

TUGAS 6 Analisis Sentiment dan Data Scrapping  
Gojek vs Grab di twitter



Di susun oleh:

Nama: Rizqy Mubarok

NPM: 23.0504.0024

Fakultas Teknik

Universitas Muhammadiyah Magelang 2025

## A. Scrapping

1. Untuk scrapping saya menggunakan Library twikit untuk mengambil data dari twitter

Import Library dan token twitter **Twikit** mengharuskan user menginputkan token sehingga tidak bisa anonim

```
1 import csv
2 import nest_asyncio
3 import asyncio
4 import random
5 import math
6 from twikit import Client, TooManyRequests
7
8 nest_asyncio.apply()
9 # SETUP CLIENT + COOKIE
10 client = Client("en-US")
11 client.set_cookies({
12     "auth_token": "d125d285011b59146b60cdcd6edd9c82dca055c9",
13     "ct0": "59e9745d43723d4cb98f25be06a03b19dc235dbb0f0410ca8b0b49f8231a184b2bc947fcc2fd77668b90979426c200556e89b9d9dbd00",
14     "twid": "u%3D1645176567038091265",
15     "gt": "1989338289653047338",
16     "guest_id": "v1%3A175121650454954896",
17     "guest_id_ads": "v1%3A175121650454954896",
18     "guest_id_marketing": "v1%3A175121650454954896",
19     "kdt": "HZdPjFR2pqsvUhxs05GasalB5x7wYSzJ8sBb5x85",
20     "lang": "en",
21     "personalization_id": "v1_JcinIxPjnV00ctgz7rKeiQ==",
22     "cf_bm": "sjuQ5SgqIGcWguc4AGCyac42n3AikDN522y1tic1cGI-1763135431.7418268-1.0.1.1-3bCnpzTDAGCtwt0vvIoOgF_Li.Jdb5N0EUXs",
23     "cuid": "f6d3d51a5a0345ee9e273cb98b2c5068",
24     "att": "1-PToudQBU2NDJJeJZeJ8k2vit0taZstJ225z7kN54g",
25     "external_referer": "padhuUp37zhD6%2F29CpQtyhGQCU105AFo|0|8e8t2xd8A2w%3D",
26     "g_state": "{\"i_l1\":0,\"i_l11\":1763130464264}",
27     "gt": "1989338289653047338",
28 })
29
```

Token yang saya pakai Adalah token yang saya ambil dari menu dev tools setelah login emnggunakan akun twitter saya, jadi datanya bisa saya pertanggung jawabkan.

2. Beberapa parameter yang saya gunakan untuk scrapping

```
1 # Config
2 BATCH_SIZE = 20
3 TARGET_PER_KEYWORD = 500
4 CHECKPOINT_EVERY = 50
5 MIN_DELAY = 10
6 MAX_DELAY = 20
7 MAX_RETRIES = 6
8 CSV_FILENAME = "tweet_gojek_grab.csv"
9
10 def write_header_if_needed(filename):
11     try:
12         with open(filename, "x", encoding="utf-8", newline="") as f:
13             writer = csv.writer(f)
14             writer.writerow(["keyword", "tweet"])
15     except FileExistsError:
16         pass
17
18 def append_rows(filename, rows):
19     with open(filename, "a", encoding="utf-8", newline="") as f:
20         writer = csv.writer(f)
21         writer.writerows(rows)
```

- Batch size = 20 agar Ketika kode melakukan pengambilan data dalam 1 kali request, maksimal hanya mengambil 20 data untuk menghindari rate limit (too many request).
  - Target per keyword = 500 per untuk keyword gojek dan grab masing-masing ada 500 data.
  - Min dan max delay untuk menentukan jarak waktu pengiriman request ke twitter, jika terlalu cepat beresiko terkena rate limit.
  - Max Retries di gunakan untuk menghentikan kode jika terjadi eror pada jaringan atau twitter menolak respon dari request.
3. Karena twitter mempunyai aturan sangat ketat terkait data pengguna, maka dari itu saya menggunakan paginasi agar Ketika mengakses data dari twitter, tidak langsung mengambil 500 data namun hanya akan mengambil per halaman

```

1 # scrapping dengan pagination, backoff, dan checkpointing
2 async def scrape_keyword(keyword, target=TARGET_PER_KEYWORD):
3     print(f"Start scrapping keyword: {keyword} (target {target})")
4     results = []
5     cursor = None
6     consecutive_empty = 0
7     write_header_if_needed(CSV_FILENAME)
8
9     while len(results) < target:
10         attempt = 0
11         while True:
12             try:
13                 response = await client.search_tweet(
14                     f"{keyword} lang:id",
15                     product="Latest",
16                     count=BATCH_SIZE,
17                     cursor=cursor
18                 )
19                 tweets = response if response is not None else []
20                 break
21             except TooManyRequests:
22                 attempt += 1
23                 backoff = min(60, 2 ** attempt + random.random() * 5)
24                 print(f"429 Rate limit. Backoff for {backoff:.1f}s (attempt {attempt})")
25                 await asyncio.sleep(backoff)
26                 if attempt >= MAX_RETRIES:
27                     print("Max retries reached for 429. Exiting this keyword early.")
28                     return results
29
30

