

Course Reminders

- Due this Wednesday (11:59 PM)
 - Pre-course survey
 - Practice Assignment
- Due this Friday (11:59 PM)
 - D1
 - #FinAid quiz on Canvas
- Due next Wednesday
 - Group signup (get in groups of 3 - 5 now!!)
 - A1

Data tidiness & intuition

Jason G. Fleischer, Ph.D
UC San Diego

• • •

Department of Cognitive Science
jfleischer@ucsd.edu
<https://jgfleischer.com>
 @jasongfleischer

Data Structures Review

Structured data

- can be stored in database
SQL
- tables with rows and columns
- requires a relational key
- 5-10% of all data

Semi-structured data

- doesn't reside in a relational database
- has organizational properties (easier to analyze)
- CSV, XML, JSON

Unstructured

- non-tabular data
- 80% of the world's data
- images, text, audio, videos

Structured Data

Databases! Programs that manage huge data so that you can find the subset of data you want with a query. DB manage the data and run analyses through queries. DB are either “relational” (aka SQL) or “non-relational” (aka noSQL). Relational DB work using tables of data with “relationships” established between tables. The next few slides on the relationships in CSV semi-structured data apply to SQL DB as well. Non-relational DB work with key-value pairs to lookup data, and that’s exactly like JSON slides coming up.

Structured Data

Examples of relational DB

- *SQLite*
- *MySQL*
- *Postgres*

Examples of non-relational DB

- *Hadoop*
- *Hive*
- *Apache CouchBase*

(Semi-)Structured Data

Data that is stored in such a way that it is easy to search and work with. These data are stored in a particular format that adheres to organization principles imposed by the file format. These are the data structures data scientists work with most often.

CSVs

Each column separated by a comma

Has the extension ".csv"

Email	First Name	Last Name	Company	Snippet 1
example1@domain.com	John	Smith	Company 1	Snippet Sentence1
example2@gmail.com	Mary	Blake	Company 2	Snippet Sentence 2
example3@outlook.com	James	Joyce	Company 3	Snippet Sentence 3

Each row is separated by a new line



Example CSV



File Edit View Insert Format Data Tools Add-ons Help All changes saved in Drive

fx

	A	B	C	D	E	F
1	Email	First Name	Last Name	Company	Snippet 1	
2	example1@domain.com	John	Smith	Company 1	Snippet Sentence1	
3	example2@gmail.com	Mary	Blake	Company 2	Snippet Sentence2	
4	example3@outlook.com	James	Joyce	Company 3	Snippet Sentence3	
5						
6	CSV file					
7						
8						

Example CSV - Sheet 1 — Notatnik

Plik Edycja Format Widok Pomoc

Email,First Name,Last Name,Company,Snippet 1

example1@domain.com,John,Smith,Company 1,Snippet Sentence1

example2@gmail.com,Mary,Blake,Company 2,Snippet Sentence2

example3@outlook.com,James,Joyce,Company 3,Snippet Sentence3

JSON: key-value pairs

nested/hierarchical data

{"Name": "Isabela"}

The diagram illustrates a JSON object consisting of a single key-value pair. The key, 'Name', is highlighted in large black font at the top left. The value, 'Isabela', is highlighted in large black font at the top right. Two pink arrows point from the words 'key' and 'value' at the bottom left and bottom right respectively, towards their corresponding parts in the JSON object above.

key

value

JSON

These are all
nested within
attributes

```
"attributes": {  
    "Take-out": true,  
    "Wi-Fi": "free",  
    "Drive-Thru": true,  
    "Good For": {  
        "dessert": false,  
        "latenight": false,  
        "lunch": false,  
        "dinner": false,  
        "breakfast": false,  
        "brunch": false  
    },
```

These are all
nested within
"Good For"

JSON



Jupyter notebooks suck to version control

<https://nextjournal.com/schmudde/how-to-version-control-jupyter>

```
{
  "cell_type": "code",
  "execution_count": null,
  "metadata": {},
  "outputs": [],
  "source": [
    "import pandas as pd\n",
    "import holoviews as hv\n",
    "hv.extension('bokeh')"
  ]
},
```

A large orange arrow pointing to the right, containing the word "DETOUR" in black capital letters, set against a black background.

DETOUR

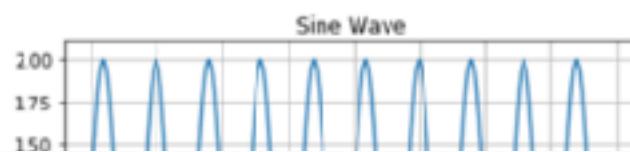
```
In [10]: import numpy as np
import matplotlib.pyplot as plt

# Data for plotting
t = np.arange(0.0, 2.0, 0.01)
s = 1 + np.sin((5 * 2)* np.pi * t)

# Note that using plt.subplots below is equivalent to using
# fig = plt.figure() and then ax = fig.add_subplot(111)
fig, ax = plt.subplots()
ax.plot(t, s)

ax.set(xlabel='time (s)', ylabel='voltage (mV)', title='Sine Wave')
ax.grid()
```

Cut[10]:



"outputs": [

{

 "data": {

 'image/png':

"iVBORw0KGgoAAAANSUhEUgAAAYWAAAECAYAAAB1xKBvAAAABHNCSVQICAgIfAhkiAAAAAlwSFzAAALEgAACxIB0t1+/AAAADl0RVh0U29mdHdcmUAbWF0cGxvdGxpYiB2ZXJzaW9uIDIuMi4yLCBudHRwCi8vbWF0cGxvdGxpYi5vcmlcvhp/UCwAAIABJREFUeJzsvXmcHNd13/s9vc4+2EgABHeQEkkVSXGGRFLembFNSPn7Wyy45i5UXh5ZjvcSy4xcr78WK5bwkzvKSeIlloqaVxZKcOJLN+FHc0dxJEVxxAgQBAiCIdbDP0tPT+80fVdXdmOnl1q17ezBm/T6f+QDdXVXnVtU996z3HFFKESNGjBgxYvRDYrkHECNGjBgxVgZigREjRowYMBQQC4wYMWLEiKGFWGDEiBEjRgwtxAIjRowYMWJoIRYYMWLEiBFDC7HAiBEDEJG/JiKPL/c4YsQ4nxELjBgfGojIXSLyoojMiMg

