

# Building of NLP Datasets

Topic 3

**Dwi Intan Af'idah, S.T., M.Kom**



# Agenda

01

Reminder about  
General Steps of NLP

02

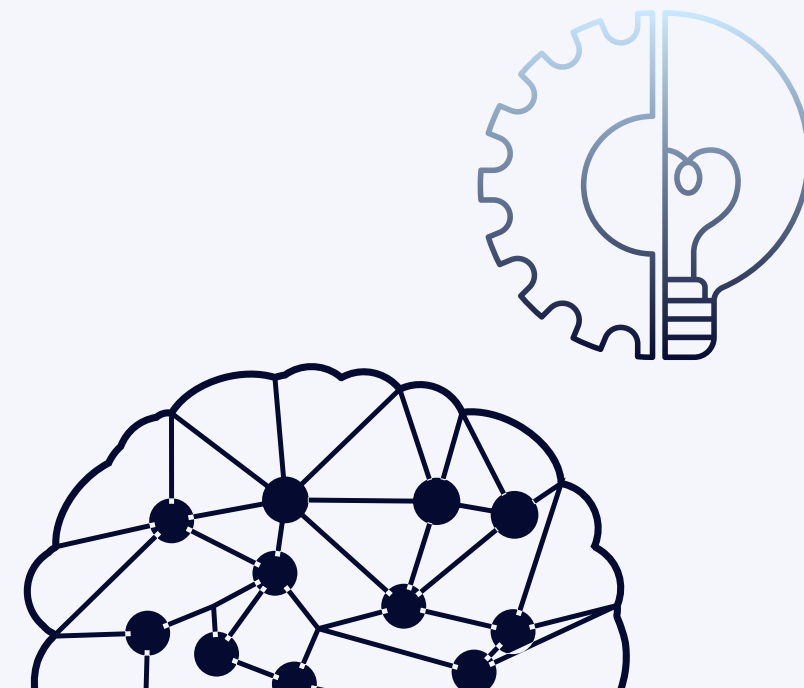
Crawling

03

Scraping

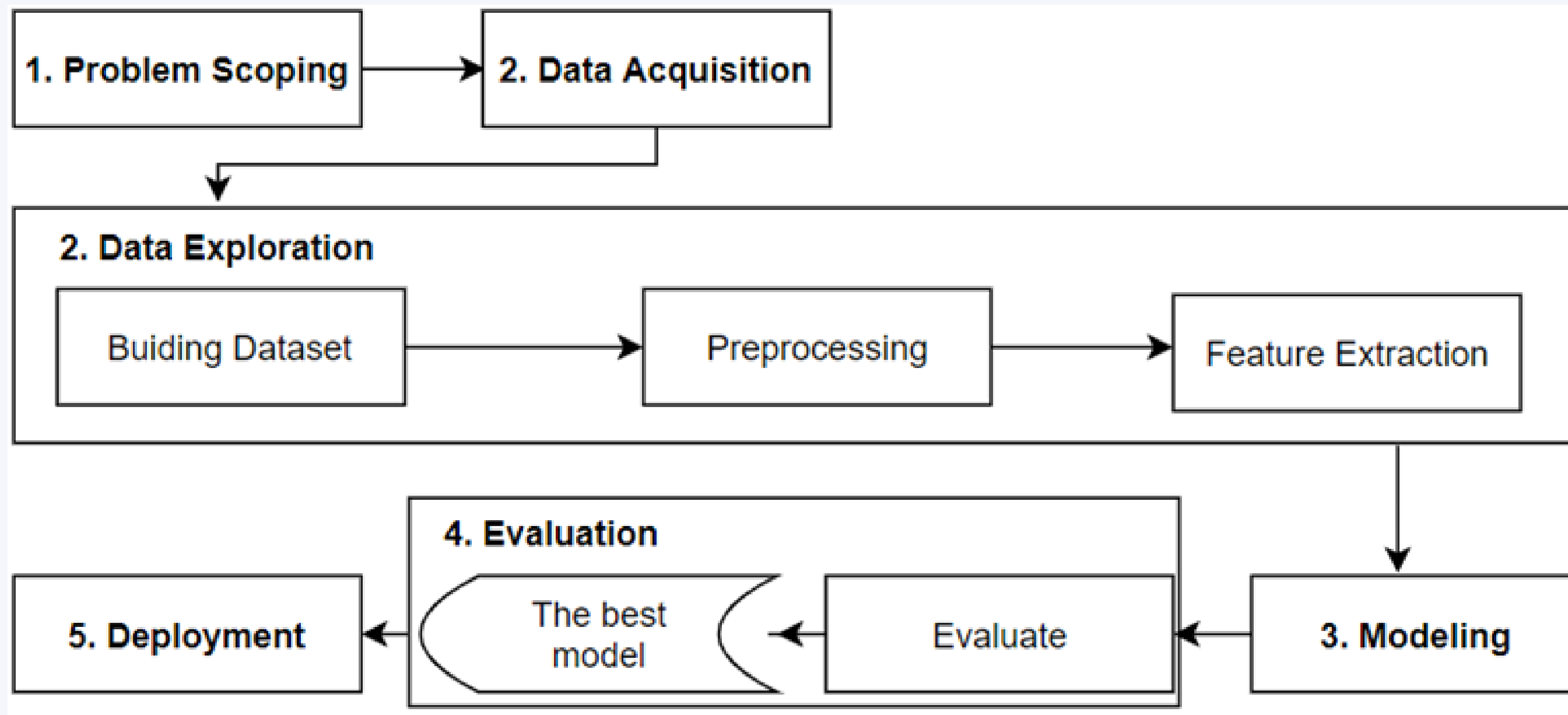
04

Basic Coding

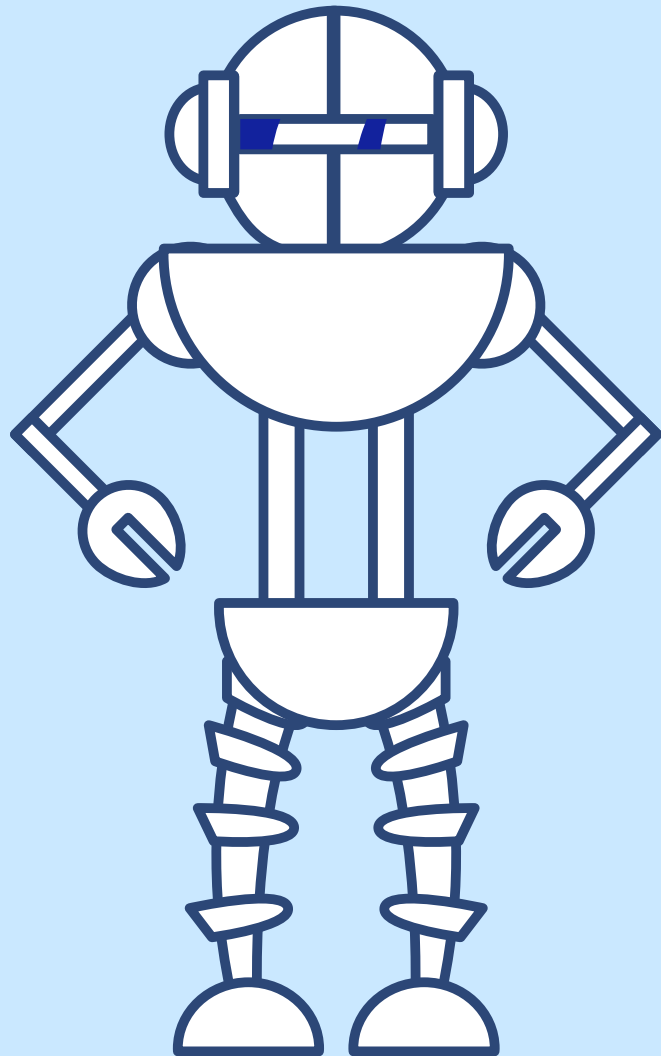




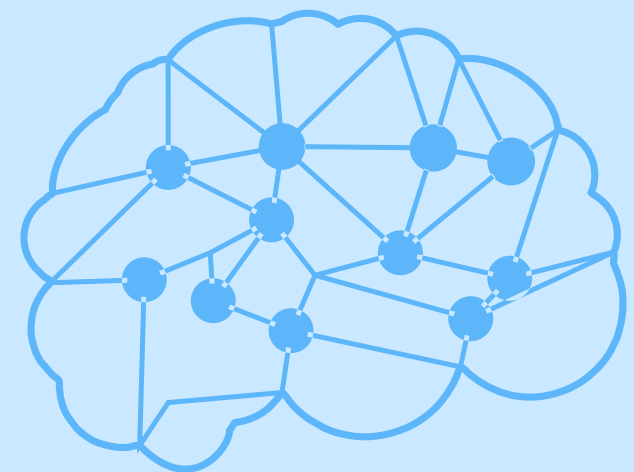
# General Steps of NLP



# File for Exercise

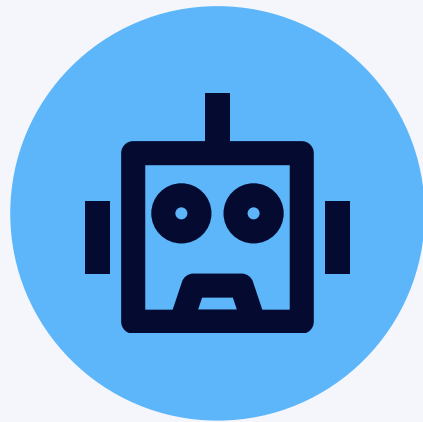
