



Beginning your own data engineering projects

Conf42: Python 2021

Hui Xiang Chua

What is data engineering?

To ensure consistent data flow for data scientists/ users, typically in the form of a data warehouse to enable large-scale data mining, analytics and reporting purposes

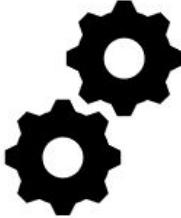


Typical process

Extract



Transform



Load



Typical process



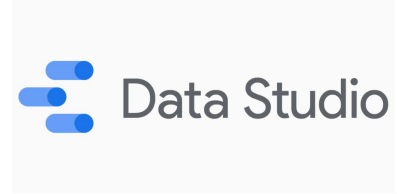
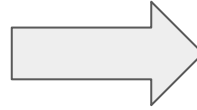
We want to automate this process as much as possible to save time and ensure consistency and accuracy.

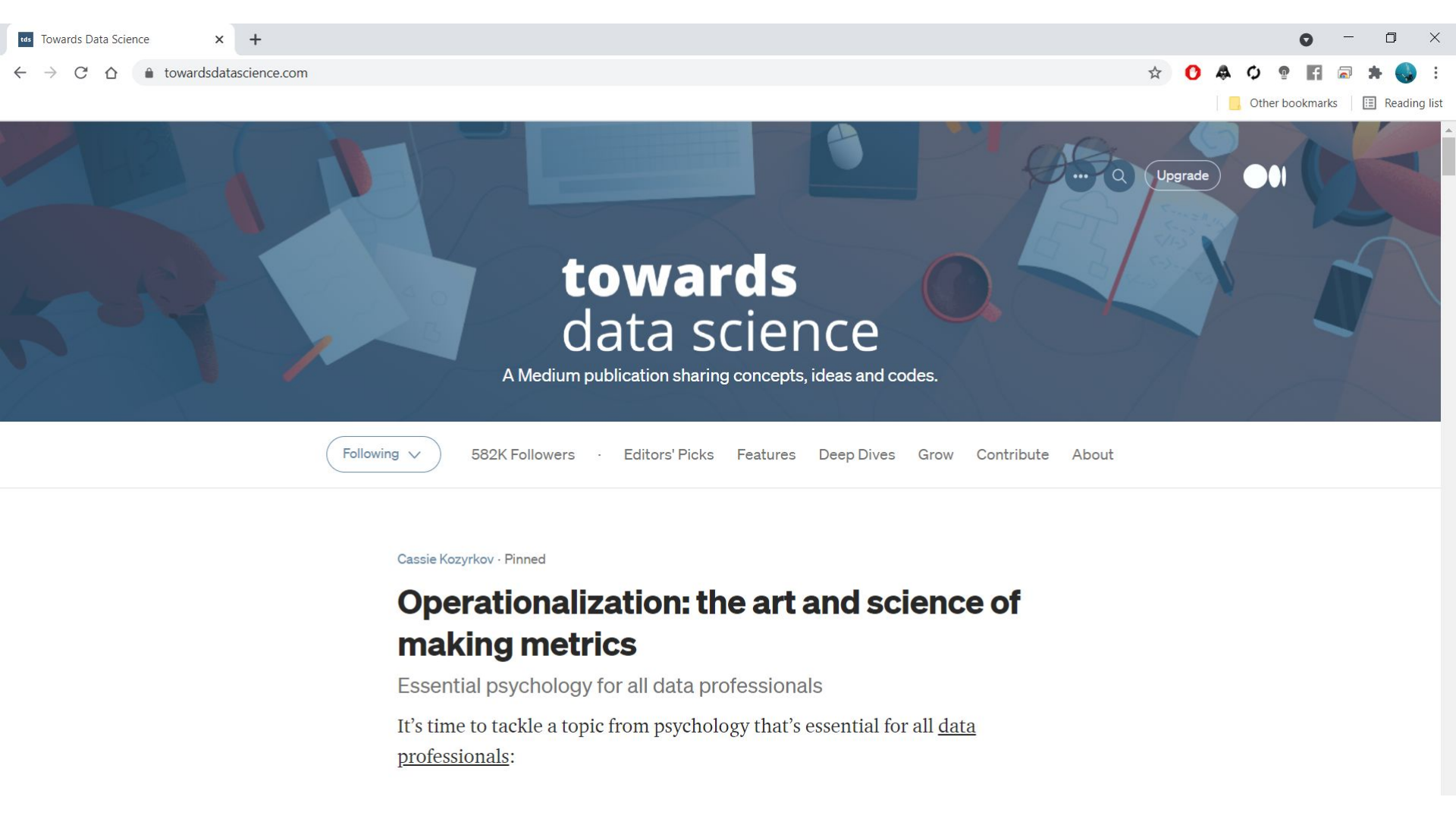
Use case #1

Towards Data Science listing



Google Sheets





towards data science

A Medium publication sharing concepts, ideas and codes.

