

Computational Analysis of Big Data

Week 1

Mostly an introduction, but also:

Coding with data in Python

About the teacher ●

About this course ● ○ ○ ○

About Big Data and Python ○ ○ ○ ○

Course overview

Course overview

Sessions

1. Coding with data in Python
2. ~~Data~~ Scientist's most fundamental tools
3. Getting data—scraping and APIs
4. Machine learning 1
5. Machine learning 2
6. Networks
7. Natural language processing
8. Crunching Big Data with MapReduce
9. Lab work on project report
10. Lab work on project report
11. Project presentations

Course overview

Sessions

- 1. Coding with data in Python**
- 2. A Data Scientist's most fundamental tools**
- 3. Getting data—scraping and APIs**
- 4. Machine learning 1**
- 5. Machine learning 2**
- 6. Networks**
- 7. Natural language processing**
- 8. Crunching Big Data with MapReduce**
- 9. Lab work on project report**
- 10. Lab work on project report**
- 11. Project presentations**

```

212 def randomize_by_edge_swaps(G, num_iterations):
213     """Randomizes the graph by swapping edges in such a way that
214     preserves the degree distribution of the original graph.
215     """
216     Source: https://gist.github.com/gotgenes/2770023
217     """
218     newgraph = G.copy()
219     edge_list = newgraph.edges()
220     num_edges = len(edge_list)
221     total_iterations = num_edges * num_iterations
222
223     for i in xrange(total_iterations):
224         rand_index1 = int(round(random.random() * (num_edges - 1)))
225         rand_index2 = int(round(random.random() * (num_edges - 1)))
226         original_edge1 = edge_list[rand_index1]
227         original_edge2 = edge_list[rand_index2]
228         head1, tail1 = original_edge1
229         head2, tail2 = original_edge2
230
231         # Flip a coin to see if we should swap head1 and tail1 for
232         # the connections
233         if random.random() >= 0.5:
234             head1, tail1 = tail1, head1
235
236         if head1 == tail2 or head2 == tail1:
237             continue
238
239         if newgraph.has_edge(head1, tail2) or newgraph.has_edge(
240             head2, tail1):
241             continue
242
243         # Succeeded checks, perform the swap
244         original_edge1_data = newgraph[head1][tail1]
245         original_edge2_data = newgraph[head2][tail2]
246
247         newgraph.remove_edges_from((original_edge1, original_edge2))
248
249         new_edge1 = (head1, tail2, original_edge1_data)
250         new_edge2 = (head2, tail1, original_edge2_data)
251
252         newgraph.add_edges_from((new_edge1, new_edge2))
253
254         # Now update the entries at the indices randomly selected
255         edge_list[rand_index1] = (head1, tail2)
256         edge_list[rand_index2] = (head2, tail1)
257
258     assert len(newgraph.edges()) == num_edges
259     return newgraph

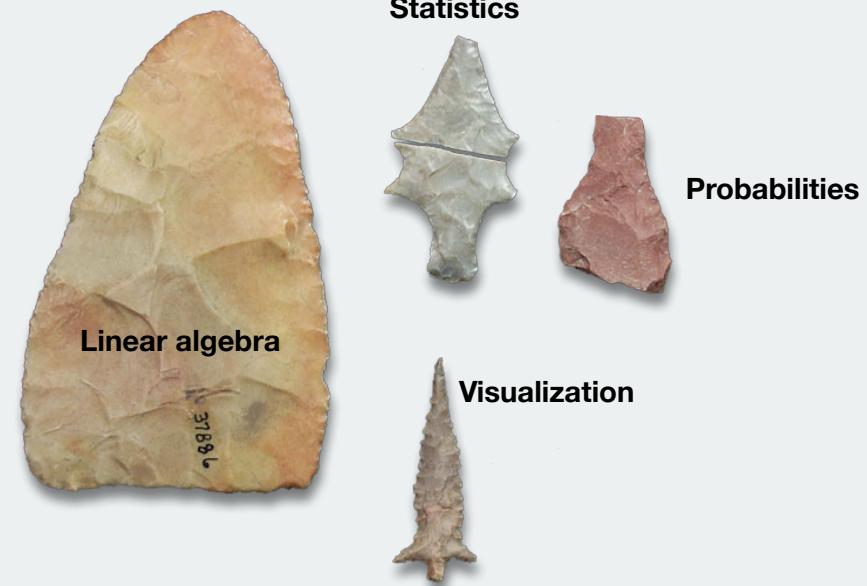
```

