

EDA Part 1: Data Engineering

Tabular Data, Pandas, and Scraping

cs109a, Fall 2017

Pavlos Protopapas, Kevin Rader, Rahul Dave, Margo Levine

It took about three years before the BellKor's Pragmatic Chaos team managed to win the prize ... The winning algorithm was ... so complex that it was never implemented by Netflix.¹

¹ <https://hbr.org/2012/10/big-data-hype-and-reality>

Machine

Data Management

Data Mining

Machine Learning

Visualization

Business Intelligence

Statistics

Data Science

Human

Human Cognition

Perception

Story Telling

Decision Making
Theory



DATA
ENGINEERING

Hacking Skills

DATA
ANALYSIS

*Math & Statistics
Knowledge*

Machine
Learning

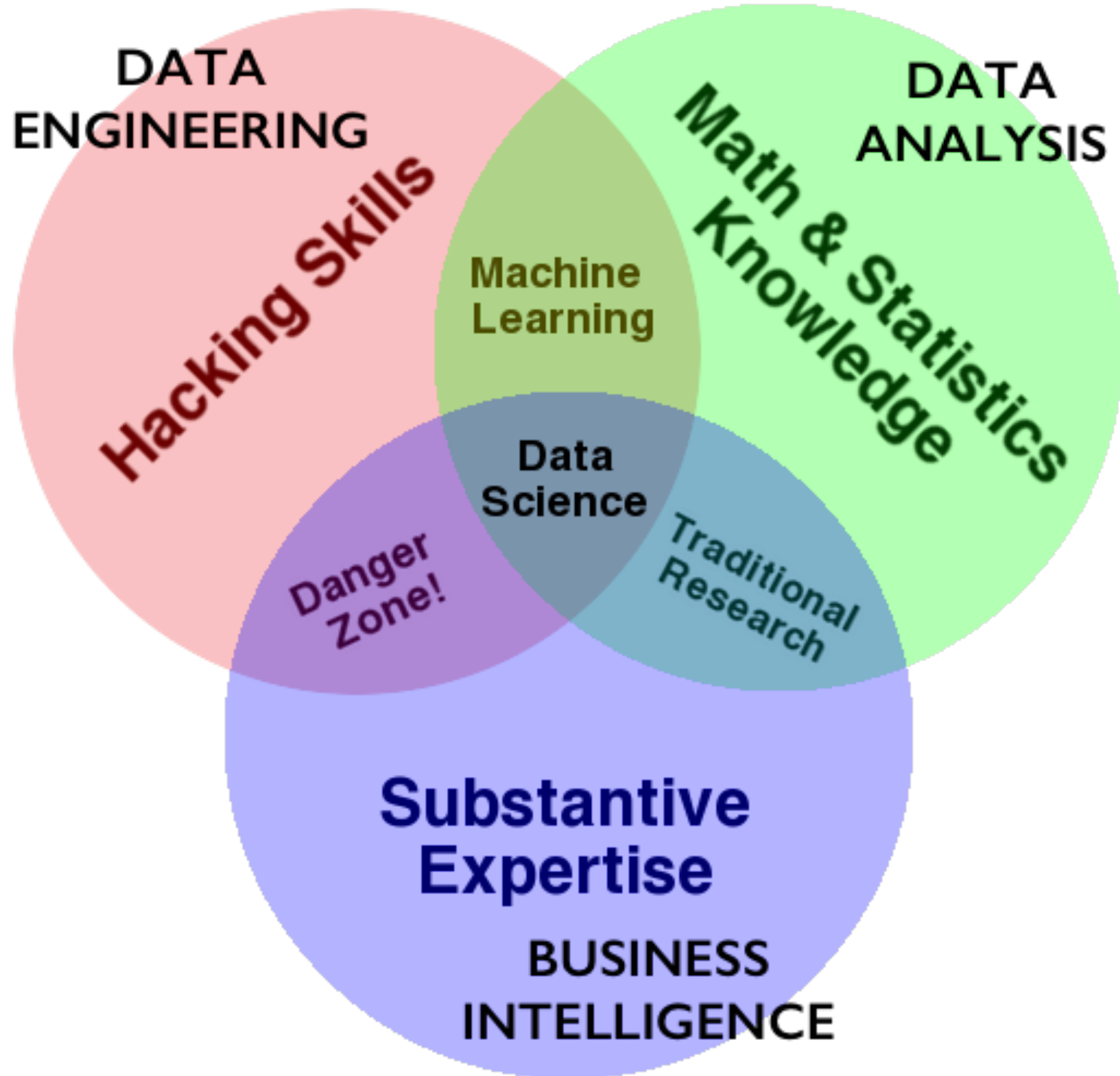
Data
Science

*Danger
Zone!*

*Traditional
Research*

**Substantive
Expertise**

BUSINESS
INTELLIGENCE



