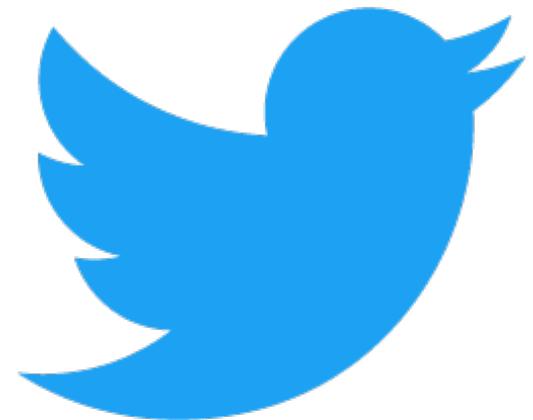
Source: GSM4 126/91

# A data science question



Question: Is 140 characters too short?

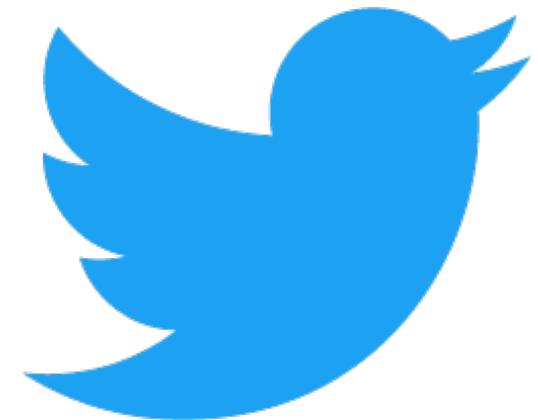
# A data science question



Question: Is 140 characters too short?

Idea: We can look at people's tweets, how long are tweets?

# A data science question

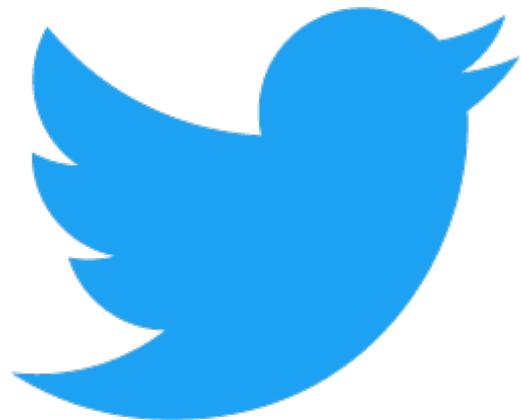


Question: Is 140 characters too short?

Idea: We can look at people's tweets, how long are tweets?

Hypothesis 1: most tweets are short,  
almost none reach 140  
→ Plenty of room

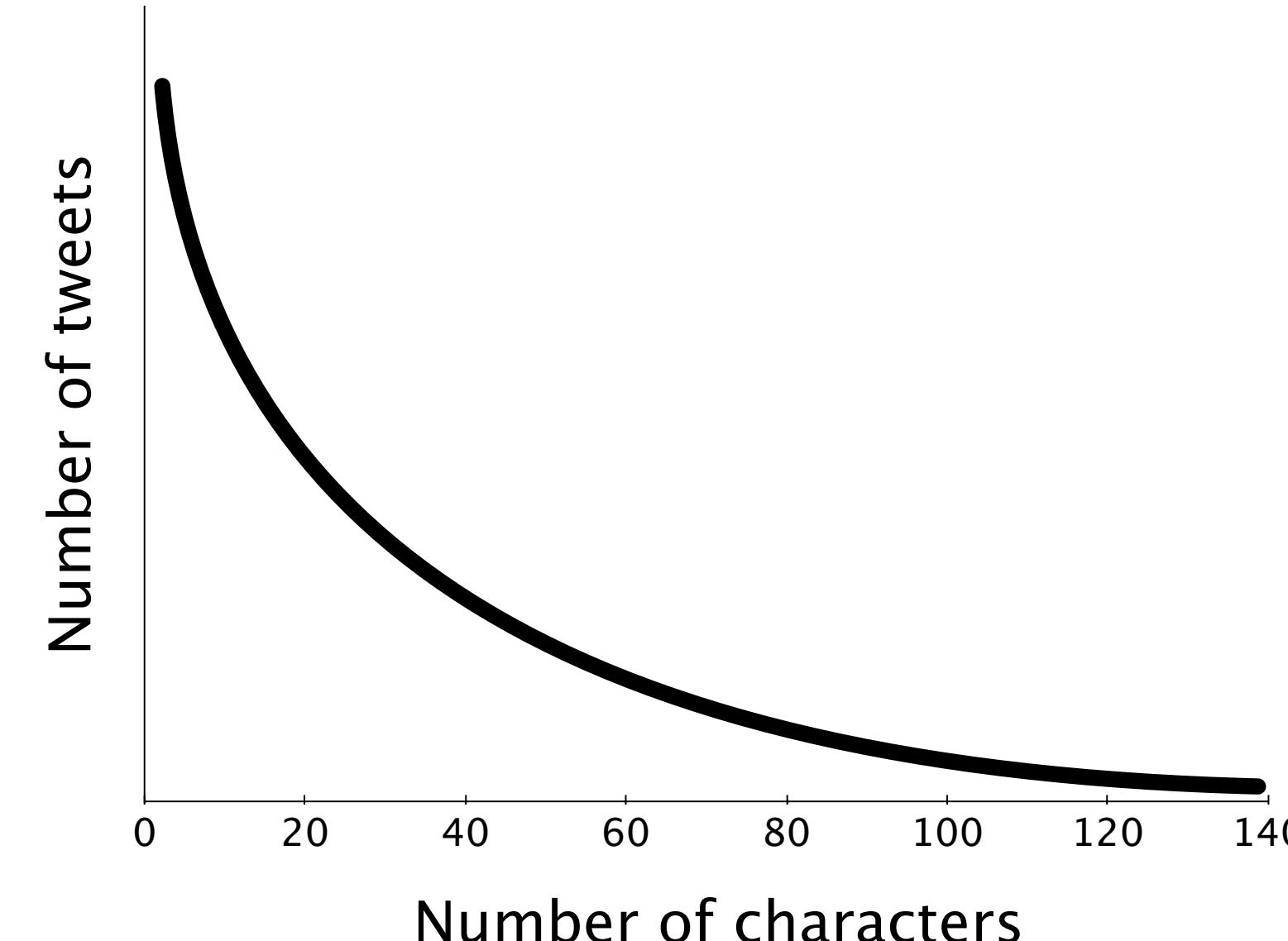
# A data science question



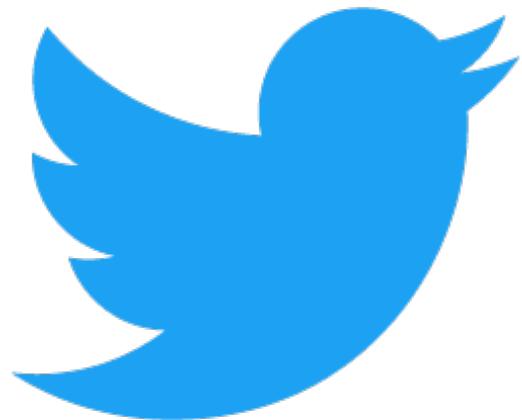
Question: Is 140 characters too short?

Idea: We can look at people's tweets, how long are tweets?

Hypothesis 1: most tweets are short,  
almost none reach 140  
→ Plenty of room



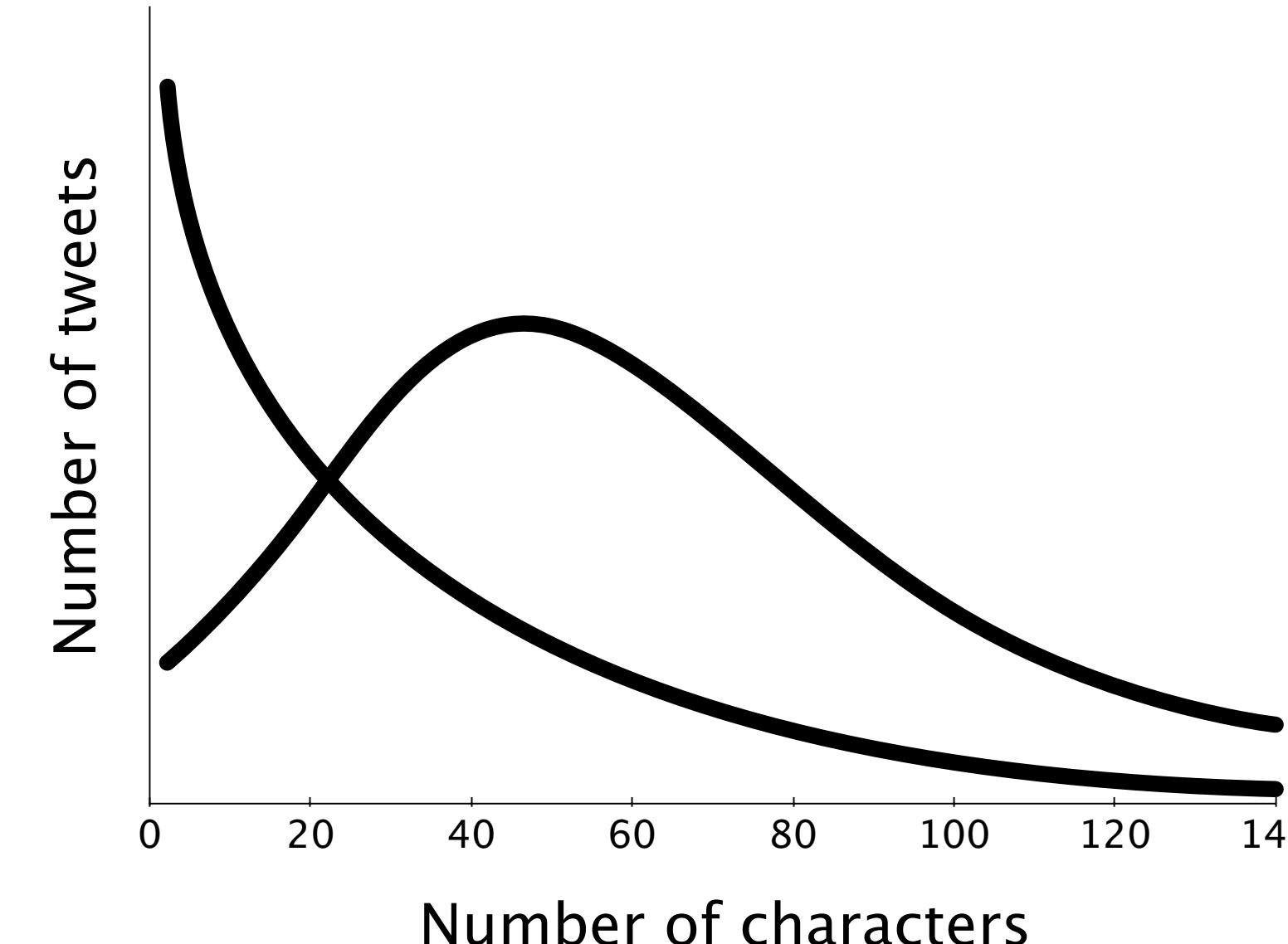
# A data science question



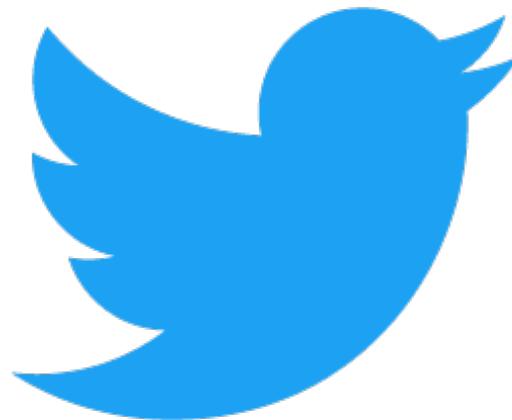
Question: Is 140 characters too short?

Idea: We can look at people's tweets, how long are tweets?