Aim: Get comfortable with basic Python

Course overview

Sessions

1. Coding with data in Python
2. ~~A Data Scientist's most fundamental tools~~
3. Getting data—scraping and APIs
4. Machine learning 1
5. Machine learning 2
6. Networks
7. Natural language processing
8. Crunching Big Data with MapReduce
9. Lab work on project report
10. Lab work on project report
11. Project presentations



Aim: Refresh math skills and learn how to use them in Python

Course overview

Sessions

1. Coding with data in Python
2. A Data Scientist's most fundamental tools
3. Getting data—scraping and APIs
4. Machine learning 1
5. Machine learning 2
6. Networks
7. Natural language processing
8. Crunching Big Data with MapReduce
9. Lab work on project report
10. Lab work on project report
11. Project presentations

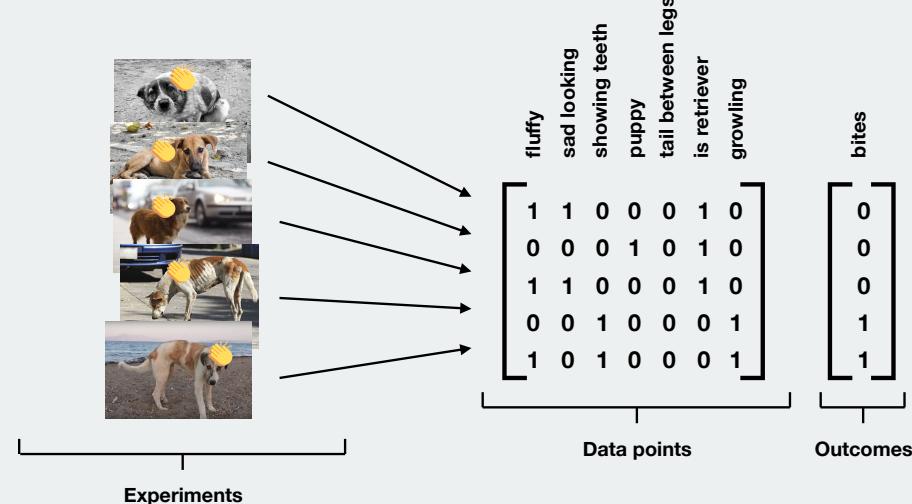
The screenshot shows the Wikipedia page for "Marvel Comics". The page content includes a brief history, key figures like Stan Lee and Martin Goodman, and various comic book series. The sidebar provides links to related articles and a navigation menu. On the right side, the developer tools (F12) are open, displaying the raw HTML code of the page. The code is heavily annotated with browser developer tool styling and script elements, illustrating how web content is rendered.

Aim: Learn how to get data from web and get some from the Wiki API

Course overview

Sessions

1. Coding with data in Python
2. ~~Data Scientist's most fundamental tools~~
3. Getting data—scraping and APIs
- 4. Machine learning 1**
- 5. Machine learning 2**
6. Networks
7. Natural language processing
8. Crunching Big Data with MapReduce
9. Lab work on project report
10. Project presentations

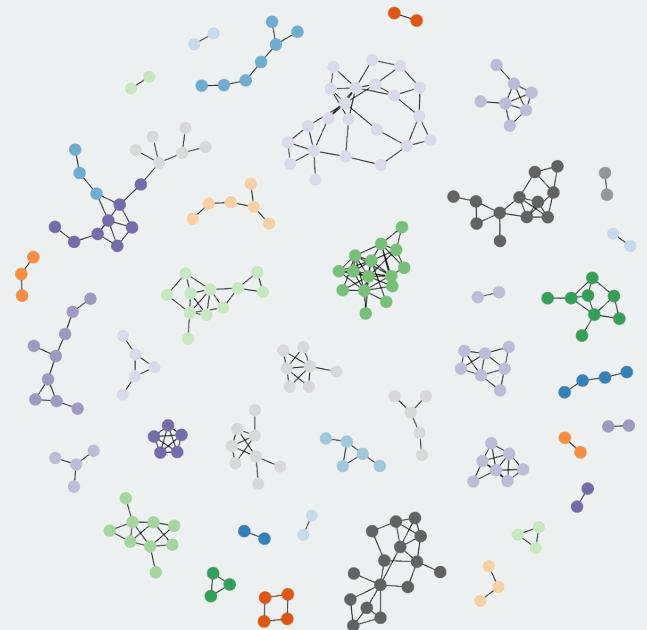


Aim: Understand the paradigm, learn how (some of) it works

Course overview

Sessions

1. Coding with data in Python
2. ~~Data Scientist's most fundamental tools~~
3. Getting data—scraping and APIs
4. Machine learning 1
5. Machine learning 2
- 6. Networks**
7. Natural language processing
8. Crunching Big Data with MapReduce
9. Lab work on project report
10. Lab work on project report
11. Project presentations