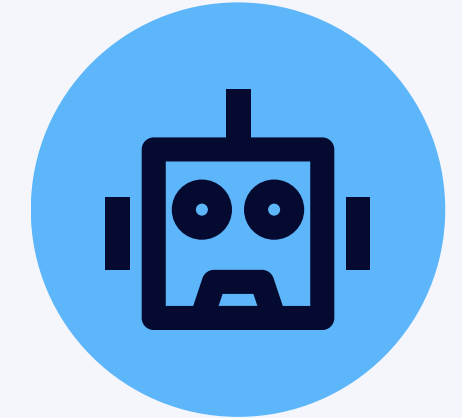


<https://bit.ly/bisaai-nlp>



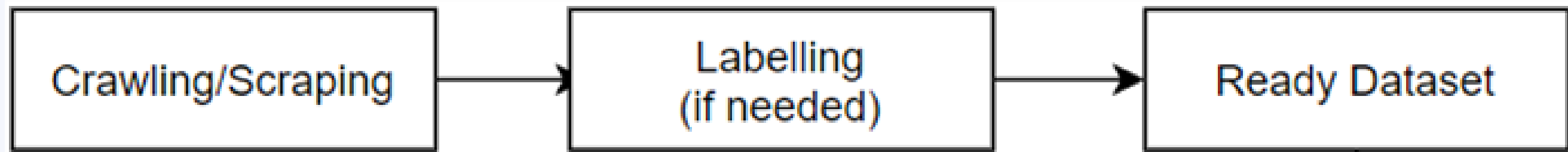
# Understanding Datasets

- Dataset adalah sekumpulan data yang disusun secara terstruktur.
- Biasanya, dataset dipresentasikan dalam bentuk tabel, alias baris dan kolom.
- Tiap baris dan kolom biasanya mewakili variable tertentu.
- Contohnya, kolom pertama merupakan data ulasan, sedangkan kolom kedua merupakan kolom label/kelas.