Hypothesis 1: most tweets are short,  
almost none reach 140  
→ Plenty of room



# A data science question

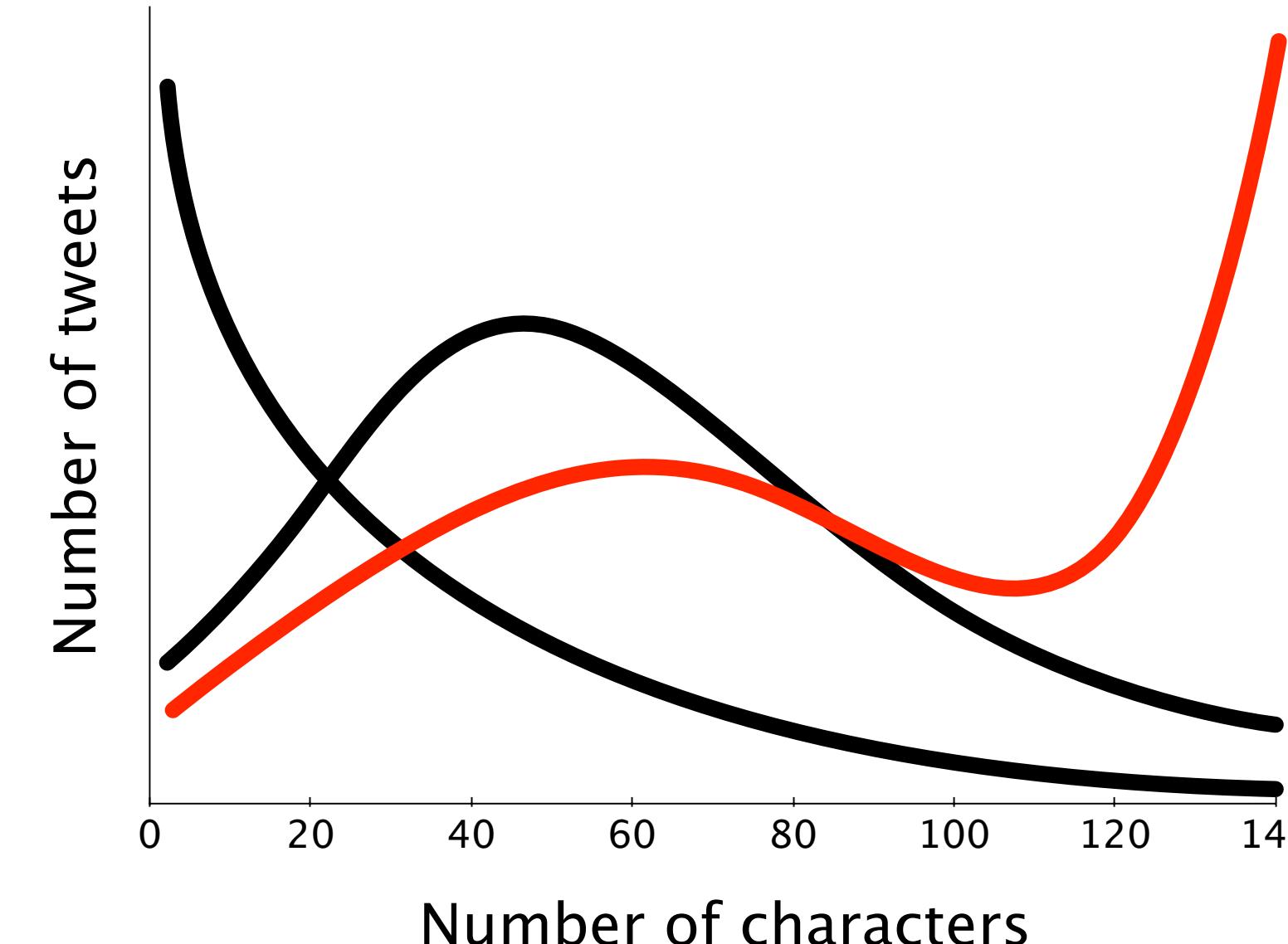


Question: Is 140 characters too short?

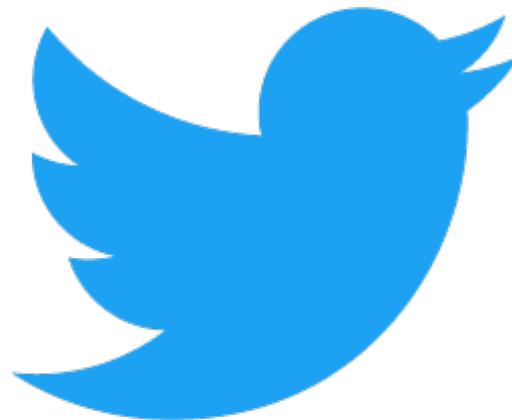
Idea: We can look at people's tweets, how long are tweets?

Hypothesis 1: most tweets are short,  
almost none reach 140  
→ Plenty of room

Hypothesis 2: many tweets hitting  
the 140 limit  
→ People are *cramming* tweets



# A data science question

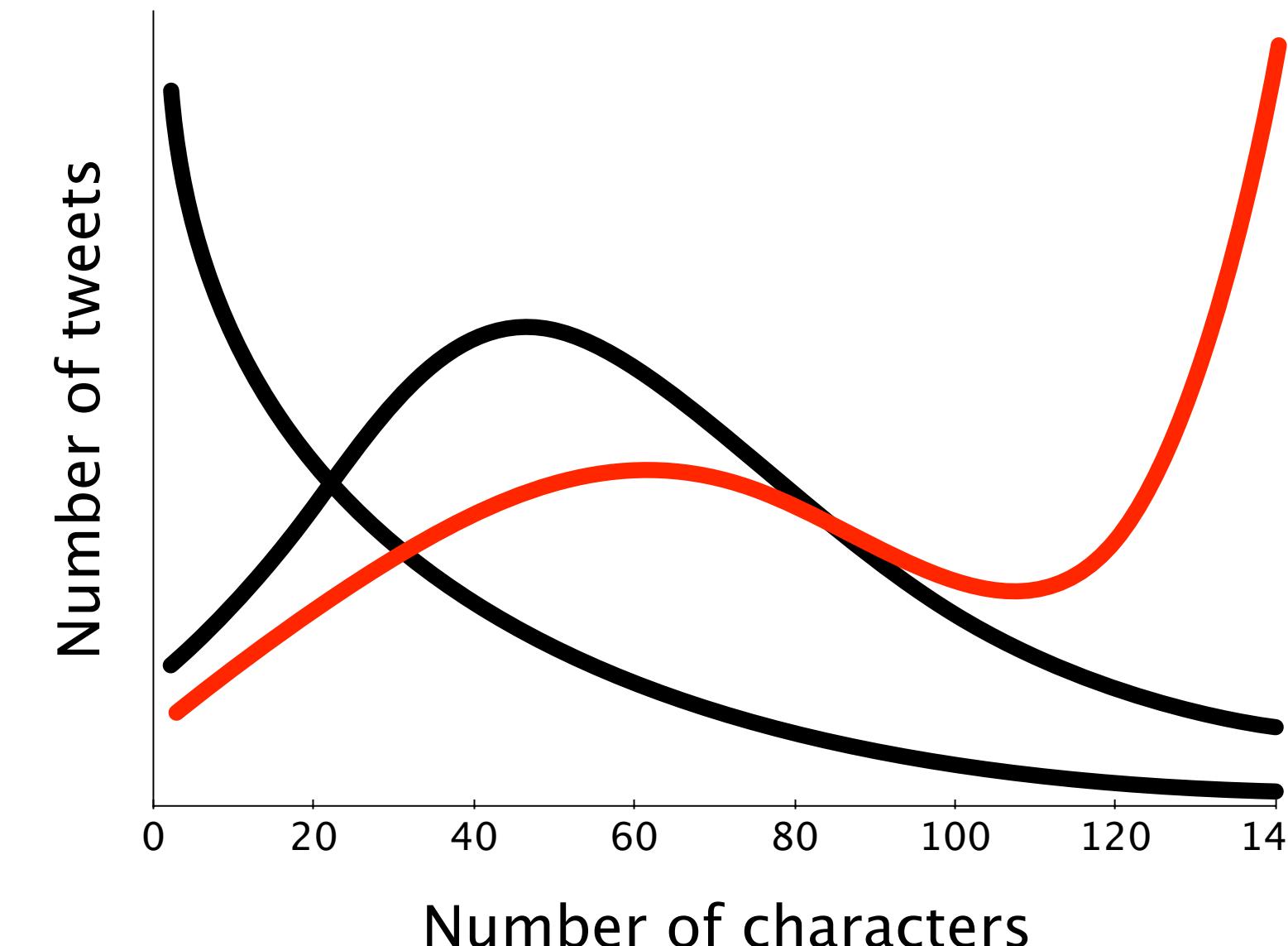


Question: Is 140 characters too short?

Idea: We can look at people's tweets, how long are tweets?

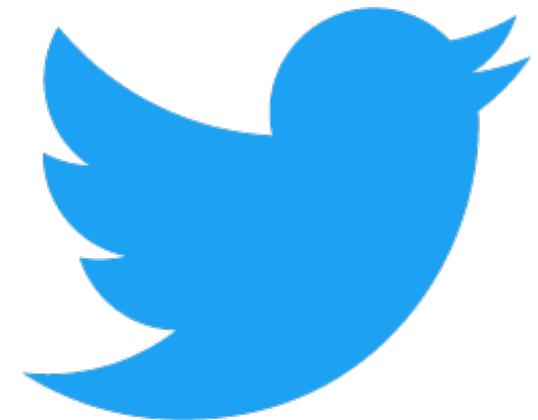
Hypothesis 1: most tweets are short,  
almost none reach 140  
→ Plenty of room

Hypothesis 2: many tweets hitting  
the 140 limit  
→ People are *cramming* tweets

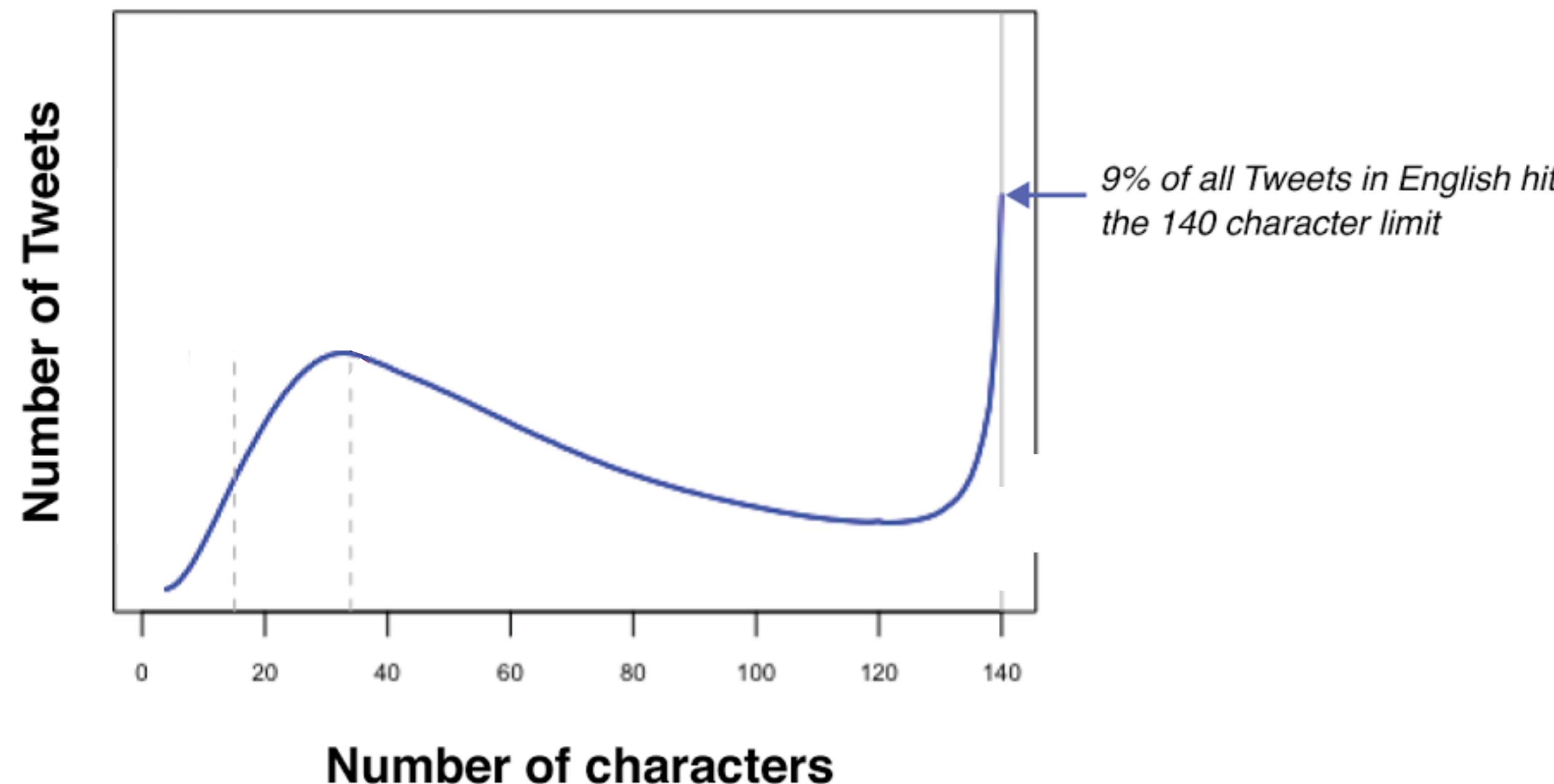


What's it really  
look like?