```

```

29         return results
30
31     except Exception as e:
32         attempt += 1
33         backoff = min(60, (2 ** attempt) + random.random() * 5)
34         print(f"transient error: {e!r}. Backoff {backoff:.1f}s (attempt {attempt})")
35         await asyncio.sleep(backoff)
36         if attempt >= MAX_RETRIES:
37             print("Max retries reached for errors. Exiting this keyword early.")
38             return results
39
40     if not tweets:
41         consecutive_empty += 1
42         print(f"No tweets returned in this batch (empty count={consecutive_empty}).")
43         if consecutive_empty >= 4:
44             print("iToo many consecutive empty batches - stopping pagination for this keyword.")
45             break
46         await asyncio.sleep(random.uniform(MIN_DELAY, MAX_DELAY))
47         continue
48     else:
49         consecutive_empty = 0
50
51     batch_texts = []
52     for t in tweets:
53         try:
54             text = t.text
55         except Exception:
56             text = str(t)
57         batch_texts.append(text)
58
59     rows = []
60     for txt in batch_texts:
61         results.append(txt)
62         rows.append([keyword.replace(" lang:id", ""), txt])
63
64     if rows:
65         append_rows(CSV_FILENAME, rows)
66
67

```

```

67
68     print(f" got {len(batch_texts)} tweets (total collected: {len(results)}/{target})")
69
70     try:
71         next_cursor = getattr(tweets, "next_cursor", None)
72         if next_cursor is None and isinstance(response, dict):
73             next_cursor = response.get("next_cursor") or response.get("cursor")
74     except Exception:
75         next_cursor = None
76
77     if not next_cursor:
78         if len(results) >= target:
79             break
80         cursor = None
81         print("no next_cursor; will try a few more small batches before stopping.")
82     else:
83         cursor = next_cursor
84         await asyncio.sleep(random.uniform(MIN_DELAY, MAX_DELAY))
85         if len(results) >= target:
86             break
87
88     print(f"Finished {keyword}: collected {len(results)} tweets.")
89     return results[:target]

```

#### 4. Eksekusi scrapping

```
1 # main driver
2 async def main():
3     write_header_if_needed(CSV_FILENAME)
4     keywords = ["gojek", "grab"]
5
6     final = {}
7     for kw in keywords:
8         tweets = await scrape_keyword(kw, TARGET_PER_KEYWORD)
9         final[kw] = tweets
10        cooldown = random.uniform(5, 12)
11        print(f"Cooldown {cooldown:.1f}s before next keyword.")
12        await asyncio.sleep(cooldown)
13
14    print("All done. Summary:")
15    for k, v in final.items():
16        print(f" - {k}: {len(v)} tweets collected")
17    await main()
```

```
... Start scraping keyword: gojek (target 500)
got 19 tweets (total collected: 19/500)
got 20 tweets (total collected: 39/500)
got 20 tweets (total collected: 59/500)
got 20 tweets (total collected: 79/500)
got 20 tweets (total collected: 99/500)
got 20 tweets (total collected: 119/500)
got 20 tweets (total collected: 139/500)
got 20 tweets (total collected: 159/500)
got 20 tweets (total collected: 179/500)
got 20 tweets (total collected: 199/500)
got 20 tweets (total collected: 219/500)
got 20 tweets (total collected: 239/500)
got 20 tweets (total collected: 259/500)
got 20 tweets (total collected: 279/500)
got 20 tweets (total collected: 299/500)
got 20 tweets (total collected: 319/500)
got 20 tweets (total collected: 339/500)
got 20 tweets (total collected: 359/500)
got 20 tweets (total collected: 379/500)
got 19 tweets (total collected: 398/500)
got 20 tweets (total collected: 418/500)
got 20 tweets (total collected: 438/500)
got 20 tweets (total collected: 458/500)
got 20 tweets (total collected: 478/500)
got 20 tweets (total collected: 498/500)
got 17 tweets (total collected: 515/500)
Finished gojek: collected 515 tweets.
Cooldown 5.1s before next keyword.
```

```
Start scraping keyword: grab (target 500)
got 20 tweets (total collected: 20/500)
got 20 tweets (total collected: 40/500)
got 20 tweets (total collected: 60/500)
got 20 tweets (total collected: 80/500)
got 20 tweets (total collected: 100/500)
got 20 tweets (total collected: 120/500)
got 20 tweets (total collected: 140/500)
got 20 tweets (total collected: 160/500)
got 20 tweets (total collected: 180/500)
got 20 tweets (total collected: 200/500)
got 20 tweets (total collected: 220/500)
got 20 tweets (total collected: 240/500)
got 18 tweets (total collected: 258/500)
got 20 tweets (total collected: 278/500)
got 20 tweets (total collected: 298/500)
got 20 tweets (total collected: 318/500)
got 20 tweets (total collected: 338/500)
got 20 tweets (total collected: 358/500)
got 20 tweets (total collected: 378/500)
got 20 tweets (total collected: 398/500)
got 20 tweets (total collected: 418/500)
got 20 tweets (total collected: 438/500)
got 20 tweets (total collected: 458/500)
got 20 tweets (total collected: 478/500)
got 20 tweets (total collected: 498/500)
got 20 tweets (total collected: 518/500)
Finished grab: collected 518 tweets.
Cooldown 10.0s before next keyword.
All done. Summary:
- gojek: 500 tweets collected
- grab: 500 tweets collected
```

## B. Analisis sentiment

### 1. Import library

```
1 import re
2 import os
3 import csv
4 import math
5 import random
6 import string
7 from collections import Counter, defaultdict
8
9 import numpy as np
10 import pandas as pd
11 from tqdm.auto import tqdm
12
13 # NLP
14 import nltk
15 from nltk.tokenize import word_tokenize
16 nltk.download('punkt')
17 nltk.download('stopwords')
18 nltk.download('punkt_tab')
19
20 from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
21
22 # sklearn
23 from sklearn.model_selection import train_test_split, cross_val_score
24 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
25 from sklearn.naive_bayes import MultinomialNB
26 from sklearn.linear_model import LogisticRegression
27 from sklearn.svm import LinearSVC
28 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, precision_recall_fscore_support
29
30 #word2vec
31 from gensim.models import Word2Vec
32
33 # plotting
34 import matplotlib.pyplot as plt
35 import seaborn as sns
36 sns.set()
37 from wordcloud import WordCloud, STOPWORDS
```