Data Scientist: Sexiest Job of the 21st Century

It's important that our data team wasn't comprised solely of mathematicians and other "data people." It's a fully integrated product group that includes people working in design, web development, engineering, product marketing, and operations. They all understand and work with data, and I consider them all data scientists... Often, an engineer can have the insight that makes it >clear how the product's design should work, or vice-versa — a designer can have the insight that helps the engineers understand how to better use the data. Or it may take someone from marketing to understand what a customer really wants to accomplish.²

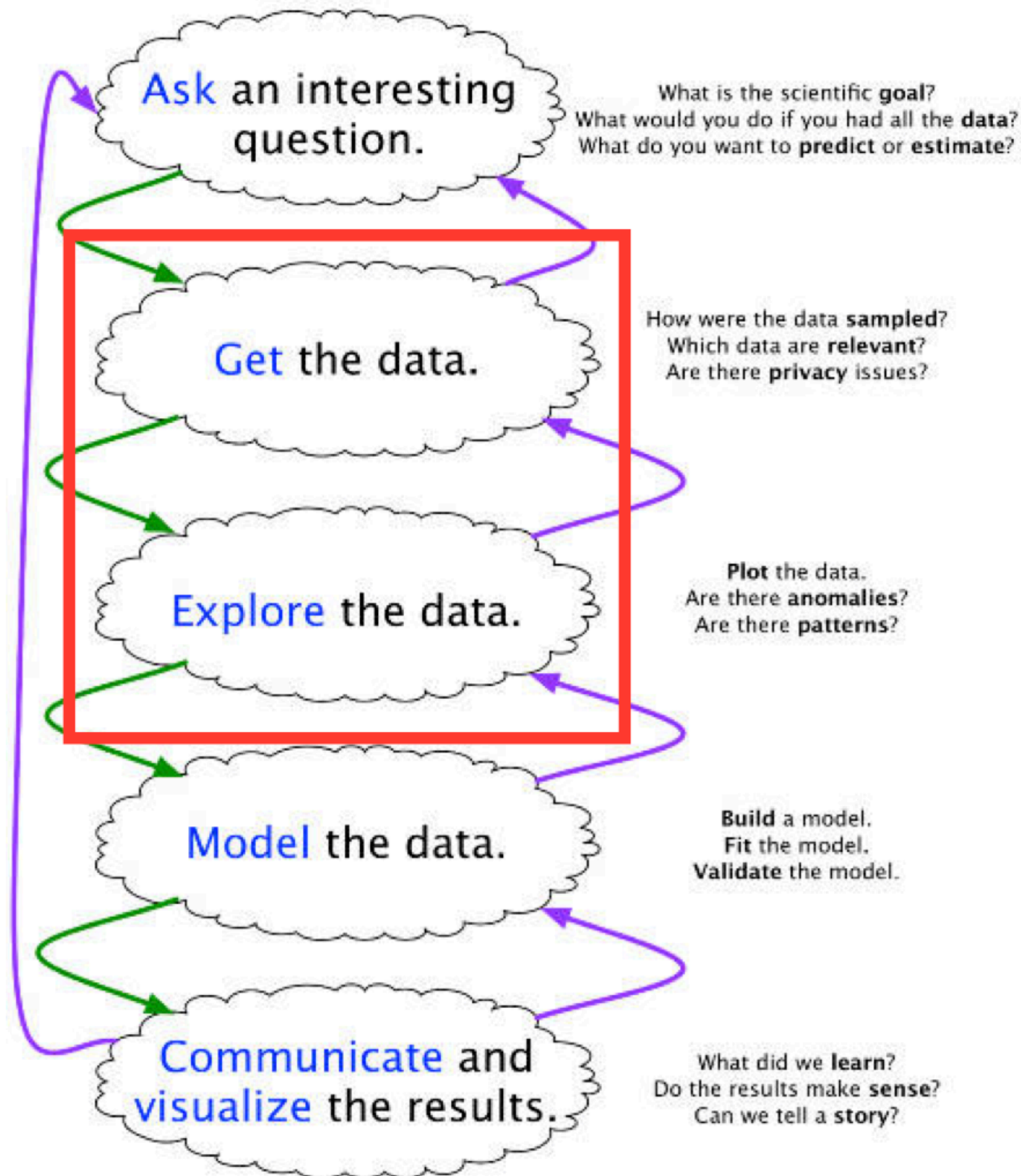
² D. J. Patil, U.S. Chief Data Scientist, Building data science teams. " O'Reilly Media, Inc.", 2011.