Aim: Learn how to describe and visualize complex data as a network

Course overview

Sessions

1. Coding with data in Python
2. ~~Data Scientist's most fundamental tools~~
3. Getting data—scraping and APIs
4. Machine learning 1
5. Machine learning 2
6. Networks
7. **Natural language processing**
8. Crunching Big Data with MapReduce
9. Lab work on project report
10. Lab work on project report
11. Project presentations

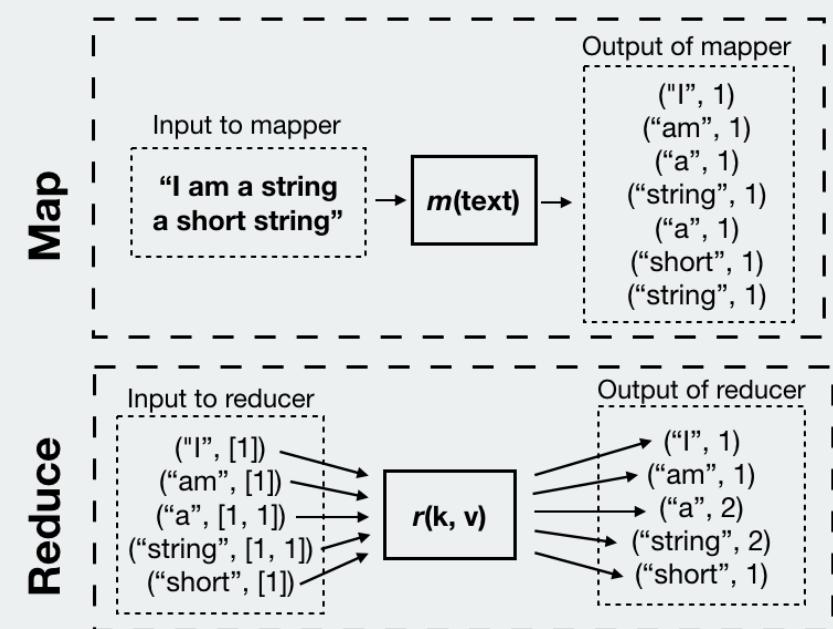


Aim: Get introduced to the (huge) field of NLP, and learn a few skills

Course overview

Sessions

1. Coding with data in Python
2. ~~Data Scientist's most fundamental tools~~
3. Getting data—scraping and APIs
4. Machine learning 1
5. Machine learning 2
6. Networks
7. Natural language processing
- 8. Crunching Big Data with MapReduce**
9. Lab work on project report
10. Lab work on project report
11. Project presentations

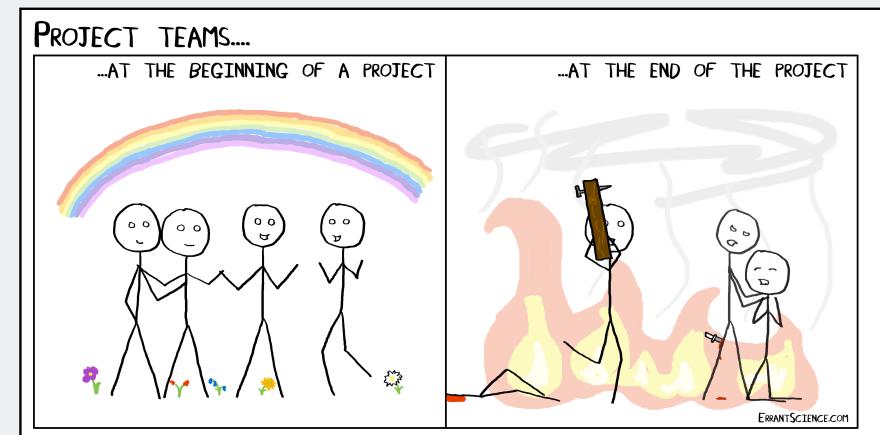


Aim: Process a massive dataset with an awesome method

Course overview

Sessions

1. Coding with data in Python
2. A Data Scientist's most fundamental tools
3. Getting data—scraping and APIs
4. Machine learning 1
5. Machine learning 2
6. Networks
7. Natural language processing
8. Crunching Big Data with MapReduce
- 9. Lab work on project report**
- 10. Lab work on project report**
- 11. Project presentations**



Aim: Synthesize all the things you have learned and make a project for your portfolio

What will you learn?