Jupyter notebooks suck to version control

<https://nextjournal.com/schmudde/how-to-version-control-jupyter>

Clear Output Manually

The simplest solution is to always clear the output before committing. **Cell → All Output → Clear → Save**. This removes any binary blobs that have been generated by the notebook. There are three main drawbacks:

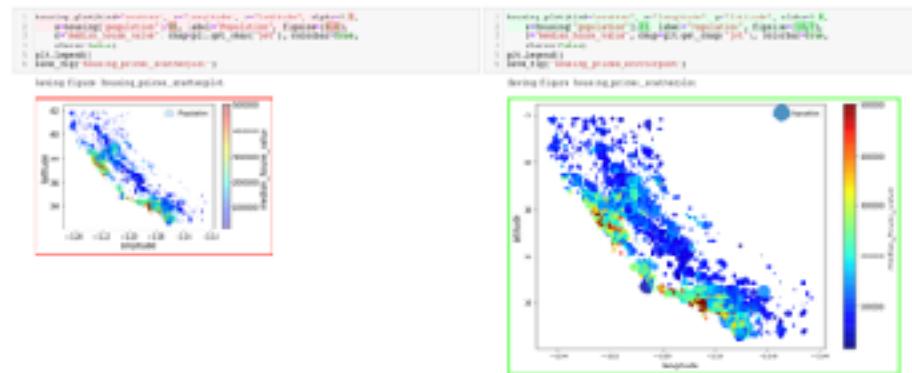
- It is a manual process.
- Collaborators on other machines will need to rerun the notebook to see the output, requiring additional time and setup.

Jupyter notebooks suck to version control

<https://nextjournal.com/schmudde/how-to-version-control-jupyter>

ReviewNB

ReviewNB is a GitHub app that also offers visual diffing with an interface that looks similar to the traditional Jupyter IDE. Because the outputs are visualized, problems associated with committing binary blobs disappear.

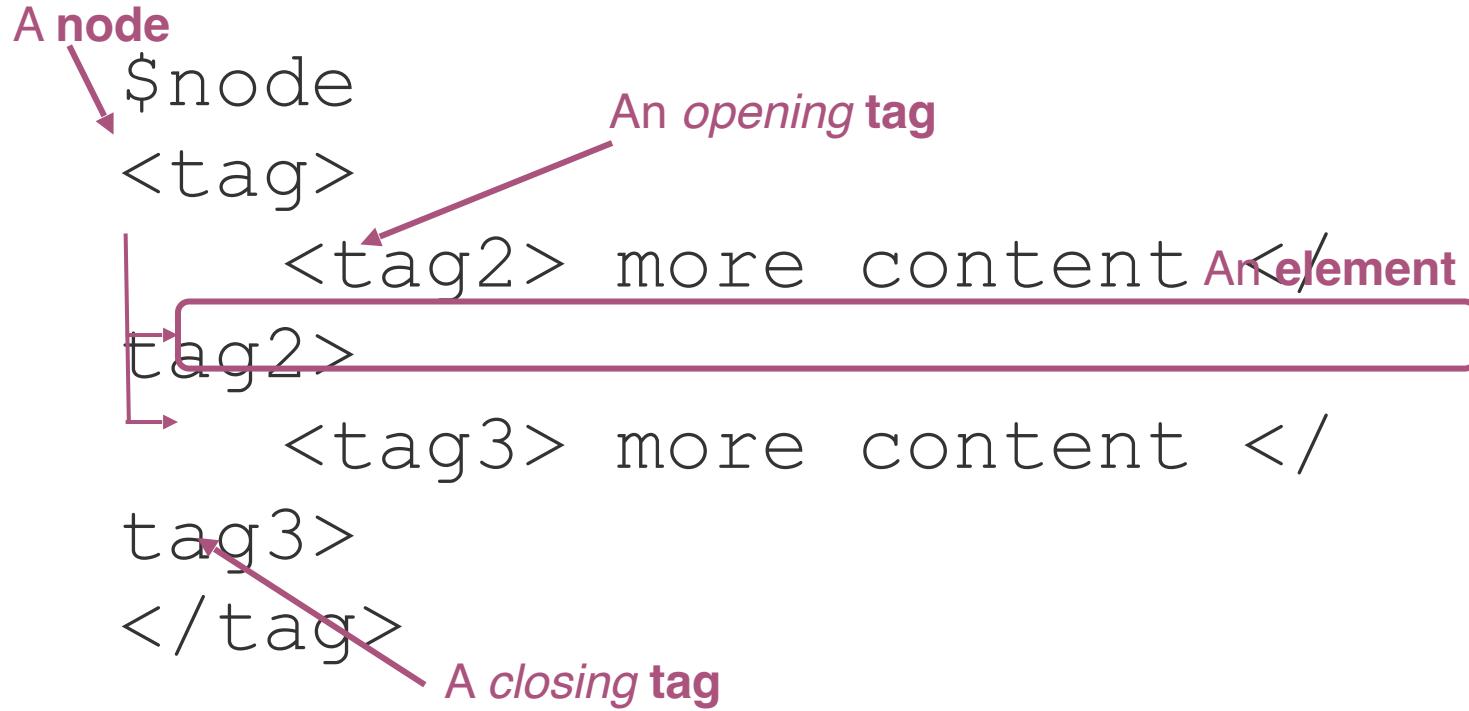


ReviewNB example courtesy of the [ReviewNB website](#)

Back to data formats...