### 2. Load Dataset

#### Load Dataset

```
1 # load dataset
2 INPUT_CSV = "tweet_gojek_grab.csv"
3 df = pd.read_csv(INPUT_CSV)
4 df = df[['keyword', 'tweet']].dropna().reset_index(drop=True)
5 print("Loaded rows:", len(df))
6 df.head()
```

Loaded rows: 1033

|   | keyword | tweet   |
|---|---------|---|
| 0 | gojek   | Lah ini gojek daerah mana nih bawa glock?         |
| 1 | gojek   | @tubirfess kalo ngga mau satu tempat sama mere... |
| 2 | gojek   | dan ongkir minimum yang diterima driver juga ...  |
| 3 | gojek   | @Ordinaryjhonny baru kmrn nyoba syopifud emg l... |
| 4 | gojek   | Bangsat, lu kalau mau nyaman ya bawa kendaraan... |

Next steps: [Generate code with df](#)

[New interactive sheet](#)

Dataset yang saya gunakan Adalah dataset hasil scrapping, dengan jumlah data sebanyak 1000 data dan masing-masing keyword ada 500 data.

### 3. Preprocessing data

```
1 # 1. Preprocessing helpers
2 NORMALIZE_MAP = {
3     "gk": "nggak", "ga": "nggak", "g": "nggak",
4     "gk2": "nggak", "tdk": "tidak", "td": "tidak",
5     "aja": "saja", "klo": "kalau", "kalo": "kalau",
6     "jg": "juga", "jd": "jadi", "dgn": "dengan",
7     "goblok": "goblok",
8     # add more slang/abbrev here
9 }
10
11 # regex patterns
12 URL_RE = re.compile(r"https?://\S+|www\.\S+")
13 MENTION_RE = re.compile(r"@w+")
14 HASHTAG_RE = re.compile(r"#w+")
15 EMOJI_RE = re.compile("[
16     u"\U0001F600-\U0001F64F"
17     u"\U0001F300-\U0001F5FF"
18     u"\U0001F680-\U0001F6FF"
19     u"\U0001F1E0-\U0001F1FF"
20     "]+", flags=re.UNICODE)
21
22 PUNCT_TABLE = str.maketrans('', '', string.punctuation)
23 from nltk.corpus import stopwords
24 try:
25     stop_words = set(stopwords.words('indonesian'))
26 except:
27     stop_words = set(stopwords.words('english'))
28
```

Preprocessing data saya mulai dengan

- membunag kata-kata yang tidak baku dengan stemming
- menghilangkan karakter-karaktaer unik yaitu emoticon dengan regex
- menormalisasi kata-kata slang menjadi kata baku
- dan melakukan tokenisasi kata

```

29
30 # stemmer
31 factory = StemmerFactory()
32 stemmer = factory.create_stemmer()
33
34 def normalize_text(s):
35     s = str(s)
36     s = s.strip()
37     s = URL_RE.sub("", s)
38     s = MENTION_RE.sub("", s)
39     s = HASHTAG_RE.sub(lambda m: m.group(0)[1:], s)
40     s = EMOJI_RE.sub("", s)
41     tokens_tmp = s.split()
42     tokens_tmp = [NORMALIZE_MAP.get(t.lower(), t) for t in tokens_tmp]
43     s = " ".join(tokens_tmp)
44     s = s.translate(PUNCT_TABLE)
45     s = s.lower()
46     s = re.sub(r'\s+', ' ', s).strip()
47     return s
48
49 def preprocess_text(s, do_tokenize=True, do_stopword=True, do_stem=True):
50     s = normalize_text(s)
51     if not do_tokenize:
52         return s
53     tokens = word_tokenize(s)
54     if do_stopword:
55         tokens = [t for t in tokens if t not in stop_words and len(t) > 1]
56     if do_stem:
57         tokens = [stemmer.stem(t) for t in tokens]
58     return tokens
59
60 print(preprocess_text("Gojek lemot bgt si update tracking drivernya #gojek 🚗 https://t.co/abc"))
61

```

```

62 tqdm.pandas()
63 df['clean_tokens'] = df['tweet'].progress_apply(lambda x: preprocess_text(x))
64 df['clean_text'] = df['clean_tokens'].apply(lambda toks: " ".join(toks))