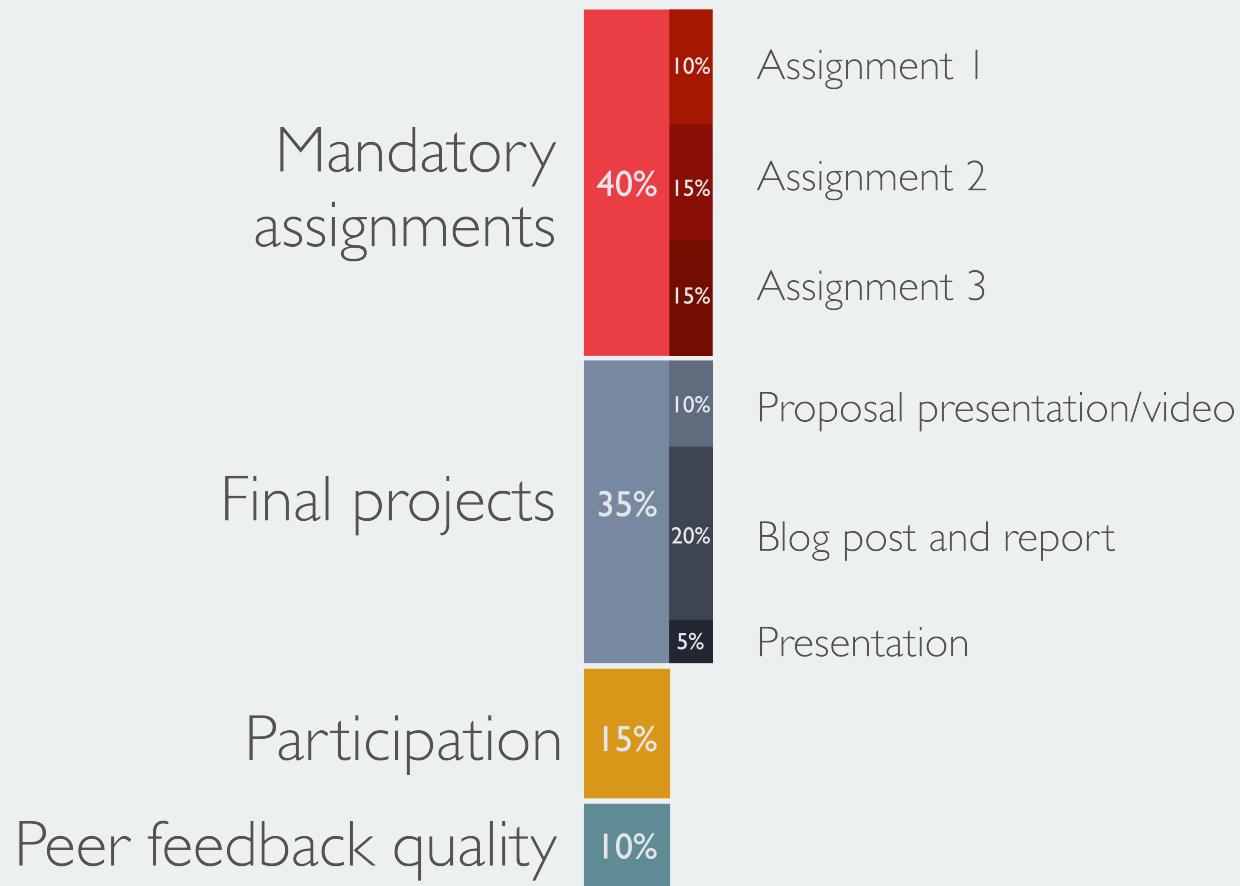
Knowledge and competences

- Get to know the landscape of problems and tools in Data Science
- Learn what Big Data is in this context
- Know where to start when you want to analyze something that requires lots of data

Concrete skills

- Juggle data in Python
- Predict outcomes from input data (machine learning)
- Visualize complex data
- Analyze Big Data using parallel methods (MapReduce)

How will you be graded?