DATA ENGINEERING

- **data:** scraping, API, feature engineering, all part of EDA
- **compute:** code, python, R, julia, spark, hadoop
- **storage/database:** pandas, SQL, NoSQL, HBase, disk, memory
- **devops:** AWS, docker, mesos, repeatability
- **product:** database, web, API, viz, UI, story

Different at different scales....

The Data Science Process



The basic EDA workflow⁵

1. **Build** a DataFrame from the data (ideally, put all data in this object)
2. **Clean** the DataFrame. It should have the following properties:
 - Each row describes a single object
 - Each column describes a property of that object
 - Columns are numeric whenever appropriate
 - Columns contain atomic properties that cannot be further decomposed
3. Explore **global properties**. Use histograms, scatter plots, and aggregation functions to summarize the data.
4. Explore **group properties**. Use groupby, queries, and small multiples to compare subsets of the data.

⁵ enunciated in this form by Chris Beaumont, the first Head TF of cs109

Today, Monday

We'll focus on data and relational storage

- How do we engineer features from the web?
- What is a relational Database?
- What Grammar of Data does it follow?
- How is this grammar implemented in Pandas?

Wednesday

We'll focus on the visualization part
of EDA.

In reality, both go together.

Relational Database

- Don't say_: seek 20 bytes onto disk and pick up from there. The next row is 50 bytes hence
- *Say*: select data from a set. I don't care where it is, just get the row to me.
- Its just the table Kevin talked about last time...

Relational Database(contd)

- A collection of tables related to each other through common data values.
- Rows represent attributes of something
- Everything in a column is values of *one* attributes
- A cell is expected to be atomic
- Tables are related to each other if they have columns called keys which represent the same values

Scales of Measurement

TABLE 1

- Quantitative (Interval and Ratio)

- Ordinal

- Nominal³

Scale	Basic Empirical Operation	Mathematical Group Structure	Permissible Statistics (invariantive)
NOMINAL	Determination of equality	Permutation group $x' = f(x)$ $f(x)$ means any one-to-one substitution	Number of cases Mode Contingency correlation
ORDINAL	Determination of greater or less	Isotonic group $x' = f(x)$ $f(x)$ means any monotonic increasing function	Median Percentiles
INTERVAL	Determination of equality of intervals or differences	General linear group $x' = ax + b$	Mean Standard deviation Rank-order correlation Product-moment correlation
RATIO	Determination of equality of ratios	Similarity group $x' = ax$	Coefficient of variation

³ S. S. Stevens, Science, New Series, Vol. 103, No. 2684 (Jun. 7, 1946), pp. 677-680

Grammar of Data

Been there for a while (SQL, Pandas), formalized in `dplyr`⁴.