Following

582K Followers · Editors' Picks Features Deep Dives Grow Contribute About

Cassie Kozyrkov · Pinned

Operationalization: the art and science of making metrics

Essential psychology for all data professionals

It's time to tackle a topic from psychology that's essential for all data professionals:

towards
data science

Medium_TDS_Titles

File Edit View Insert Format Data Tools Add-ons Help

Last edit was on January 2

100%

\$ % .0 .00 123

Default (Ari...

10

B I S A

A1

fx

author

	A	B	C	D	E
1	author	recency	post_title	time_extracted	
2	Youness Mansar	-1 day ago	Learning to Play CartPole and LunarLander with Proximal Policy Optimization	2020112321292	
3	Rose Day	-1 day ago	15 Topics to Consider as You Review Code in Data Science	2020112321292	
4	Sara A. Metwalli	-1 day ago	NLP 101: Towards Natural Language Processing	2020112321292	
5	Richmond Alake	-1 day ago	Interesting AI/ML Articles On Medium This Week (Nov 22)	2020112321292	
6	Arthur Kakande	-1 day ago	Analysis of Uganda's Social Media Data Regarding the 2021 General Presidential Elections	2020112321292	
7	Sanjay Singh	-1 day ago	Spark	2020112321292	
8	Florent Poux, Ph.D.	-1 day ago	How to automate LiDAR point cloud sub-sampling with Python	2020112321292	
9	Audhi Apriliant	-1 day ago	Hands-on Tutorial	2020112321292	
10	Alan Jones	-1 day ago	How to Scrape Dynamic Web pages with Selenium and BeautifulSoup	2020112321292	
11	Idil Ismiguzel	-1 day ago	Linear Regression Model with Python	2020112321292	
12	Sebastian Poliak	-1 day ago	1 to 5 Star Ratings — Classification or Regression?	2020112321292	
13	Dimitris Pouloupoulos	-1 day ago	Mini Kubeflow on AWS is your new ML workstation	2020112321292	
14	Vicky Yu	-1 day ago	Advice for New Data Analysts	2020112321292	
15	Acusio Bivona	-1 day ago	A Tutorial on Scraping Images from the Web Using BeautifulSoup	2020112321292	
16	Sam Watts	-1 day ago	Building a Custom Semantic Segmentation Model	2020112321292	
17	Tan Pengshi Alvin	-1 day ago	A Comprehensive Guide to Metis Data Science Bootcamp	2020112321292	
18	Mahnoor Javed	-1 day ago	Building a Facial Recognition Model using PCA & SVM Algorithms	2020112321292	
19	Brian Mwangi.	-1 day ago	Analysis of Teen Pregnancy in Kenya Using R Shiny	2020112321292	
20	mugoh mwaura	-1 day ago	9 Promising Applications of Reinforcement Learning in 2021	2020112321292	
21	Yash Indulkar	-1 day ago	Twitter Sentimental Analysis & Algorithm Comparison for Uber & Ola Using 'R'	2020112321292	
22	Alberta Odamea Anim-Ayeko	-1 day ago	Python Tweet Deleter.	2020112321292	
23	Paul Torres	-1 day ago	Querying Databases	2020112321292	
24	Terence S	-1 day ago	All Machine Learning Algorithms You Should Know in 2021	2020112321292	
25	Joshua Yeung	-1 day ago	Pass CKAD (Certified Kubernetes Application Developer) with a Score of 97!	2020112321292	