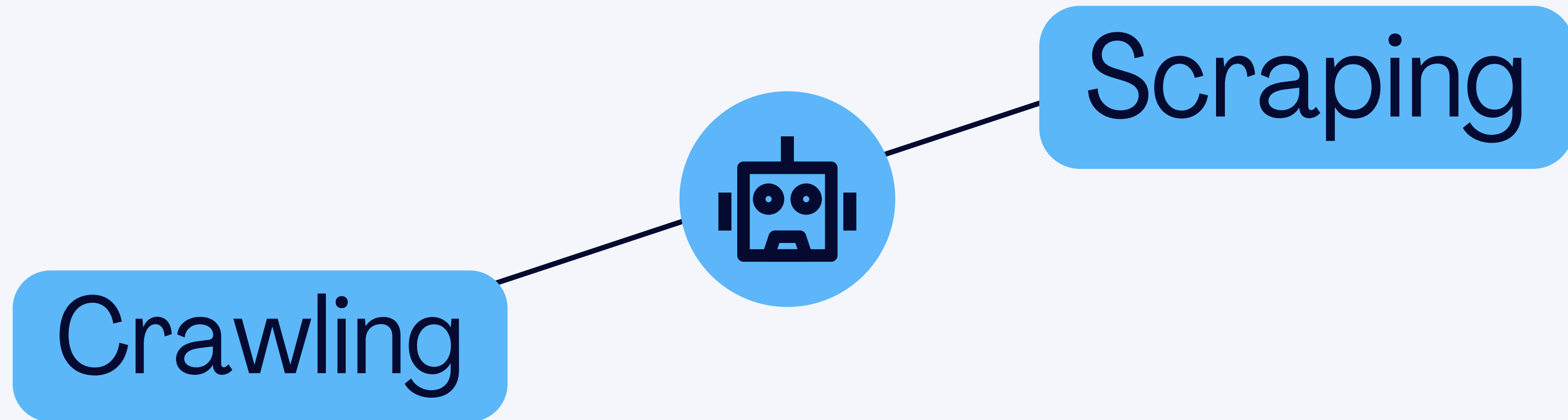


- Dataset pada penelitian NLP biasanya berupa data teks.
- Dataset pada riset NLP dapat menggunakan dataset dari:
  1. Penelitian sebelumnya yang sejenis,
  2. Membangun dataset sendiri.

# Step of Building NLP Datasets



## Collect Text Datasets



# Crawling

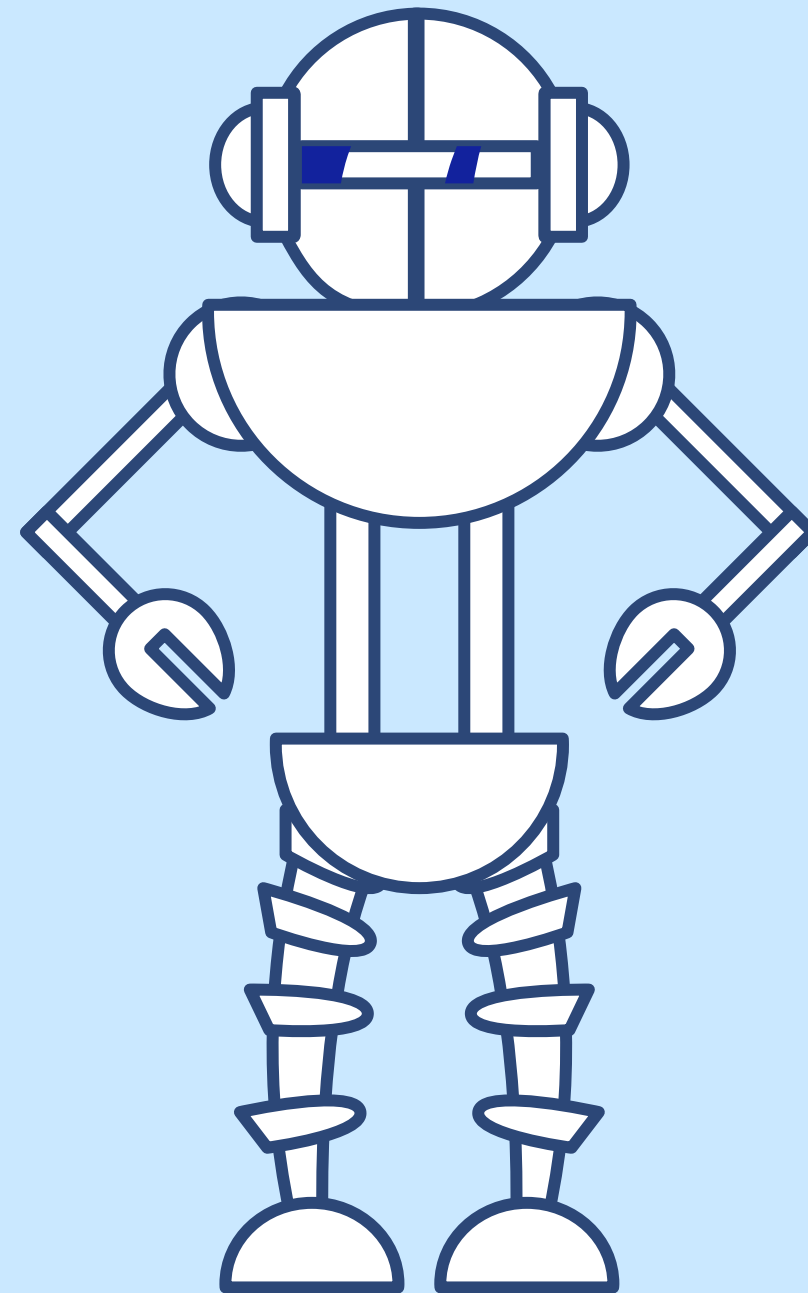
- Penerapan crawling menggunakan automation program dan menggunakan Application Programming Interface (API) sebagai jalur komunikasi dalam mendapatkan data.
- Dengan API dapat dikumpulkan data lebih spesifik sesuai dengan link URL yang ada tanpa harus mengetahui element HTML pada sebuah website.



