



Having Fun with Python – the best libraries for EDA and personal portfolios

**Ned Stratton @
PyData London
– Jan 2023**



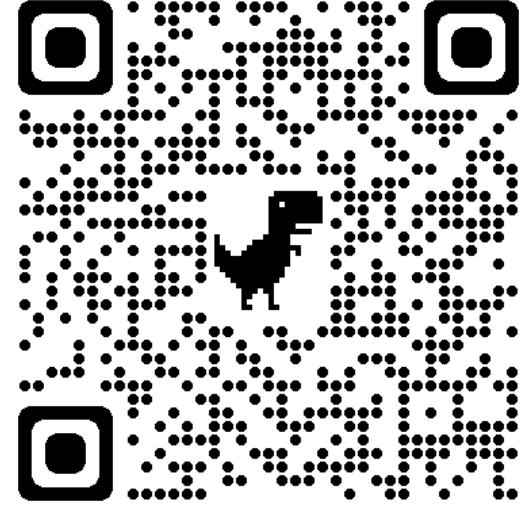
ACTION
BUSINESS INTELLIGENCE

Ned Stratton – Action BI Ltd

Data Analyst & Rehabilitated Humanities Graduate

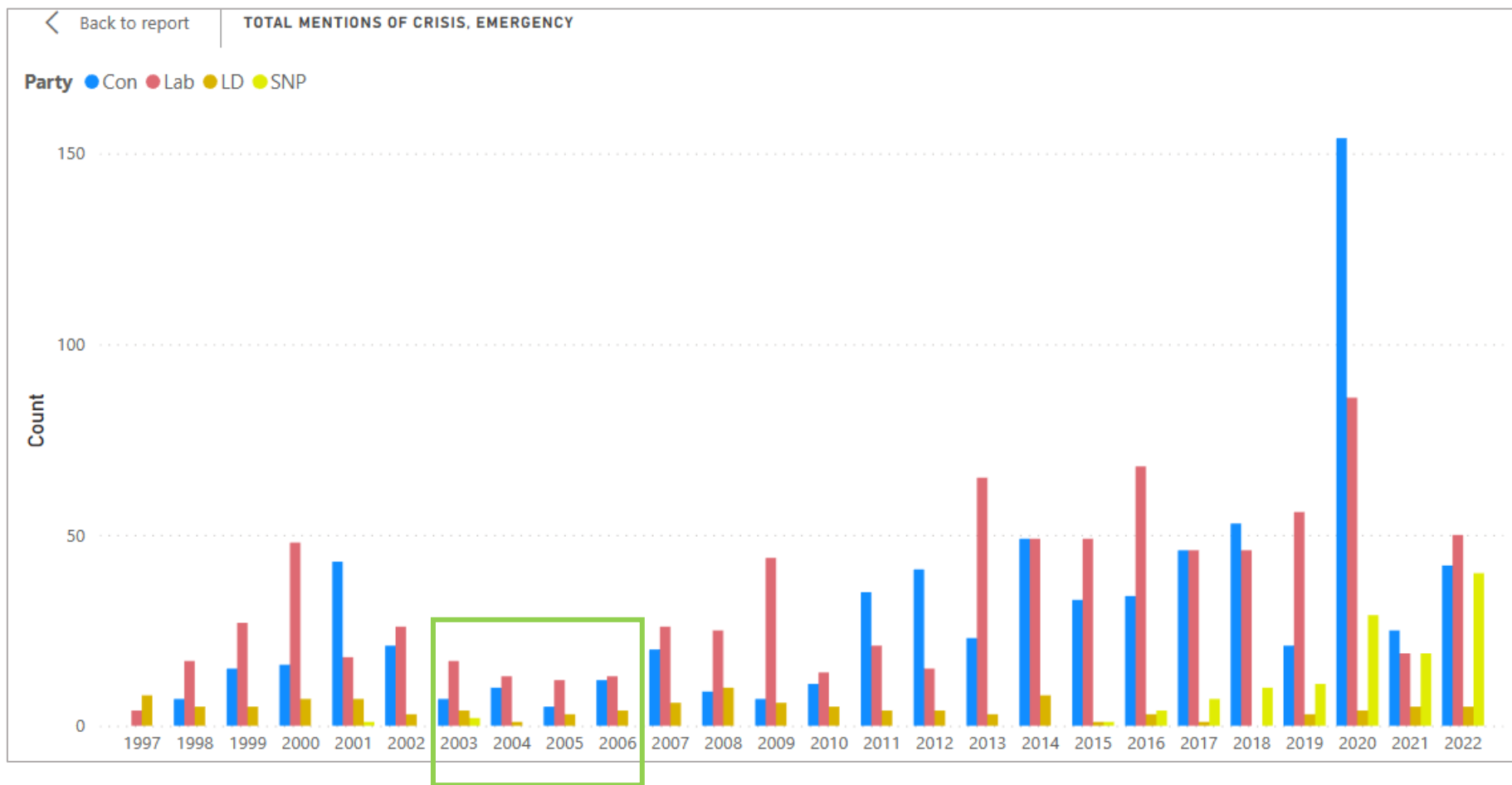