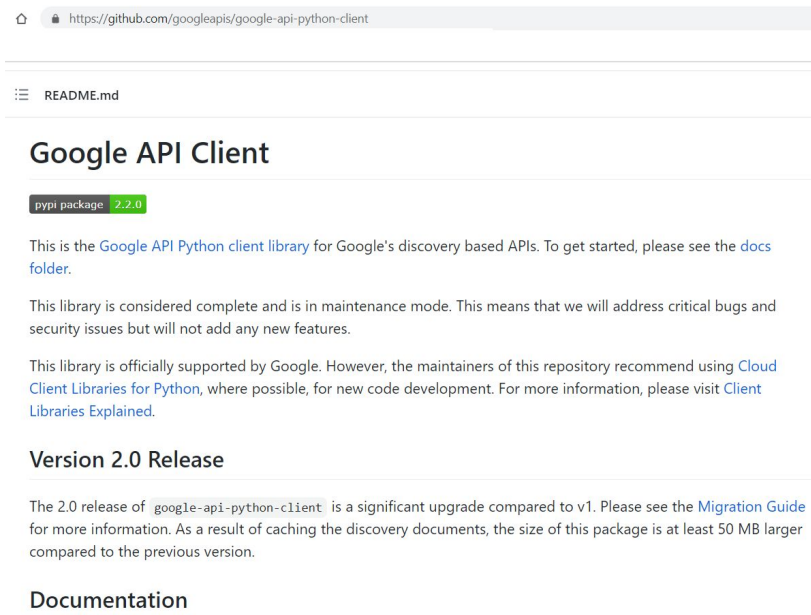
+

≡

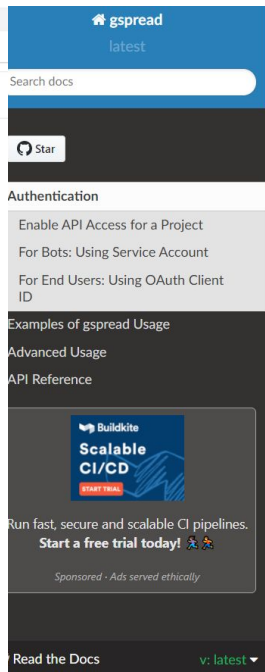
Sheet1


```
10 from __future__ import print_function
11 import numpy as np
12 import pandas as pd
13 import gspread_dataframe as gd
14 import gspread
15 from google.oauth2.service_account import Credentials
16 from bs4 import BeautifulSoup
17 from datetime import datetime
18 from selenium import webdriver
19 from selenium.webdriver.chrome.options import Options
20 import time
21
```

To use Google API, follow the instructions [here](#) for installation and [here](#) for authentication. To use gspread and gspread_dataframe, please run `pip install gspread` and `pip install gspread_dataframe`.



The screenshot shows the GitHub repository for the Google API Client. The URL in the browser is `https://github.com/googleapis/google-api-python-client`. The page title is "Google API Client". Below the title, it says "pypi package 2.2.0". The main text describes the library as the "Google API Python client library" for Google's discovery based APIs. It mentions that the library is in maintenance mode and that the maintainers recommend using "Cloud Client Libraries for Python" for new code development. There is a section for "Version 2.0 Release" which states that the 2.0 release of `google-api-python-client` is a significant upgrade compared to v1. The "Documentation" section points to the `docs` folder for more detailed instructions.



The screenshot shows the documentation page for the gspread library. The header is blue with the "gspread" logo and the word "latest". There is a search bar and a "Star" button. The "Authentication" section is highlighted, with sub-sections for "Enable API Access for a Project" (for bots and end users) and "Examples of gspread Usage". There is also a "Scalable CI/CD" banner for Buildkite.

Docs » Authentication

[Edit on GitHub](#)

Authentication

To access spreadsheets via Google Sheets API you need to authenticate and authorize your application.

- If you plan to access spreadsheets on behalf of a bot account use [Service Account](#).
- If you'd like to access spreadsheets on behalf of end users (including yourself) use [OAuth Client ID](#).

Enable API Access for a Project

1. Head to [Google Developers Console](#) and create a new project (or select the one you already have).
2. In the box labeled "Search for APIs and Services", search for "Google Drive API" and enable it.
3. In the box labeled "Search for APIs and Services", search for "Google Sheets API" and enable it.