- provide simple verbs for simple things. These are functions corresponding to common data manipulation tasks
- second idea is that backend does not matter. Here we constrain ourselves to Pandas.
- multiple backends implemented in Pandas, Spark, Impala, Pig, dplyr, ibis, blaze

⁴ Hadley Wickham: <https://cran.rstudio.com/web/packages/dplyr/vignettes/introduction.html>

Why bother

- learn how to do core data manipulations, no matter what the system
- relational databases critical for non-memory fits. Big installed base.
- one off questions: google, stack-overflow, <http://chrisalbon.com>

The grammar of data

For cleaning and for transformation:

VERB	dplyr	pandas	SQL
QUERY/SELECTION	filter() (and slice())	query() (and loc[], iloc[])	SELECT WHERE
SORT	arrange()	sort()	ORDER BY
SELECT-COLUMNS/PROJECTION	select() (and rename())	(and rename())	SELECT COLUMN
SELECT-DISTINCT	distinct()	unique(),drop_duplicates()	SELECT DISTINCT COLUMN
ASSIGN	mutate() (and transmute())	assign	ALTER/UPDATE
AGGREGATE	summarise()	describe(), mean(), max()	None, AVG(),MAX()
SAMPLE	sample_n() and sample_frac()	sample()	implementation dep, use RAND()
GROUP-AGG	group_by/summarize	groupby/agg, count, mean	GROUP BY
DELETE	?	drop/masking	DELETE/WHERE

Example: Candidates

	id	first_name	last_name	middle_name	party
	Filter	Filter	Filter	Filter	Filter
1	16	Mike	Huckabee		R
2	20	Barack	Obama		D
3	22	Rudolph	Giuliani		R
4	24	Mike	Gravel		D
5	26	John	Edwards		D
6	29	Bill	Richardson		D
7	30	Duncan	Hunter		R
8	31	Dennis	Kucinich		D
9	32	Ron	Paul		R

unique integers

strings

string w subset of options (categorical variables)

Contributors

Table: contributors												
	id	last_name	first_name	middle_name	street_1	street_2	city	state	zip	amount	date	candidate_id
	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter
1	1	Agee	Steven	NULL	549 Laurel ...	NULL	Floyd	VA	24091	500	2007-06-30	16
2	5	Akin	Charles	NULL	10187 Suga...	NULL	Bentonville	AR	72712	100	2007-06-16	16
3	6	Akin	Mike	NULL	181 Baywo...	NULL	Monticello	AR	71655	1500	2007-05-18	16
4	7	Akin	Rebecca	NULL	181 Baywo...	NULL	Monticello	AR	71655	500	2007-05-18	16
5	8	Aldridge	Brittni	NULL	808 Capitol...	NULL	Washington	DC	20024	250	2007-06-06	16
6	9	Allen	John D.	NULL	1052 Cann...	NULL	North Augu...	SC	29860	1000	2007-06-11	16
7	10	Allen	John D.	NULL	1052 Cann...	NULL	North Augu...	SC	29860	1300	2007-06-29	16
8	11	Allison	John W.	NULL	P.O. Box 10...	NULL	Conway	AR	72033	1000	2007-05-18	16
9	12	Allison	Rebecca	NULL	3206 Sum...	NULL	Little Rock	AR	72227	1000	2007-04-25	16

Operations

- QUERY: `dfcwc_i[(dfcwc_i.state=='VA') & (dfcwc_i.amount < 400)]` series (pandas) —> like an array
- SORT: `dfcwc_i.sort_values(by="amount", ascending=False)`
- SELECT-COLUMNS: `dfcwc_i[['first_name', 'amount']]`
- SELECT-DISTINCT:
`dfcwc_i[['last_name', 'first_name']].drop_duplicates()`

- **ASSIGN:**
`dfcwci['name']=dfcwci['last_name']+"",
"+dfcwci['first_name']`
- **ASSIGN(in-place):**
`dfcwci.loc[dfcwci.state=='VA',
'name']="junk"`
- **AGGREGATE:** `dfcwci.amount.max(),
dfcwci.describe()`
- **DELETE:** `del dfcwci['name']` (DROP-COLUMN)