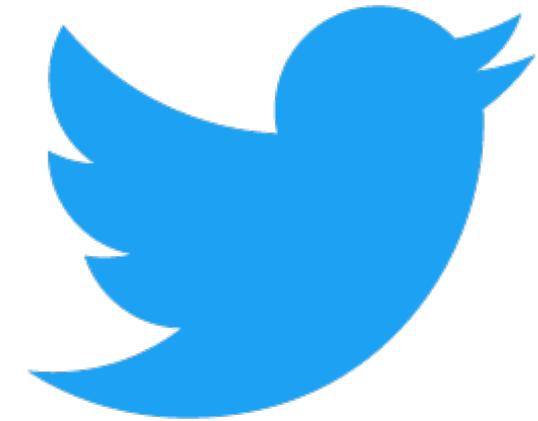
# A data science question



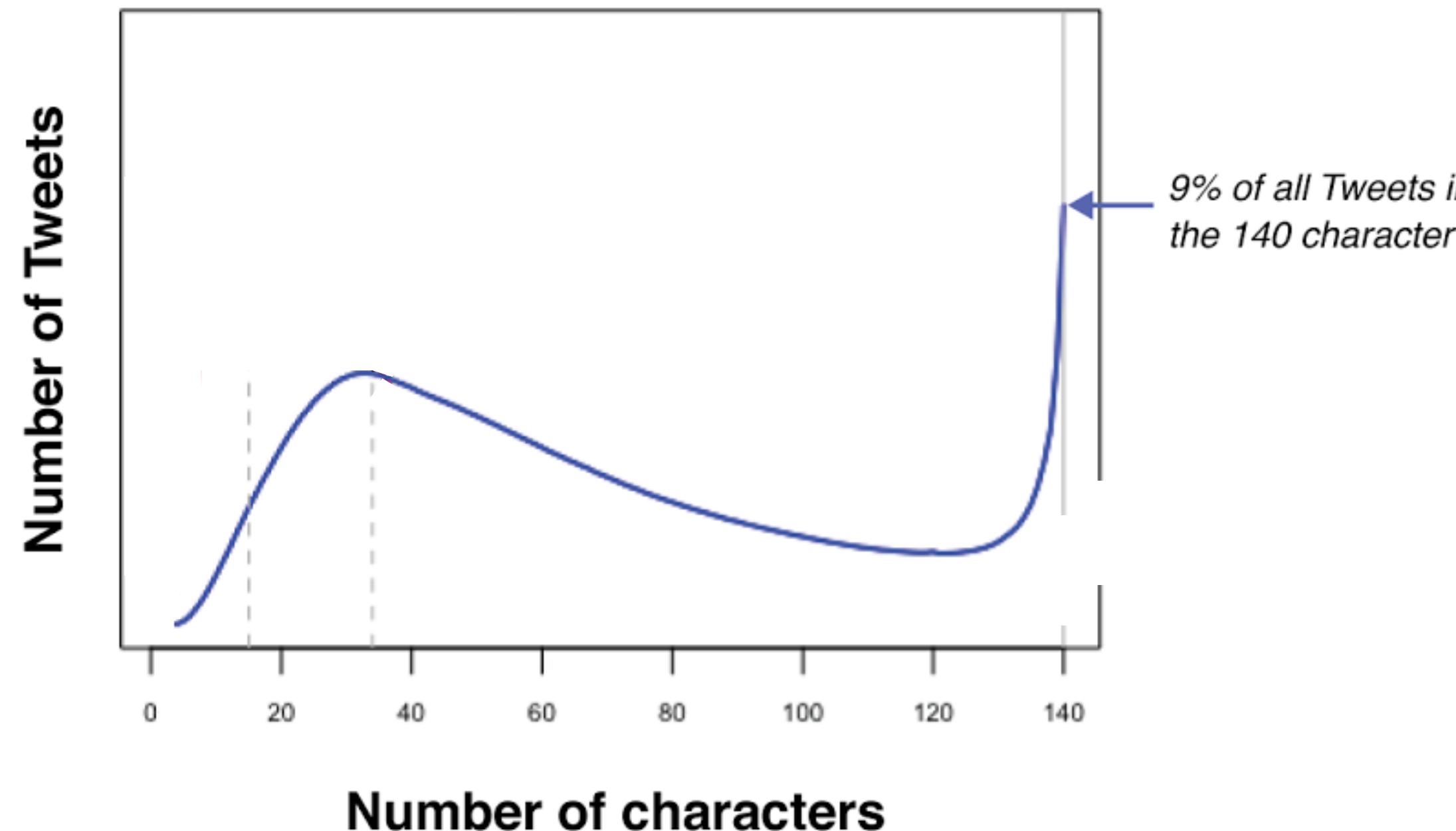
Question: Is 140 characters too short?



# A data science question

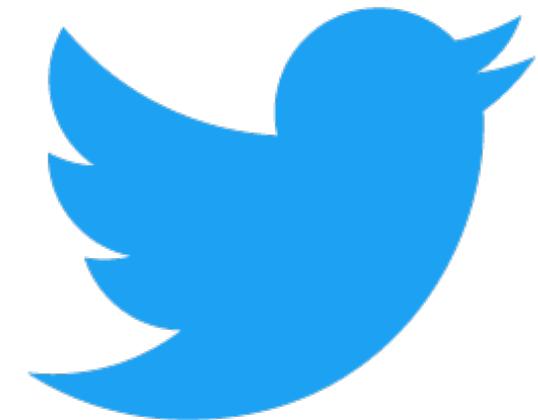


Question: Is 140 characters too short?

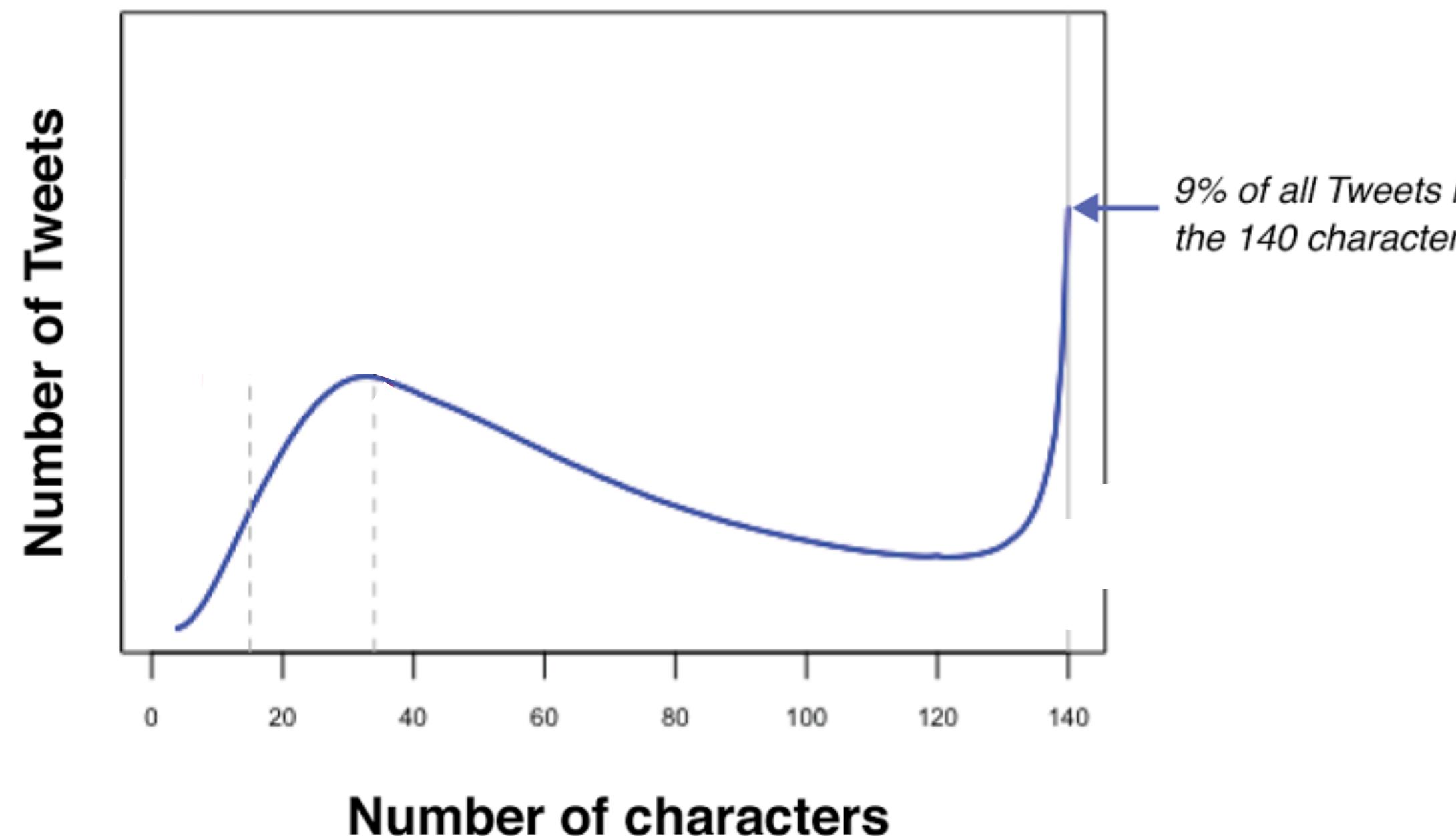


Answer: Yes, probably too short

# A data science question



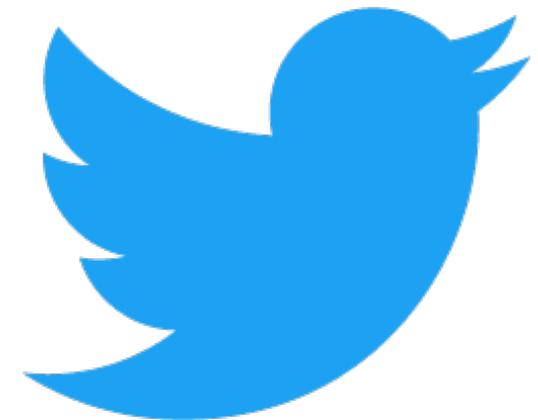
Question: Is 140 characters too short?