For Bots: Using Service Account

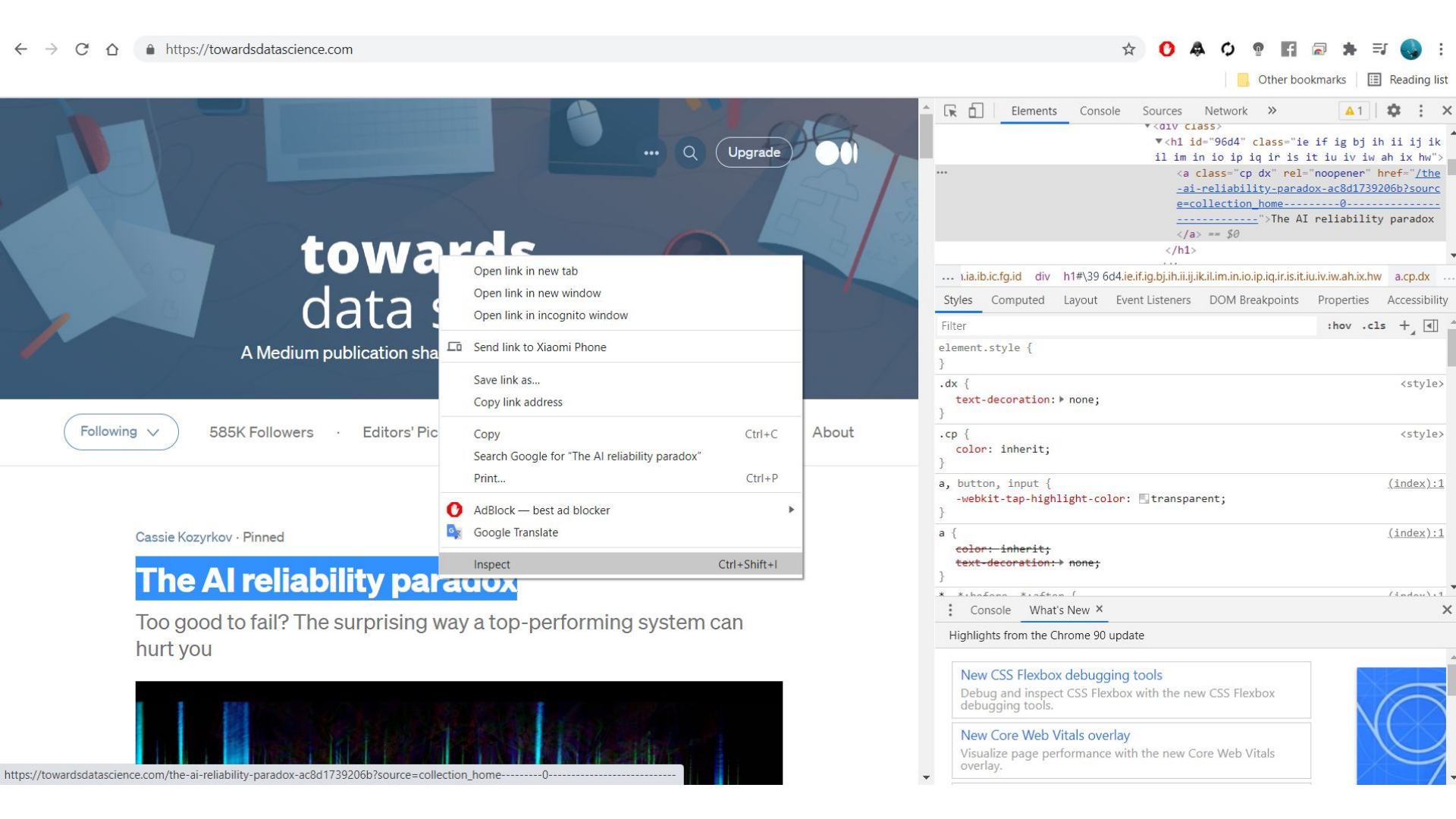
A service account is a special type of Google account intended to represent a non-human user that needs to authenticate and be authorized to access data in Google APIs [sic].

Since it's a separate account, by default it does not have access to any spreadsheet until you share it with this account. Just like any other Google account.

Here's how to get one:

See the [docs](#) folder for more detailed instructions and additional documentation.

```
22 # If modifying these scopes, delete the file token.pickle.
23 SCOPES = ['https://www.googleapis.com/auth/spreadsheets',
24           'https://www.googleapis.com/auth/drive']
25
26 # The ID and range of a sample spreadsheet.
27 SAMPLE_SPREADSHEET_ID = '17e0aX_IcL6fF1GpUDaaaiGqCFijYorWtFrnKy7vUcUI' ##change this
28 #SAMPLE_RANGE_NAME = 'Sheet1!A:E' ##change this
29
30
31
32 credentials = Credentials.from_service_account_file('service_account.json', scopes=SCOPES)
33 gc = gspread.authorize(credentials)
34
35 # Connecting with `gspread` here
36
37 ws = gc.open_by_key(SAMPLE_SPREADSHEET_ID).worksheet("Sheet1")
38 existing = gd.get_as_dataframe(ws)
39 #print(existing)
40
41 #scraping the site
42 options = Options()
43 options.add_argument('--headless')
44 options.add_argument('--disable-gpu')
45 driver = webdriver.Chrome('C:/Users/hxchu/env/Lib/site-packages/chromedriver',options=options)
46 site = "https://towardsdatascience.com/"
47 driver.get(site)
48 #time.sleep(3)
49
```