Split-Apply-Combine

```
In [28]: dfewci.groupby("state").sum()
```

```
Out[28]:
```

	zip	amount	candidate_id
state			
AK	2985459621	1210.00	111
AR	864790	14200.00	192
AZ	860011121	120.00	37
CA	14736360720	-5013.73	600
CO	2405477834	-5823.00	111
CT	68901376	2300.00	35
DC	800341853	-1549.91	102
FL	8970626520	-4050.00	803

- GROUP-AGG
- splitting the data into groups based on some criteria
- applying a function to each group independently
- combining the results into a data structure

RELATIONSHIPS (in addition to rubric)

- we usually need to combine data from multiple sources
- different systems have different ways, most copy SQL (pandas)
- sub-select:

```
obamaid=dfcand.query("last_name=='Obama'")['id'].values[0]  
obamacontrib=dfcwcic.query("candidate_id==%i" % obamaid)
```

JOINS

- combine tables on a common key-value
- 90% of the time, EXPLICIT INNER JOIN

```
In [40]: cols_wanted=['last_name_x', 'first_name_x', 'candidate_id', 'id', 'last_name_y']  
dfcwc.merge(dfcand, left_on="candidate_id", right_on="id")[cols_wanted]
```

```
Out[40]:
```

	last_name_x	first_name_x	candidate_id	id	last_name_y
0	Agee	Steven	16	16	Huckabee
1	Akin	Charles	16	16	Huckabee
2	Akin	Mike	16	16	Huckabee
3	Akin	Rebecca	16	16	Huckabee
4	Aldridge	Brittni	16	16	Huckabee

Web Servers

- A server is a long running process (also called daemon) which listens on a pre-specified port
- and responds to a request, which is sent using a protocol called HTTP
- A browser must first we must parse the url.
Everything after a # is a fragment. Until then its the DNS name or ip address, followed by the URL.

Example

Our notebooks also talk to a local web server on our machines:
`http://localhost:8888/Documents/cs109/
BLA.ipynb#something`

- protocol is `http`, hostname is `localhost`, port is `8888`
- url is `/Documents/cs109/BLA.ipynb`
- url fragment is `#something`

Request is sent to localhost on port 8888. It says:

Request:

```
GET /request-URI HTTP/version
```

Example with Response: Google

GET / HTTP/1.0
Host: www.google.com

HTTP/1.0 200 OK
Date: Mon, 14 Nov 2016 04:49:02 GMT
Expires: -1
Cache-Control: private, max-age=0
Content-Type: text/html; charset=ISO-8859-1
P3P: CP="This is ..."
Server: gws
X-XSS-Protection: 1; mode=block
X-Frame-Options: SAMEORIGIN
Set-Cookie: NID=90=gb5q7b0...; expires=Tue, 16-May-2017 04:49:02 GMT; path=/; domain=.google.com; HttpOnly
Accept-Ranges: none
Vary: Accept-Encoding

<!doctype html><html itemscope=""
itemtype="http://schema.org/WebPage" lang="en">
<head><meta content="Search the world's information,

HTTP Status Codes⁶

- 200 OK:
Means that the server did whatever the client wanted it to, and all is well.
- 201 Created:
The request has been fulfilled and resulted in a new resource being created. The newly created resource can be referenced by the URI(s) returned in the entity of the response, with the most specific URI for the resource given by a Location header field.
- 400: Bad request
The request sent by the client didn't have the correct syntax.
- 401: Unauthorized
Means that the client is not allowed to access the resource. This may change if the client retries with an authorization header.
- 403: Forbidden
The client is not allowed to access the resource and authorization will not help.
- 404: Not found
Seen this one before? :) It means that the server has not heard of the resource and has no further clues as to what the client should do about it. In other words: dead link.
- 500: Internal server error
Something went wrong inside the server.
- 501: Not implemented
The request method is not supported by the server.