Extensible Markup Language (XML): nodes, tags, and elements

nested/hierarchical data



XML

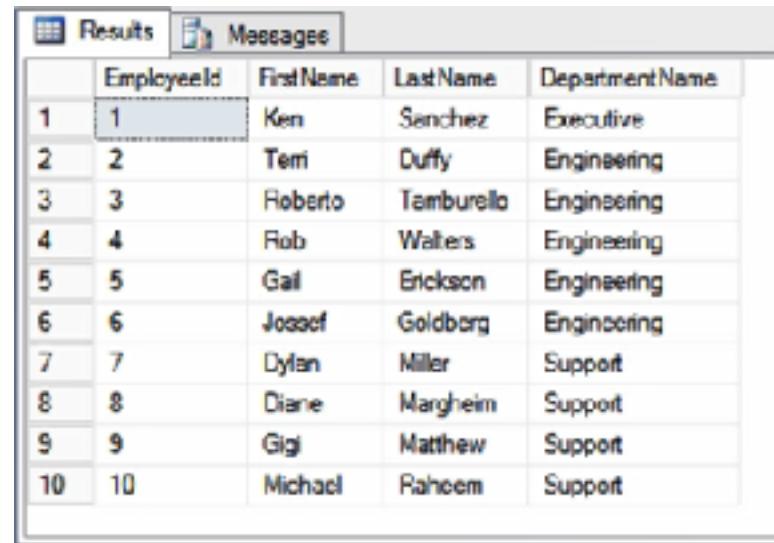
```
<?xml version="1.0" encoding="UTF-8"?>
<customers>
    <customer>
        <customer_id>1</customer_id>
        <first_name>John</first_name>
        <last_name>Doe</last_name>
        <email>john.doe@example.com</email>
    </customer>
    <customer>
        <customer_id>2</customer_id>
        <first_name>Sam</first_name>
        <last_name>Smith</last_name>
        <email>sam.smith@example.com</email>
    </customer>
    <customer>
        <customer_id>3</customer_id>
        <first_name>Jane</first_name>
        <last_name>Doe</last_name>
        <email>jane.doe@example.com</email>
    </customer>
</customers>
```

XML

adapted from Chris Keown

Relational Databases: A set of interdependent tables

1. Efficient Data Storage
2. Avoid Ambiguity
3. Increase Data Privacy



The screenshot shows a Microsoft SQL Server Management Studio (SSMS) interface with the 'Results' tab selected. The results grid displays a table of employee data with the following columns: EmployeeId, FirstName, LastName, and DepartmentName. The data consists of 10 rows, each representing an employee with a unique EmployeeId from 1 to 10, and corresponding FirstName, LastName, and DepartmentName.

	EmployeeId	FirstName	LastName	DepartmentName
1	1	Ken	Sanchez	Executive
2	2	Terri	Duffy	Engineering
3	3	Roberto	Tamburello	Engineering
4	4	Rob	Walters	Engineering
5	5	Gail	Erickson	Engineering
6	6	José	Goldberg	Engineering
7	7	Dylan	Miller	Support
8	8	Diane	Margheim	Support
9	9	Gigi	Matthew	Support
10	10	Michael	Rahiem	Support

relational database

Information is stored across tables

unique_identifier
AH13JK
JJ29JJ
CI21AA

unique_identifier
AH13JK
JJ29JJ
JJ29JJ
XJ11AS
CI21AA

unique_identifier
AH13JK
SE92FE
CI21AA

entries are *related* to one another by their unique identifier

relational database

restaurant

name	id	address	type
Taco Stand	AH13JK	1 Main St.	Mexican
Pho Place	JJ29JJ	192 Street Rd.	Vietnamese
Taco Stand	XJ11AS	18 W. East St.	Fusion
Pizza Heaven	CI21AA	711 K Ave.	Italian

health inspections

id	inspection_date	inspector	score
AH13JK	2018-08-21	Sheila	97
JJ29JJ	2018-03-12	D'eonte	98
JJ29JJ	2018-01-02	Monica	66
XJ11AS	2018-12-16	Mark	43
CI21AA	2018-08-21	Anh	99

rating

id	stars
AH13JK	4.9
JJ29JJ	4.8
XJ11AS	4.2
CI21AA	4.7

relational database

restaurant

name	id	address	type
Taco Stand	AH13JK	1 Main St.	Mexican
Pho Place	JJ29JJ	192 Street Rd.	Vietnamese
Taco Stand	XJ11AS	18 W. East St.	Fusion
Pizza Heaven	CI21AA	711 K Ave.	Italian

Two different restaurants with
the same name will have
different unique identifiers

health inspections

id	inspection_date	inspector	score
AH13JK	2018-08-21	Sheila	97
JJ29JJ	2018-03-12	D'eonte	98
JJ29JJ	2018-01-02	Monica	66
XJ11AS	2018-12-16	Mark	43
CI21AA	2018-08-21	Anh	99

rating

id	stars
AH13JK	4.9
JJ29JJ	4.8
XJ11AS	4.2
CI21AA	4.7

relational database

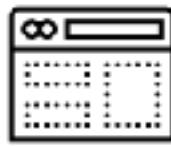
Unstructured Data

Some datasets record information about the state of the world, but in a more heterogeneous way. Perhaps it is a large text corpus with images and links like Wikipedia, or the complicated mix of notes and test results appearing in personal medical records.