# Crawling



Tweepy



**Tweepy** adalah *library* Python yang berguna untuk mengakses data Twitter melalui akses API (*Application Programming Interface*). Ini biasanya digunakan oleh para pengembang aplikasi untuk merancang atau menyempurnakan aplikasi yang membutuhkan data dari Twitter di dalamnya.

*Library* ini memberikan akses otomatis untuk menarik data dari Twitter melalui Python. Tidak hanya itu, kamu juga bisa mengatur akun Twitter mu menggunakan *library* ini seperti memberikan perintah untuk mengirimkan *tweet* baru, menghapus *tweet*, mengikuti ataupun berhenti mengikuti suatu akun.

Saat ini, Tweepy cenderung dimanfaatkan untuk membuat *bot* Twitter dengan tujuan tertentu ataupun melakukan analisis terhadap pengguna Twitter sendiri, utamanya sentimen yang terbentuk pada suatu kasus tertentu untuk mengetahui trennya.

## Coding of Tweepy (1)

```
1 pip install tweepy
```

```
1 import tweepy
2 access_token = "1101429342041604096-brplhFporydhLgep9zE2JJD63Yehmc"
3 access_token_secret = "m0LewkNeSpBTH03m9jMF1xkE03Vg8e36ReHnRN1Pj0MA0"
4 consumer_key = "N5qi9C8J4c1MSmL7mAeH2JW4v"
5 consumer_secret = "jjsz4HTeI1vUHTCxdp3PdnKD17REBioqkhz9ossKiXj5AekuJJ"
```

```
1 tweets_for_csv = []
2
3 def get_tweets(username):
4     auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
5     auth.set_access_token(access_token, access_token_secret)
6     api = tweepy.API(auth)
7
8     limit = 100
9     print('- Username : @'+username)
10
11     for tweet in tweepy.Cursor(api.user_timeline,
12                                screen_name=username, include_rts=False,
13                                tweet_mode='extended').items(limit):
14         actualTweet = re.sub(r'\s+', ' ', tweet.full_text)
15         tweets_for_csv.append(
16             [actualTweet])
```

## Coding of Tweepy (2)

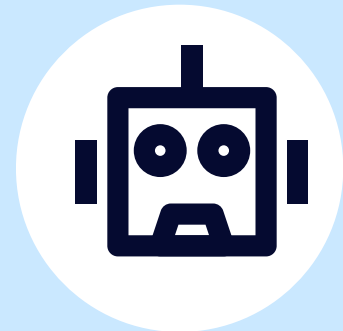
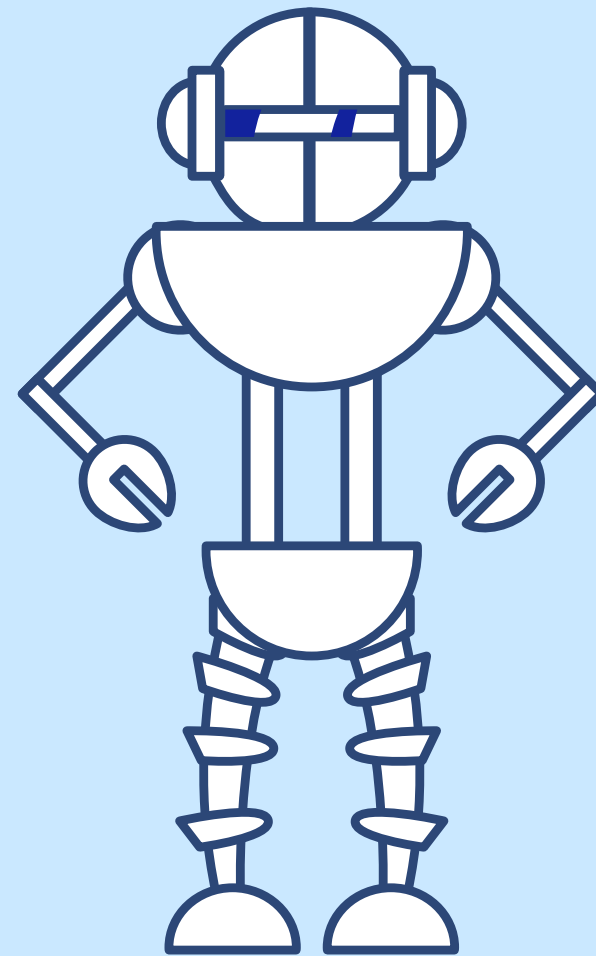
```
1 import re
2 import csv
3
4 if __name__ == '__main__':
5     users = ['CNNIndonesia']
6
7     print("\nGet tweet from username ...")
8     for user in users:
9         get_tweets(user)
10
11     outfile = "Topic 3-tweet-cnn-indonesia.csv"
12     with open(outfile, 'w', newline='', encoding='utf-8') as csvfile:
13         csvwriter = csv.writer(csvfile)
14         csvwriter.writerow(["tweet"])
15         csvwriter.writerows(tweets_for_csv)
16
17     print("\nwriting to '" + outfile + "' complete.")
```

```
1 import pandas as pd
2
3 dataset = pd.read_csv('Topic 3-tweet-cnn-indonesia.csv')
4 dataset[:5]
```