The image below shows how a DAG is a unidirectional, acyclic graph, where each node in the graph is a task and edges define dependencies among tasks. ...

[Read more](#) · 4 min read



Tirthajyoti Sarkar · 10 hours ago ★

PyTest for Machine Learning — a simple

Show more ▾

```
50 for i in range(20):
51     try:
52         loadMoreButton = driver.find_element_by_xpath("//button[text()='Show more']")
53         #find_element_by_xpath("//button[contains(@aria-label,'Show more')]")
54         time.sleep(2)
55         loadMoreButton.click()
56         time.sleep(2)
57     except Exception as e:
58         print(e)
59         break
60 print("Complete")
61 time.sleep(10)
62 page = driver.page_source
63 driver.quit()
64
65 soup = BeautifulSoup(page, 'html.parser')
66
67 text_ls = []
68 for i in range(len(soup.find_all('a'))):
69     text = soup.find_all('a')[i].get_text()
70     text_ls.append((text))
71
72 # size of the list
73 size = len(text_ls)
74 index_ls = []
75 # looping till length - 2
76 for i in range(size - 2):
77
78     # checking the conditions
79     if text_ls[i] == text_ls[i + 1] and text_ls[i + 1] == text_ls[i + 2]:
80
81         # printing the element as the
82         # conditions are satisfied
83         #print(i)
84         index_ls.append((i))
85
```



```
text_ls
```

```
''  
,  
'Paul May',  
'·9 hours ago',  
'Named Tuples: A Little Known Machine Learning Helper',  
'Glenn Carstens-Peters',  
'Unsplash',  
'Read more · 8 min read',  
''  
,  
''  
,  
''  
,  
'Dean McGrath',  
'·11 hours ago',  
'How to Combine Python, Pandas & XlsxWriter',  
'Safar Safarov',  
'Unsplash',  
'XlsxWriter',  
'DataFrame',  
'pip',  
'Read more · 3 min read',  
''  
,
```

```

86 data_ls = []
87 for i in index_ls[:-2]:
88     name = text_ls[i+3]
89     recency = text_ls[i+4]
90     if(len(text_ls[i+5]) <20):
91         title = text_ls[i+7]
92     else:
93         title = text_ls[i+5]
94     data_ls.append((name,recency,title))
95
96 df = pd.DataFrame(np.array(data_ls))
97 df.columns = ['author','recency','post_title']
98
99 dt = datetime.now()
100 seq = int(dt.strftime("%Y%m%d%H%M%S"))
101 df['time_extracted']=seq
102
103
104 #appending with latest data
105 updated = existing.append(df)
106 print(updated)
107 updated = updated[['author','recency','post_title','time_extracted']]
108 updated = updated.drop_duplicates(['author','post_title'])
109 updated.fillna('', inplace=True)
110 updated.drop(updated.index[0], inplace=True)
111 #updated.dropna(axis=0,how='all',inplace=True)
112 ws.update([updated.columns.values.tolist()] + updated.values.tolist())
113
114
115
116 ###References:
117 #https://stackoverflow.com/questions/23377533/python-beautifulsoup-parsing-table
118 #https://stackoverflow.com/questions/31328861/python-pandas-replacing-header-with-top-row
119 #https://www.graspdata.tech/solved-the-caller-does-not-have-permission-using-the-api-with-a-private-goog
120 #https://gspread.readthedocs.io/en/latest/