ucovi-data.com

Power BI | SQL | Python | PowerShell | Excel





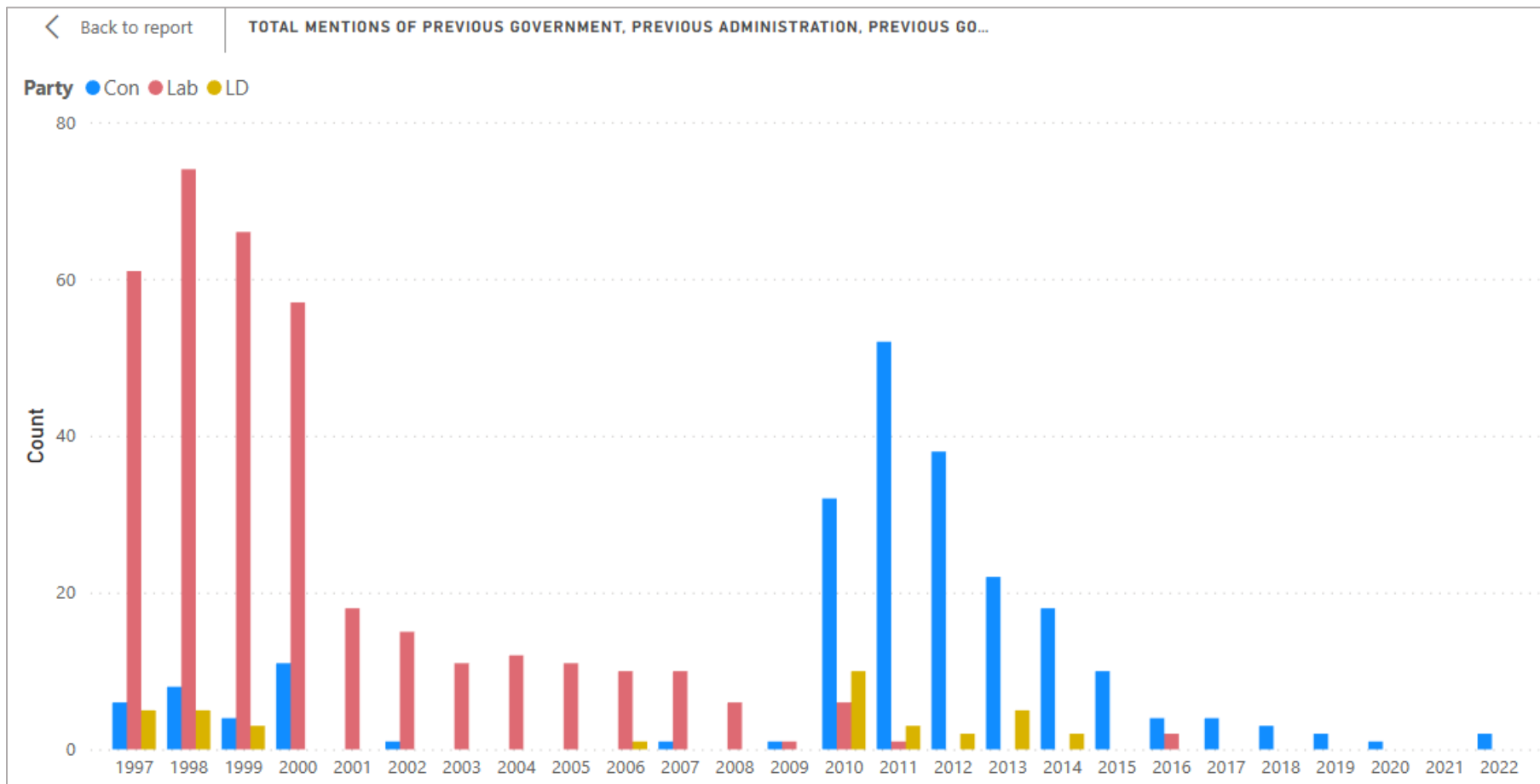
**When was UK politics last
not in a crisis?**



ucovi-data.com/PMQs



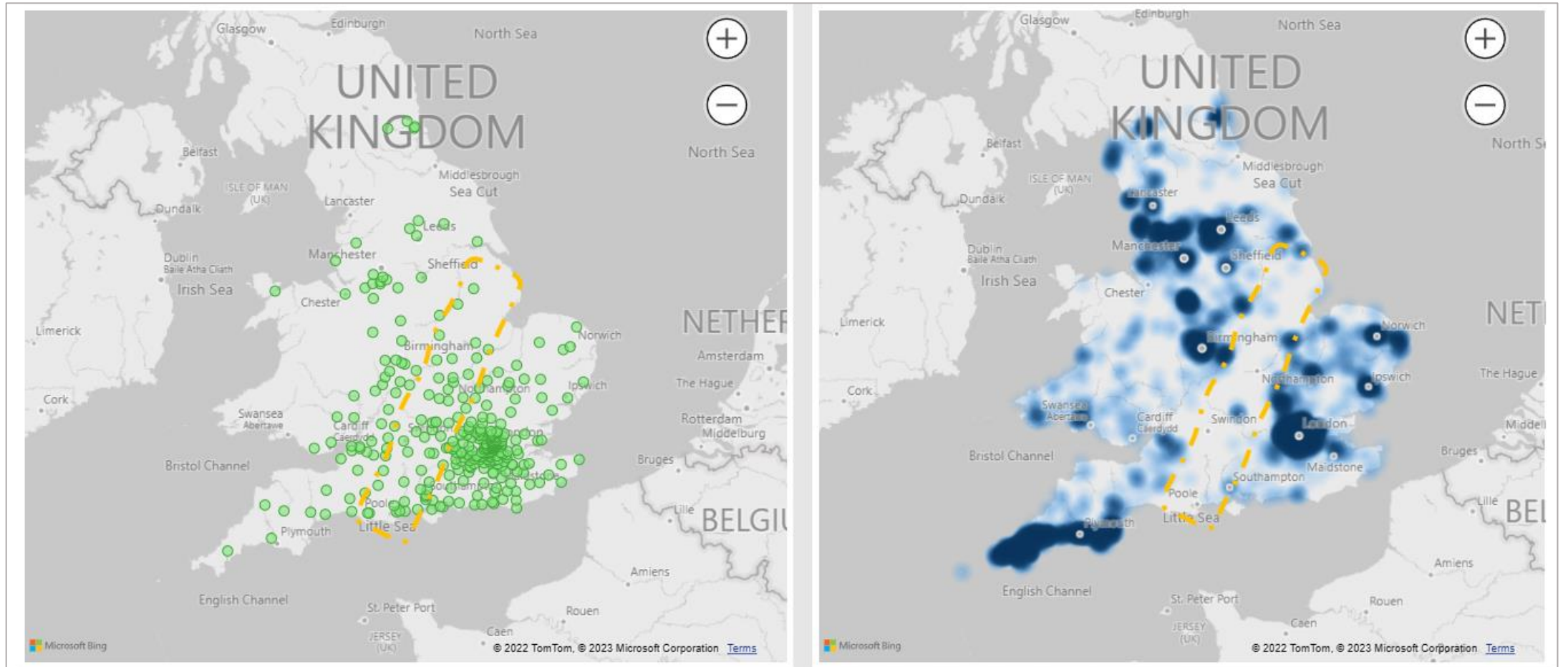
**How long do politicians
blame the previous
government for?**



