

HW#1: Data Processing (13 points)

Cloud Computing and Big Data Analytics

Dataset



- **Statistics for trending YouTube videos**

USvideos.csv

- video_id, trending_date, title, channel_title, category_id, publish_time, tags, views, likes, dislikes, comment_count, thumbnail_link, comments_disabled, ratings_disabled, video_error_or_removed, description

If you open with Excel, you will find a new world...

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y
video_id	trending_d	title	channel_title	category_id	publish_time	tags	views	likes	dislikes	comment_count	thumbnail	comments	ratings_dis	video_error	description									
RxGQe4Ee	17.14.11	鮑 ? by 潤 譚		22	2017-11-13	潤 譚	156130	1422	40	272	https://i.ytimg.com/vi/RxGQe4Ee/default.jpg,FALSE,FALSE,FALSE,鮑 ? by 潤 譚	FALSE	FALSE	FALSE	?木???鮑 ????虛? '鮑 ? 甄 ? 惡虛?\n蓮??科????渥 勞圍未 鵝勞?譚 ? 輪渥 ?? \n鮑 ?, ??? 尊 鮑									
hH7wVE8	17.14.11	JSA 蓮??	25	2017-11-13	JSA"蓮??	76533	211	28	113	https://i.ytimg.com/vi/hH7wVE8/default.jpg,FALSE,FALSE,FALSE,JSA 蓮??	FALSE	FALSE	FALSE	[麋 ?A?刺]??輪 ? ? ? 暉 ???? ? ?到破鵝??鵝?\n[麋 ?A?刺]?圍收鮑?鮑 ? 50m ? ?? ?? ? ? 圍 ?										
9V8bnWU	17.14.11	? 狐?潤 ? 狐?潤		22	2017-11-13	??木	421409	5112	166	459	https://i.ytimg.com/vi/9V8bnWU/default.jpg,FALSE,FALSE,FALSE,? 狐?潤 ? 狐?潤	FALSE	FALSE	FALSE	?澎??月? 業?粉 ? 暉? ??木?									
Q_8py-t5R	17.14.11	? ? 諱賈??渠	2017-11-13	渠?諱 "	222850	2093	173	1219	https://i.ytimg.com/vi/Q_8py-t5R/default.jpg,FALSE,FALSE,FALSE,? ? 諱賈??渠	FALSE	FALSE	FALSE	?木 ??? 筋寢 ? 簪???? ?											
bk55RbxiC	17.14.11	簪?伉疏? iNocutV	25	2017-11-13	nocutV"賄遠V"l"CBS"l"mbc"l"簪?伉疏"l"渥?"l" ??"	84466,1094,109,450	https://i.ytimg.com/vi/bk55RbxiCQdI/default.jpg,FALSE,FALSE,FALSE,簪?伉疏 MBC ?科 ??寢國筏?渥? ?\nMBC ?駭潰 ??諱																	
AmP0ryzD	17.14.11	簪? ??科? ??25	2017-11-13	簪? ?"簪	188707	545	311	357	https://i.ytimg.com/vi/AmP0ryzD/default.jpg,FALSE,FALSE,FALSE,簪? ??科? ??25	FALSE	FALSE	FALSE	簪? ??科 ????麋 ? 勞?諱賄 ????蛟 駭潺邪? ?? ? ????? ? 賈??木????蛟 ??簪 ?\nKevin MacLeod??Easy J											
4Nxb_nQL	17.14.11	? ???麋 旭?渥	2017-11-13	Positive"En	114858	252	40	36	https://i.ytimg.com/vi/4Nxb_nQL/default.jpg,FALSE,FALSE,FALSE,? ???麋 旭?渥	FALSE	FALSE	FALSE	譚木 2麋?? ? ? ????拘???^n鵝 ? 寢?輪渥 ? 月底 蓀禹?惡?鮑 ?? 未? 駭潰 ??^^											
cplEUy1zk	17.14.11	[?潺零諱 輪湊? ?塑	22	2017-11-13	潺零"l"	70166	301	37	352	https://i.ytimg.com/vi/cplEUy1zk/default.jpg,FALSE,FALSE,FALSE,[?潺零諱 輪湊? ?塑	FALSE	FALSE	FALSE	Kevin MacLeod??At Launch?(?? Creative Commons Attribution ?潰 ? (https://creativecommons.org/licenses/										
					? ?" ? 綿 ??"	65547,91,207,154,https://i.ytimg.com/vi/ToRdbxuMtg/default.jpg,FALSE,FALSE,FALSE,?																		
iToRdbxuM	17.14.11	? ?潺零 籀 ????2	2017-11-13	?? 鵝 ? 麋? ? ?	22	2017-11-13	※?渥	17868	312	0	4	https://i.ytimg.com/vi/iToRdbxuMtg/default.jpg,FALSE,FALSE,FALSE,? ?潺零 籀 ????2	FALSE	FALSE	FALSE	? 偏??譚 ??諱??圍收? ??蛙???n 麋 溢?圍采 簪? ???? ? 綿 ? 寢?? 麋 ? 勞??科 暉?諱賈								