How assignments work

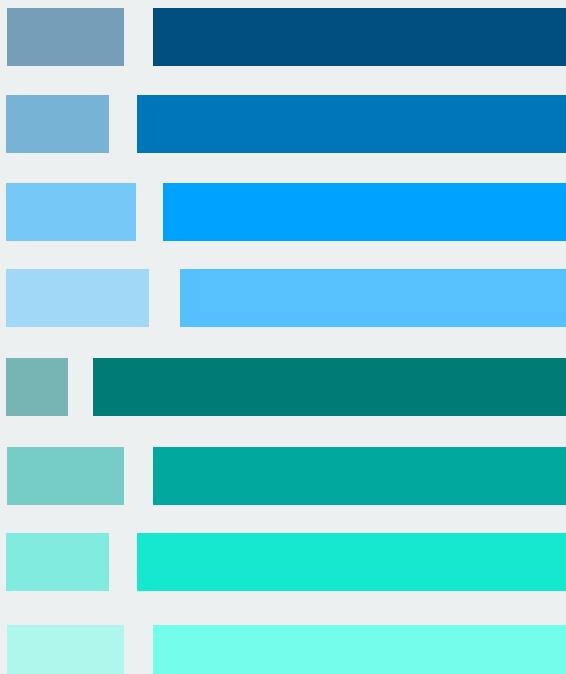
Sessions

Talk

Solve exercises

How assignments work

Sessions



How assignments work



The final project

Feb. 17

I validate your
project idea

Part A Feb. 22

You present a
proposal video

Deliverable

Part B March 18

You deliver a blog
post and your code

Deliverable

March 19

You give a
presentation

Previous students projects

“What makes us happy?”

“The secrets hidden in your Venmo and Instagram data”

“Pizza Makes the World go Round”

How to do well in this course

Best strategy:

1. Complete the *preparation goals* for each session (see calendar on Canvas)
2. Be inquisitive. Ask lots of questions to your neighbors and me, and up your googling-game

Everything else is on Canvas

About the teacher ●

About this course ● ● ● ●

About Big Data and Python ○ ○ ○ ○

What is Big Data?

What is Big Data?

big data

noun COMPUTING

extremely large data sets that may be analysed computationally to reveal patterns, trends, and associations, especially relating to human behaviour and interactions.

"much IT investment is going towards managing and maintaining big data"



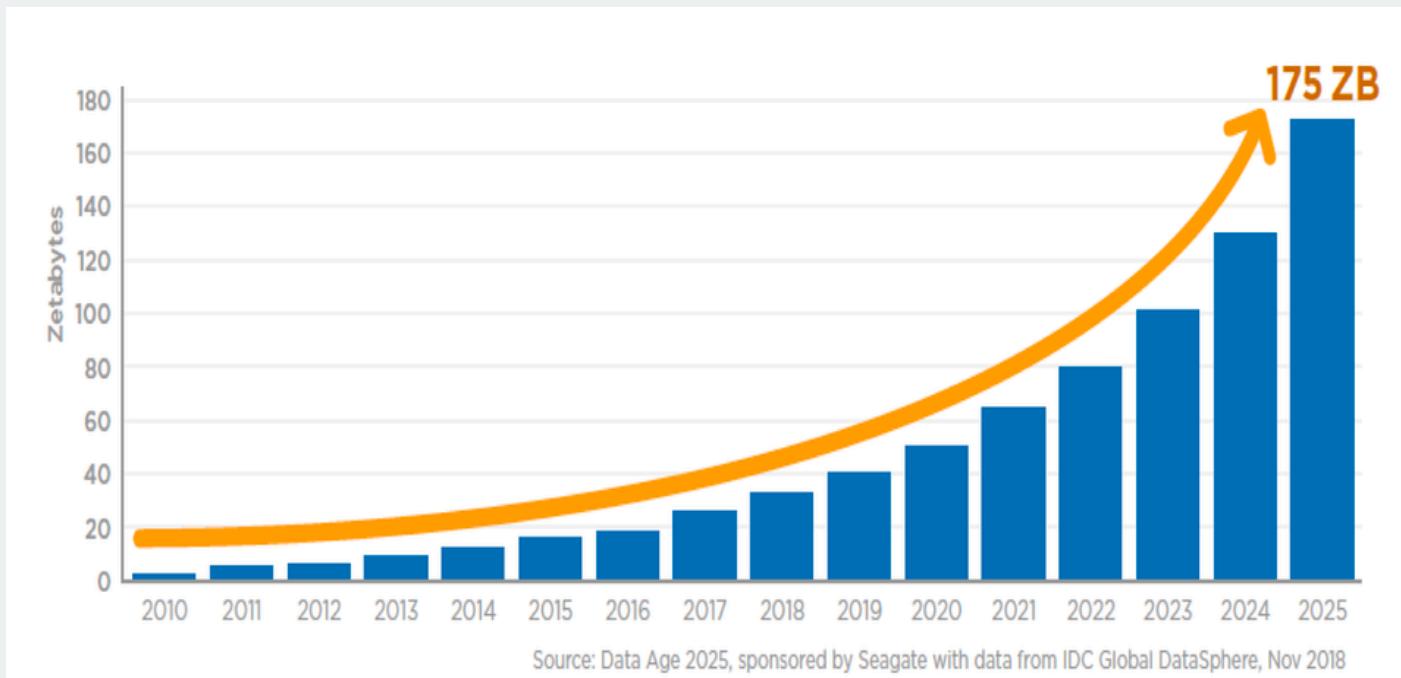
Translations, word origin, and more definitions