```

#### 4. Membuat wordcloud

```

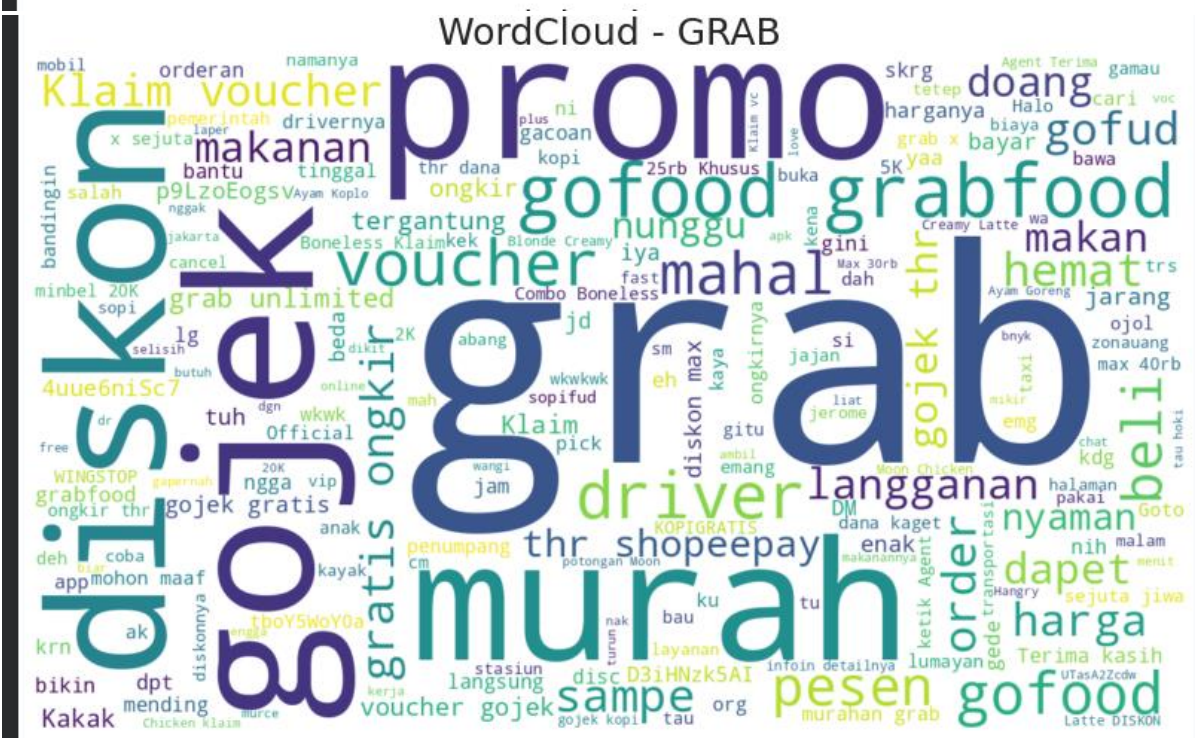
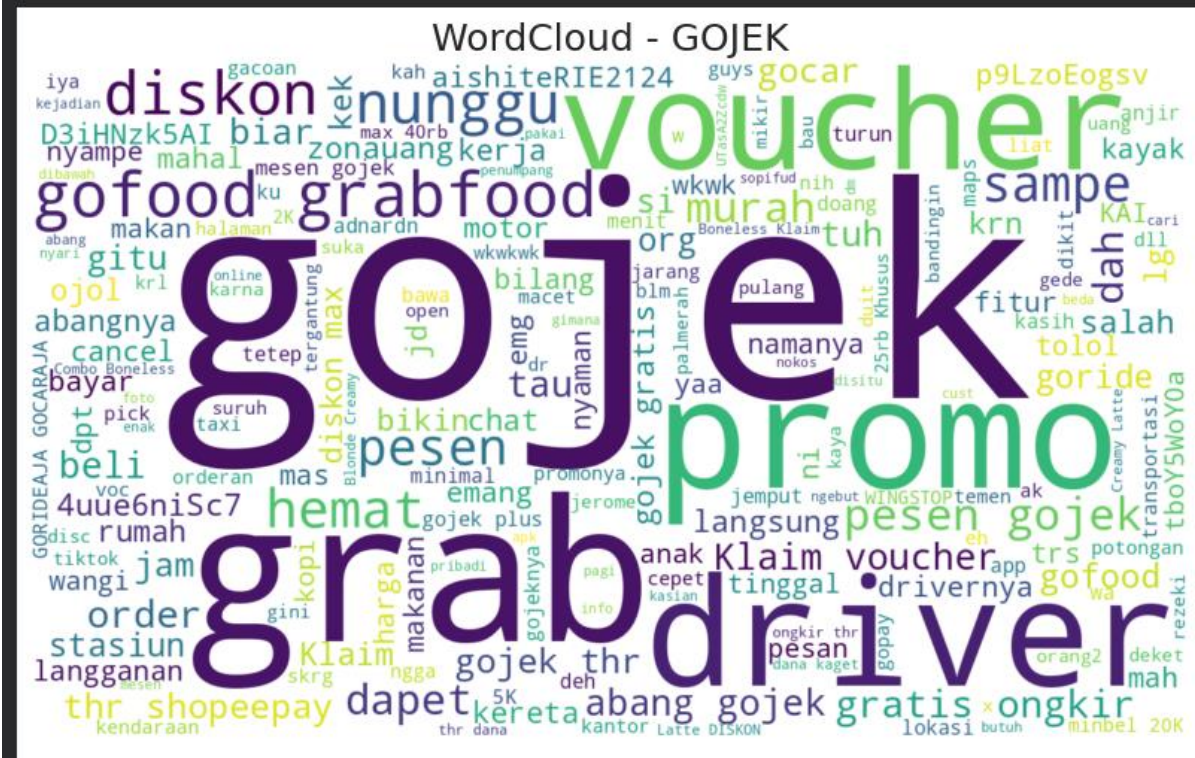
(38)
✓ 6s
1 # Stopwords dasar
2 stop_words = set(stopwords.words("indonesian"))
3 stop_words = stop_words.union(set(stopwords.words("english")))
4 stop_words = stop_words.union(STOPWORDS)
5
6 # Tambahkan stopwords manual
7 extra_sw = {"rt", "via", "amp", "co", "https", "t", "nya", "gk", "ga", "https", "shoppe", "co", "kalo", "yang", "aja", "pake", "Ordinaryjhu",
8             "shoppefood", "kalo", "aja", "pake", "yang", "udah", "ya", "shopee", "bgt", "go", "Cetepbebek", "kadang", "sih", "banget",
9             "gue", "gw", "lu", "klo", "gua", "tp", "pas", "sfood", "food", "resto", "udh", "udah", "Shopeefood", "yg", "gak", "daget",
10            "cek", "jg", "abis", "habis", "kl", "kak", "titik", "kode", "jalan", "orang", "lo", "kali", "tubirfess", "akun", "utk",
11            "gt", "dh", "seenaknya", "knp"}
12 stop_words = stop_words.union(extra_sw)
13
14 # Fungsi bikin wordcloud
15 def generate_wordcloud(text, title):
16     wc = WordCloud(
17         width=1200,
18         height=700,
19         background_color="white",
20         stopwords=stop_words,
21         collocations=True
22     ).generate(text)
23
24     plt.figure(figsize=(10, 6))
25     plt.imshow(wc, interpolation='bilinear')
26     plt.axis("off")
27     plt.title(title, fontsize=18)
28     plt.show()
29