Unstructured Data Types



Text files
and
documents



Websites
and
applications



Sensor
data



Image
files



Audio
files



Video
files



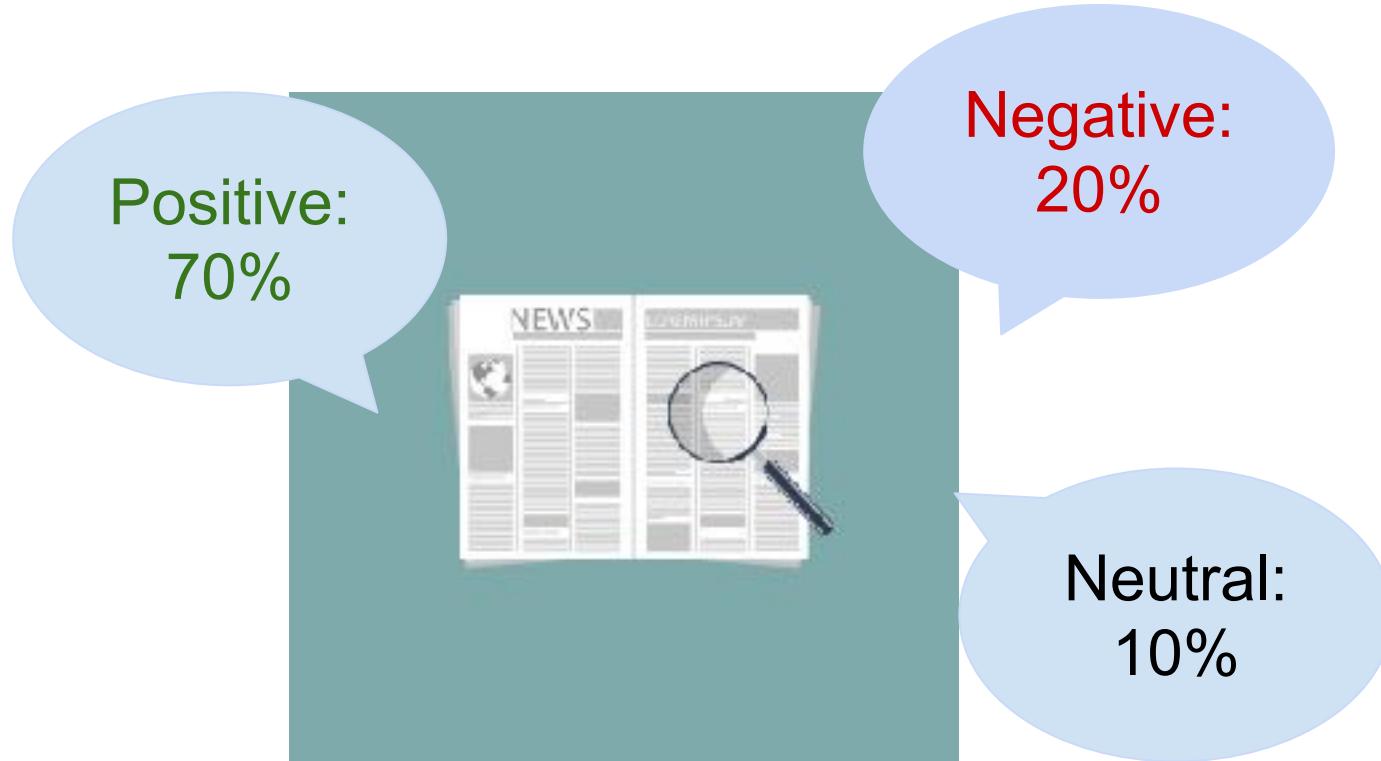
Email
data

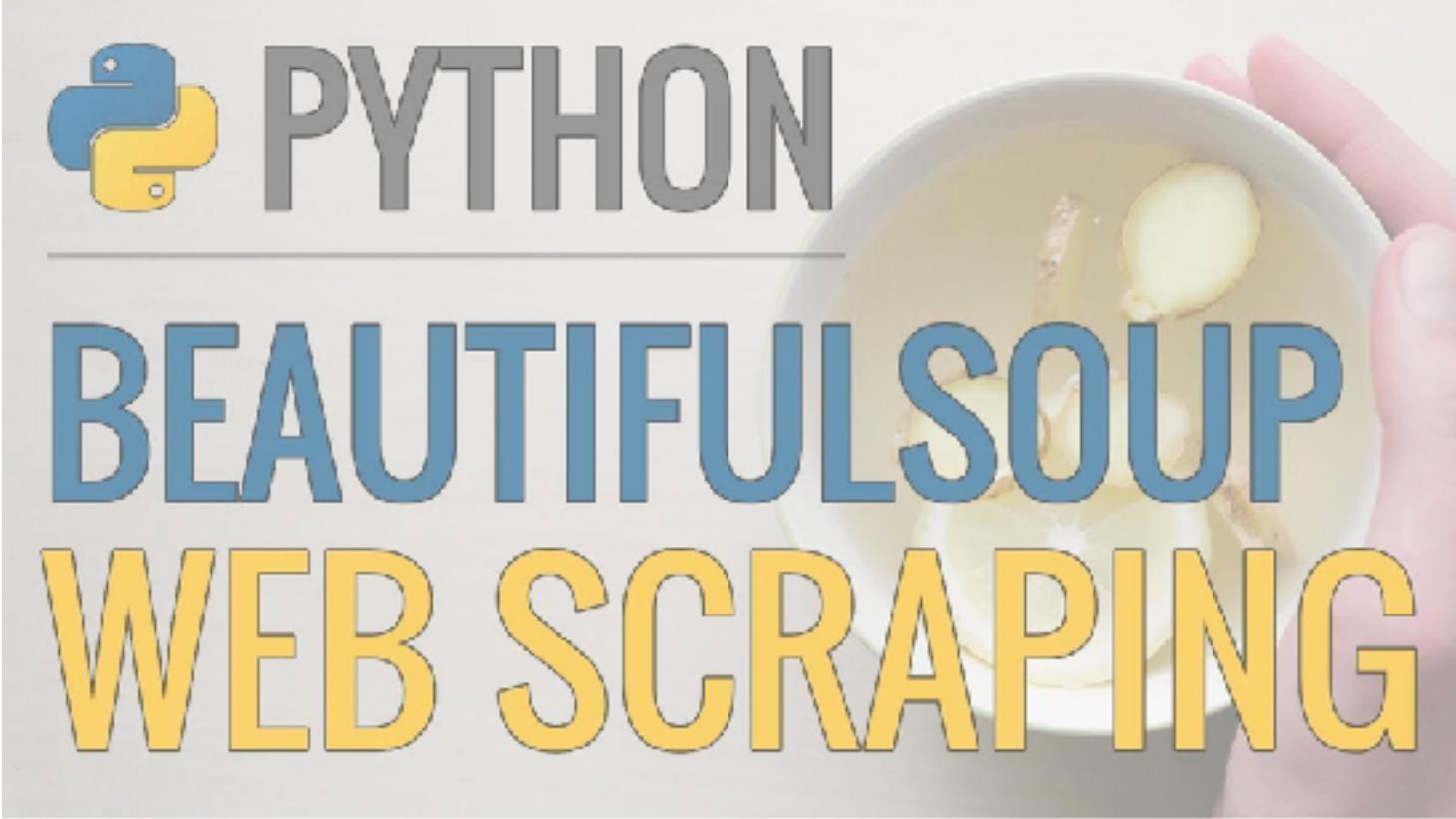


Social
media
data



Text: Sentiment Analysis

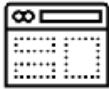




PYTHON

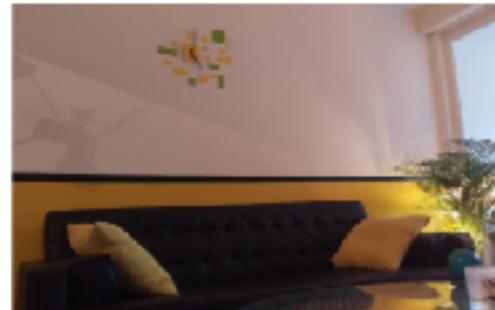
BEAUTIFULSOUP

WEB SCRAPING





Bedroom Or Not?



"The left two photos were correctly predicted as bedrooms; The right two photos were correctly predicted NOT as bedrooms."

Tidy Data

"Good data scientists understand, in a deep way, that the heavy lifting of cleanup and preparation isn't something that gets in the way of solving the problem: it is the problem."

- DJ Patil



Australian Bureau of Statistics

1800.0 Australian Marriage Law Postal Survey, 2017

Released on 15 November 2017

TABLE 5 Participation by Federal Electoral Division(a), Males and Age

(a) The Federal Electoral Roll was current as at 24 August 2017.

(b) Includes those whose age is

(c) Includes Christmas Island and the Cocos (Keeling) Islands

Includes Madalyn Island

(d) Includes ~~Revolving~~ ~~new~~

Return of the table junk

untidy data

Australian Bureau of Statistics										
2010-11 Australian Marriage Laws Postal Survey, 2011										
Assessment on 15% subsample (n=170)										
Table 1: Participation by response location, residence status and age										
Demographic variables										
Regions										
	10-19 years	20-29 years	30-39 years	40-49 years	50-59 years	60-69 years	70-79 years	80-89 years	90-99 years	All ages
Anglo-Celtic	592	2056	2163	2163	2155	2158	2153	2158	2153	2154
Capital cities	272	2350	2350	2350	2350	2350	2350	2350	2350	2350
Non-capital cities	31.8	36.7	36.7	36.7	36.7	36.7	36.7	36.7	36.7	36.7
Primary key variables										
Household										
	Total households	847	3463	3069	2387	2324	1109	1124	1173	1173