The 5Vs of Big Data



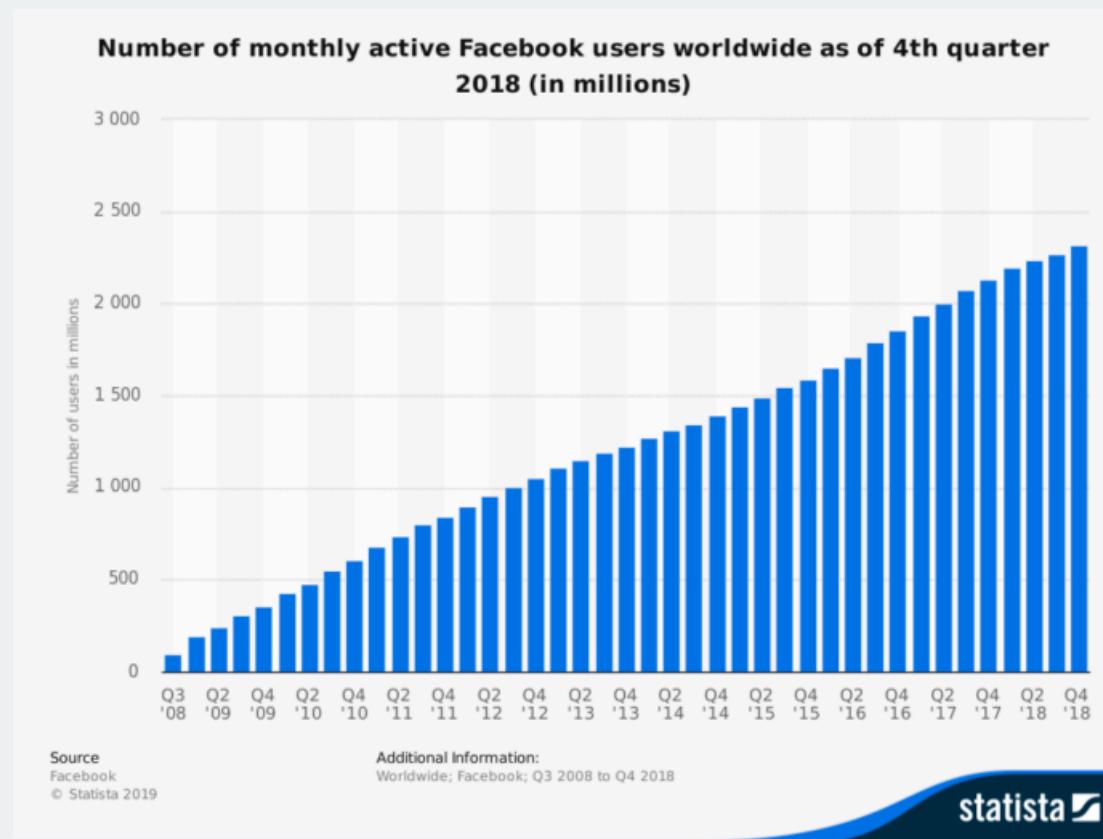
<https://suryagutta.medium.com/the-5-vs-of-big-data-2758bfcc51d>