```



```
data_ls = []
for i in index_ls[:-2]:
    name = text_ls[i+3]
    recency = text_ls[i+4]
    if(len(text_ls[i+5]) <20):
        title = text_ls[i+7]
    else:
        title = text_ls[i+5]
    data_ls.append((name,recency,title))
```

data_ls

```
[('Paul May',
  '.9 hours ago',
  'Named Tuples: A Little Known Machine Learning Helper'),
 ('Dean McGrath',
  '.11 hours ago',
  'How to Combine Python, Pandas & XlsxWriter'),
 ('Dan Baker',
  '.12 hours ago',
  'Deploying your Dash App to Heroku – THE MAGICAL GUIDE'),
 ('Shivangi Sareen', '.13 hours ago', 'Terminal Multiplexers to the Rescue'),
 ('Rashida Nasrin Sucky',
  '.14 hours ago',
  'A Complete Neural Network Algorithm from Scratch in Python'),
 ('Erdem Isbilen', '.14 hours ago', 'A Complete Guide to Python Lists'),
 ('Travis Tang (Voon Hao)',
  '.14 hours ago',
  'Dangers of a Pizza-Making Robot'),
 ('Brandon Smith',
  '.14 hours ago',
  'Type I and Type II Errors in COVID-19 Serology Testing')]
```

```
Anaconda Prompt (Anaconda3) - python tds_script.py

(base) C:\Users\hxchu>cd Documents

(base) C:\Users\hxchu\Documents>cd BACKUP2020

(base) C:\Users\hxchu\Documents\BACKUP2020>goog\Scripts\activate.bat

(goog) (base) C:\Users\hxchu\Documents\BACKUP2020>python tds_script.py
_
```

Home

PUBLIC

Questions

Tags

Users

FIND A JOB

Jobs

Companies

TEAMS

Stack Overflow for Teams – Collaborate and share knowledge with a private group.



Create a free Team

38



Creating the exe should be the best method. But if you want to run it v
can do it in this way:

1. Launch Window's Task Scheduler

2.

3.

To e
uses
(64b
Pyth

After determining the location of python.exe, this is what is entered in the Action panel of the task

scheduler:

Program/script:

C:\Python27\ArcGIS10.2\python.exe

Browse...

Add arguments (optional):

"E:\My script.py"

Start in (optional):

First, open the Control Panel and then click on the **Administrative Tools**:



Next, double-click on the **Task Scheduler**, and then choose the option to 'Create Basic Task...'



Type a name for your task (you can also type a description if needed), and then press Next.

Suppose the script you want to run is E:\My script.py. Instead of
instruct the task scheduler to run python.exe with the script as an argument. For example:

C:\Python27\ArcGIS10.2\python.exe "E:\My script.py"

The location of python.exe depends on your install. If you don't know where it is, you can discover its location; copy and paste the following code into a new Python script then execute the script. The script will print the location of python.exe as well as other information about your Python environment.

towards
data science

Medium_TDS_Titles ☆ ↻ ⌂

File Edit View Insert Format Data Tools Add-ons Help *Last edit was on January 2*

↶ ↷ 🖨️ 🔍 100% ▼ \$ % .0 .00 123 ▾ Default (Ari... ▾ 10 ▾ B I S A 🔗 📐 📊 ▾ ☰ ▴ ▾ ▸ ▹ ► ▸ ▹ ► 🔍 🔗 + 📄 🔍 ▾

A1	fx	author			
		A	B	C	D
1		author	recency	post_title	time_extracted
2		Youness Mansar	:1 day ago	Learning to Play CartPole and LunarLander with Proximal Policy Optimization	2020112321292
3		Rose Day	:1 day ago	15 Topics to Consider as You Review Code in Data Science	2020112321292
4		Sara A. Metwalli	:1 day ago	NLP 101: Towards Natural Language Processing	2020112321292
5		Richmond Alake	:1 day ago	Interesting AI/ML Articles On Medium This Week (Nov 22)	2020112321292
6		Arthur Kakande	:1 day ago	Analysis of Uganda's Social Media Data Regarding the 2021 General Presidential Elections	2020112321292
7		Sanjay Singh	:1 day ago	Spark	2020112321292
8		Florent Poux, Ph.D.	:1 day ago	How to automate LiDAR point cloud sub-sampling with Python	2020112321292
9		Audhi Apriliant	:1 day ago	Hands-on Tutorial	2020112321292
10		Alan Jones	:1 day ago	How to Scrape Dynamic Web pages with Selenium and BeautifulSoup	2020112321292
11		Idil Ismiguzel	:1 day ago	Linear Regression Model with Python	2020112321292
12		Sebastian Poliak	:1 day ago	1 to 5 Star Ratings — Classification or Regression?	2020112321292
13		Dimitris Pouloupoulos	:1 day ago	Mini Kubeflow on AWS is your new ML workstation	2020112321292
14		Vicky Yu	:1 day ago	Advice for New Data Analysts	2020112321292
15		Acusio Bivona	:1 day ago	A Tutorial on Scraping Images from the Web Using BeautifulSoup	2020112321292
16		Sam Watts	:1 day ago	Building a Custom Semantic Segmentation Model	2020112321292
17		Tan Pengshi Alvin	:1 day ago	A Comprehensive Guide to Metis Data Science Bootcamp	2020112321292
18		Mahnoor Javed	:1 day ago	Building a Facial Recognition Model using PCA & SVM Algorithms	2020112321292
19		Brian Mwangi.	:1 day ago	Analysis of Teen Pregnancy in Kenya Using R Shiny	2020112321292
20		mugoh mwaura	:1 day ago	9 Promising Applications of Reinforcement Learning in 2021	2020112321292
21		Yash Indulkar	:1 day ago	Twitter Sentimental Analysis & Algorithm Comparison for Uber & Ola Using 'R'	2020112321292
22		Alberta Odamea Anim-Ayeko	:1 day ago	Python Tweet Deleter.	2020112321292
23		Paul Torres	:1 day ago	Querying Databases	2020112321292
24		Terence S	:1 day ago	All Machine Learning Algorithms You Should Know in 2021	2020112321292
25		Joshua Yeung	:1 day ago	Pass CKAD (Certified Kubernetes Application Developer) with a Score of 971	2020112321292