ucovi-data.com/PMQs



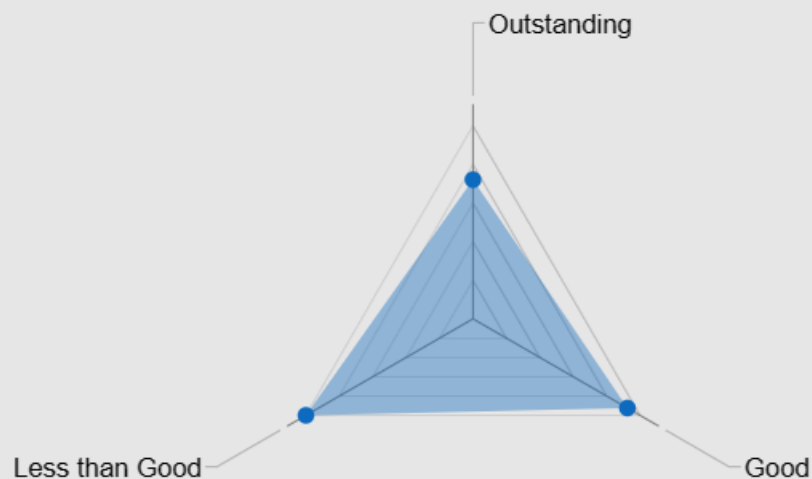
**Is Waitrose a positive
indicator for car accidents
and outstanding schools?**



ucovi-data.com/WaitroseIndex



Average distance to nearest Waitrose by Ofsted rating

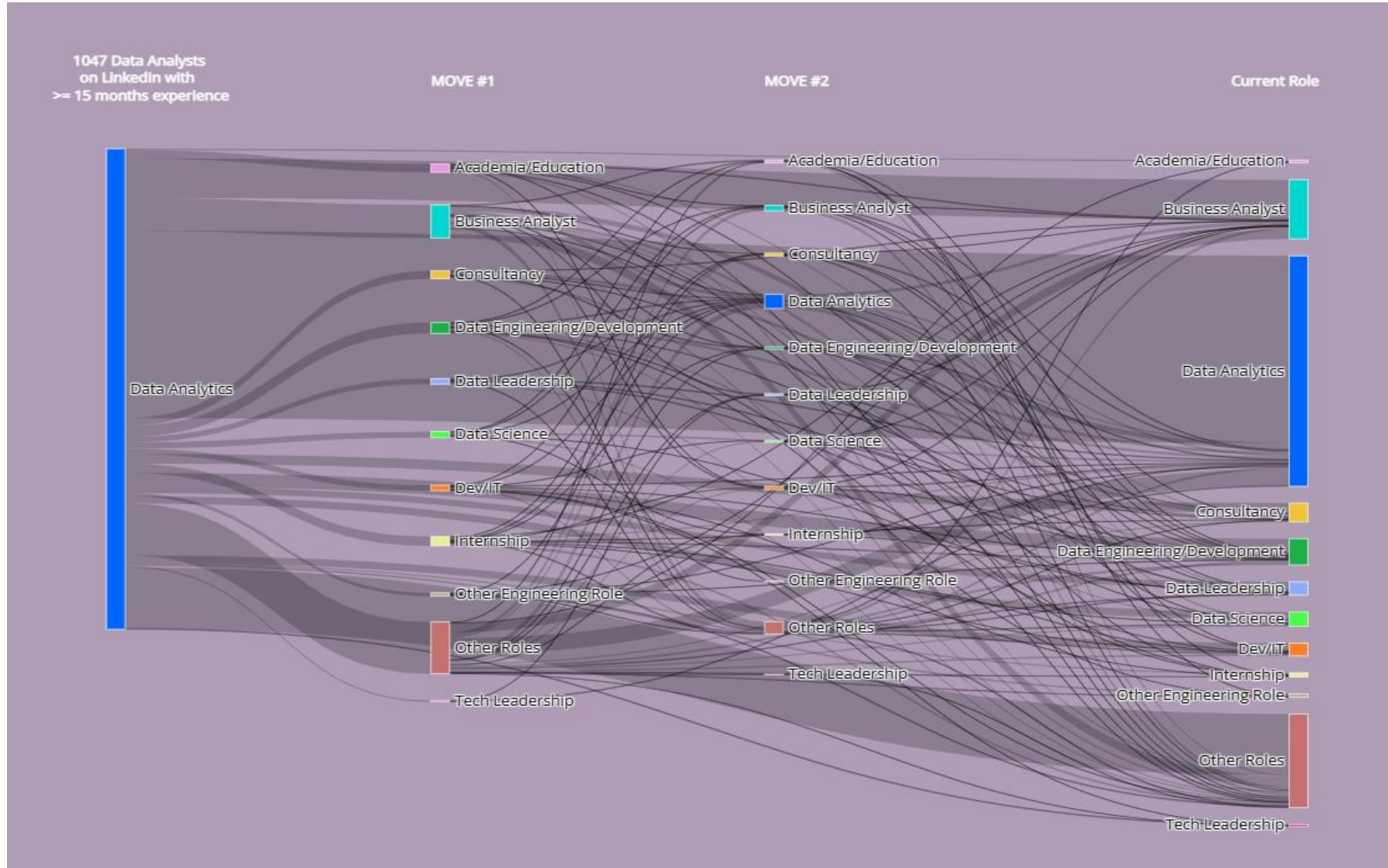


Ofsted Rating	# of Schools	Average KM to nearest Waitrose branch
Outstanding	2,634	9.57
Good	13,189	12.27
Less than Good	1,908	13.28
Total	17,731	11.98