```

Wordcloud saya buat agar memudahkan saya mengidentifikasi kata mana saja yang lebih sering muncul yang nantinya saya gunakan Ketika merancang stopwords.



```
30 # Filter tweet berdasarkan keyword
31
32 # Tweet yang mengandung "gojek"
33 df_gojek = df[df["keyword"].str.lower() == "gojek"]
34 # Tweet yang mengandung "grab"
35 df_grab = df[df["keyword"].str.lower() == "grab"]
36
37 # Gabungkan jadi satu string
38 text_gojek = " ".join(df_gojek["tweet"].astype(str).tolist())
39 text_grab = " ".join(df_grab["tweet"].astype(str).tolist())
40
41 generate_wordcloud(text_gojek, "WordCloud - GOJEK")
42 generate_wordcloud(text_grab, "WordCloud - GRAB")
43
```



## 5. Labelling dataset

```
1 # 2. Lexicon-based sentiment (simple)
2 LEXICON_CSV = "indonesian_lexicon.csv"
3 lexicon = {}
4
5 if os.path.exists(LEXICON_CSV):
6     lex = pd.read_csv(LEXICON_CSV)
7     for _, r in lex.iterrows():
8         w = str(r['word']).strip().lower()
9         s = float(r.get('score', 0))
10        lexicon[w] = s
11 else:
12     # fallback minimal lexicon (expand for better results)
13     lexicon = {
14         "bagus": 1.0, "mantap": 1.0, "mantapjiwa": 1.0, "puas": 1.0, "murah": 1.0, "baguslah": 1.0,
15         "jelek": -1.0, "buruk": -1.0, "payah": -1.0, "kurang": -0.5, "telat": -1.0, "mal": -0.5,
16         "lemot": -1.0, "nggak": -0.5, "ga": -0.5, "benci": -1.0, "suka": 1.0, "senang": 1.0,
17         "marah": -1.0, "tolol": -1.0, "nyaman": 1.0, "murah": -1.0, "bagus": 1.0,
18         "jelek": -1.0, "buruk": -1.0, "payah": -1.0, "kurang": -0.5, "telat": -1.0, "mal": -0.5,
19         "lemot": -1.0, "nggak": -0.5, "ga": -0.5, "benci": -1.0, "suka": 1.0, "senang": 1.0,
20         "nunggu": -1.0, "tolol": -1.0, "nyaman": 1.0, "hemat": 1.0, "promo": 1.0, "diskon": 1.0,
21         "voucher": 1.0, "promosi": 1.0, "ramah": 1.0, "goblok": -1.0, "jahat": -1.0, "bangsat": -1.0,
22         "tolol": -1.0, "sehat": 1.0, "enak": 1.0, "cape": -1.0, "gaada": -1.0, "gamau": -1.0,
23         "fiktif": -1.0, "pukimak": -1.0, "telat": -1.0, "komunikatif": 1.0, "baik": 1.0, "anjir": -1.0,
24         "tidak puas": -1.0, "salah": -1.0, "nyolot": -1.0, "setia": 1.0, "lama": -1.0, "jarang": -1.0,
25         "matiiin": -1.0, "mati": -1.0, "lambat": -1.0, "parah": -1.0, "parah": -1.0, "mengecewakan": -1.0,
26         "mahal": -1.0, "top": 1.0, "terbaik": 1.0, "keren": 1.0, "aman": 1.0, "cepat": 1.0,
27         "thr": 1.0, "ragu": -1.0, "jarang": -1.0, "mual": -1.0, "pernah": 1.0, "tidak pernah": -1.0,
28         "nggak pernah": -1.0, "ga pernah": -1.0, "ga ada": -1.0, "klaim": 1.0, "gratis": 1.0, "gratislah": 1.0,
29         "potongan": 1.0, "potong": 1.0, "potongin": 1.0,
30     }
31
```

Sama seperti stop words, kata-kata yang saya masukan untuk labelling manual juga berdasarkan kata yang sering muncul di wordcloud.

```
32 def lexicon_score(tokens):
33     score = 0.0
34     for t in tokens:
35         score += lexicon.get(t, 0.0)
36     return score
37
38 def lexicon_label(score, pos_thr=0.5, neg_thr=-0.5):
39     if score >= pos_thr:
40         return "positive"
41     elif score <= neg_thr:
42         return "negative"
43     else:
44         return "neutral"
45
46 df['lex_score'] = df['clean_tokens'].apply(lambda toks: lexicon_score(toks))
47 df['lex_label'] = df['lex_score'].apply(lexicon_label)
48 print("Lexicon distribution:\n", df['lex_label'].value_counts())
```

Untuk libray yang saya pakai yaitu lexicon.

```
*** Lexicon distribution:
lex_label
neutral      436
negative     312
positive     285
Name: count, dtype: int64
```

Hasil akhir labellingnya seperti pada output, namun ini masih mengacu pada 2 keyword tadi.