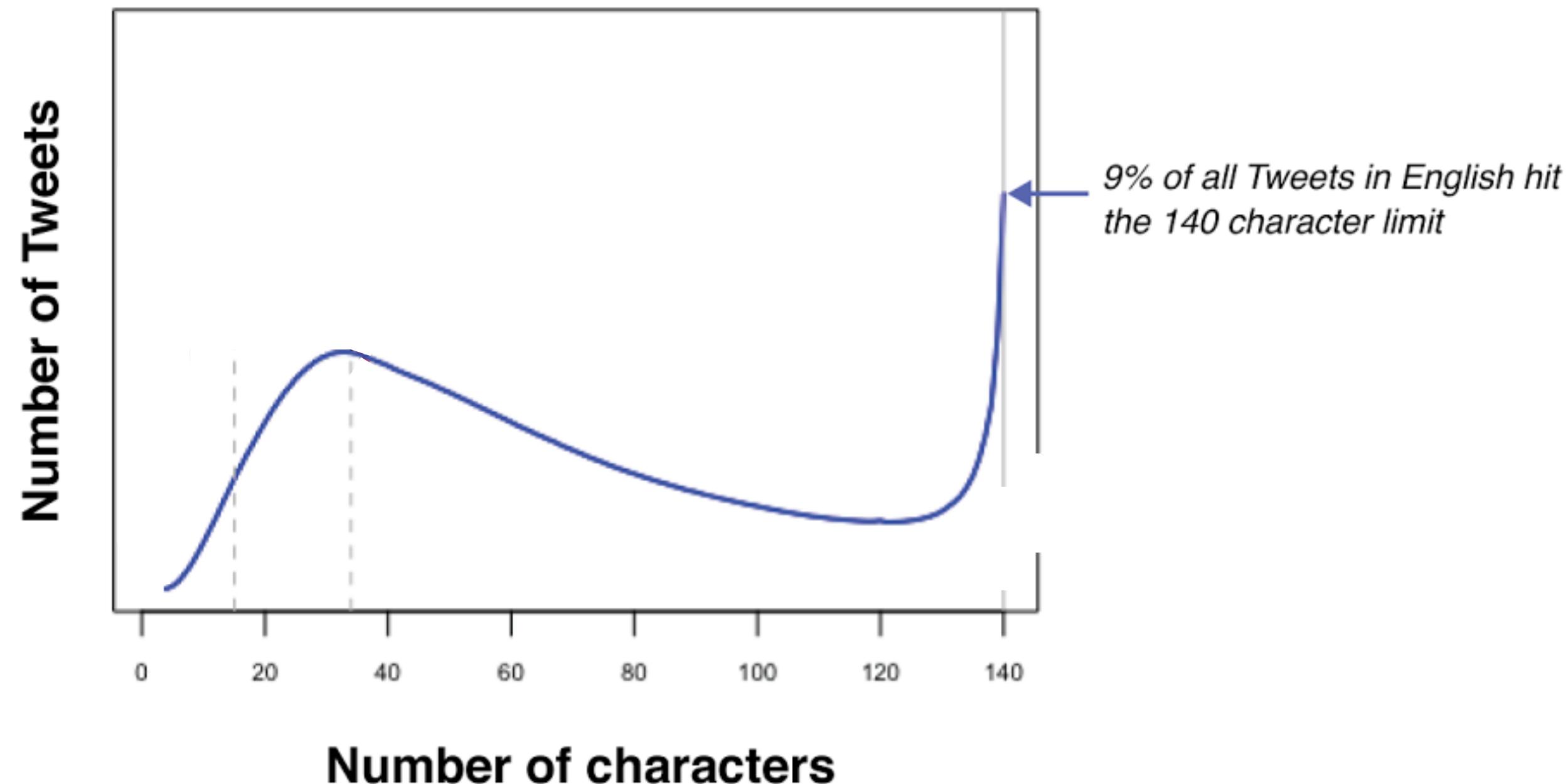
Answer: Yes, probably too short

New question: if the limit is raised, will people just cram into that new length? What should the limit be?

# A data science question

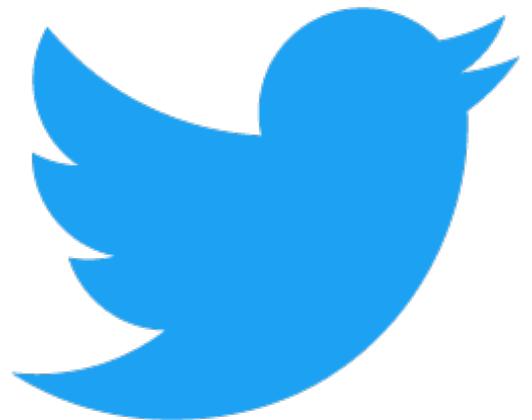


New question: if the limit is raised, will people just cram into that new length? What should the limit be?

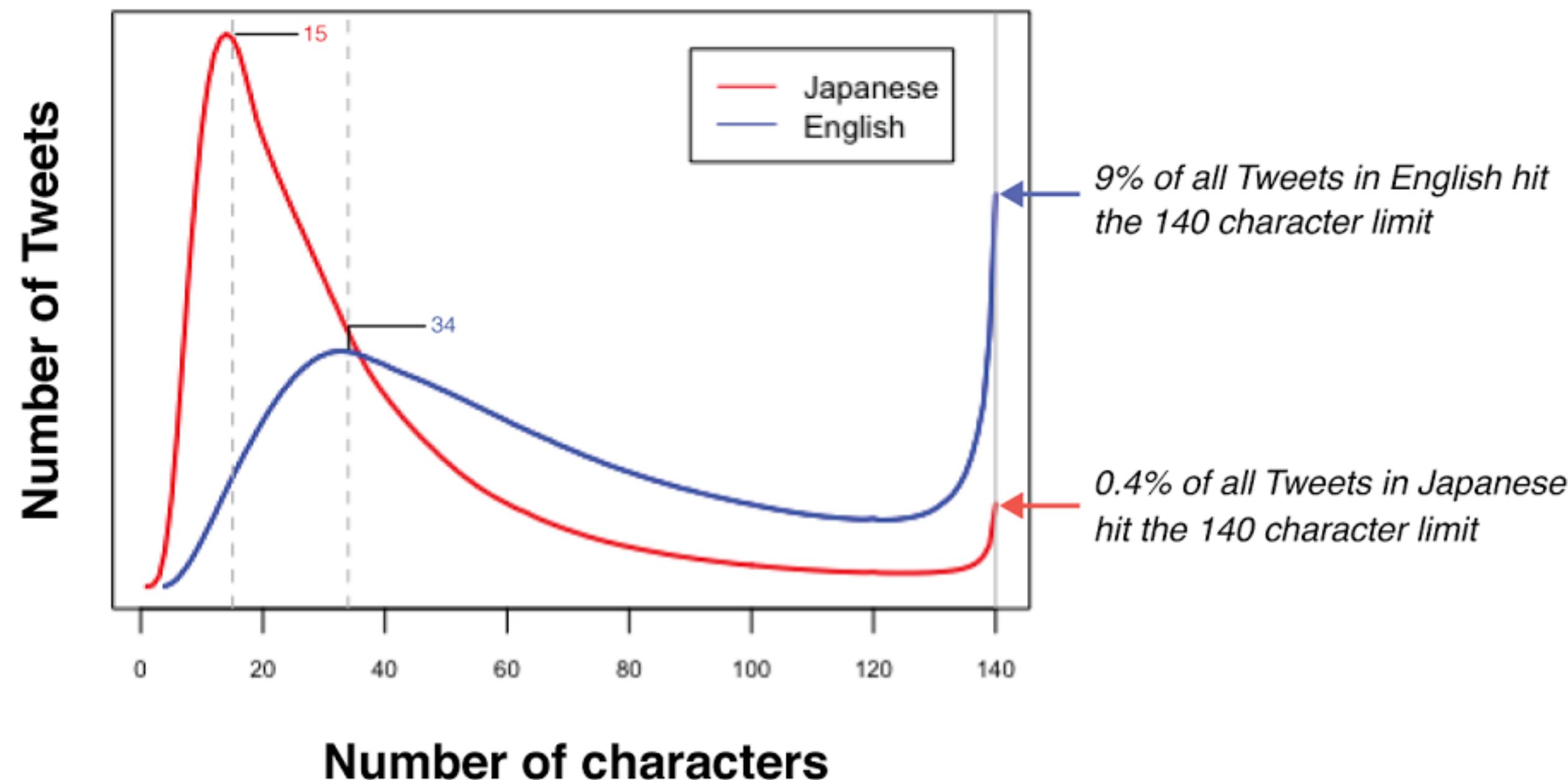


Nice "natural experiment"  
built into Twitter

# A data science question

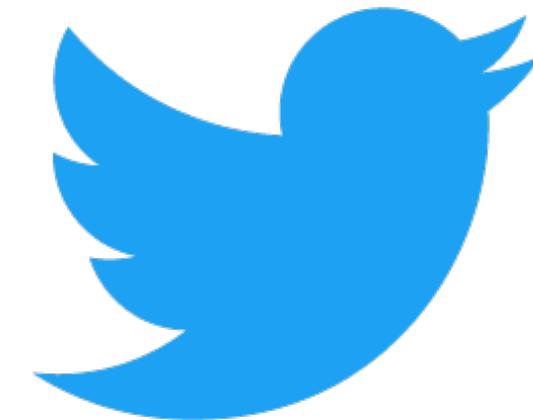


New question: if the limit is raised, will people just cram into that new length? What should the limit be?



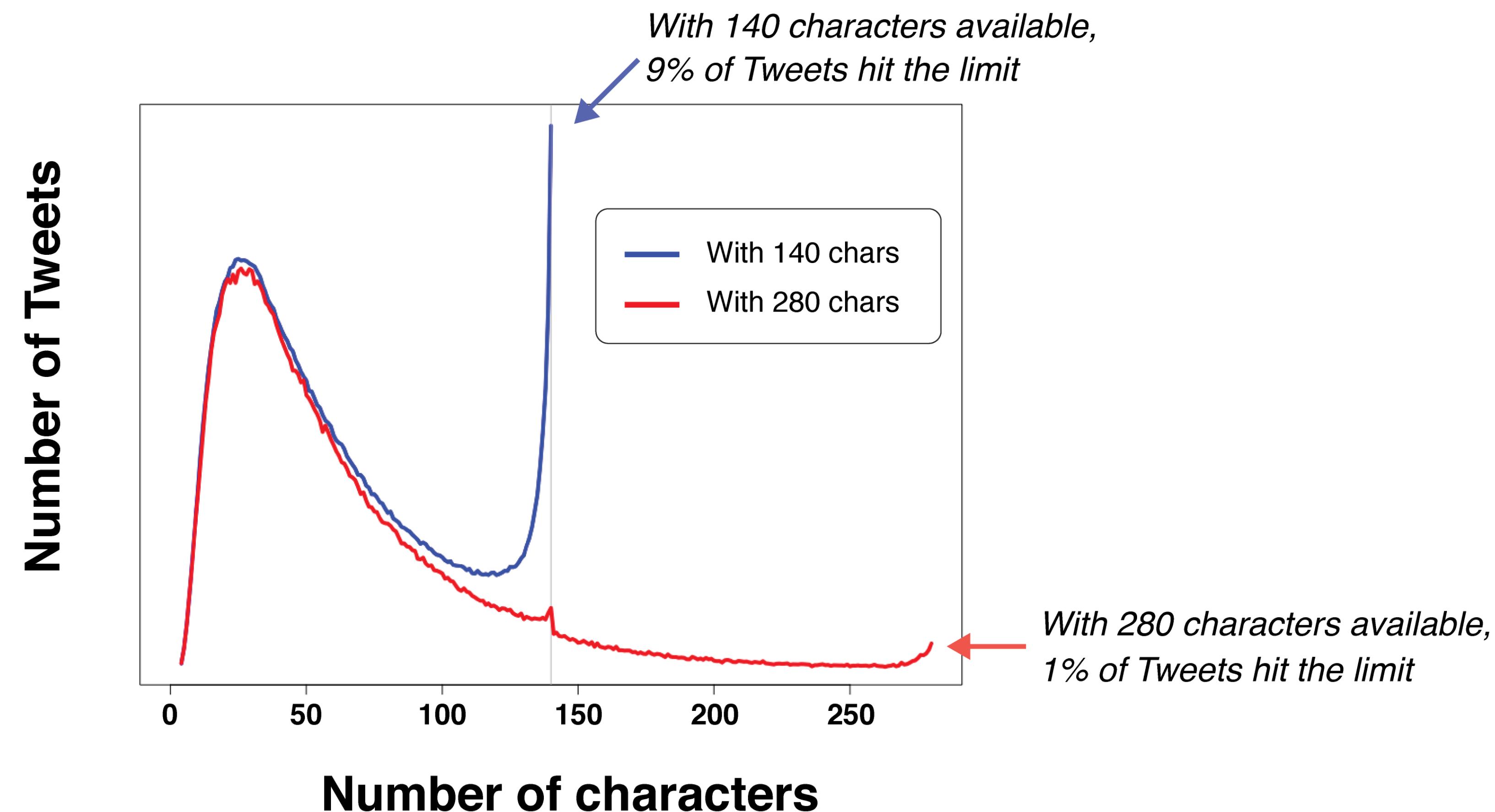
Nice "natural experiment"  
built into Twitter

# A data science question



Twitter decided to double the limit to 280

The change was rolled out in 2017



# A data science question

Notice the process

1. Ask question: *140 limit correct?*
  - i. Why ask this question? Is it important? History/context?
  - ii. Can we answer this? With data?
2. Hypothesize about what the data will tell us
3. Test the hypotheses
4. Insight → *English-vs-Japanese*
5. Decision: *double limit to 280*