+ ☰ Sheet1 ▾

This dashboard shows titles of articles published in [Towards Data Science](#) on Medium between a six-week period, 21 Nov 2020 and 2 Jan 2021. Data is scraped into a [Google Sheet](#) and then populated into tables here. This allows us to get a sense of all articles published within the day quickly. We can make use of this list to self-select articles that are of interest to us instead of solely depending on Medium's recommendations. And of course, we can also do text mining to identify popular entities/ topics/ keywords! Get to the article easily by copying the title and searching for it on Towards Data Science.

Prepared by Hui Xiang Chua. For feedback/ comments/ questions, pls write to datadoubleconfirm@gmail.com.

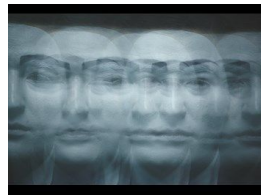
	author	Record Count ▾
1.	Soner Yildirim	32
2.	Dimitris Pouloupoulos	17
3.	Richmond Alake	15
4.	Terence S	14
5.	Dario Radečić	14
6.	Matt Przybyla	12
7.	Rose Day	11
8.	Emmett Boudreau	11
9.	Bharath K	10
10.	Kurtis Pykes	8

1 - 100 / 973 < >

	author	post_title	time_extracted
1.	Salma Elshahawy, MSc.	How to Configure Github Actions the Easy Way.	20210102141454
2.	Sophia Yang	Jupyter workflow for data scientists	20210102141454
3.	Terence Shin	7 Most Recommended Skills to Learn in 2021 to be a Data Scientist	20210102141454
4.	Lauren Holzbauer	SQL Interview Prep: The Next Level	20210102141454
5.	Vinay Prabhu	A machine learning practitioner's tour of 10 under-appreciated PyPi pack...	20210102141454
6.	Yefeng Xia	Data analytics helps warehouse management	20210102141454
7.	Amanda West	Association Analysis Explained	20210102141454
8.	Renato Boemer	Trends in Data Science That Will Change Business Strategies	20210102141454
9.	Lev Maximov	Broadcasting in NumPy	20210102141454
10.	Samarth Agrawal	Understanding the Confusion Matrix from Scikit learn	20210102141454
11.	Jeroen van Zeeland	Shakespeare versus Eminem—who's the better lyricist?	20210102141454
12.	Vicky Yu	How to Present Machine Learning Results to Non-Technical People	20210102141454
13.	Yash Prakash	5 Step Guide to Setting Up a New Python Environment For Data Science	20210102141454
14.	Sun Weiran	Real-time Age, Gender and Emotion Prediction from Webcam with Keras...	20210102141454
15.	Mathew Datta	The Data Science Game	20210102141454

Use case #2

Creating a Blended face



Battleground Singapore: Get to know your GE2020 candidates

The slate for the July 10 election is clear after Nomination Day. Find out who will be standing in your constituency and who the new faces are.