ucovi-data.com/WaitroseIndex



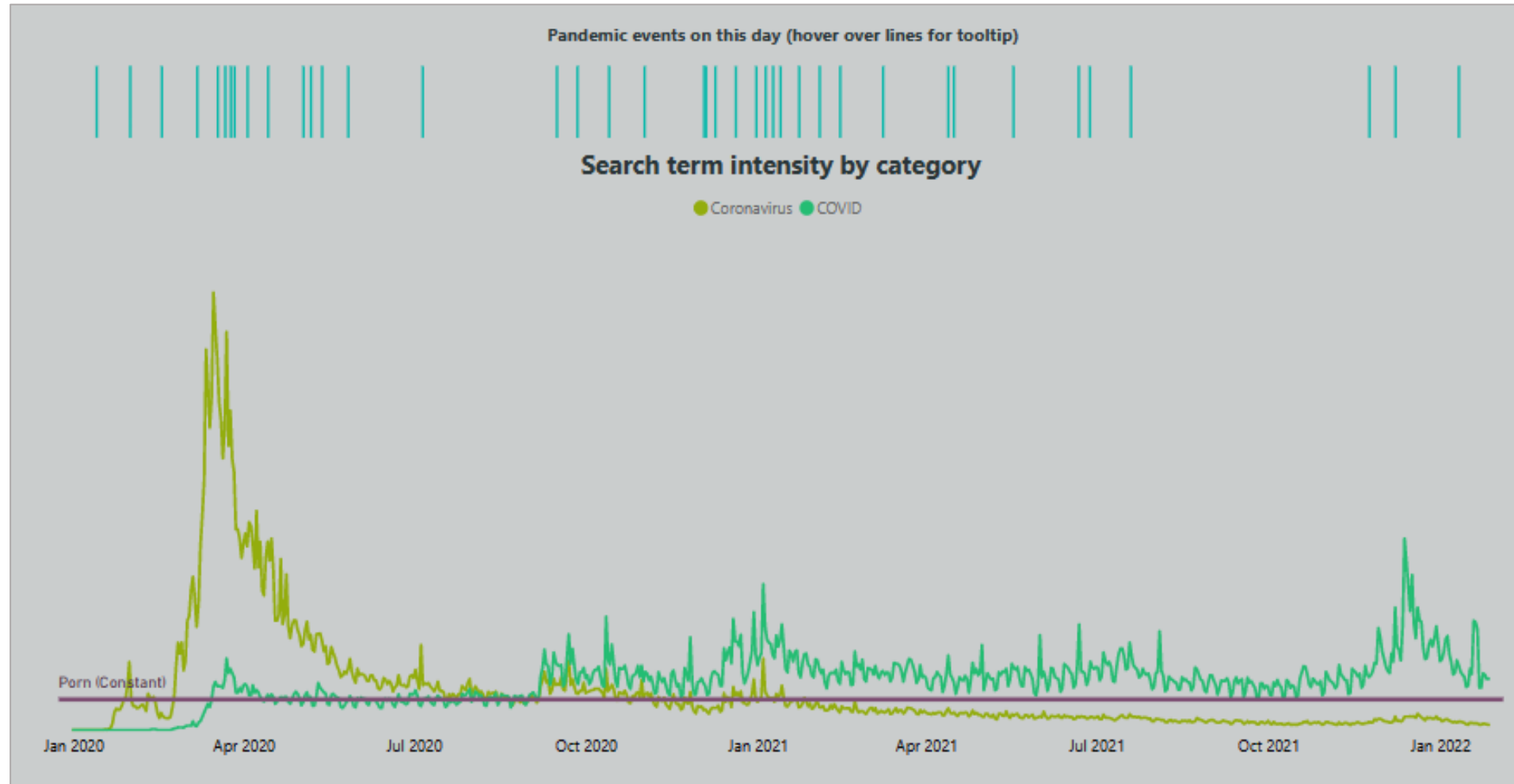
**Is data analytics just a
feeder role for data science,
software dev and IT?**