# Low-tech or High-tech?

# Low-tech or High-tech?

Neural language representations predict outcomes of scientific research

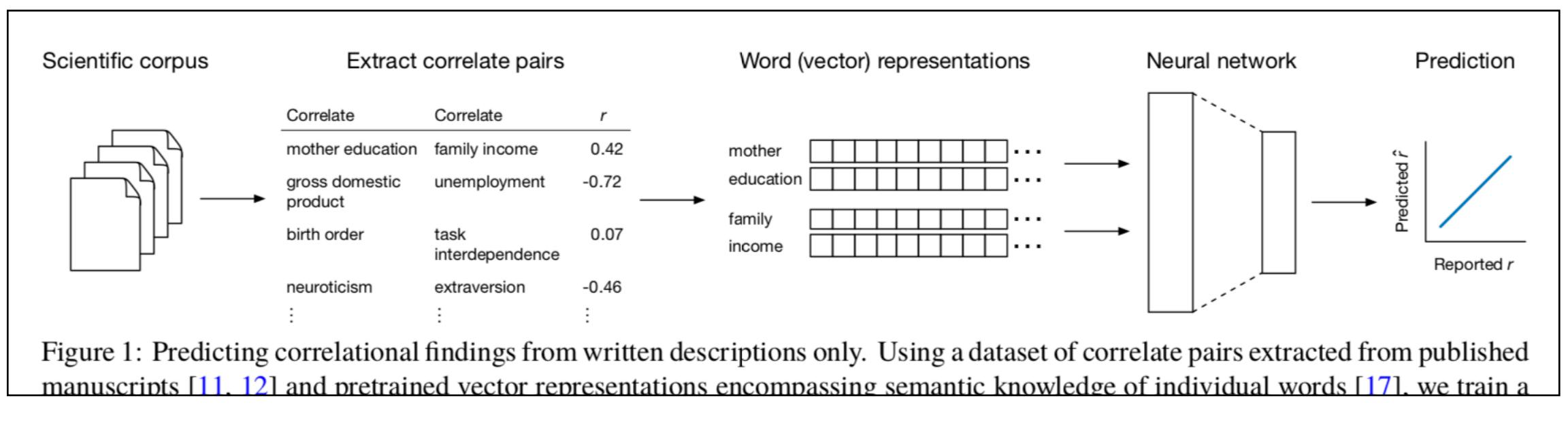
James P. Bagrow<sup>1,2,\*</sup>, Daniel Berenberg<sup>3,2</sup>, and Joshua Bongard<sup>3,2</sup>

<sup>1</sup>Department of Mathematics & Statistics, University of Vermont, Burlington, VT, United States

<sup>2</sup>Vermont Complex Systems Center, University of Vermont, Burlington, VT, United States

<sup>3</sup>Department of Computer Science, University of Vermont, Burlington, VT, United States

\*Corresponding author. Email: [james.bagrow@uvm.edu](mailto:james.bagrow@uvm.edu), Homepage: [bagrow.com](http://bagrow.com)



# Low-tech or High-tech?

Google Forms

## Neural language representations predict outcomes of scientific research

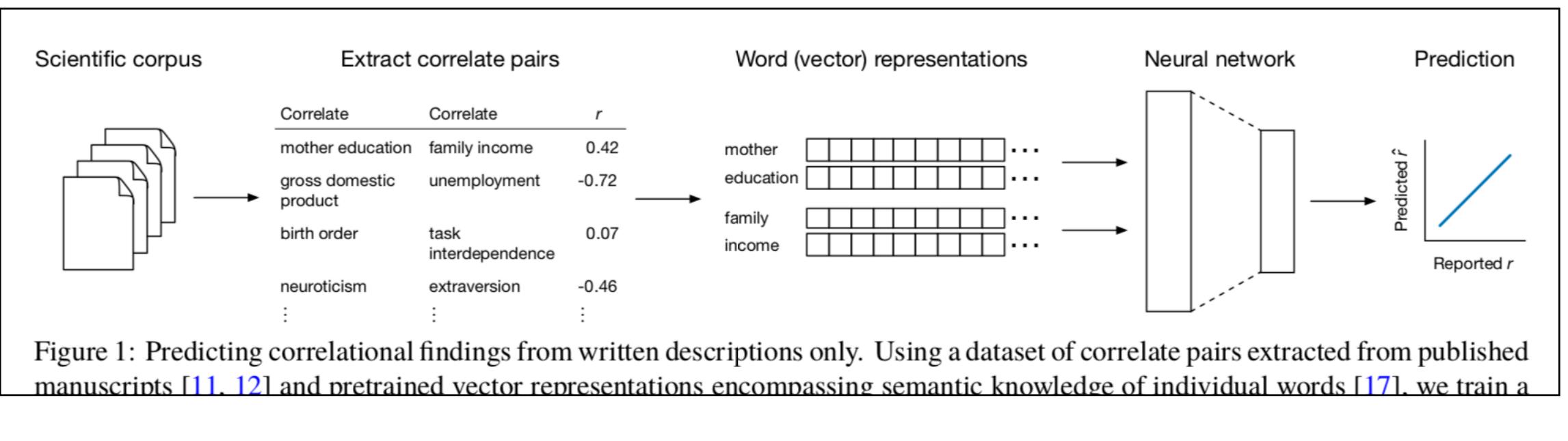
James P. Bagrow<sup>1,2,\*</sup>, Daniel Berenberg<sup>3,2</sup>, and Joshua Bongard<sup>3,2</sup>

<sup>1</sup>Department of Mathematics & Statistics, University of Vermont, Burlington, VT, United States

<sup>2</sup>Vermont Complex Systems Center, University of Vermont, Burlington, VT, United States

<sup>3</sup>Department of Computer Science, University of Vermont, Burlington, VT, United States

\*Corresponding author. Email: [james.bagrow@uvm.edu](mailto:james.bagrow@uvm.edu), Homepage: [bagrow.com](http://bagrow.com)



# Low-tech or High-tech?

Google Forms

Neural language representations predict outcomes of scientific research

James P. Bagrow<sup>1,2,\*</sup>, Daniel Berenberg<sup>3,2</sup>, and Joshua Bongard<sup>3,2</sup>

<sup>1</sup>Department of Mathematics & Statistics, University of Vermont, Burlington, VT, United States

<sup>2</sup>Vermont Complex Systems Center, University of Vermont, Burlington, VT, United States

<sup>3</sup>Department of Computer Science, University of Vermont, Burlington, VT, United States

\*Corresponding author. Email: [james.bagrow@uvm.edu](mailto:james.bagrow@uvm.edu), Homepage: [bagrow.com](http://bagrow.com)

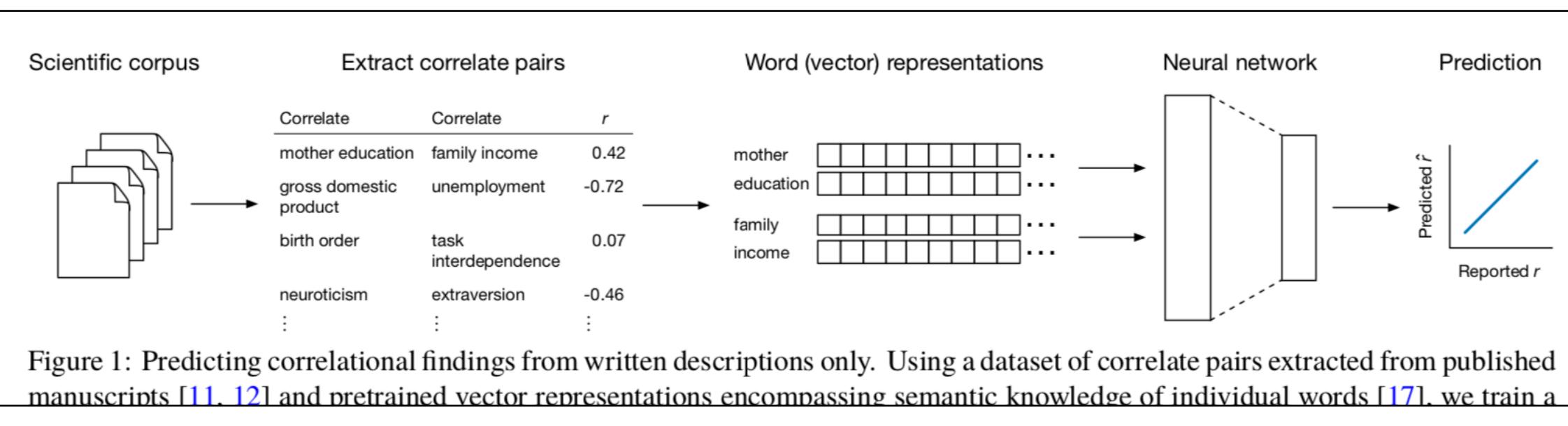


Figure 1: Predicting correlational findings from written descriptions only. Using a dataset of correlate pairs extracted from published manuscripts [11, 12] and pretrained vector representations encompassing semantic knowledge of individual words [17], we train a

A screenshot of a Google Form titled "Amazing survey". The form has a purple header with "QUESTIONS" and "RESPONSES" tabs. Below the title, there is a text input field with the placeholder "Enter to win a free car". On the right side, there is a vertical toolbar with icons for adding questions, opening responses, and other form settings.

# Low-tech or High-tech?

QUESTIONS    RESPONSES

## Amazing survey

Enter to win a free car

Do you have a driver license? \*

Yes  
 No

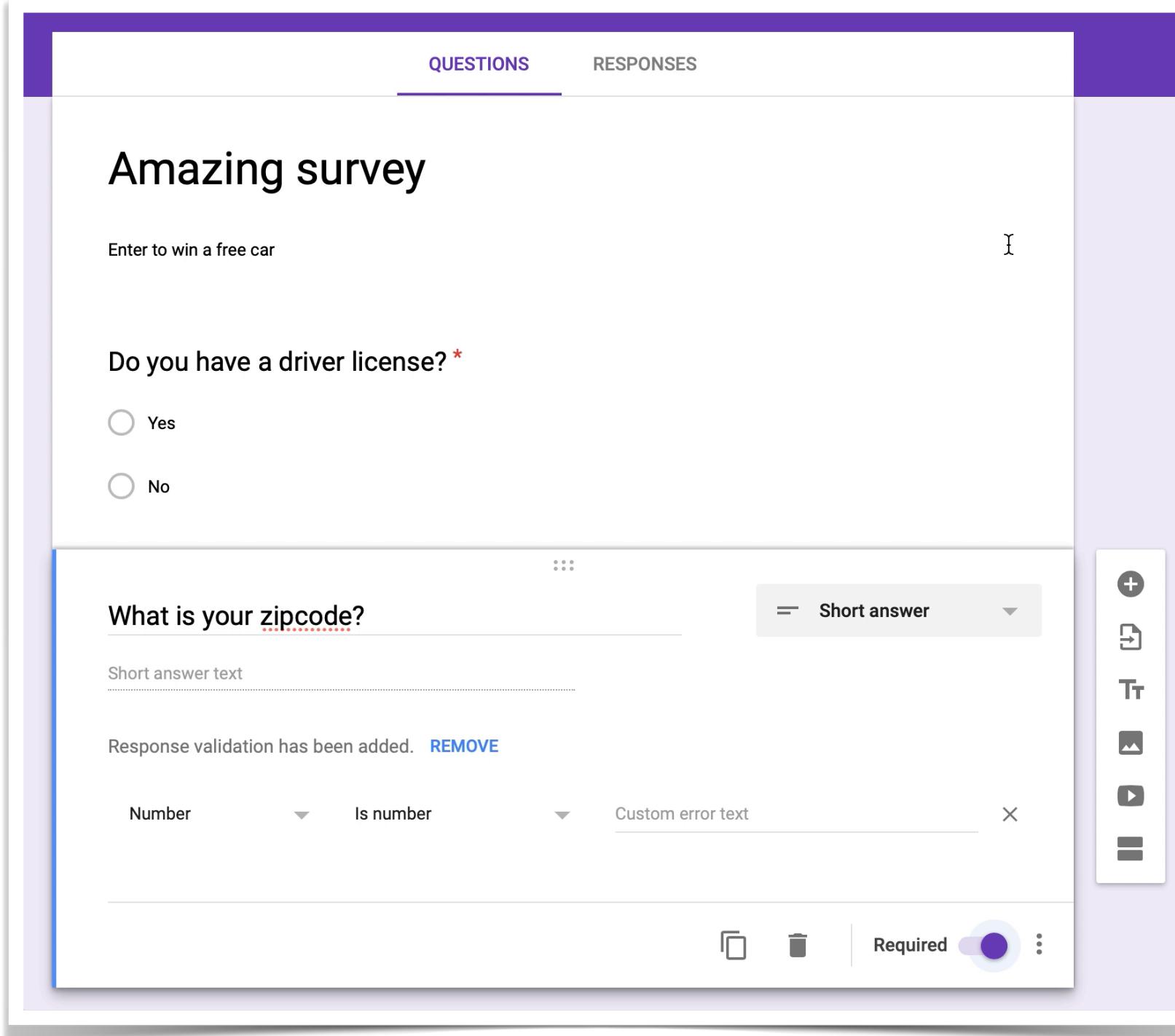
What is your zipcode?