## 6. Ekstraksi fitur dengan TF-IDF

```

1 # 3. Feature extraction
2 # BoW and TF-IDF
3 tfidf = TfidfVectorizer(max_features=8000, ngram_range=(1,2))
4 X_tfidf = tfidf.fit_transform(df['clean_text'])
5
6 # optionally Word2Vec (train on our token lists)
7 w2v_model = Word2Vec(sentences=df['clean_tokens'].tolist(), vector_size=100, window=5, min_count=2, workers=2, epochs=10)
8 def doc_vector(tokens):
9     vecs = []
10    for t in tokens:
11        if t in w2v_model.wv:
12            vecs.append(w2v_model.wv[t])
13    if not vecs:
14        return np.zeros(w2v_model.vector_size)
15    return np.mean(vecs, axis=0)
16 X_w2v = np.vstack([doc_vector(toks) for toks in df['clean_tokens']])

```

## 7. Split data dan training dengan beberapa model

```

1 # 4. Supervised ML
2 df_ml = df.copy()
3 df_ml = df_ml[df_ml['lex_label'].isin(['positive','negative','neutral'])].reset_index(drop=True)
4 y = df_ml['lex_label'].map({'negative':0,'neutral':1,'positive':2}).values
5 X = tfidf.transform(df_ml['clean_text']) # use TF-IDF features; or use X_w2v
6
7 # train/test split
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
9
10 # models to try
11 models = {
12     "MultinomialNB": MultinomialNB(),
13     "LogisticRegression": LogisticRegression(max_iter=200),
14     "LinearSVC": LinearSVC(max_iter=2000)
15 }
16
17 results = {}
18 for name, model in models.items():
19     print("Training", name)
20     model.fit(X_train, y_train)
21     y_pred = model.predict(X_test)
22     acc = accuracy_score(y_test, y_pred)
23     report = classification_report(y_test, y_pred, target_names=['neg','neu','pos'])
24     print(name, "accuracy:", acc)
25     print(report)
26     results[name] = {
27         "model": model,
28         "accuracy": acc,
29         "report": report,
30         "y_pred": y_pred
31     }
32 best_name = max(results.keys(), key=lambda k: results[k]['accuracy'])
33 best_model = results[best_name]['model']
34 print("Best model:", best_name, results[best_name]['accuracy'])

```

Hasil:

### I. Model multinomial

|  |           |        |          |         |
|--|-----------|--------|----------|---------|
| Training MultinomialNB                     |           |        |          |         |
| MultinomialNB accuracy: 0.7439613526570048 |           |        |          |         |
|  | precision | recall | f1-score | support |
| neg  | 0.76      | 0.67   | 0.71     | 63      |
| neu  | 0.65      | 0.85   | 0.74     | 87      |
| pos  | 0.97      | 0.67   | 0.79     | 57      |
| accuracy                                   |           |        | 0.74     | 207     |
| macro avg                                  | 0.80      | 0.73   | 0.75     | 207     |
| weighted avg                               | 0.78      | 0.74   | 0.75     | 207     |

## II. Model logistic regression

|   |           |        |          |         |  |
|---|-----------|--------|----------|---------|--|
| Training LogisticRegression                     |           |        |          |         |  |
| LogisticRegression accuracy: 0.7729468599033816 |           |        |          |         |  |
|   | precision | recall | f1-score | support |  |
| neg   | 0.86      | 0.67   | 0.75     | 63      |  |
| neu   | 0.67      | 0.93   | 0.78     | 87      |  |
| pos   | 1.00      | 0.65   | 0.79     | 57      |  |
| accuracy  |           |        | 0.77     | 207     |  |
| macro avg                                       | 0.84      | 0.75   | 0.77     | 207     |  |
| weighted avg                                    | 0.82      | 0.77   | 0.77     | 207     |  |

## III. Model linearSVC

|  |           |        |          |         |  |
|--|-----------|--------|----------|---------|--|
| Training LinearSVC                       |           |        |          |         |  |
| LinearSVC accuracy: 0.8115942028985508   |           |        |          |         |  |
|  | precision | recall | f1-score | support |  |
| neg                                      | 0.84      | 0.76   | 0.80     | 63      |  |
| neu                                      | 0.74      | 0.89   | 0.81     | 87      |  |
| pos                                      | 0.93      | 0.75   | 0.83     | 57      |  |
| accuracy                                 |           |        | 0.81     | 207     |  |
| macro avg                                | 0.84      | 0.80   | 0.81     | 207     |  |
| weighted avg                             | 0.82      | 0.81   | 0.81     | 207     |  |
| Best model: LinearSVC 0.8115942028985508 |           |        |          |         |  |

## 8. Training dengan keseluruhan data menggunakan model terbaik

```
1 # menggunakan model terbaik untuk prediksi ke seluruh dataset
2 X_all = tfidf.transform(df['clean_text'])
3 y_all_pred = best_model.predict(X_all)
4 inv_map = {0:'negative',1:'neutral',2:'positive'}
5 df['ml_label'] = [inv_map[int(i)] for i in y_all_pred]
6
7 #membandingkan lexicon dengan ml counts
8 print("Lexicon counts:\n", df.groupby('keyword')['lex_label'].value_counts().unstack(fill_value=0))
9 print("\nML counts:\n", df.groupby('keyword')['ml_label'].value_counts().unstack(fill_value=0))
10
11 #Menyimpan csv dengan label
12 OUT_CSV = "tweet_gojek_grab_labeled.csv"
13 df.to_csv(OUT_CSV, index=False)
14 print("Saved labeled CSV:", OUT_CSV)
```

Hasil:

```
*** Lexicon counts:
      lex_label  negative  neutral  positive
keyword
gojek          146       232       137
grab           166       204       148

ML counts:
      ml_label  negative  neutral  positive
keyword
gojek          140       246       129
grab           166       207       145
Saved labeled CSV: tweet_gojek_grab_labeled.csv
```