Anglo-Celtic	847	3463	3069	2387	2324	1109	1124	1173	1173	1173
Capital cities	760	2993	2993	2364	2364	1107	1107	1107	1107	1107
Non-capital cities	85.9	85.9	85.9	85.9	85.9	85.9	85.9	85.9	85.9	85.9
Geography										
	Total households	744	3039	2635	2053	2053	1054	1054	1054	1054
Anglo-Celtic	744	3039	2635	2053	2053	1054	1054	1054	1054	1054
Capital cities	653	2958	2958	2155	2155	1037	1037	1037	1037	1037
Non-capital cities	85.5	85.5	85.5	85.5	85.5	85.5	85.5	85.5	85.5	85.5
Demographic Subtables										
Covariate and Subtables										
Summary of data in the data										
Variables										
Demographic										
	Total households	3,764	4,788	1,827	6,023	6,406	4,162	5,274	1,168	4,356
Anglo-Celtic	3,764	4,788	1,827	6,023	6,406	4,162	5,274	1,168	4,356	4,356
Capital cities	2,960	4,671	1,840	5,529	5,963	3,606	5,062	1,064	3,637	3,637
Non-capital cities	86.5	86.5	86.5	86.5	86.5	86.5	86.5	86.5	86.5	86.5
Participation										
	Total households	1,473	4,667	1,316	5,766	6,096	4,143	5,264	1,144	4,346
Anglo-Celtic	1,473	4,667	1,316	5,766	6,096	4,143	5,264	1,144	4,346	4,346
Capital cities	1,064	4,554	1,317	5,762	5,960	3,605	5,061	1,062	3,634	3,634
Non-capital cities	71.6	71.6	71.6	71.6	71.6	71.6	71.6	71.6	71.6	71.6
WA, WA, WA										
Australian Capital Territory (ACT)										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Queensland										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
South Australia										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Western Australia										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Tasmania										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327
Anglo-Celtic	3,442	4,474	1,676	5,826	6,192	4,047	5,169	1,127	4,327	4,327
Capital cities	2,638	3,349	1,693	5,822	6,188	3,993	5,155	1,125	4,321	4,321
Non-capital cities	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8	71.8
Victoria										
	Total households	3,442	4,474	1,676	5,826	6,192	4			