ucovi-data.com/DataAnalysts



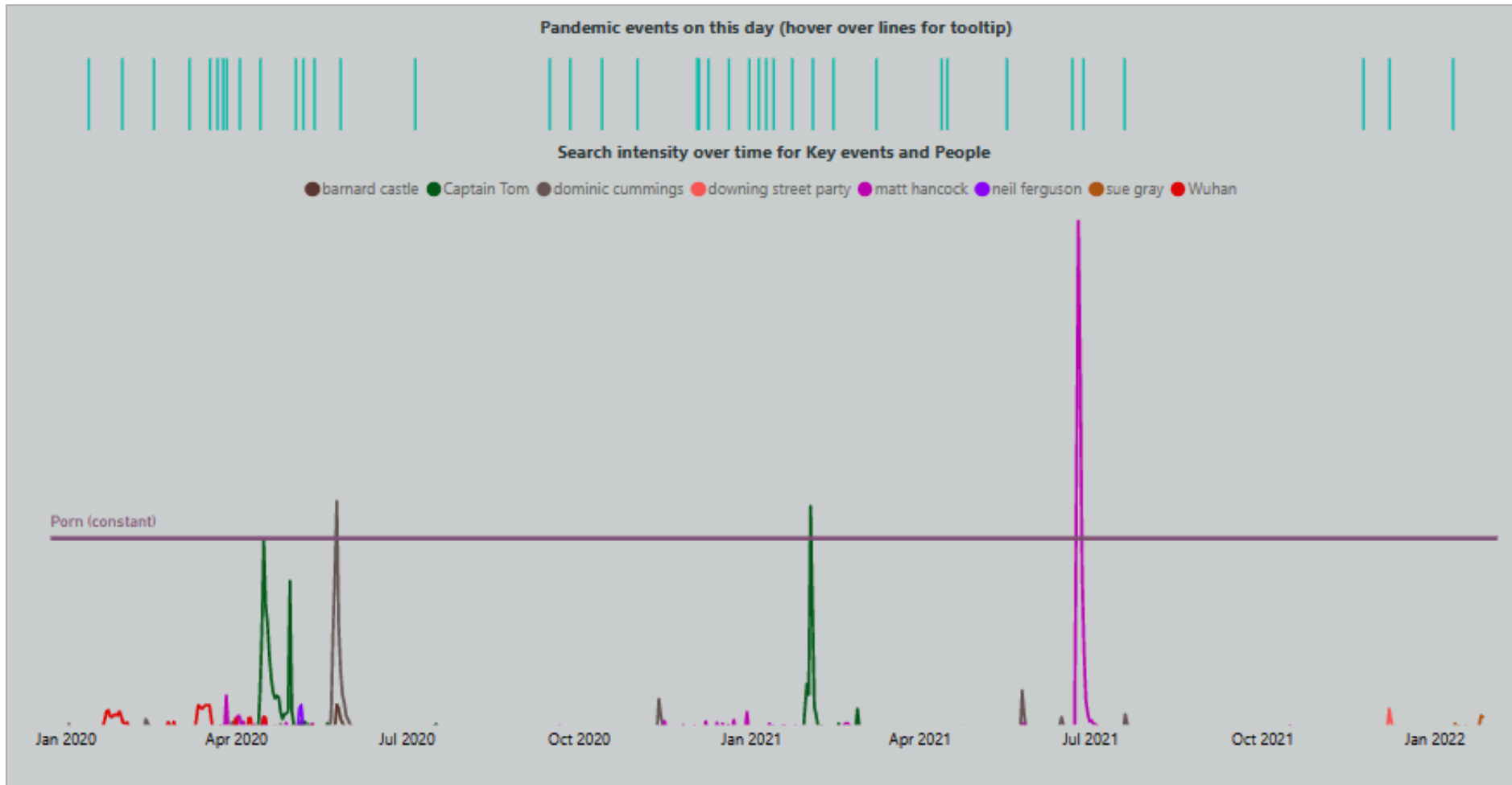
**When did we stop saying
Coronavirus and start
saying COVID?**



ucovi-data.com/CoronaVersusCovid



**Who was the biggest
pandemic villain: Dominic
Cummings or Matt
Hancock?**



ucovi-data.com/CoronaVersusCovid

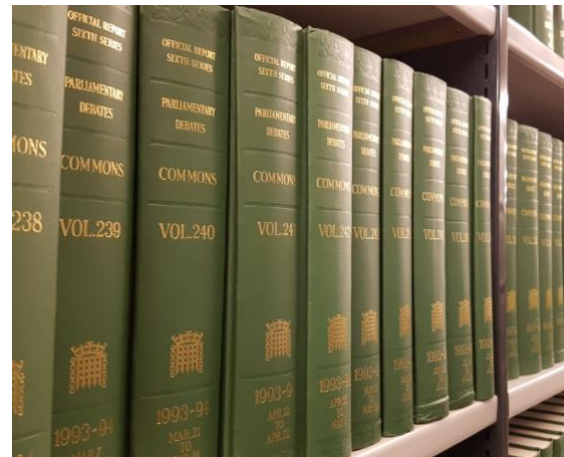


LinkedIn®

Google Trends



Waitrose





Text Analysis: NLTK | re

Web Scraping: selenium | bs4 | requests

Geo Distance Calculation: geopy

APIs: requests | pytrends (relies on pandas)

File Reading and Management: csv | open



PMQs Part 1 - APIs

```
1 #FINAL
2
3 debateq = "https://hansard-api.parliament.uk/debates/debate/{0}.json"

1 #FINAL
2
3 #tosrape is a list of lists, with each inner list being a debate ID
4 appendlist = []
5 for t in tosrape:
6     if len(t[0]) == 0:
7         print('No debates')
8     elif len(t[0]) >= 2:
9         print('More than 1 debate with Questions in title - investigate manually')
10    else:
11        datx = datetime.strptime(t[1], '%d/%m/%Y')
12        datf = datetime.strftime(datx, "%Y-%m-%d")
13        debx = requests.get(debateq.format(t[0][0]))
14        debj = json.loads(debx.text)
15        appendlist.append([t[1], debj])
```