## 9. Visualisasi

```
1 # 6. Evaluation helper
2 def evaluate_with_gold(y_true, y_pred, labels=['negative','neutral','positive']):
3     print(classification_report(y_true, y_pred, labels=labels))
4     cm = confusion_matrix(y_true, y_pred, labels=[0,1,2])
5     sns.heatmap(cm, annot=True, fmt="d", xticklabels=labels, yticklabels=labels)
6     plt.xlabel("pred")
7     plt.ylabel("true")
8     plt.show()
```

```
1 #Hitung jumlah
2 count_df = df.groupby(['keyword', 'ml_label']).size().unstack(fill_value=0)
3
4 #Hitung persentase
5 percent_df = count_df.div(count_df.sum(axis=1), axis=0) * 100
6
7 print("COUNT PER SENTIMENT")
8 display(count_df)
9
10 print("\n PERCENTAGE PER SENTIMENT (%) ")
11 display(percent_df.round(2))
12
13 # --- Plot ---
14 plt.figure(figsize=(10,4))
15 ax = count_df.plot(kind='bar', stacked=False)
16
17 plt.title("Sentiment Counts per Keyword (ML Labels)")
18 plt.xlabel("Keyword")
19 plt.ylabel("Count")
20
21 # tampilkan nilai di atas bar
22 for p in ax.patches:
23     ax.annotate(
24         str(p.get_height()),
25         (p.get_x() + p.get_width()/2, p.get_height()),
26         ha='center', va='bottom', fontsize=8
27     )
28
29 plt.tight_layout()
30 plt.show()
```

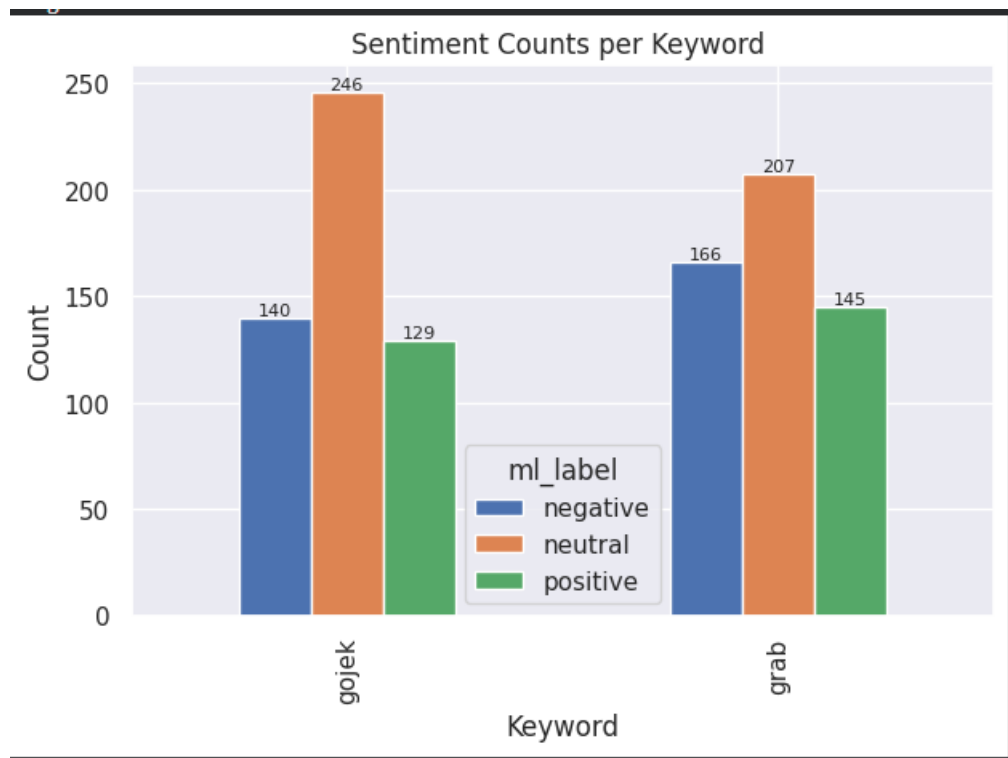
```
31
32 # Plot Persentase
33 plt.figure(figsize=(10,4))
34 ax2 = percent_df.plot(kind='bar', stacked=False)
35
36 plt.title("Sentiment Percentage per Keyword (ML Labels)")
37 plt.xlabel("Keyword")
38 plt.ylabel("Percentage (%)")
39
40 # tampilkan angka persentasenya
41 for p in ax2.patches:
42     ax2.annotate(
43         f"{p.get_height():.1f}%",
44         (p.get_x() + p.get_width()/2, p.get_height()),
45         ha='center', va='bottom', fontsize=8
46     )
47
48 plt.tight_layout()
49 plt.show()
50
```

Hasil:

| COUNT PER SENTIMENT |          |         |          |
|---------------------|----------|---------|----------|
| ml_label            | negative | neutral | positive |
| keyword             |          |         |          |
| gojek               | 140      | 246     | 129      |
| grab                | 166      | 207     | 145      |