Short answer text

Response validation has been added. [REMOVE](#)

Number    Is number    Custom error text

Required



## Cloud AutoML BETA

Train high-quality custom machine learning models with minimal effort and machine learning expertise.

[Try AutoML](#) ▼ [View documentation](#)

[cloud.google.com/automl/](http://cloud.google.com/automl/)



## Democratizing AI: Non-expert design of prediction tasks

James P. Bagrow<sup>1,2</sup>

<sup>1</sup>Department of Mathematics & Statistics, University of Vermont

<sup>2</sup>Vermont Complex Systems Center

Corresponding author:  
James P. Bagrow

Email address: [james.bagrow@uvm.edu](mailto:james.bagrow@uvm.edu)

[bagrow.com/#pub50](http://bagrow.com/#pub50)

# Low-tech or High-tech?

QUESTIONS    RESPONSES

## Amazing survey

Enter to win a free car

Do you have a driver license? \*

Yes  
 No

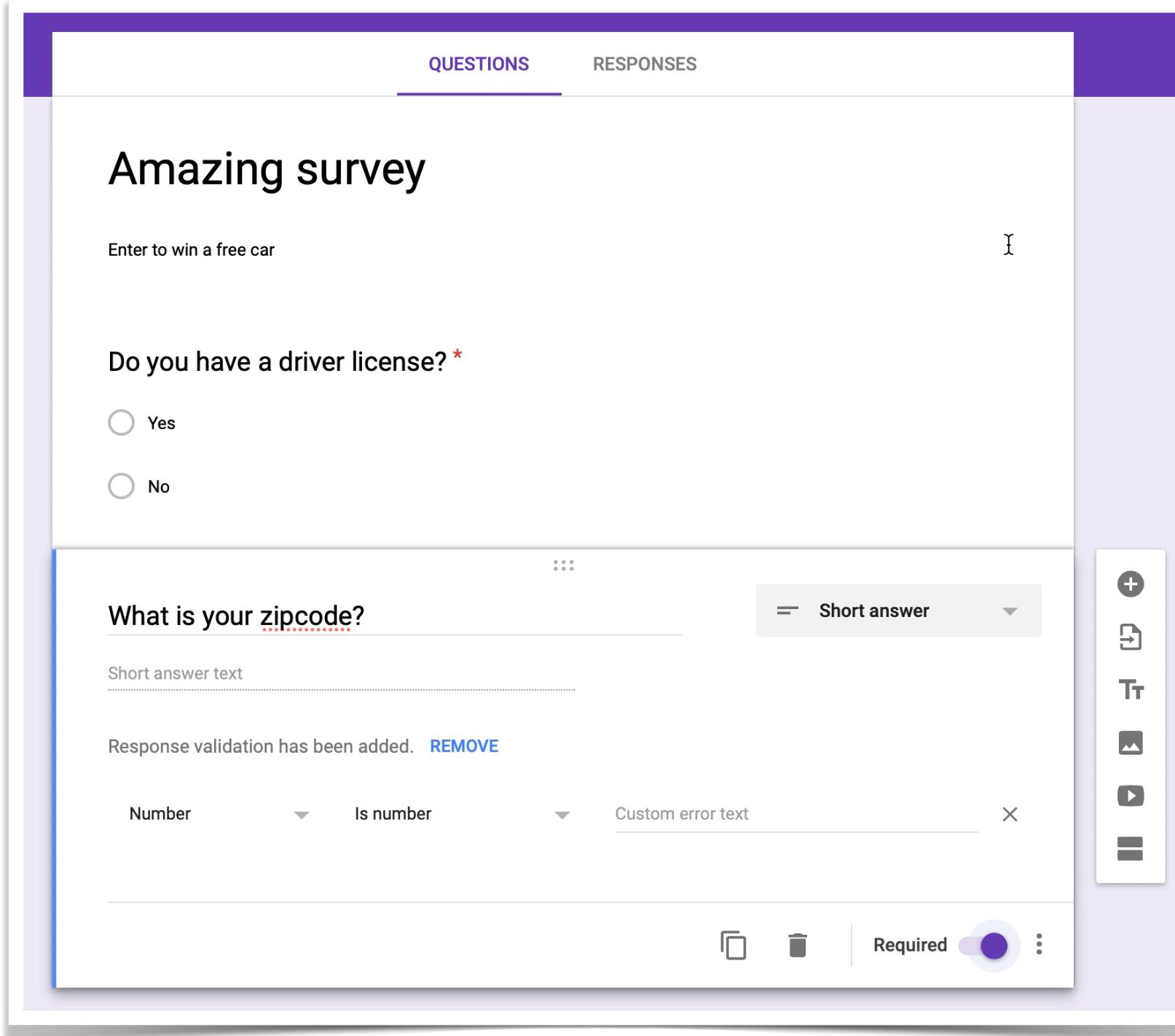
What is your zipcode?

Short answer text

Response validation has been added. [REMOVE](#)

Number    Is number    Custom error text

Required



## Cloud AutoML BETA

Train high-quality custom machine learning models with minimal effort and machine learning expertise.

[Try AutoML](#) ▼ [View documentation](#)

[cloud.google.com/automl/](http://cloud.google.com/automl/)



## Democratizing AI: Non-expert design of prediction tasks

James P. Bagrow<sup>1,2</sup>

<sup>1</sup>Department of Mathematics & Statistics, University of Vermont

<sup>2</sup>Vermont Complex Systems Center

Corresponding author:  
James P. Bagrow

Email address: [james.bagrow@uvm.edu](mailto:james.bagrow@uvm.edu)

[bagrow.com/#pub50](http://bagrow.com/#pub50)

# Today's Schedule

- ✓ Course logistics
- ✓ Intro + Motivation

Intro/Review of programming

# Intro/Review of programming



This is a programming-intensive  
class. The language will be **Python**.

# Why Python?

Why not: R      Perl      STATA  
MATLAB      Julia      :  
IDL      SAS      ?

## Answers:

1. Popularity
2. Personal Preference
3. Plasticity

# Why Python?

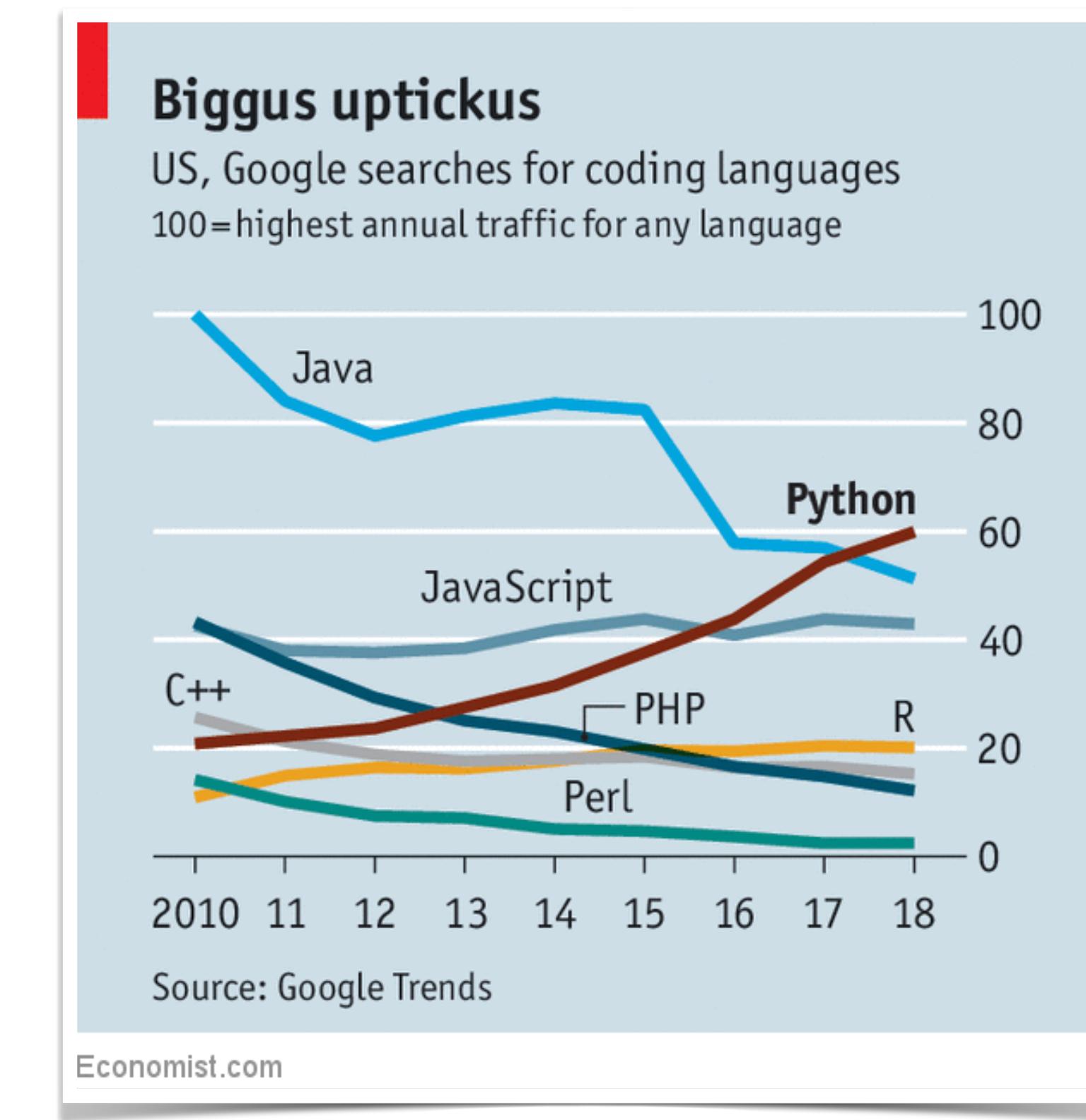
Answers:

1. Popularity
2. Personal Preference
3. Plasticity

Why not: R  
MATLAB  
IDL

Perl  
Julia  
SAS

STATA  
⋮  
?

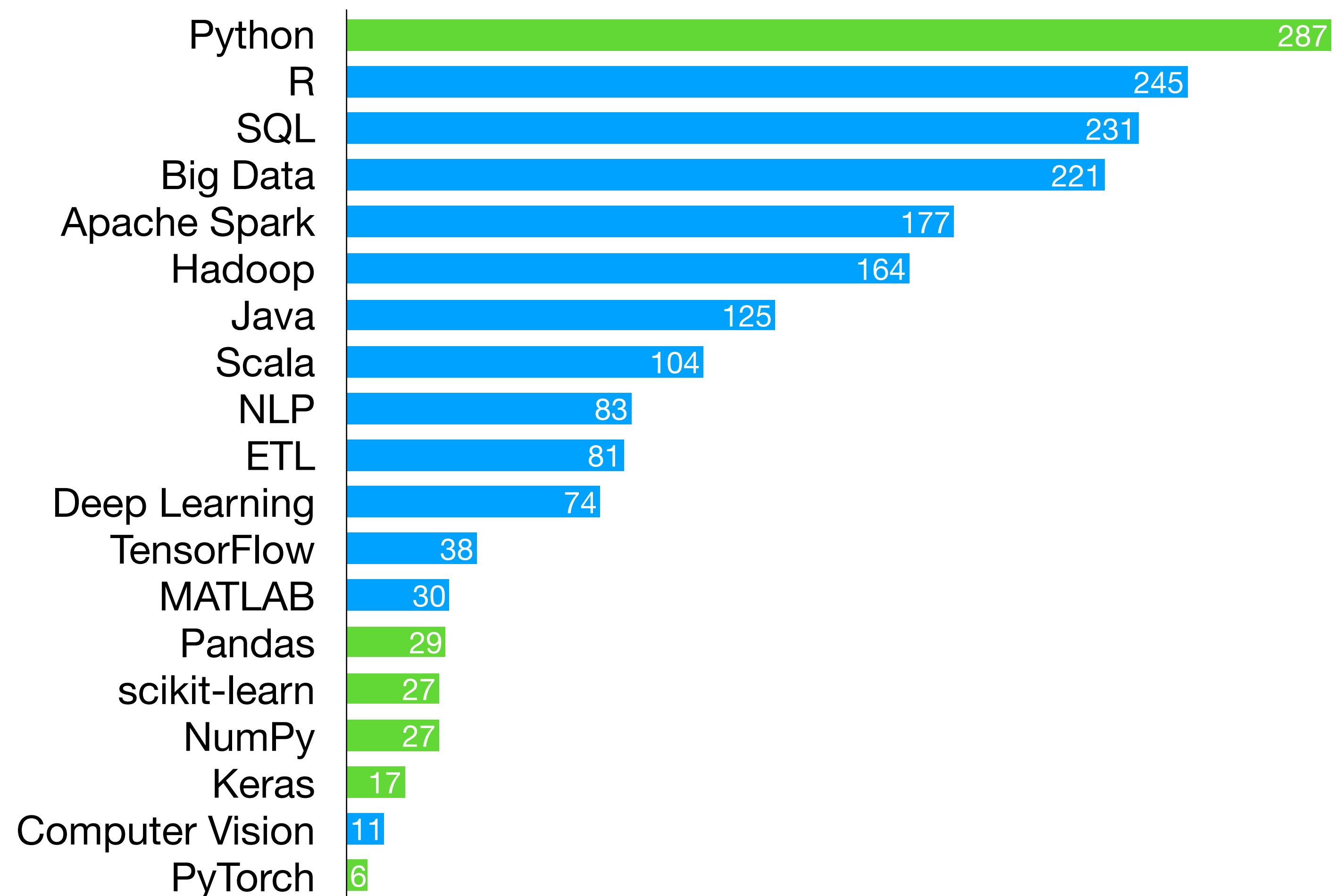


# Why Python?

Answers:

1. Popularity
2. Personal Preference
3. Plasticity

Skills mentioned in 300 data scientist job postings (June 2019)



# Programming review

```
print("Hello world")
```

This is a programming-intensive class. The language will be Python.