1	area	gender	age	State	Area (sq km)	Eligible participants	Participation rate (%)	Total participants	Total Participants
2	Adelaide	Female	18-19 years	SA	76	1341	83.5	1120	1120
3	Adelaide	Female	20-24 years	SA	76	4820	81.2	3750	3750
4	Adelaide	Female	25-29 years	SA	76	4897	81.0	4004	4004
5	Adelaide	Female	30-34 years	SA	76	4784	79.8	3820	3820
6	Adelaide	Female	35-39 years	SA	76	4319	79	3411	3411
7	Adelaide	Female	40-44 years	SA	76	4310	80.6	3472	3472
8	Adelaide	Female	45-49 years	SA	76	4579	81.4	3728	3728
9	Adelaide	Female	50-54 years	SA	76	4475	84.7	3791	3791
10	Adelaide	Female	55-59 years	SA	76	4622	87.3	4033	4033
				SA	76	4342	89.3	3879	3879
				SA	76	3870	90.7	3602	3602
				SA	76	3009	90.3	2716	2716
				SA	76	2158	88.5	1908	1908
				SA	76	1673	85.1	1423	1423

Australian Bureau of Statistics

Table Book

ABBRD Australian Bureau of Statistics Postal Survey, 2017

Released on 05 December 2017

Name & participation by recent electronic messaging status and age

Gender, age, sex

Primary respondents

Mobile phones

Secondary respondents

Mobile phones

Household members

Mobile phones

Non-household members

Mobile phones

Business contacts

Mobile phones

Other contacts

Mobile phones

Business contacts

Tidy Data

1. Each **variable** you measure should be in a single column

	A	B	C	D	E	F	G
1	ID	LastName	FirstName	Sex	City	State	Occupation
2	1004	Smith	Jane	female	Frederick	MD	Welder
3	4587	Nayef	Mohammed	male	Upper Darby	PA	Nurse
4	1727	Doe	Janice	female	San Diego	CA	Doctor
5	6879	Jordan	Alex	male	Birmingham	AL	Teacher

2. Every **observation** of a variable should be in a different row

	A	B	C	D	E	F	G
1	ID	LastName	FirstName	Sex	City	State	Occupation
2	1004	Smith	Jane	female	Frederick	MD	Welder
3	4587	Nayef	Mohammed	male	Upper Darby	PA	Nurse
4	1727	Doe	Janice	female	San Diego	CA	Doctor
5	6879	Jordan	Alex	male	Birmingham	AL	Teacher

3. There should be one table for each type of data

Demographic Survey Data

	A	B	C	D	E	F	G
1	ID	LastName	FirstName	Sex	City	State	Occupation
2		1004	Smith	Jane	female	Frederick	MD
3		4587	Nayef	Mohammed	male	Upper Darby	PA
4		1727	Doe	Janice	female	San Diego	CA
5		6879	Jordan	Alex	male	Birmingham	AL

Doctor's Office Measurements Data

	A	D	E	F	G
1	ID	Height_Inches	Weight_lbs	Insulin	Glucose
2		1004	65	180	0.60
3		4587	75	215	1.46
4		1727	62	124	0.72
5		6879	77	150	1.23

4. If you have multiple tables, they should include a column in each *with the same column label* that allows them to be joined or merged

	A	B	C	D	E	F	G
1	ID	LastName	FirstName	Sex	City	State	Occupation
2	1004	Smith	Jane	female	Frederick	MD	Welder
3	4587	Nayef	Mohammed	male	Upper Darby	PA	Nurse
4	1727	Doe	Janice	female	San Diego	CA	Doctor
5	6879	Jordan	Alex	male	Birmingham	AL	Teacher

	A	D	E	F	G
1	ID	Height_inches	Weight_lbs	Insulin	Glucose
2	1004	65	180	0.60	163
3	4587	75	215	1.46	150
4	1727	62	124	0.72	177
5	6879	77	160	1.23	205

Tidy data == rectangular data

A

	A	B	C	D	E
1	Id	sex	glucose	insulin	triglyc
2	101	Male	134.1	0.60	273.4
3	102	Female	120.0	1.18	243.6
4	103	Male	124.8	1.23	297.6
5	104	Male	83.1	1.16	142.4
6	105	Male	105.2	0.73	215.7

Tidy Data Benefits

1. consistent data structure
2. foster tool development
3. require only a small set of tools to be learned
4. allow for datasets to be combined

TIDY data is **NOT** the same as **CLEAN** data

Stopped here for time



Group Signup is on Canvas!

[Signup link](#)

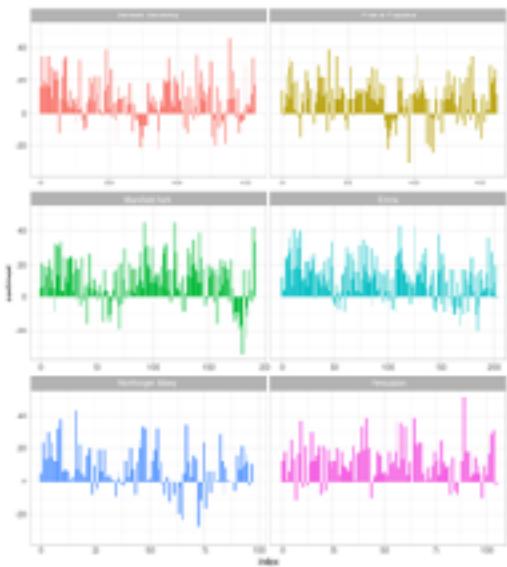
[GITHUB username signup link \(MUST DO!\)](#)

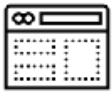
BOTH due Wed Apr 19

results

tidy dataset

Word	Novel	Frequency
good	Emma	359
young	Emma	192
friend	Emma	166





website

1. The first sentence of the following text is from the book *How to Win Friends and Influence People* by Dale Carnegie. The second sentence is from the book *The Art of Persuasion* by Robert H. Schuller. Both sentences are taken from the beginning of their respective books. In each sentence, there is one word that is underlined. In the first sentence, the underlined word is "you". In the second sentence, the underlined word is "you".
2. The first sentence of the following text is from the book *How to Win Friends and Influence People* by Dale Carnegie. The second sentence is from the book *The Art of Persuasion* by Robert H. Schuller. Both sentences are taken from the beginning of their respective books. In each sentence, there is one word that is underlined. In the first sentence, the underlined word is "you". In the second sentence, the underlined word is "you".