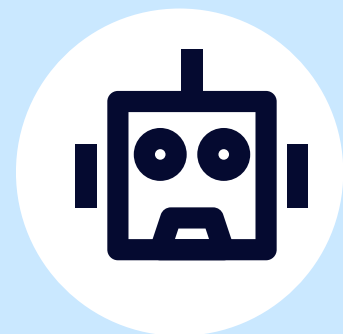
- Scraping merupakan teknik mengumpulkan data pada sebuah website melalui proses ekstraksi informasi menggunakan Hypertext Transfer Protocol (HTTP).
- Scraping dapat digunakan secara manual ataupun secara automation program.
- Namun untuk mendapatkan data kita perlu mengetahui element HTML ataupun XML pada sebuah website.
- Kemudian kita masukkan ke dalam program yang dibuat untuk mencari data sesuai nama id atau nama class dari element HTML tersebut.



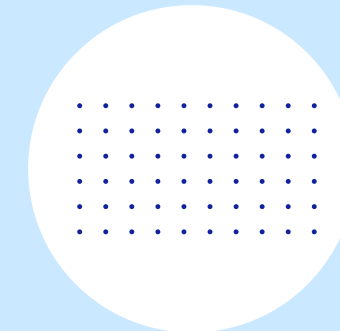
# Scraping



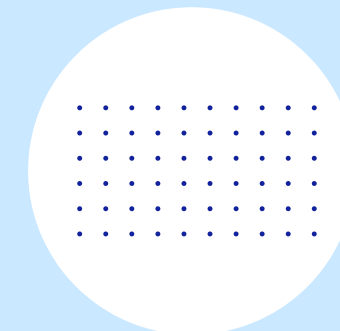
**Beautiful Soup**



**Twint**



**Selenium**



**Googleplay Scraper**

- **Twint is an advanced Twitter scraping tool written in Python that allows for scraping Tweets from Twitter.**
- **The advantage of Twint is that you don't need Twitter's API to make TWINT work. Twint utilizes Twitter's search operators to let you :**
  - 1.scrape Tweets from specific users
  - 2.scrape Tweets relating to certain topics
  - 3.hashtags & trends
  - 4.or sort out sensitive information from Tweets like e-mail and phone numbers.



```
1 pip install --upgrade git+https://github.com/twintproject/twint.git@origin/master#egg=twint
```



```
1 !pip install nest_asyncio
```

## Coding of Twint (1)

```
1 import twint
2 import nest_asyncio
3
4 nest_asyncio.apply()
5 c = twint.Config()
6
7 c.Search = "Jokowi"
8 c.Since = "2022-1-1"
9 c.Until = "2022-1-31"
10 c.Limit = 100
11 c.Pandas = True
12
13 twint.run.Search(c)
```

## Coding of Twint (2)

```
1 def column_names():
2     return twint.output.panda.Tweets_df.columns
3
4 def twint_to_pd(columns):
5     return twint.output.panda.Tweets_df[columns]
6
7 dataset = twint_to_pd(['tweet'])
8 print(dataset)
9
10 dataset.to_csv("Topic 3-tweet-jokowi.csv", index=False)
```