# Programming review

## Data types

4      3.45

“a”    True

## Data structures

vectors, arrays, hashes,  
matrices, sets, etc.

## Code

keywords    for

+   -   %   #       while

## Code structures

functions

libraries  
modules

# Programming review

## Data types

4      3.45

“a”      True

## Data structures

vectors, arrays, hashes,  
matrices, sets, etc.

## Code

keywords      for

+    -    %    #      while

## Code structures

functions

libraries

modules

# Programming review

## Data types

4      3.45  
“a”    True

## Data structures

vectors, arrays, hashes,  
matrices, sets, etc.

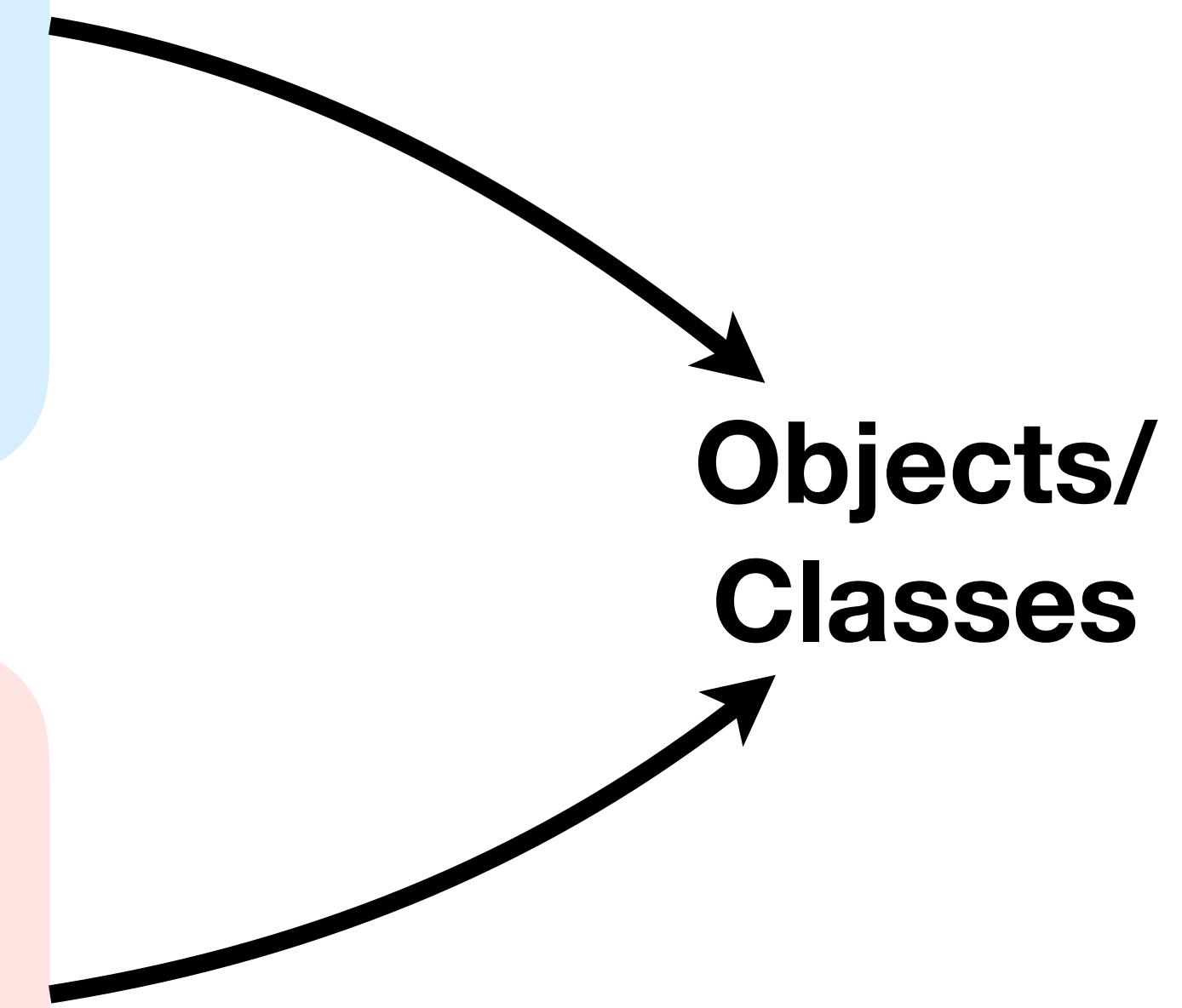
## Code

keywords    for  
+   -   %   #      while

## Code structures

functions  
libraries  
modules

## Objects/ Classes



# Use Python interactively

```
Python 3.10.2 |Anaconda custom (x86_64)| (default, Aug 31 2021, 17:52:12)
[GCC 4.2.1 Compatible Apple LLVM 5.0 (clang-500.2.79)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

# Use Python interactively

```
Python 3.10.2 |Anaconda custom (x86_64)| (default, Aug 31 2021, 17:52:12)  
[GCC 4.2.1 Compatible Apple LLVM 5.0 (clang-500.2.79)] on darwin  
Type "help", "copyright", "credits" or "license" for more information.
```

```
>>>
```



Prompt

# Use Python interactively

```
Python 3.10.2 |Anaconda custom (x86_64)| (default, Aug 31 2021, 17:52:12)
[GCC 4.2.1 Compatible Apple LLVM 5.0 (clang-500.2.79)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

# Use Python interactively

```
Python 3.10.2 |Anaconda custom (x86_64)| (default, Aug 31 2021, 17:52:12)
[GCC 4.2.1 Compatible Apple LLVM 5.0 (clang-500.2.79)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>> print("hello world")
```

# Use Python interactively

```
Python 3.10.2 |Anaconda custom (x86_64)| (default, Aug 31 2021, 17:52:12)
[GCC 4.2.1 Compatible Apple LLVM 5.0 (clang-500.2.79)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>> print("hello world")
hello world
```

# Use Python interactively

```
Python 3.10.2 |Anaconda custom (x86_64)| (default, Aug 31 2021, 17:52:12)
[GCC 4.2.1 Compatible Apple LLVM 5.0 (clang-500.2.79)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>> print("hello world")
hello world
```

As part of your homework and setup,  
you will learn about a much more  
advanced Python prompt called **IPython**

# Use Python files

awesome\_science.py:



```
import sys # not used  
  
def product(x,y):  
    return x*y  
  
print("x =", product(5,7))
```

x = 35

# Simple syntax

```
x = 34          # a comment
y = "Hello"
z = 3.45

if z == 3.45 or y == "Hello":
    x = x+1
    y = y + " World" # String concat.

print(x)
print(y)
```

# Data types

integers	x = 5
floats	pi = 3.14
strings	s = "hi\namy"
booleans	found = True

# Data structures

## Lists

```
nums = [4,7,12,9]
names = ["John", "Alice", "Bob"]
[5, [1,2,3], "nested and mixed"]
```

# Data structures

Example:  
Using lists

```
>>> nums = [4,7,12,9]
>>> print(nums[0], nums[-1])
4 9
>>> print(nums[:3])
[4,7,12]
>>> print(nums[2:])
[12,9]
>>> nums.append(15)
>>> print(nums)
[4,7,12,9,15]
```

index access

slicing

modifying

# Data structures

OK, but so what?

- Suppose I want to study a social network
- For each person, I have a list of their friends:

# Data structures

OK, but so what?

- Suppose I want to study a social network
- For each person, I have a list of their friends:

$$\mathbf{LF} = [\mathbf{F}_1, \mathbf{F}_2, \dots]$$

# Data structures

OK, but so what?

- Suppose I want to study a social network
- For each person, I have a list of their friends:

$$\mathbf{LF} = [\mathbf{F}_1, \mathbf{F}_2, \dots]$$

```
>>> LF = [["Alice", "John", "Bob"],  
           ["John", "Bob", "Stacy"],...  
           ]  
>>> print(LF[1])  
["John", "Bob", "Stacy"]
```

# Data structures

OK, but so what?

- Suppose I want to study a social network
- For each person, I have a list of their friends:

$$\mathbf{LF} = [\mathbf{F}_1, \mathbf{F}_2, \dots]$$

```
>>> LF = [["Alice", "John", "Bob"],  
           ["John", "Bob", "Stacy"],...  
           ]  
>>> print(LF[1])  
["John", "Bob", "Stacy"]
```

But whose friends are these? Who is person "1"?

# Data structures

Dictionaries

A list **maps** the integers 0, 1, ... to some values.  
*Why limit ourselves to integers?*

# Data structures

## Dictionaries

A list **maps** the integers 0, 1, ... to some values.  
*Why limit ourselves to integers?*

```
>>> D = {"Gwen": ["Alice", "John", "Bob"],  
         "Peter": ["John", "Bob", "Stacy"]}
```

# Data structures

## Dictionaries

A list **maps** the integers 0, 1, ... to some values.  
*Why limit ourselves to integers?*

```
>>> D = {"Gwen": ["Alice", "John", "Bob"],  
           "Peter": ["John", "Bob", "Stacy"]}  
>>> print(D['Gwen'])  
['Alice', 'John', 'Bob']
```

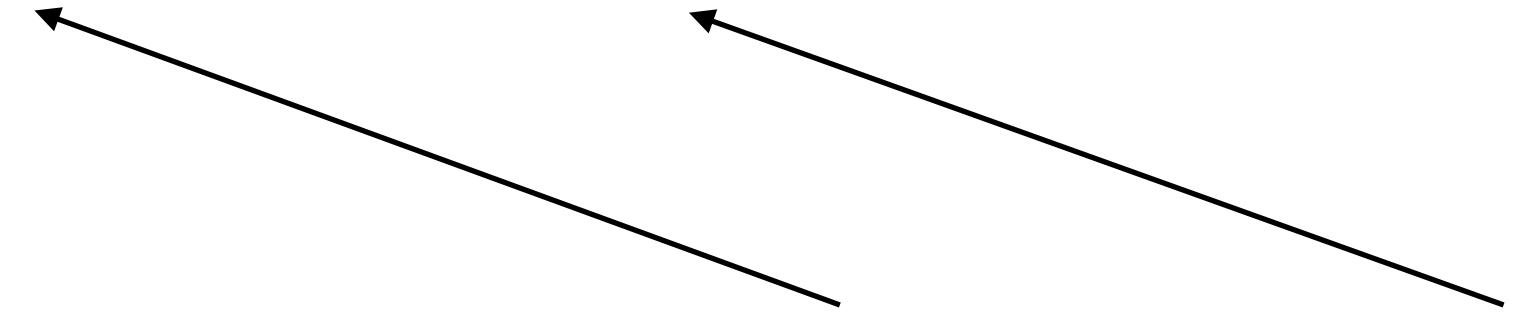
# Data structures

## Dictionaries

A list **maps** the integers 0, 1, ... to some values.  
*Why limit ourselves to integers?*

```
>>> D = {"Gwen": ["Alice", "John", "Bob"],  
         "Peter": ["John", "Bob", "Stacy"]}
```