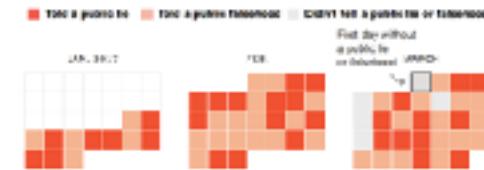


tidy dataset

date	list	description	ref
8 Jan 2011	I wanted another place I didn't want to go to it.	He made an observation where he was against it.	https://archive.is/0mJZL E...
1 Jan 2011 2011	A response to the magazine and I wasn't happy.	Trump was on the cover of Time and News week.	https://archive.is/0mJZL D... https://archive.is/0mJZL D... https://archive.is/0mJZL D...
2 Jan 2011 2011	Because I wanted to be a citizen again after I left.	TRUMP TO ENDURE OF legal ruling.	https://archive.is/0mJZL D...
4 Jan 2011 2011	"Now, the government has the biggest ever bill in..."	Official press photos show Obama 2011 press.	https://archive.is/0mJZL D...
8 Jan 2011 2011	"See a lot of the Iran agents were showed."	The most recent meeting on Iran.	https://archive.is/0mJZL D...

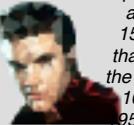


results



text (lyrics)

"I'll be analyzing the repetitiveness of a dataset of 15,000 songs that charted on the Billboard Hot 100 between 1958 and 2017."

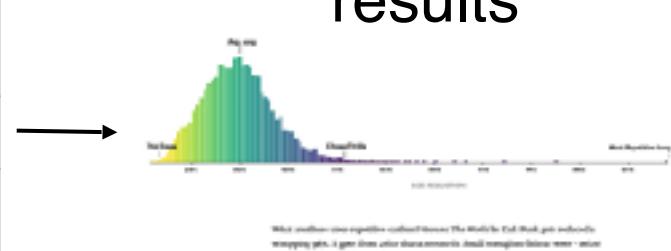


Are Pop Lyrics Getting More Repetitive?

tidy dataset

song	Artist	Released	Reduction
Cheap Thrills	Sia	2016	76
Around The World	Daft Punk	1997	98
Everybody Dies	J. Cole	2018	27

results



Data Intuition

In today's pattern recognition class my professor talked about PCA, eigenvectors and eigenvalues.

1011

I understood the mathematics of it. If I'm asked to find eigenvalues etc. I'll do it correctly like a machine. But I didn't **understand** it. I didn't get the purpose of it. I didn't get the feel of it.

I strongly believe in the following quote:



1375

You do not really understand something unless you can explain it to your grandmother. -- Albert Einstein



Well, I can't explain these concepts to a layman or grandma.

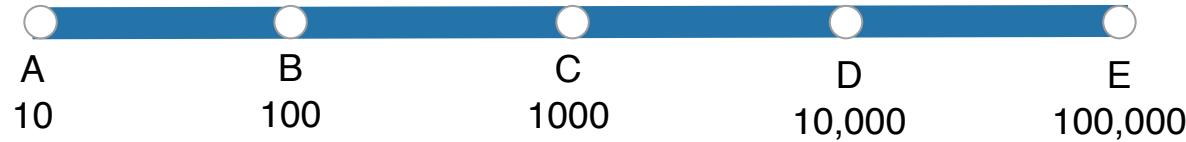
1. Why PCA, eigenvectors & eigenvalues? What was the *need* for these concepts?
2. How would you explain these to a layman?



Fermi Estimation

<https://forms.gle/C982naWtU9RvHqAb7>

Approximately how many piano tuners do you think there are
in the city of Chicago?







<https://www.youtube.com/watch?v=0YzvupOX8ls>

**Has humanity produced enough
paint to cover the entire land area of
the Earth?**

—Josh (Bolton, MA)

Fermi Estimation

<https://forms.gle/shS84W1tai4SDrVF9>



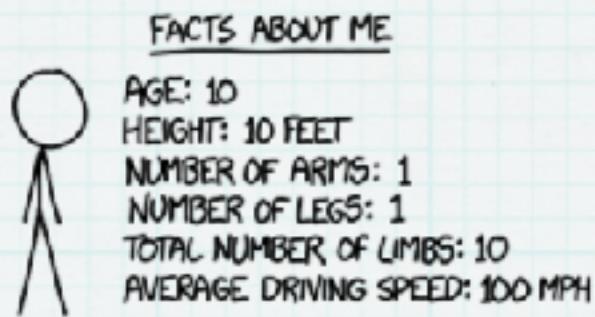
Has humanity produced enough paint to cover the entire land area of the Earth?



This answer is pretty straightforward. We can look up the size of the world's paint industry, extrapolate backward to figure out the total amount of paint produced. We'd also need to make some assumptions about how we're painting the ground. Note: When we get to the Sahara desert, I recommend not using a brush.



But first, let's think about different ways we might come up with a guess for what the answer will be. In this kind of thinking—often called Fermi estimation—all that matters is getting in the right ballpark; that is, the answer should have about the right number of digits. In Fermi estimation, you can round [1] all your answers to the nearest order of magnitude:



Let's suppose that, on average, everyone in the world is responsible for the existence of two rooms, and they're both painted. My living room has about 50 square meters of paintable area, and two of those would be 100 square meters. 7.15 billion people times 100 square meters per person is a little under a trillion square meters – an area smaller than Egypt.

NOT ENOUGH	EXACTLY ENOUGH	MORE THAN ENOUGH
/		

Let's make a wild guess that, on average, one person out of every thousand spends their working life painting things. If I assume it would take me three hours to paint the room I'm in,^[2] and 100 billion people have ever lived, and each of them spent 30 years painting things for 8 hours a day, we come up with 150 trillion square meters ... just about exactly the land area of the Earth.

NOT ENOUGH	EXACTLY ENOUGH	MORE THAN ENOUGH
/	/	

How much paint does it take to paint a house? I'm not enough of an adult to have any idea, so let's take another Fermi guess.