The 5Vs of Big Data: Volume



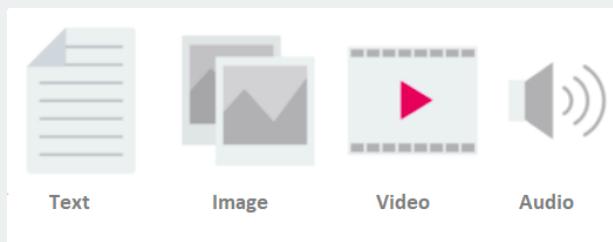
<https://suryagutta.medium.com/the-5-vs-of-big-data-2758bfcc51d>

The 5Vs of Big Data: Velocity



<https://suryagutta.medium.com/the-5-vs-of-big-data-2758bfcc51d>

The 5Vs of Big Data: Variety



<https://suryagutta.medium.com/the-5-vs-of-big-data-2758bfcc51d>

The 5Vs of Big Data: Veracity

- Quality
- Integrity
- Credibility
- Accuracy



<https://suryagutta.medium.com/the-5-vs-of-big-data-2758bfcc51d>

The 5Vs of Big Data: Value



<https://suryagutta.medium.com/the-5-vs-of-big-data-2758bfcc51d>

Other Vs

- **Viscosity** (complexity or degree of correlation)
- **Variability** (inconsistency in data flow)
- **Volatility** (durability or how long time data is valid and how long it should be stored)
- **Viability** (capability to be live and active)
- **Validity** (understandable to find the hidden relationships)

<https://suryagutta.medium.com/the-5-vs-of-big-data-2758bfcc51d>

Where and for what is Big Data being used?

Where and for what is Big Data being used?

Ads



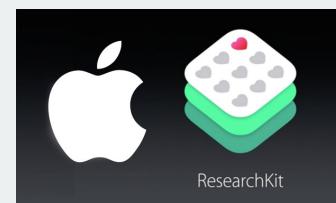
 **theTradeDesk**[®]
Omnicom WPP

Ted Talk: “Is Big Data Killing Creativity?”