⁶ (from <http://www.garshol.priv.no/download/text/http-tut.htm>)

requests

- great module built into python for http requests

```
req = requests.get("https://en.wikipedia.org/wiki/Harvard_University")
```

```
<Response [200]>
```

```
page = req.text
```

```
'<!DOCTYPE html>\n<html class="client-nojs" lang="en" dir="ltr">\n<head>\n<meta charset="UTF-8"/>\n<title>Harvard University -\nWikipedia</title>\n<script>document.documentElement.className =\ndocument.documentElement.className.replace( /(^\|\\s)client-nojs(\\s|$)/,\n"$1client-js$2"\n);</script>\n<script>(window.RLQ=window.RLQ||[]).push(function(){mw.config.set({\n"wgCanonicalNamespace":"","wgCanonicalSpecialPageName":false,"wgNamespaceNumber"\n:0,"wgPageName":"Harvard_University","wgTitle":"Harva...'`
```



WIKIPEDIA
The Free Encyclopedia

[Main page](#)
[Contents](#)
[Featured content](#)
[Current events](#)
[Random article](#)
[Donate to Wikipedia](#)
[Wikipedia store](#)

[Interaction](#)

[Help](#)
[About Wikipedia](#)
[Community portal](#)
[Recent changes](#)
[Contact page](#)

[Tools](#)

[What links here](#)
[Related changes](#)
[Upload file](#)
[Special pages](#)
[Permanent link](#)
[Page information](#)
[Wikidata item](#)
[Cite this page](#)

[Print/export](#)

Not logged in [Talk](#) [Contributions](#) [Create account](#) [Log in](#)

Article

Talk

Read

[View source](#)

[View history](#)



**Wiki Loves Monuments: The world's largest
photography competition is now open!**



Photograph a historic site, learn more about our history, and win prizes.

Harvard University



From Wikipedia, the free encyclopedia

Coordinates: 42°22′28″N 71°07′01″W﻿ / ﻿

"Harvard" redirects here. For other uses, see [Harvard \(disambiguation\)](#).

Harvard University is a private Ivy League research university in Cambridge, Massachusetts, established in 1636, whose history, influence, and wealth have made it one of the world's most prestigious universities.^[7]

Established originally by the Massachusetts legislature and soon thereafter named for John Harvard (its first benefactor), Harvard is the United States' oldest institution of higher learning,^[8] and the Harvard Corporation (formally, the *President and Fellows of Harvard College*) is its first chartered corporation. Although

Harvard University



Latin: *Universitas Harvardiana*

Former names	Harvard College
Motto	<i>Veritas</i> ^[1]
Motto in English	Truth
Type	Private research
Established	1636 ^[2]
Endowment	\$24.541 billion (2016) ^[3]

Python data scraping

- Why scrape the web?
- vast source of information, combine with other data sets
- companies have not provided APIs
- automate tasks
- keep up with sites
- fun!

copyrights and permission:

- be careful and polite
- give credit
- care about media law
- don't be evil (no spam, overloading sites, etc.)

put sleep commands

Robots.txt

- specified by web site owner
- gives instructions to web robots (aka your script)
- is located at the top-level directory of the web server

e.g.: <http://google.com/robots.txt>

HTML

- angle brackets
- should be in pairs, eg `<p>Hello</p>`
- maybe in implicit bears, such as `
`

```
<!DOCTYPE html>
```

```
<html>
```

```
  <head>
```

```
    <title>Ttle</title>
```

```
  </head>
```

```
  <body>
```

```
    <h1>Body Title</h1>
```

```
    <p>Body Content</p>
```

```
  </body>
```

```
</html>
```

Developer Tools

- ctrl/cmd shift i in chrome
- cmd-option-i in safari
- look for "inspect element"
- locate details of tags

Beautiful Soup

- will normalize dirty html
- basic usage

```
import bs4
## get bs4 object
soup = bs4.BeautifulSoup(source)
## all a tags
soup.findAll('a')
## first a
soup.find('a')
## get all links in the page
link_list = [l.get('href') for l in soup.findAll('a')]
```