# Beautiful soup (1)

## Tentang BeautifulSoup

Kita sebenarnya bisa melakukan web scraping secara manual. Namun, dengan request url sederhana, komputer akan memberikan data HTML yang nampak membingungkan. Lihat contoh tampilannya di bawah ini.

```
'<!DOCTYPE html>\n<html class="client-nojs" lang="en" dir="ltr">\n<head>\n<meta charset="UTF-8"/>\n<title>Joko Widodo - Wiki-
pedia</title>\n<script>document.documentElement.className="client-js";RLCONF={"wgBreakFrames":!1,"wgSeparatorTransformTabl
e":["",""],"wgDigitTransformTable":["",""],"wgDefaultDateFormat":"dmy","wgMonthNames":["","January","February","March","Apri
l","May","June","July","August","September","October","November","December"],"wgRequestId":"704ea0fd-86b8-4984-8de2-6376c01a
4a14","wgCSPNonce":!1,"wgCanonicalNamespace":"","wgCanonicalSpecialPageName":!1,"wgNamespaceNumber":0,"wgPageName":"Joko_Wid
odo","wgTitle":"Joko Widodo","wgCurRevisionId":1030025186,"wgRevisionId":1030025186,"wgArticleId":29410367,"wgIsArticle":!
0,"wgIsRedirect":!1,"wgAction":"view","wgUserName":null,"wgUserGroups":[""],"wgCategories":["CS1 Indonesian-language source
s (id)","CS1 maint: multiple names: authors list","Articles with short description","Short description is different from Wik
idata","Wikipedia pending changes protected pages","Use British English from August 2020","Use dmy dates from May 2018","A
rticles containing Indonesian-language text","Commons category link from Wikidata","Wikipedia articles with GND identifier
s","Wikipedia articles with ISNI identifiers","Wikipedia articles with VIAF identifiers","Wikipedia articles with LCCN ident
ifiers","Wikipedia articles with PLWABN identifiers","Wikipedia articles with FAST identifiers","Wikipedia articles with SUD
OC identifiers","Wikipedia articles with WORLDCATID identifiers","Joko Widodo","1961 births","Gadjah Mada University alumn
i","Governors of Jakarta","Indonesian businesspeople","Indonesian Democratic Party of Struggle politicians","Indonesian engi
neers","Indonesian Muslims","Javanese people","Living people","Mayors of Surakarta","People from Surakarta","Presidents of I
ndonesia","Mayors of places in Indonesia"],"wgPageContentLanguage":"en","wgPageContentModel":"wikitext","wgRelevantPageNam
e":"Joko_Widodo","wgRelevantArticleId":29410367,"wgIsProbablyEditable":!0,"wgRelevantPageIsProbablyEditable":!0,"wgRestricti
onEdit":\n[{"wgRestrictionMove":[""],"wgFlaggedRevsParams":{"tags":{"status":{"levels":-1}}},"wgStableRevisionId":103002518
6,"wgMediaViewerOnClick":!0,"wgMediaViewerEnabledByDefault":!0,"wgPopupsFlags":10,"wgVisualEditor":{"pageLanguageCode":"e
n","pageLanguageDir":"ltr","pageVariantFallbacks":"en"},"wgMFDisplayWikibaseDescriptions":{"search":!0,"nearby":!0,"watchlis
t":!0,"tagline":!1},"wgWMESchemaEditAttemptStepOversample":!1,"wgULSCurrentAutonym":"English","wgNoticeProject":"wikipedi
a","wgCentralAuthMobileDomain":!1,"wgEditSubmitButtonLabelPublish":!0,"wgULSPosition":"interlanguage","wgULSisCompactLinksEn
abled":!0,"wgGNewcomerTasksGuidanceEnabled":!0,"wgGEAskQuestionEnabled":!1,"wgGELinkRecommendationsFrontendEnabled":!1,"wgW
ikibaseItemId":"Q3318231"}];RLSTATE={"ext.globalCssJs.user.styles":"ready","site.styles":"ready","noscript":"ready","user.sty
les":"ready","ext.globalCssJs.user":"ready","user":"ready","user.options":"loading","ext.flaggedRevs.icons":"ready","\n"oojs-
ui-core.styles":"ready","oojs-ui.styles.indicators":"ready","mediawiki.widgets.styles":"ready","oojs-ui-core.icons":"read
y","ext.cite.styles":"ready","skins.vector.styles.legacy":"ready","jquery.makeCollapsible.styles":"ready","ext.flaggedRevs.b
asic":"ready","ext.visualEditor.desktopArticleTarget.noscript":"ready","ext.uls.interlanguage":"ready","ext.wikimediaBadge
s":"ready","wikibase.client.init":"ready"};RLPAGEMODULES=["ext.cite.ux-enhancements","ext.scribunto.logs","site","mediawiki.
page.readv","jquery.makeCollapsible","mediawiki.toc","skins.vector.legacy.js","ext.flaggedRevs.advanced","ext.pageret.Referen
```