PMQs Part 2 – Regexes to clean API data

(Removing web tags
and other non-natural
language from
transcripts of PMQs
debates speeches on
Hansard.com)

```
3
4 def cleanup(x):
5
6     cleanups_log = []
7     straight_with_space = ['\n', '\t', '\r', '(LD):', '(Con):', '(Lab):', '(Ind):', '(Lab/Co-op):', '\x96',
8                             '\x97', '\u2060', '\u07de', '\u2004', '\xa0', '\u2002', '\u2003', '\u2014']
9     for s in straight_with_space:
10         x = x.replace(s, ' ')
11     straight_w_custom = [['\u1d9c', 'c'],
12                           ['\u011f', 'g'], ['\u0177', 'y'], ['\u0175', 'w'], ['\x93', ''], ['\u201c', ''], ['\u201d', ''],
13                           ['\u2019', ''], ['\x92', ''], ['\x94', ''], ['\u0103', 'a']]
14     for s in straight_w_custom:
15         x = x.replace(s[0], s[1])
16         cleanups_log.append([s[0], s[1]])
17
18     #regex replacements
19
20     # mps names in brackets
21     membreg = ['Member for [A-Za-z\s, \-]{1,40} \([A-Z][A-Za-z]{1,20} [A-Z][A-Za-z]{1,20} \)',
22                'Member for [A-Za-z\s, \-]{1,40} \([A-Z][A-Za-z]{1,14} [A-Z][A-Za-z]{1,15} [A-Z][A-Za-z]{1,15} \)']
23     xmembreg = '|'.join(membreg)
24     if re.search(xmembreg, x):
25         mbmr = re.finditer(xmembreg, x)
26
27         for d in mbmr:
28             new_d = re.sub('\s\([A-Za-z\s]{1,}\)$', '', d[0])
29             x = x.replace(d[0], new_d)
30             cleanups_log.append([d[0], new_d])
31
32     complexr = ['<span ', 'data-house-id="[a-zA-Z0-9\-\-]{1,}"',
33                 'data-volume-number="[a-zA-Z0-9\-\-]{1,}"',
34                 'id="[a-zA-Z0-9\-\-]{1,}"', 'class="[a-zA-Z0-9\-\-]{1,}"',
35                 'data-column-number="[a-zA-Z0-9\-\-]{1,}"', '</span>',
36                 '</Question>', '</QuestionText>', 'Q[1-9]{1,2}\. ',
37                 '<Question HRSContentId="\{[A-Z0-9\-\-]{1,}\}">', '<QuestionText HRSContentId="\{[A-Z0-9\-\-]{1,}\}">',
38                 '<em>Interruption.</em>', '<em>Interruption</em>', '\[Interruption\]', '\[Interruption.\]',
39                 '<em>', '</em>', r'\[Hon\.. Members: "[\Ww]{1,28}"\]', '\[[A-Za-z0-9\.;\s, \-]{1,65}\]', '\[R\] ',
40                 , #for training - remember that the first version didn't have a '-'
41                 '\[', '\]']
42
43     for cm in complexr:
44         x = re.sub(cm, ' ', x) #repeated space to space
45         cleanups_log.append([cm, ' '])
46
47     x = re.sub('\s\s+', ' ', x) #repeated space to space
48     x = re.sub('^s+', '', x) #Leading space
49     x = re.sub('^Q[0-9]+\.', '', x) #Leading Q with number for the questions
50     x = re.sub(' \([ [0-9]+\)\s*', '', x) #question entries brackets
51     x = re.sub(' \([ [0-9]+\]\s*', '', x) #question entries square brackets
52     x = re.sub('\s+$', '', x) #trailing space
53     x = re.sub('<>', '', x) #triangular brackets
54     x = re.sub('^ [0-9]{1,2}\.', '', x) #Leading number at start of question with no Q
55     return x, cleanups_log
56
```



PMQs Part 3 – N-Grams and Lemmatization

Stemming: “Criticized” > “Critic”

Lemmatization: “Criticized” > “Criticize”

N-grams (2): “The dog barked” >
“The dog” | “dog barked”

```
1 import csv
2 from nltk import word_tokenize, ngrams, WordNetLemmatizer, PorterStemmer, sent_tokenize # had to download punkt
3 from nltk.corpus import wordnet as wordn # show if the word is an adjective, noun, verb or adverb
4 from collections import Counter
5 from nltk.corpus import stopwords
6 from nltk.sentiment import SentimentIntensityAnalyzer
7 import pandas as pd
8 import re
9 import os
10 import datetime
11 from textblob import TextBlob
12 import statistics
13 english_sw = stopwords.words('english') # same as comment above for stopwords corpus
```

```
37 lemmatizer = WordNetLemmatizer()
38 ps = PorterStemmer()
```

```
93 def smart_lemma(x):
94     try:
95         wordtype = wordn.synsets(x)[0].pos() # get the first pos tag of the word. 'a' and 's' are adjectives, 'v'=verb
96         if wordtype == 'n':
97             out = lemmatizer.lemmatize(x) #noun is default for lemmatizer. This should standardise plurals
98         elif wordtype == 'v':
99             out = lemmatizer.lemmatize(x, pos='v') #standardize verbs
100         else:
101             #Leave alone
102             out = x
103     except:
104         out = x
105     return out
106
107 def ngramlist(x,y,swmax,bgswitch=False):
108     #return list of distinct ngrams with counts, where x is a list of speeches, y is n grams, and swmax is the limit.
109     #... to the number of REDACTED stopwords that can appear in the gram for it to be counted and returned.
110     #bgswitch tells the function that it is running on bigrams - so the ngrams_useful call should factor this in and
111     #...both words in bigrams against the extended stopwords
112     gramslst = []
113     for spch in x:
114
115         if y > 1:
116             toks = tokensprep(list(word_tokenize(spch)))
117             rawgrams = [list(b) for b in list(ngrams(toks,y))]
118             cleangrams = [re.sub('\.$',''," ".join(c).lower()) for c in rawgrams if ngrams_useful(c,swmax,bgswitch)]
119             gramslst.extend(cleangrams)
120         else:
121             toks = tokensprep(list(word_tokenize(spch.replace('...',' '))))
122             cleangrams = [re.sub('\.$',''," ".join(c).lower()) for c in rawgrams if not re.search('[^a-z]',
123                                                         tko.lower()) and tko.lower() not in english_sw]
124             if len(cleangrams) > 0:
125                 cleangramslemma = [smart_lemma(cl) for cl in cleangrams]
126                 gramslst.extend(cleangramslemma)
127
128     ngram_cts = Counter(gramslst)
129     return [[i[0],i[1]] for i in ngram_cts.items()]
```