HTML is a tree

```
tree = bs4.BeautifulSoup(source)
```

```
## get html root node  
root_node = tree.html
```

```
## get head from root using contents  
head = root_node.contents[0]
```

```
## get body from root  
body = root_node.contents[1]
```

```
## could directly access body  
tree.body
```

Demographics table we want

Student life

Demographics of student body^{[124][125][126]}

	Undergraduate	Graduate and professional	U.S. census
Asian/Pacific Islander	17%	11%	5%
Black/non-Hispanic	6%	4%	12%
Hispanics of any race	9%	5%	16%
White/non-Hispanic	46%	43%	64%
Mixed race/other	10%	8%	9%
International students	11%	27%	N/A

Student body

In the last six years, Harvard's student body has grown to 21,000, across all programs.^[127] Harvard has 10,722 students in undergraduate programs, 3,738 students in professional programs, and 6,540 students in graduate and professional programs. The undergraduate population is 51% female, the graduate and professional population is 49% female.

Athletics

Main article: [Harvard Crimson](#)

The [Harvard Crimson](#) competes in 42 intercollegiate sports in the [NCAA Division I Ivy League](#). Harvard has an intense athletic rivalry with [Yale University](#) culminating in *The Game*, although the [Harvard–Yale Regatta](#) predates the football game. This rivalry is put aside every two years when the Harvard and Yale

Table with sole class wikitable

United States, both for students and parents.^[122] *College ROI Report: Best Value Colleges* by [PayScale](#) puts Harvard 22nd nationwide in the most recent 2016 edition.^[123]

Student life

Demographics of student body^{[124][125][126]}

	Undergraduate	Graduate and professional	U.S. census
Asian/Pacific Islander	17%	11%	5%
Black/non-Hispanic	6%	4%	12%
Hispanics of any race	9%	5%	16%
White/non-Hispanic	46%	43%	64%
Mixed race/other	10%	8%	9%
International students	11%	27%	N/A

Student body

In the last six years, Harvard's student population ranged from 19,000 to 21,000, across all programs.^[127] Harvard enrolled 6,655 students in undergraduate programs, 3,738 students in graduate programs, and 10,722 students in professional programs.^[124] The undergraduate population is 51% female, the graduate population is 48% female, and the professional population is 49% female.^[124]

Athletics

Main article: [Harvard Crimson](#)

The [Harvard Crimson](#) competes in 42 intercollegiate sports in the [NCAA Division I Ivy League](#). Harvard has an intense athletic rivalry with [Yale University](#) culminating in *The Game*, although the [Harvard–Yale](#)

Elements Console Sources Network Performance Memory Application Security Audits CoffeeConsole 1

<p>...</p>

<p>...</p>

<h2>...</h2>

<table style="text-align:center; float:left; font-size:85%; margin-right:2em;" class="wikitable">

<caption>

<i>Demographics of student body</i> = \$0

^{...}

^{...}

^{...}

</caption>

<tbody>...</tbody>

</table>

<h3>...</h3>

<p>...</p>

<h3>...</h3>

Styles Computed >>

Filter :hov .cls +

element.style {

}

i, user agent styleshee

ci

te, em, var, address,

dfn {

font-style: italic;

}

Inherited from caption

load.php?debug=...tor.des

table.wikitable >

caption <

Beautiful Soup Code

```
dfinder = lambda tag: tag.name=='table' and tag.get('class') == ['wikitable']
table_demographics = soup.find_all(dfinder)
rows = [row for row in table_demographics[0].find_all("tr")]
header_row = rows[0]
columns = [col.get_text() for col in header_row.find_all("th") if col.get_text()]
columns = [rem_nl(c) for c in columns]
indexes = [row.find("th").get_text() for row in rows[1:]]
values = []
for row in rows[1:]:
    for value in row.find_all("td"):
        values.append(to_num(value.get_text()))
stacked_values_lists = [values[i::3] for i in range(len(columns))]
stacked_values_iterator = zip(*stacked_values_lists)
df = pd.DataFrame(list(stacked_values_iterator), columns=columns, index=indexes)
```