Untuk memudahkan scraping, kita bisa menggunakan BeautifulSoup. BeautifulSoup adalah library Python yang digunakan untuk mengambil data HTML dan XML. BeautifulSoup berfungsi sebagai parser untuk memisahkan komponen-komponen HTML menjadi rangkain elemen yang mudah dibaca.

```
<!DOCTYPE html>

<html class="client-nojs" dir="ltr" lang="en">
<head>
<meta charset="utf-8"/>
<title>Joko Widodo - Wikipedia</title>
<script>document.documentElement.className="client-js";RLCONF={"wgBreakFrames":!1,"wgSeparatorTransformTable":["",""],"wgDigitTransformTable":["",""],"wgDefaultDateFormat":"dmy","wgMonthNames":["","January","February","March","April","May","June","July","August","September","October","November","December"],"wgRequestId":"704ea0fd-86b8-4984-8de2-6376c01a4a14","wgCSPNonce":!1,"wgCanonicalNamespace":"","wgCanonicalSpecialPageName":!1,"wgNamespaceNumber":0,"wgPageName":"Joko_Widodo","wgTitle":"Joko Widodo","wgCurRevisionId":1030025186,"wgRevisionId":1030025186,"wgArticleId":29410367,"wgIsArticle":!0,"wgIsRedirect":!1,"wgAction":"view","wgUserName":null,"wgUserGroups":["*"],"wgCategories":["CS1 Indonesian-language sources (id)","CS1 maint: multiple names: authors list","Articles with short description","Short description is different from Wikidata","Wikipedia pending changes protected pages","Use British English from August 2020","Use dmy dates from May 2018","Articles containing Indonesian-language text","Commons category link from Wikidata","Wikipedia articles with GND identifiers","Wikipedia articles with ISNI identifiers","Wikipedia articles with VIAF identifiers","Wikipedia articles with LCCN identifiers","Wikipedia articles with PLWABN identifiers","Wikipedia articles with FAST identifiers","Wikipedia articles with SUDOC identifiers","Wikipedia articles with WORLDCATID identifiers","Joko Widodo","1961 births","Gadjah Mada University alumni","Governors of Jakarta","Indonesian businesspeople","Indonesian Democratic Party of Struggle politicians","Indonesian engineers","Indonesian Muslims","Javanese people","Living people","Mayors of Surakarta","People from Surakarta","Presidents of Indonesia","Mayors of places in Indonesia"],"wgPageContentLanguage":"en","wgPageContentModel":"wikitext","wgRelevantPageName":"Joko_Widodo","wgRelevantArticleId":29410367,"wgIsProbablyEditable":!0,"wgRelevantPageIsProbablyEditable":!0,"wgRestrictionEdit":
```

```
1 # SCRAPING URL ADDRESS
2 import requests
3 import csv
4 from bs4 import BeautifulSoup
5
6 # Write ke file csv
7 csv_output = csv.writer(open('Topic 3-buku-url.csv', 'w', newline=''))
8
9 pages = []
10
11 # Collecting & parsing konten Web hlm 1-6
12 for i in range(1, 7):
13     url = 'https://www.goodreads.com/list/show/2405.Buku_Non_Fiksi_Indonesia_Terbaik_Sepanjang_Masa' + str(i)
14     pages.append(url)
15
16 for item in pages:
17     page = requests.get(item)
18     soup = BeautifulSoup(page.text, 'html.parser')
19
20     # Find elemen class bookTitle di dalam class tableList
21     novel_title_list = soup.find(class_='tableList js-dataTooltip')
22     novel_title_list_items = novel_title_list.find_all(class_='bookTitle')
23
24     # Get masing-masing title dan url dari class tableList
25     for novel_title in novel_title_list_items:
26         link = 'https://www.goodreads.com' + novel_title.get('href') + '?language_code=id'
27
28         csv_output.writerow([link])
```

```

1 # SCRAPING REVIEW DATA FROM CSV FILE
2 import requests
3 from bs4 import BeautifulSoup
4 import csv
5
6 #Read input url dari file csv & write output review
7 with open('Topic 3-buku-url.csv', newline='') as f_urls, open('Topic 3-buku-ulasan.csv', 'w', newline='', encoding="utf-8") as f_output:
8     csv_urls = csv.reader(f_urls)
9     csv_output = csv.writer(f_output)
10    csv_output.writerow(['Nama', 'Ulasan', 'Rating'])
11
12    # Collecting & parsing konten web
13    for line in csv_urls:
14        r = requests.get(line[0]).text
15        soup = BeautifulSoup(r, 'lxml')
16
17        # Find elemen class review di dalam id bookReviews
18        #novel_review_list = soup.find('div', {'id': 'bookReviews'})
19        #novel_review_list_items = novel_review_list.find_all(class_='review')
20        novel_review_list2 = soup.find_all('div', class_="friendReviews elementListBrown")
21
22        # Get masing-masing nama & review dari class bookReviews
23        for novel_review in novel_review_list2:
24            #name = novel_review.find(class_='user').get_text()
25            review = novel_review.find(class_='readable').find('span', recursive=False).get_text()
26            rating_element = novel_review.find('span', {'size': '15x15'})
27            # Skip item review ketika elemen rating tidak ditemukan
28            if rating_element == None:
29                continue
30            else :
31                rating = rating_element.get_text()
32
33            csv_output.writerow([review, rating])

```

## Googleplay Scraper

Ada jutaan aplikasi, buku, dan film di Google Play Store, dan jumlahnya terus bertambah setiap hari. Menurut [Penelitian AppBrain](#) ada sekitar 3 juta aplikasi pada kuartal pertama 2021. Miliaran komentar dimasukkan ke dalam jutaan aplikasi, buku, dan film ini setiap hari.

Peneliti, pengembang aplikasi, pakar pemasaran, dan pakar yang bekerja di berbagai bidang ingin mengorek dan memeriksa komentar yang dibuat di Google Play karena berbagai alasan.

Meskipun dimungkinkan untuk mengunduh komentar aplikasi yang Anda kembangkan dari Google Dashboard, [Outscraper Pengikis Ulasan Google Play](#) adalah alat yang tepat jika Anda mencari sistem di mana Anda dapat mengikis komentar untuk aplikasi lain tanpa batasan apa pun.