Waitrose Branches 1 – Simple Web Scraping

```
1 import requests
2 from bs4 import BeautifulSoup
3 import time
4 import re
5 waitstart = 'https://www.waitrose.com/content/waitrose/en/bf_home/bf/616.html' #Clapham Common for demo
6
7 reqstart = requests.get(waitstart).text
8 optionsoup = BeautifulSoup(reqstart, 'lxml').find('div', class_="branch-finder-form")
9
10 options = []
11 for o in optionsoup.find_all('option'):
12     if not o['value'] == "":
13         options.append([o['value'], o.text])
14
15
16
17 templatelink = 'https://www.waitrose.com/content/waitrose/en/bf_home/bf/{0}.html'
18
19 for o in options:
20     pcodes = []
21     storelink = templatelink.format(o[0])
22     storereq = requests.get(storelink).text
23     storesoup = BeautifulSoup(storereq, 'lxml').find('div', class_="col branch-details").find_all('p')
24     for ss in storesoup:
25         pc = parse_postcode(ss.text) #PARSE POSTCODE is a SIMPLE REGEX func to extract a UK postcode
26         if pc:
27             pcodes.append(pc)
28     o.extend(pcodes)
29     time.sleep(0.5)
30
```



Waitrose Branches 2 – Distance Calculation from Lat/Long Coordinate pairs

```
1 import geopy.distance
2 def distance_km (waitr,sch):
3     return geopy.distance.geodesic(waitr, sch).km

1 for sc in schools_ll: # a csv file of UK schools with Long/lat coordinates read into a list of lists
2
3     sccords = (float(sc[11]),float(sc[12]))
4     distances = []
5     for wt in waitrose_ll: #waitrose branches with long/lat coords
6         sname = wt[1]
7         spc = wt[2]
8         coords = (float(wt[6]),float(wt[7]))
9         dista = round(distance_km(sccords,coords),3)
10        distances.append([sname,spc,dista])
11    nearest = sorted(distances,key=lambda k: k[2])[0]
12    sc.extend(nearest)
```

UK Postcodes & Lats/Longs: <https://geoportal.statistics.gov.uk>

UK Schools Data (has postcode, lat and long for each school: <https://www.get-information-schools.service.gov.uk/Downloads>)



LinkedIn Data Analysts – Advanced Scraping

[LinkedIn Group Members Export | PhantomBuster](#)

<https://phantombuster.com/automations/linkedin/2852/linkedin-group-members-export>

Experience



Microsoft Platform Data Solutions Architect

ACTION BI LTD · Full-time

Sep 2022 - Present · 5 mos

London, England, United Kingdom

Skills: Powershell



Founder and Contributor

UCOVI · Part-time

Aug 2019 - Present · 3 yrs 6 mos



Senior Data Analyst

Fastmarkets · Full-time

Nov 2020 - Sep 2022 · 1 yr 11 mos

London, England, United Kingdom

Data analytics role within the marketing division of a global commodities e agency. I worked principally with the Python-based data blending tool Dat



LinkedIn Data Analysts – Advanced Scraping

```
9 pth = 'C:\\Users\\nedst\\chromedriver.exe'
10
11 randsleeper = {1:11.2,2:13.7,3:14.5,4:12.3,
12               5:14.9,6:13.1,7:11.8,8:15.3,
13               9:15.7,10:14.11:16.7,12:32,13:19,14:21,15:29,16:18.9,17:25,
14               18:21,19:23,20:15.7,21:14.4,22:22.4,23:16.3,24:8.8,25:20.1}
15 chrome = webdriver.Chrome(pth)
16 chrome.get("https://www.linkedin.com/login")
17 sleep(2)
18 chrome.find_element_by_id('username').send_keys(email)
19 chrome.find_element_by_id('password').send_keys(password)
20 chrome.find_element_by_id('password').send_keys(Keys.RETURN)
21 sleep(10)
22
23 for no in notdone[:50]:
24     chrome.get(gurl + no)
25     sleep(1)
26     try:
27         css_sel = 'section.artdeco-card.ember-view.break-words.pb3'
28         experience = [x for x in bs(chrome.page_source,
29                                'lxml').select(css_sel) if x.find('div',{'id':'experience'})]
30         # css selector for when you need to select an element with a space-separated class
31         career = []
32         #classes
33         jtitle = 'mr1 t-bold'
34         timeframe = 't-14 t-normal t-black--light'
35         jtitlei = "mr1 hoverable-link-text t-bold"
36         newcss = 'div.display-flex.flex-row.justify-space-between'
37         #newcss = 'div.display-flex.flex-column.full-width'
38         jobslist = experience[0].find('div',{'class':'pvs-list__outer-container'}).select(newcss)
39         jobslist = [x for x in jobslist if not x.find('a',{'data-field':'experience_company_logo'})]
40         for j in jobslist:
41             templ = []
42             for sp in j.find_all('span'):
43                 if 'class' in sp.attrs:
44                     if ' '.join(sp['class']) in [jtitle,timeframe,jtitlei]:
45                         txtsp = sp.find('span',{'aria-hidden':'true'}).text
46                         templ.append(txtsp)
47             career.append(templ)
48         except:
49             career = ['ERROR']
50         profcareer.append([no,career])
51         sleep(1)
52         chrome.get("https://www.linkedin.com/feed")
53         sleep(randsleeper[random.randint(1,25)])
```



```
1 from time import sleep
2 from selenium import webdriver
3 from selenium.webdriver.common.keys import Keys
4 from selenium.webdriver.common.proxy import Proxy, ProxyType
5 from selenium.webdriver.common.by import By
6 #https://medium.com/nerd-for-tech/linkedin-web-scraper-using
7
8 #https://www.us-proxy.org/
9 from bs4 import BeautifulSoup as bs
```

Tech-Ninja...

```
1 #AUTOMATED CHROME DRIVER DOWNLOAD/INSTALL
2 from selenium import webdriver
3 from webdriver_manager.chrome import ChromeDriverManager
4 chrome = webdriver.Chrome(ChromeDriverManager().install())
5
6 #RUNNING WITHOUT CHROME WINDOW OPENING
7
8 options = webdriver.ChromeOptions()
9 op.add_argument('headless')
10
11 #ONLINE PDF DOWNLOAD AND SAVE
12
13 options.add_experimental_option('prefs', {
14     "download.default_directory": "C:\\Users\\username\\ScriptPDFs\\",
15     "download.prompt_for_download": False,
16     "download.directory_upgrade": True,
17     "plugins.always_open_pdf_externally": True
18 })
```




United Kingdom 9/12/19 - 9/1/22 All categories Web Search

Interest over time

The chart displays search interest for five terms: covid (blue), coronavirus (red), vaccines (yellow), lockdown (green), and omicron (purple). The y-axis represents interest level from 0 to 100. The x-axis shows dates from Sep 15, 2019, to Oct 10, 2021. A vertical line labeled 'Note' is positioned at approximately Oct 10, 2021. A small bar chart in the bottom left corner shows the average interest for each term: covid (blue), coronavirus (red), and lockdown (green).