Modules we can use...

```
import os
```

```
import numpy as np
```

```
import pandas as pd
```

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt
```

Answer the following questions

1. How many empty entries we have for each attributes (5%)?

- Hint: Google-isnull()

2. What are the average values of "likes", "dislikes", "views", "comment_count" in 2017? (10%)

- Hint: Google-groupby

3. Plot the boxplot of #dislikes for each month in 2017. (10%)

- Hint: `Google-sns.boxplot`

4. Plot the histogram of #views for each category in 2017 and 2018 in one figure. (10%)

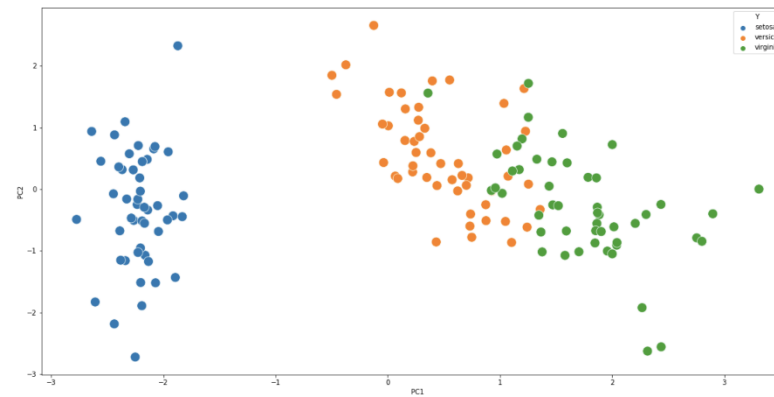
- Hint: `Google-sns.barplot`

5. Write a myPCA.py from scratch. (50%)

- Hint: <https://towardsdatascience.com/principal-component-analysis-pca-from-scratch-in-python-7f3e2a540c51>

Details

- myPCA.py is the function for PCA
 - Input: dataset path, target #dim Output: new data object
- In the main jupyter notebook,
 - select 50 samples from each category as the dataset, use “views, likes, dislikes, comment_count, comments_disabled, ratings_disabled, video_error_or_removed” as the features
 - output visualization of data points in 2-dim with color representing the category_id



6. Plot the word cloud of “title”. (15%)

Word Segmentation > Count > Cutoff > Visualization

- Hint: You can use “wordcloud”
- <https://pypi.org/project/wordcloud/>

Submission

- File Name: [studentID].ipynb, myPCA.py
- Deadline: 23:55. April 3rd, 2021

Lab 1: Data Cleaning

Cloud Computing and Big Data Analytics

Jupyter Notebook

- Blog
 - <https://www.dataquest.io/blog/jupyter-notebook-tutorial/>
- Tutorial Video
 - <https://www.youtube.com/watch?v=HW29067qVWk>
- Cloud Services
 - <https://www.dataschool.io/cloud-services-for-jupyter-notebook/>

Search in Drive



Google

My Drive

Cloud Computing and Big Data
Lecture 3: Philosophy and impact of data mining

Hong-Han Shuai
National Chiao Tung University

2021 Spring

Lecture3.pptx
You opened today

New folder

Upload files

Upload folder

Google Docs >

Google Sheets >

Google Slides >

Google Forms >

More >

Date	Presenter 1	Presenter 2	Presenter 3	Presenter 4
3/10				
3/17	蔡明儀	李俊傑		
3/24	林凱人	陳國治		
3/31	陳志仁	吳思聰		
4/7	林國華	陳智強		
4/14	Midterm			
4/21	黃清海	林國宇		
4/28	黃清海	陳國治		
5/5	陳國治	陳志仁		
5/12	陳志仁	李俊傑		
5/19	黃清海			

BASIC LAB

Edited in the past month by Tom

Google Drawings

Google My Maps

Google Sites

Google Apps Script

Google Colaboratory

Google Jamboard

Name

Owner

MAC_Backup

me

data science 2020

me

Colab Test

me

Lab1.ipynb

- https://colab.research.google.com/drive/1J9NfyYvM0Ti4TEojGQtX8dmjsz_As3RR?usp=sharing