Generalization: maps **KEYS** to **VALUES**



# Data structures

Dictionaries

A list **maps** the integers 0, 1, ... to some values.  
*Why limit ourselves to integers?*

```
>>> user2friends = {"Gwen": ["Alice", "John", "Bob"],  
                    "Peter": ["John", "Bob", "Stacy"]}  
>>> print(user2friends['Gwen'])  
['Alice', 'John', 'Bob']
```

# Data structures

Dictionaries

A list **maps** the integers 0, 1, ... to some values.  
*Why limit ourselves to integers?*

```
>>> user2friends = {"Gwen": ["Alice", "John", "Bob"],  
                     "Peter": ["John", "Bob", "Stacy"]}  
>>> print(user2friends['Gwen'])  
['Alice', 'John', 'Bob']
```



Less cryptic with  
a good name!

```
>>> print(D['Gwen'])
```

# Data structures

## Dictionaries

A list **maps** the integers 0, 1, ... to some values.  
*Why limit ourselves to integers?*

```
>>> {"Gwen" : ["Alice", "John", "Bob"],  
      "Peter": ["John", "Bob", "Stacy"]}
```

Keys must be **unique**  
and **immutable**

Finding a key in a dict is **very fast**  
(make great lookup tables)

# Data structures

Dictionaries  
cont.

```
>>> D = {} # empty
>>> D[5] = 9
>>> D[3] = 12
>>> print(D)
{3: 12, 5: 9}
>>> D[5] = 8
>>> del D[3]
>>> print(D)
{5: 8}
```

dicts are unordered

replaced a value

deleted a (key,value)

# Data structures

Dictionaries  
cont.

```
>>> D = {} # empty
>>> D[5] = 9
>>> D[3] = 12
>>> print(D)
{3: 12, 5: 9}
```

dicts are unordered

More data structures to come

```
>>> print(D)
{5: 9, 3: 12}
```

# Control flow | for loops

A Python for-loop is like a "for each" loop:

The looping variable (*x*) becomes *each* item in the list

```
L = ["Alice", "Bob"]  
for x in L:  
    print("x =", x)  
...
```

Also: white space is  
**significant!**

# Control flow | for loops

A Python for-loop is like a "for each" loop:

The looping variable (*x*) becomes *each* item in the list

```
L = ["Alice", "Bob"]  
for x in L:  
    print("x =", x)
```

Output:

```
x = Alice  
x = Bob
```

# Control flow | for loops

A Python for-loop is like a "for each" loop:

The looping variable (*x*) becomes *each* item in the list

Compare to a C-style for loop:

```
L = ["Alice", "Bob"]  
for x in L:  
    print("x =", x)
```

```
for (i=0; i<2; i++) {  
    printf("%s\n", L[i]);  
}
```

# Control flow | for loops

```
L = ["Alice", "Bob"]
```

compare: `for x in L:  
 print("x =", x)`

versus: `for i in range(len(L)):  
 print("x =", L[i])`

# Control flow | for loops

```
L = ["Alice", "Bob"]
```

✓ compare: `for x in L:  
 print("x =", x)`

✗ versus: `for i in range(len(L)):  
 print("x =", L[i])`

Becoming *comfortable* with "for each" loops:  
→ less and more readable code!  
→ fewer mistakes!

# Control flow | for loops

One temptation is multiple  
"iterables":

```
X = [...]
Y = [...]
for i in range(len(X)):
    print("s =", X[i]+Y[i])
```

# Control flow | for loops

One temptation is multiple "iterables":

```
X = [...]  
Y = [...]  
for i in range(len(X)):  
    print("s =", X[i]+Y[i])
```



Instead, embrace new tools.  
For example, use `zip`:

```
for x,y in zip(X,Y):  
    print("s =", x+y)
```



More "Pythonic" code

# Control flow

## If statements

```
x = 5  
y = "yes"  
z = True  
p = 0.333
```

```
if x < 10 and y == "yes" and not z:  
    print("won't happen")  
print("if statement over")
```

```
if 0 < p < 1:  
    print("yep")
```

← **nice!** same as:  
 $p > 0$  and  $p < 1$

# Control flow

## While loops

```
x = 0
while True:
    x += 1
    if x >= 100:
        break
print(x) →
x = 0
while x < 100:
    x += 1
    print(x)
```

( $x += 1$  is a **shortcut** for  $x = x + 1$ )

# Functions

```
def my_function(x, y):
    """This is the docstring.

    This is my function. There are many
    like it but this one is mine.

    """
# The code would go here..
pass

def second_func(username, password=None):
    if password is not None:
        print("found a password")
```

Colors!

(syntax highlighting)

# Functions

```
def my_function(x, y):  
    """This is the docstring.
```

This is my function. There are many like it but this one is mine.

```
"""
```

```
# The code would go here..
```

```
pass
```

```
def second_func(username, password=None):  
    if password is not None:  
        print("found a password")
```

Colors!

(syntax highlighting)

default value for optional argument

# "Batteries included"

Want to write a [web crawler](#)? No problem!

To download web pages:

# "Batteries included"

Want to write a [web crawler](#)? No problem!

To download web pages:

```
>>> import urllib  
>>> web = urllib.request.urlopen("http://bagrow.com/ds1")  
>>> text = web.read() # read?  
>>> print(text[2800:3400])
```

# "Batteries included"

Want to write a [web crawler](#)? No problem!

To download web pages:

```
>>> import urllib  
>>> web = urllib.request.urlopen("http://bagrow.com/ds1")  
>>> text = web.read() # read?  
>>> print(text[2800:3400])
```

b'data,\nso it is more important than ever to understand how to collect, process, and\nanalyze these data. A picture is worth a thousand words, so visualizations,\nfrom scientific plots and infographics to interactive data explorers, are\n crucial to summarize and communicate new discoveries.\n</p>\n\n</div><!--row-->\n\n<div class="row">\n <h3>Resources</h3>\n <div class="col-md-7">\n <ul>\n <li>The <b><a href="syllabus\_stat287\_fall2016.pdf">course syllabus</a></b>. Be sure to check this out. </li>\n <li>The course introductory reading, <a href="whirlwindtourpython/00-Title.html">A W'

# "Batteries included"

Want to write a [web crawler](#)? No problem!

Programmatic power:

```
import urllib  
  
website = "http://uvm.edu/"  
names = ["alice", "bob", "john"]  
  
for name in names:  
    url = website + name + ".html"  
    print(url)
```

<http://uvm.edu/alice.html>  
<http://uvm.edu/bob.html>  
<http://uvm.edu/john.html>

# "Batteries included"

Want to write a [web crawler](#)? No problem!

Programmatic power:

```
import urllib  
  
names = ["alice", "bob", "john"]  
  
for name in names:  
    url = "http://uvm.edu/{}.html".format(name)
```

Handy reference:  
<https://pyformat.info>

# "Batteries included"

Want to write a [web crawler](#)? No problem!

Programmatic power:

```
import urllib  
  
names = ["alice", "bob", "john"]  
  
for name in names:  
    url = "http://uvm.edu/{}.html".format(name)
```



string substitution

[Better than using + to build strings](#)

Handy reference:  
<https://pyformat.info>

\* Even cooler: [f-strings](#)! → `url = f"http://uvm.edu/{name}.html"`

# Learning resources

## A Whirlwind Tour of Python for STAT/CS 287

James Bagrow

This mini-text serves as introductory reading for my course, [Data Science I](#), held at the University of Vermont. It has been very lightly adapted and specialized from the original "A Whirlwind Tour of Python" by [Jake VanderPlas](#), and he deserves most of the credit for this work.

*A Whirlwind Tour of Python* is a fast-paced introduction to essential components of the [Python language](#) for researchers and developers who are already familiar with programming in another language.

The material is particularly aimed at those who wish to use Python for data science and/or scientific programming, and in this capacity serves as an introduction to Jake VanderPlas' book, [The Python Data Science Handbook](#).

### Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	About the authors . . . . .	4
1.2	Installation and Practical Considerations . . . . .	4
<b>2</b>	<b>How to Run Python Code</b>	<b>5</b>
<b>3</b>	<b>A Quick Tour of Python Language Syntax</b>	<b>7</b>
3.1	Comments Are Marked by # . . . . .	8
3.2	End-of-Line Terminates a Statement . . . . .	8
3.3	Semicolons Can Optionally Terminate a Statement . . . . .	9
3.4	Indentation: Whitespace Matters! . . . . .	9
3.5	Whitespace Within Lines Does Not Matter . . . . .	10

## DS1 Lecture 01 supplement

Jim Bagrow

This [notebook](#) acts as an additional supplement on programming in Python. You are expected to learn this material during the lecture. A quiz of the semester will assess this.

This document complements the course's "[Whirlwind Tour](#)" reading. You are also expected to complete this reading as part of the supplement after reading the whirlwind tour.

### About this document

- This is a [Jupyter Notebook](#). I'll be using these for code-heavy lectures and lecture supplements. Notebooks contain code in a document, interleaved with narrative (non-code) text.
- Your [projects and assignments](#) will be submitted as Python scripts (`.py` files).

### Contents:

- A bit more on data structures and control flow (see also the Whirlwind Tour)
  - tuples vs. lists
    - multiassignment
    - sorting
    - zipping
- Working with [files](#)!
- Exploring the [Python standard library](#) (builtin modules and packages)
- A tour of useful third-party packages
  - [numpy and scipy](#)
  - [matplotlib](#)
- The [random](#) library
- Conclusion

# Computer setup

We're using a scientific Python environment called **Anaconda**



Supports Windows/Mac/Linux

Free for academic use

Includes **package manager**, text editor, and  
interactive prompt

<https://www.anaconda.com/download/>

# Today's Schedule

- ✓ Intro + Motivation
- ✓ Course logistics
- ✓ Intro/Review of programming

**Next time:** Intro/Review of prob/stat

Homework →

# Homework (1 of 2)

1. Register and install Anaconda Python on your machine (see Python Setup slides)
  - Complete **HW01**
  - HW01 includes additional setup reading
2. Work through **reading**—**there's a lot** but the more prepared you are now, the faster you can work later

(Quiz 01 will assess your preparation following this reading)

con't →

# Homework (2 of 2)

## Python setup

Jim Bagrow  
STAT/CS 287 — Data Science 1

## Preparing and submitting assignments

Jim Bagrow  
STAT/CS 287 — Data Science 1

# Homework (2 of 2)

## Python setup

Jim Bagrow  
STAT/CS 287 — Data Science 1

## Preparing and submitting assignments

Jim Bagrow  
STAT/CS 287 — Data Science 1

## A Whirlwind Tour of Python for STAT/CS 287

James Bagrow

This mini-text serves as introductory reading for my course, [Data Science I](#), held at the University of Vermont. It has been very lightly adapted and specialized from the original "A Whirlwind Tour of Python" by [Jake VanderPlas](#), and he deserves most of the credit for this work.

*A Whirlwind Tour of Python* is a fast-paced introduction to essential components of the [Python language](#) for researchers and developers who are already familiar with programming in another language.

The material is particularly aimed at those who wish to use Python for data science and/or scientific programming, and in this capacity serves as an introduction to Jake VanderPlas' book, *The Python Data Science Handbook*.

### Contents

1	Introduction	3
1.1	About the authors	4
1.2	Installation and Practical Considerations	4
2	How to Run Python Code	5
3	A Quick Tour of Python Language Syntax	7
3.1	Comments Are Marked by #	8
3.2	End-of-Line Terminates a Statement	8
3.3	Semicolons Can <i>Optionally</i> Terminate a Statement	9
3.4	Indentation: Whitespace Matters!	9
3.5	Whitespace <i>Within</i> Lines Does Not Matter	10

## DS1 Lecture 01 supplement

Jim Bagrow

This [notebook](#) acts as an additional supplement on programming in Python. You are expected to learn this material during the lecture, and a quiz of the semester will assess this.

This document complements the course's "Whirlwind Tour" reading. You are also expected to complete this reading as part of the assignment after reading the whirlwind tour.

### About this document

- This is a [Jupyter Notebook](#). I'll be using these for code-heavy lectures and lecture supplements. Notebooks contain code in a document, interleaved with narrative (non-code) text.
- Your [projects and assignments](#) will be submitted as Python scripts (`.py` files).

### Contents:

- A bit more on data structures and control flow (see also the Whirlwind Tour)
  - tuples vs. lists
    - multiassignment
    - sorting
    - zipping
  - Working with [files](#)!
  - Exploring the [Python standard library](#) (builtin modules and packages)
  - A tour of useful third-party packages
    - [numpy and scipy](#)
    - [matplotlib](#)
  - The [random](#) library
  - Conclusion