User customization

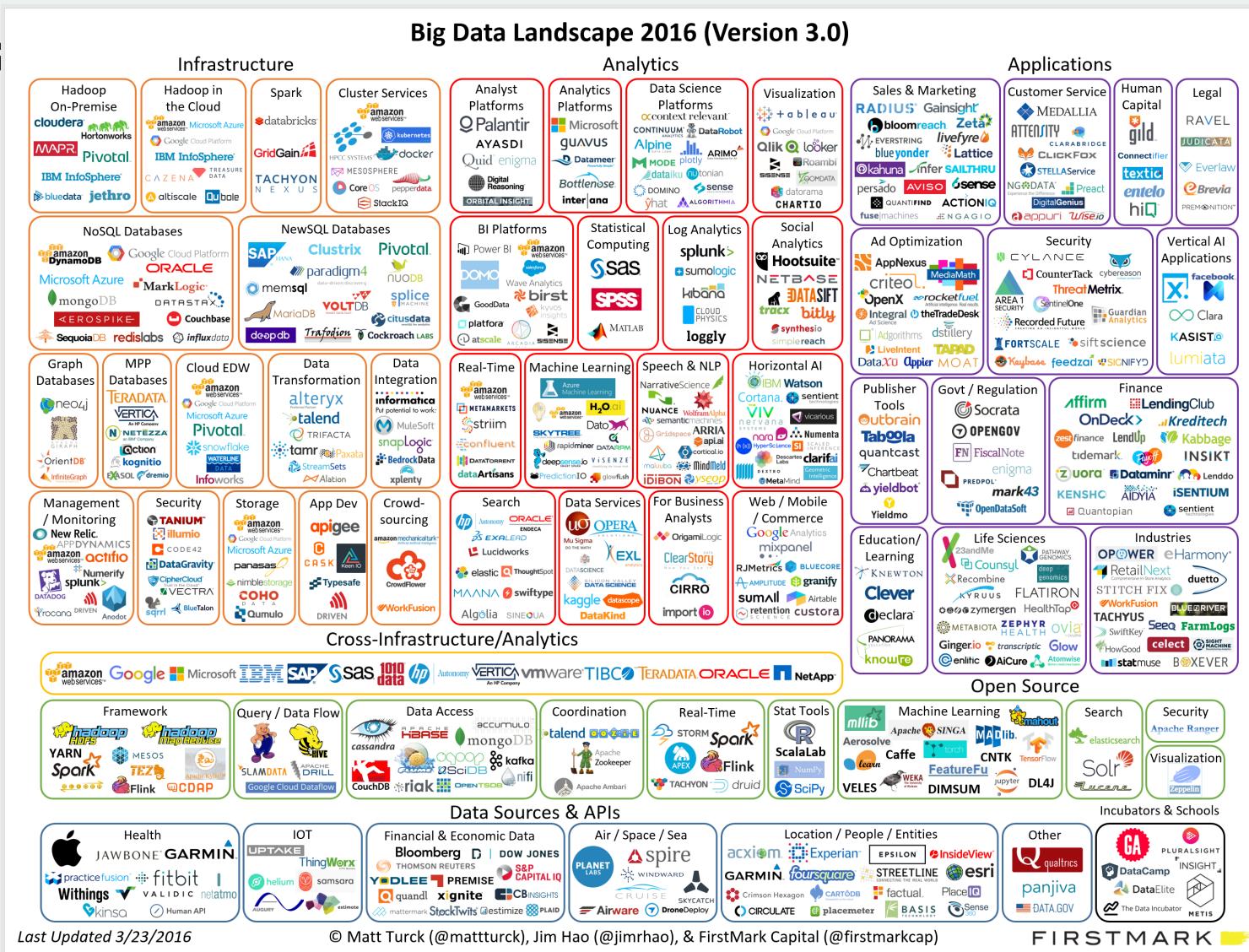


Research



How does one work with Big Data?

How does



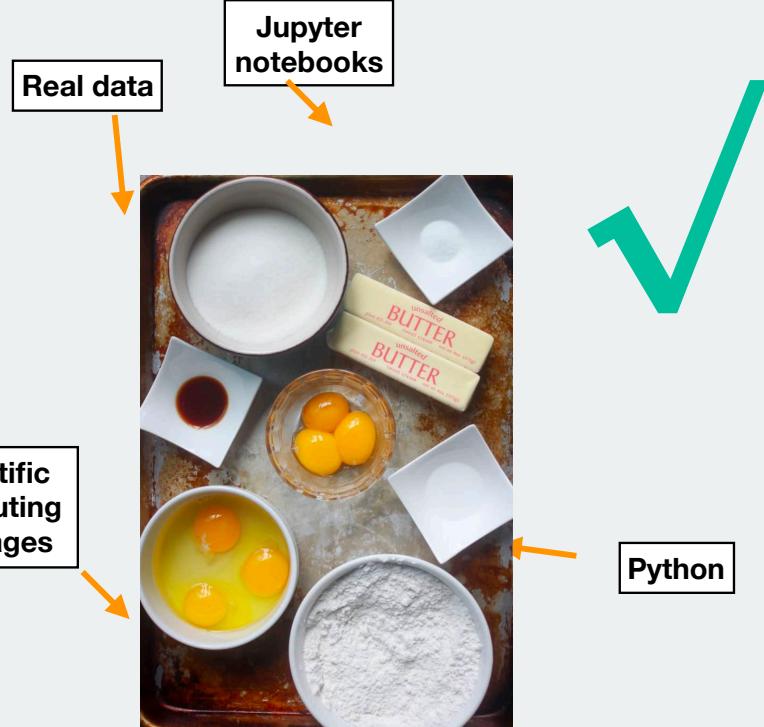
Last Updated 3/23/2016

© Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark Capital (@firstmarkcap)

FIRSTMARK

How does one work with Big Data?

In this course: Everything from scratch



Let's get started!

1. If you haven't already: download and install Anaconda (Python 3.8 version) <https://www.continuum.io/downloads>.
2. Make a folder for this course. Open a terminal (Mac) or a console (PC) and navigate to that folder.
3. Run the following command in your terminal/console:

```
git clone https://github.com/lucian979/CarletonBD
```

4. In your terminal/console, navigate into the exercises folder inside of the newly created **CarletonBD** directory, and run the following command:

```
jupyter notebook exercises_week1.ipynb
```