Based on my impressions from walking down the aisles, home improvement stores stock about as many light bulbs as cans of paint. A normal house might have about 20 light bulbs, so let's assume a house needs about 20 gallons of paint.^[3] Sure, that sounds about right.

The average US home costs about \$200,000. Assuming each gallon of paint covers about 300 square feet, that's a square meter of paint per \$300 of real estate. I vaguely remember that the world's real estate has a combined value of something like \$100 trillion,^[4] which suggests there's about 300 billion square meters of paint on the world's real estate. That's about one New Mexico.

NOT ENOUGH	EXACTLY ENOUGH	MORE THAN ENOUGH
//	/	

Of course, both of the building-related guesses could be overestimates (lots of buildings are not painted) or underestimates (lots of things that are not buildings [5] are painted) But from these wild Fermi estimates, my guess would be that there probably isn't enough paint to cover all the land.

So, how did Fermi do?

According to the report **The State of the Global Coatings Industry**, the world produced 34 billion liters of paints and coatings in 2012.

There's a neat trick that can help us here. If some quantity—say, the world economy—has been growing for a while at an annual rate of n —say, 3% (0.03)—then the most recent year's share of the whole total so far is $1 - \frac{1}{1+n}$, and the whole total so far is the most recent year's amount times $1 + \frac{1}{n}$.

If we assume paint production has, in recent decades, followed the economy and grown at about 3% per year, that means the total amount of paint produced equals the current yearly production times 34.^[6] That comes out to a little over a trillion liters of paint. At 30 square meters per gallon,^[7] that's enough to cover 9 trillion square meters—about the area of the United States.

So the answer is no; there's not enough paint to cover the Earth's land, and—at this rate—probably won't be enough until the year 2100.

Stopped here for time

Data Intuition

1. Think about your question and your expectations
2. Do some Fermi calculations (back of the envelope calculations)
3. Write code & look at outputs <- think about those outputs
4. Use your gut instinct / background knowledge to guide you
5. Review code & fix bugs

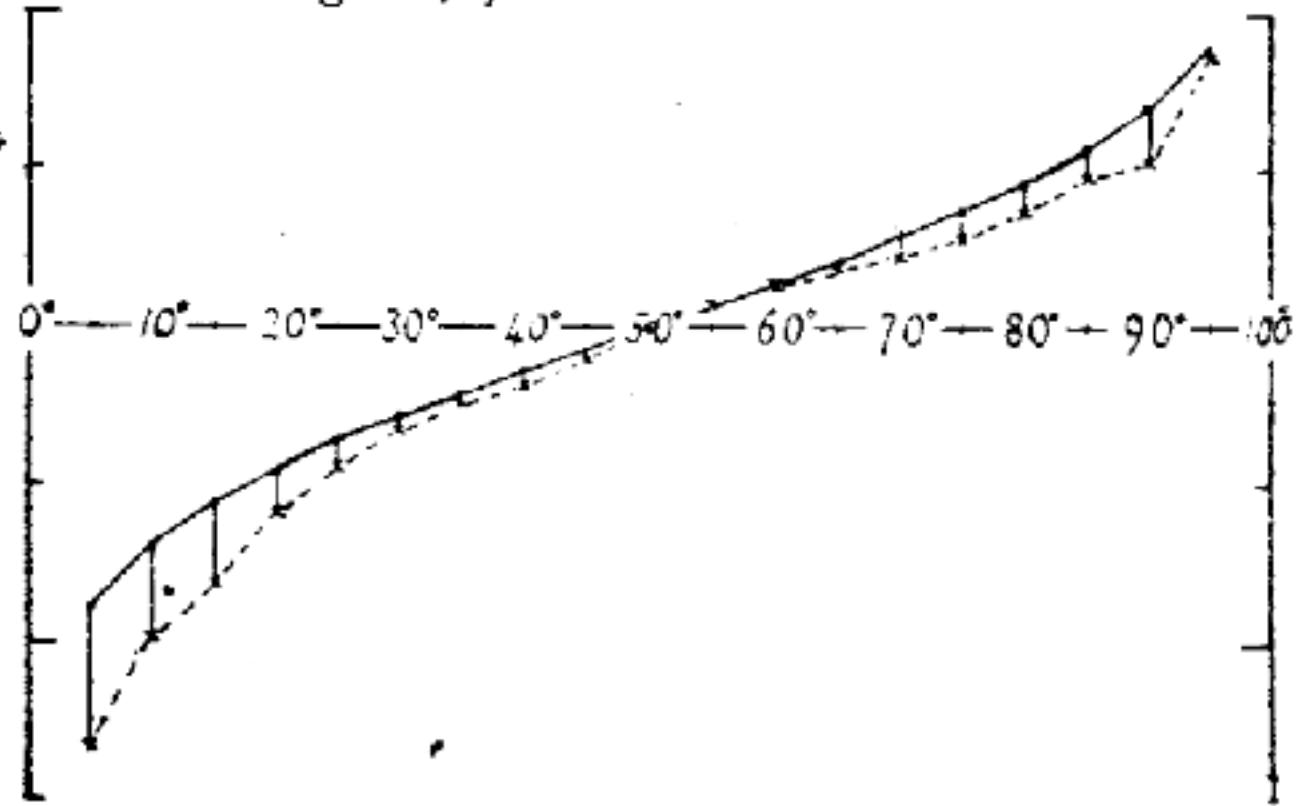
On your own (meaning w/o Googling), please fill out quickly:

<https://forms.gle/CREcpMkYDLYTUp2s6>

Other kinds of
guessing and
intuitions

Diagram, from the tabular values.

Vox Populi



The Wisdom of the Crowds

- Diversity of opinion: Each person should have private information....even if it's just an eccentric interpretation of the known facts
- Independence: People's opinions aren't determined by the opinions of those around them
- Decentralization: People are able to specialize and draw on local knowledge
- Aggregation: Some mechanism exists for turning private judgements into a collective decision