PUBLISHED: JUNE 30, 2020

Which party is standing where



Meet the candidates

Aljunied 5-member GRC



WP



Pritam Singh



Sylvia Lim



Faisal Manap



Gerald Giam



Leon Perera



PAP



Victor Lye



Chua Eng
Leong



Shamsul
Kamar



Alex Yeo



Chan Hui Yuh

Ang Mo Kio 5-member GRC

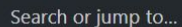


PAP



```
from urllib.request import Request, urlopen
import requests
from bs4 import BeautifulSoup
import csv
import shutil
import re
import time
from selenium import webdriver
from selenium.webdriver.chrome.options import Options
```

```
url =  
"https://www.straitstimes.com/multimedia/graphics/2020/06/singapore-general-election-ge2020-  
candidates/index.html"  
options = Options()  
options.add_argument('--headless')  
options.add_argument('--disable-gpu')  
  
driver = webdriver.Chrome(options=options)  
driver.get(url)  
time.sleep(3)  
page = driver.page_source  
driver.quit()  
soup = BeautifulSoup(page, 'html.parser')  
profile_pics = soup.find_all('img')  
for i in range(31,369):  
    if(profile_pics[i]['alt']!=''):  
        pic_url = profile_pics[i]['src']  
        pic_name = profile_pics[i]['alt'].replace(' ','_')  
        pic_name = pic_name.replace('"','_')  
        r = requests.get(pic_url, stream=True, headers={'User-agent': 'Mozilla/5.0'})  
        if r.status_code == 200:  
            with open(pic_name+pic_url[-4:], 'wb') as f:  
                r.raw.decode_content = True  
                shutil.copyfileobj(r.raw, f)  
        print(i)
```



Watch 4 Star 50 Fork 9

[Code](#)
[Issues](#) 2
 [Pull requests](#)
[Actions](#)
[Projects](#)
[Wiki](#)
[Security](#)
[Insights](#)

Go to file

Add file ▾

↓ Code

johnwmillr Merge pull request #7 from veggiedefender/master ... 73ac53b on Sep 27, 2020 30 commits

📁	facer	Only add image to warped/incremental if saving gif	7 months ago
📁	tests	Add preliminary support for GIF creation	2 years ago
📄	.gitignore	Add preliminary support for GIF creation	2 years ago
📄	AverageFaces_RapRockCountry.png	Add example image to README	12 months ago
📄	README.md	Add example image to README	12 months ago
📄	requirements.txt	Update requirements.txt	12 months ago
📄	setup.py	Add setup.py	2 years ago

☰ README.md

About

Simple (🐡) face averaging (😊) in Python (🐍)

opencv image-processing
face-detection dlib face-average
facerecognition

 [Readme](#)

Releases

No releases published

Packages

No packages published

 johnwmillr / Facer

👁 Watch 4 ☆ Star 50 🍴 Fork 9

<> Code ! Issues 2 🔗 Pull requests 🎬 Actions 📁 Projects 📖 Wiki 🛡 Security 📈 Insights

🔗 master ▾ Facer / facer / facer.py / <> Jump to ▾

Go to file ...

 veggiedefender Only add image to warped/incremental if saving gif ...

Latest commit 1b862eb on Sep 13, 2020 🕒 History

👤 2 contributors  

387 lines (326 sloc) 13.9 KB

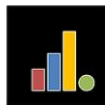
Raw Blame 🖨 ✎ 🗑

```
1 import cv2
2 import dlib
3 import matplotlib.pyplot as plt
4 from matplotlib import animation
5 import numpy as np
6 import math
7 import os
8 import glob
9
10 from facer.utils import similarityTransform, constrainPoint, calculateDelaunayTriangles, warpTriangle
11
```



Wrapping up

- Think of what analysis you/ your data science team want(s) to do
- Think about what data is necessary
- Write a script to automate the ETL (or ELT) process



DATA DOUBLE CONFIRM

[Home](#)

[About](#)

[Accolades](#)

[Collaborations](#)

[Contact](#)

[Subscribe](#)

DATA DIVE DAYS

**Process and product of various data science tasks—
from data collection, data preparation, data
visualization, to basic statistical analysis and
modelling. Datasets for practice available.**

*Selected as Top 100 Data Science Resources for 2018/2019
on [MastersInDataScience.com](#)*

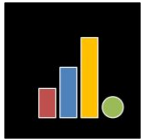
Connect with me



linkedin.com/in/hui-xiang-chua/



[@hxchuaruns](https://twitter.com/hxchuaruns)



projectosyo.wixsite.com/datadoubleconfirm