Search term	Color	Approximate Peak Interest	Approximate Peak Date
covid	Blue	45	Oct 10, 2021
coronavirus	Red	100	Mar 2020
vaccines	Yellow	5	Oct 2021
lockdown	Green	15	Mar 2020
omicron	Purple	10	Jan 2022



Google Trends Search Data – Custom API Library

<https://github.com/GeneralMills/pytrends>

```
1 # library imports
2 from pytrends.request import TrendReq
3 import datetime
4 import csv
5 from time import sleep

1 # pytrends object (function depends on this)
2 pytrends = TrendReq(hl='en-GB', tz=0)

1 #main API call/parse function to get searchterms as % of constant search term within time frames (_dates)
2 def prepare_data(_const,_srch,_dates):
3
4     outlist = []
5
6     def normalize(x,y):
7         x = float(x)
8         y = float(y)
9         return round(x/y,4)
10
11     query_list = _const
12     query_list.extend(_srch) #search terms as list with the constant value the first item
13     pytrends.build_payload(query_list, cat=0, timeframe=_dates, geo='GB', gprop='')
14     querydf = pytrends.interest_over_time()
15
16     for _term in _srch:
17         querydf[_term] = querydf.apply(lambda x: normalize(x[_term],x[_const][0]),axis=1)
18
19     querydf.drop('isPartial',axis=1,inplace=True)
20     querydf.drop(_const[0],axis=1,inplace=True)
21     querydict = querydf.to_dict()
22
23     for i in querydict:
24         for _date in querydict[i]:
25             d = _date.date()
26             outlist.append([i,d,querydict[i][_date]])
27     return outlist
```



```
1 searchterm_batches = [
2     ['Coronavirus','COVID','Wuhan'],
3     ['vaccine','variant','furlough'],
4     ['long COVID','lockdown','NHS','Captain Tom'],
5     ['Tier 1','Tier 2','Tier 3','Tier 4'],
6     ['scotch egg','herd immunity','barnard castle'],
7     ['matt hancock','neil ferguson','dominic cummings'],
8     ['Pfizer','AstraZeneca','Moderna','travel corridor'],
9     ['red list','amber list','green list','South African variant'],
10    ['Brazilian variant','Indian variant','Delta variant'],
11    ['R number','infection rates','COVID symptoms','quarantine'],
12    ['self isolation','NHS app','test and trace','5G'],
13    ['booster jab','vaccine passport','plan b','omicron'],
14    ['downing street party','sue gray']
15 ]
16
17
18 mainlist = []
19 constant = "ryanair"
20 daterng = '2021-08-15 2021-08-21' #range of dates to get daily data
21 for _batch in searchterm_batches:
22     #for _batch in newbatch: # custom for captain tom
23
24     batchdata = prepare_data([constant],_batch,daterng)
25     sleep(4)
26     mainlist.append(batchdata)
```



Cross-Project – File Read and Write Using CSV/Open

```
1 import csv
2
3 #WRITING (to append, switch the 'w' in open with 'a' and it will start the writing...
4 #...from after the last line of the input file var)
5
6 _pmqmaster_rawtd = "C:\\***\\***\\OneDrive\\Documents\\***\\***\\***\\***\\Files\\rawcontribs.txt"
7 with open(_pmqmaster_rawtd, 'w', newline='', encoding='utf-8') as fl:
8     wrt = csv.writer(fl, delimiter='\\t')
9     wrt.writerow(['Name', 'Link', 'RawText', 'Date', 'Role', 'Party', 'Name'])
10
11     for cn in contribs:
12         if cn[1] is not None:
13
14             cnrow = [cn[2], #unrefined text name
15                     '',
16                     '/search/MemberContributions?house=Commons&memberId={0}'.format(cn[1]), # link for mp
17                     cn[4], #raw text
18                     cn[0], #date
19                     mp_data_list[cn[1]][3], #role
20                     mp_data_list[cn[1]][1], #party
21                     mp_data_list[cn[1]][0], #name
22                     ]
23             wrt.writerow(cnrow)
24
25 #READING
26
27 with open(_pmqmaster_rawtd, 'r') as flp:
28     flp.readline() #skips headers
29     rdr = csv.reader(flp, delimiter = '\\t')
30     for row in rdr:
31         pcodes.append(row)
```



Text Analysis: NLTK | re

Web Scraping: selenium | bs4 | requests

Geo Distance Calculation: geopy

APIs: requests | pytrend (relies on pandas)

File Reading and Management: csv | open



Creativity with data:

Dataset-driven or idea-driven?



Why have a personal portfolio anyway?

- ❖ Learning and development
- ❖ 'Keeping your hand in'
- ❖ Display of keenness and interest
- ❖ Sanity – 'Man is born free and yet is everywhere in chains' (Rousseau)



And why Python?

- ❖ Open source
- ❖ Powerful, well-supported, and widely-used
- ❖ Feeder into PowerShell & JavaScript (if that's your thing)
- ❖ OOP & very versatile



Thanks for listening! (Any questions?)

ucovi-data.com

My LinkedIn: <https://www.linkedin.com/in/nedstratton1/>