Outscraper memiliki Layanan Aplikasi Web dan API untuk Google Play Reviews scraping. Anda dapat langsung menggunakan Web App Scraper tanpa pengkodean apa pun, atau Anda dapat menggunakan API kami untuk menggunakannya di aplikasi / layanan Anda sendiri.

```
1 from google_play_scraper import app
2
3 import pandas as pd
4
5 import numpy as np
```

```
1 #Scrape desired number of reviews
2
3 from google_play_scraper import Sort, reviews
4
5 result, continuation_token = reviews(
6     'com.telkom.mwallet',
7     lang='id', # defaults to 'en'
8     country='id', # defaults to 'us'
9     sort=Sort.MOST_RELEVANT, # defaults to Sort.MOST_RELEVANT you can use Sort.NEWEST to get newst reviews
10    count=5000, # defaults to 100
11    filter_score_with= 1 # defaults to None(means all score) Use 1 or 2 or 3 or 4 or 5 to select certain score
12 )
```



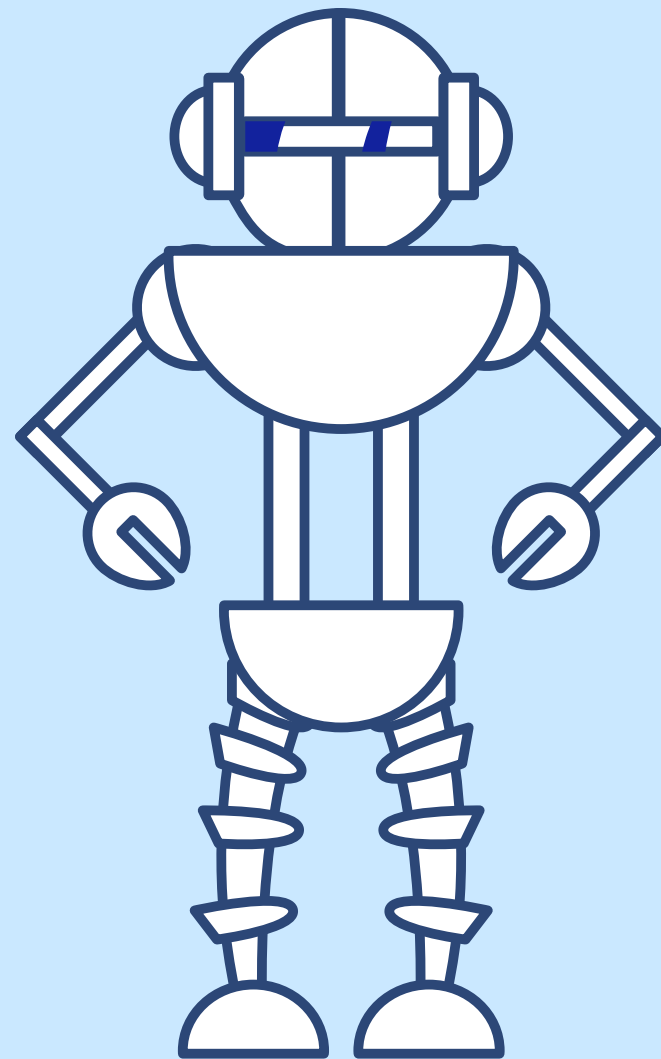
```
1 # Dataframe dengan nama
2 df_twt = pd.DataFrame(np.array(result), columns=['review'])
3
4 df_twt = df_twt.join(pd.DataFrame(df_twt.pop('review').tolist()))
5
6 df_twt.head()
```

```
1 len(df_twt.index) #count the number of data we got
```

```
1 my_df = df_twt[['userName', 'score', 'at', 'content']] #get userName, rating, date-time, and reviews only
```

```
1 my_df.to_csv("linkaja-rating1.csv", index = False) #Save the file as CSV , to download: click the folder icon on the left.
2 #the csv file should be there.
```

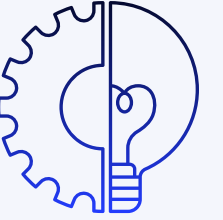
# Labelling of Dataset



Ulasan	Sentimen
Sistemnya harus lebih dibenahi,karna sering terjadi transaksi berhasil namun no tdk terdaftar	0
Jaringan udah kenceng, udah hapus cache, hapus data, install ulang, ganti provider, ganti p	0
Saya mengalami masalah pada akun saya,tapi respon dana sangat2 lambat.. Akun saya di	0
Bagaimana sih dr smalam sampai skrng dana tidak bisa digunakan?? Saya sangat kecewa.	0
Saya jujur saja ...jgn pernah pakai aplikasi e-money ini... Ini aplikasi e-money terburuk mau	0
Buruk bener akun dana , tgl 10-12-2021 saya transfer dari akun dana ke mandiri belum mas	0
tambah ribet...padahal cuma mau upgread ke premium.. susah nya minta ampun.. tapi saya	0
Customer service dana care terburuk yg pernah ada!!!! Kirim uang engga sampai,dana ditah	0
Sangat memudahkan bnget	1
lumayan membantu..meskipun promo2 berkurang	1
Tingkat kn trus	1
Kamis, 23-09-2021, 16.00wib Saya transfer ke rek BCA dengan nominal ;Rp. 275.000 (TIDA	1
Mantap saya suka sekali aplikasinya	1
Gokiiiii buangeett nih apk OVO ayo download buruann gak bakal nyeseelll	1
Sangat membantu sekali	1
Tagihan begitu mudah	1
Sangat membantu sekali	1
Sangat membantu sekali	1
Bagus mudah dan aman	1
Tolong jaga keamanan	1

## Reference

- Kedia, A., dan Rasu, M., 2020, Hands-On Python Natural Language Processing, Packt Publishing Ltd., Brimingham, UK, 35-41
- <https://www.tweepy.org/>
- <https://github.com/twintproject/twint>
- <https://pypi.org/project/beautifulsoup4/>
- <https://pypi.org/project/google-play-scraper/>



**Thank You**