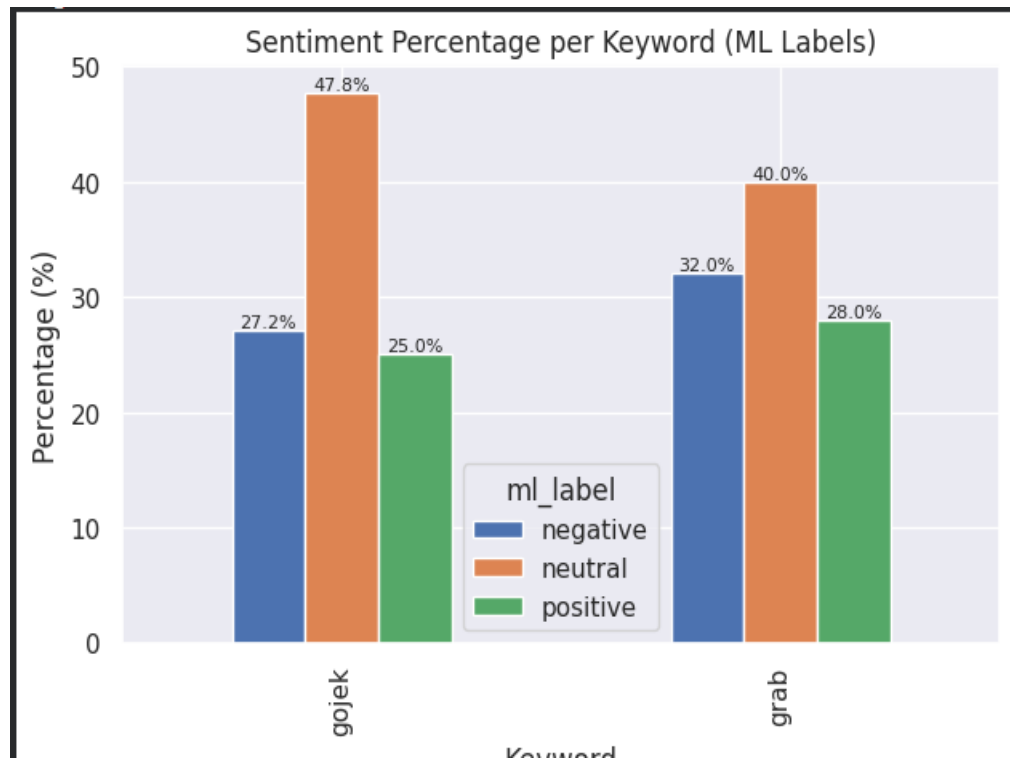
  

| PERCENTAGE PER SENTIMENT (%) |          |         |          |
|------------------------------|----------|---------|----------|
| ml_label                     | negative | neutral | positive |
| keyword                      |          |         |          |
| gojek                        | 27.18    | 47.77   | 25.05    |
| grab                         | 32.05    | 39.96   | 27.99    |



Gojek memiliki sentiment netral tertinggi, ini terjadi karena banyak kata-kata pada keyword gojek jarang muncul di wordcloud sehingga tidak terinput di dalam list labelling manual.

Sentiment negative dan positif dua-duanya paling banyak di sentiment grab.



Hasil prediksi antara model machine learning dan lexicon tidak terlalu signifikan walaupun menggunakan pendekatan yang berbeda, sentiment netral masih berada di keyword gojek.

Analisis hasil

| No | Tweet asli   | Setelah preprocessing   | Prediksi model | Analisa manual | Cocok/tidak |
|----|--|---|----------------|----------------|-------------|
| 1. | Lah ini gojek daerah mana nih bawa glock?  | gokjek daerah nih bawa glock  | netral         | netral         | Tidak cocok |
| 2. | bangun pagi ga sudi naik gojek juga ga sudi huf  | bangun pagi nggak sudi gokjek nggak sudi huf  | negative       | negative       | cocok       |
| 3. | hallo!! aku masih open nokos ridi banyak yak, gokjek, grab, wa, shopee, tokopedia yukk sung ke tele/ wa fastrespon yak 🐼 #zonauang                         | hallo open nokos ridi yak gokjek grab wa shopee tokopedia yukk sung tele wa fastrespon yak zonauang | netral         | netral         | cocok       |
| 4. | @laalaafun @cassiocpeia @gtauadasiapadmn Rating akun kita as costumer ka Kalo di gokjek driver bisa ngasih rating ke penumpang dan cuma bisa diliat driver | rating akun costumer ka gokjek driver ngasih rating tumpang liat driver                             | netral         | netral         | Tidak cocok |
| 5. | INI AMANG GOJEK NYA SPEK RACING KAHHH  | amang gokjek spek racing kahhh  | netral         | netral         | cocok       |
| 6. | gaada dan belum ada sih yag grab   | gaada bom sih yag grab  | negativ        | negativ        | cocok       |

|     |  |   |         |         |       |
|-----|--|---|---------|---------|-------|
| 7.  | for me to grab it sambil bilang<br>""istriku"  | grab bilang istri   | netral  | netral  | cocok |
| 8.  | Mari beli makanan enak sekaligus<br>memajukan UMKM<br>😊😊😊😊https://t.co/ox1dSFWiZV                      | mari beli makan<br>enak maju umkm                                     | positiv | positiv | cocok |
| 9.  | @tebakaksyapa bisa order disini,<br>ayam goreng soto jeniper by hangry<br>&lt;3https://t.co/Rh0adWiyIL | order ayam goreng<br>soto jeniper hangry<br>lt3                       | Netral  | netral  | cocok |
| 10. | knp grab di makassar kasar banget<br>kah, i literally cried bcs sekasar itu<br>tadi                    | knp grab makassar<br>kasar banget kah<br>literally cried bcs<br>kasar | negativ | negativ | cocok |

## . Refleksi dan Diskusi

1. Apakah model lebih baik menggunakan data mentah atau yang sudah distemming?  
Model lebih baik menggunakan data yang sudah di stemming, karena data mentah lebih sulit di kenali oleh model
2. Apa yang menyebabkan model salah memprediksi tweet netral?  
Karena kata-kata pada tweet netral belum di definisikan menjadi sentiment positif atau negative.
3. Jika dataset diperbesar 10x lipat, bagaimana pengaruhnya terhadap akurasi model?  
Akurasi model belum tentu lebih baik karena tantangan terbesarnya ada di preprocessing
4. Bagaimana jika digunakan model *deep learning* seperti IndoBERT?  
Mungkin bisa jadi opsi yang lebih baik jika preprocessing benar-benar di lakukan dengan teliti dan cermat.

Link repository dan dataset [Klik Di sini](#)