```
    Maftuh Mashuri (11160940000076)

    Zahrotul Aulia (11160940000024)

    Alif Tito SA (11160940000052)

    Khairul Umam (11160940000073)

    Fathi Syuhada (1115094000001)

         Import modul
In [43]: import json
         import pandas as pd
         import nltk
         from nltk.tokenize import TweetTokenizer as tw tokenizer
         from unidecode import unidecode
         import time
         Pengoperasian file
         Membuka file dan memasukkan data tweet json dan menyimpannya kedalam variabel list twetts_data
In [44]: tweets data = [] # Membuat list kosong untuk menyimpan data json perbaris
         tweets file = open('data/dataset.txt', "r") # membuka file
         for line in tweets file:
                 tweet = json.loads(line) # Membaca data dalam format json dari file perbaris
                 tweets data.append(tweet) # Menambahkan data dari file ke dalam list
                 continue
         tweets file.close()

    Mapping data json kedalam bentuk dataframe

In [45]: tweets = pd.DataFrame()
         # Informasi tweet
         tweets['created at'] = list(map(lambda tweet: tweet['created at'], tweets data))
         tweets['id str'] = list(map(lambda tweet: tweet['id str'], tweets data))
         tweets['text'] = list(map(lambda tweet: tweet['text'], tweets data))
         tweets['source'] = list(map(lambda tweet: tweet['source'], tweets data))
         tweets['lang'] = list(map(lambda tweet: tweet['lang'], tweets_data))
         # Informasi user
         tweets['user id'] = list(map(lambda tweet: tweet['user']['id str'], tweets data))
         tweets['user_name'] = list(map(lambda tweet: tweet['user']['name'], tweets_data))
         tweets['user_screen_name'] = list(map(lambda tweet: unidecode(tweet['user']['screen name']), tweets
         data))
         tweets['user location'] = list(map(lambda tweet: tweet['user']['location'], tweets data))
         tweets['user_url'] = list(map(lambda tweet: tweet['user']['url'], tweets_data))
         tweets['user description'] = list(map(lambda tweet: tweet['user']['description'], tweets data))
         tweets['user_followers_count'] = list(map(lambda tweet: tweet['user']['followers_count'], tweets_dat
         tweets['user friends count'] = list(map(lambda tweet: tweet['user']['friends count'], tweets data))
         tweets['user_favourites_count'] = list(map(lambda tweet: tweet['user']['favourites_count'], tweets_d
         tweets['user statuses count'] = list(map(lambda tweet: tweet['user']['statuses count'], tweets data
         tweets['user created at'] = list(map(lambda tweet: tweet['user']['created at'], tweets data))
         # Informasi Tempat
         tweets['address'] = list(map(lambda tweet: tweet['place']['full_name'] if tweet['place'] != None els
         e None, tweets data))
         tweets['country'] = list(map(lambda tweet: tweet['place']['country'] if tweet['place'] != None else
         None, tweets data))
         # Entities
         hashtags = list(map(lambda tweet: tweet['entities']['hashtags'] if tweet['entities'] != None else No
         ne, tweets _data))
         tweets['hashtags'] = list(map(lambda tweet: ', '.join(list(map(lambda tw: tw['text'], tweet))), ha
         shtags))
         urls = list(map(lambda tweet: tweet['entities']['urls'] if tweet['entities'] != None else None, twee
         tweets['url'] = list(map(lambda tweet : ', '.join(list(map(lambda tw : tw['url'], tweet))), urls))
         mentions = list(map(lambda tweet: tweet['entities']['user mentions'] if tweet['entities'] != None el
         se None, tweets data))
         tweets['mentions'] = list(map(lambda tweet : ', '.join(list(map(lambda tw : tw['screen name'], tweet
         ))), mentions))
         # Informasi quoted status
         tweets['quoted_status_id'] = list(map(lambda tweet : tweet['quoted_status']['id_str'] if 'quoted_sta
         tus' in tweet.keys() else None, tweets data))
         tweets['quoted_status_name'] = list(map(lambda tweet : tweet['quoted_status']['user']['name'] if 'qu
         oted status' in tweet.keys() else None, tweets data))
         tweets['quoted_status_screen_name'] = list(map(lambda tweet : tweet['quoted_status']['user']['screen
          name'] if 'quoted status' in tweet.keys() else None, tweets data))
         tweets['quoted_status_text'] = list(map(lambda tweet : tweet['quoted_status']['text'] if 'quoted_sta
         tus' in tweet.keys() else None, tweets data))
         tweets['quoted_status_quote_count'] = list(map(lambda tweet : tweet['quoted_status']['quote_count']
         if 'quoted status' in tweet.keys() else None, tweets data))
         tweets['quoted_status_reply_count'] = list(map(lambda tweet : tweet['quoted_status']['reply_count']
         if 'quoted status' in tweet.keys() else None, tweets data))
         tweets['quoted_status_retweet_count'] = list(map(lambda tweet : tweet['quoted_status']['retweet_coun
         t'] if 'quoted status' in tweet.keys() else None, tweets data))
         tweets['quoted_status_favorite_count'] = list(map(lambda tweet : tweet['quoted_status']['favorite_co
         unt'] if 'quoted_status' in tweet.keys() else None, tweets_data))
         tweetsIn = tweets[tweets.lang == 'in']
         tweetsIn.to csv("data/" + str(time.time()) + " export dataframe.csv") # Eksport dataframe ke csv
         # tweetsIn.to_excel("data/" + str(time.time()) + ' export dataframe.xlsx') # Eksport dataframe ke ex
         Pengoperasian stopword

    Membuka file stopword(indonesia, inggris, noise)

    kemudian menggabungkan semua kata yang ada di ketiga file sehingga menjadi satu teks panjang dan

             menyimpannya kedalam variabel stopword_file_all
           • kemudian teks panjang di tokenize(dipisah perkata) dan menyimpannya ke dalam variabel list stopwords
In [46]: stopword_file1 = open('stopword_id.txt', "r").read() # Membuka file stopword bahasa indones
         ia dan menjadikan isi file tersebut sebagai string
         stopword file2 = open('stopword en/stopwords en.txt', "r").read() # Membuka file stopword bahasa in
         ggris dan menjadikan isi file tersebut sebagai string
         stopword file3 = open('stopword noise/stopword noise.txt', "r").read() # Membuka file stopword noise
         e dan menjadikan isi file tersebut sebagai string
         stopword file all = stopword file1 + stopword file2 + stopword file3 # Menggabungkan ketiga string
          stopword sebelumnya kedalam satu string
         stopwords = stopword file all.split('\n') # Memisahkan kata dalam string yang sudah digambungkan ber
         dasarkan baris
         # print(stopwords)
         Pengoperasian slangwords

    Membuka file slangword(colloquial-indonesian-lexicon.csv dan 20190327_slangword.txt)

    File yang pertama membuka dengan modul pandas dan mengkonfersinya kedalam dataframe, kemudian menyimpan

             masing-masing kata kedalam variabel dictionary slangwords

    File yang kedua sama halnya pengoperasian pada file stopwords kemudian menyimpan masing-masing kata

             kedalam variabel dictionary slangwords
In [47]: slangwords = dict() # Membuat dictionary kosong untuk menyimpan kata slang dan formal sebagai key da
         n value
         slangwords dataframe = pd.read csv('slangword/colloquial-indonesian-lexicon.csv') # Membuka file csv
         yang berisi kata slang dan formal dan mengkonversi kedalam dataframe
         for slang, formal in zip(slangwords_dataframe['slang'], slangwords_dataframe['formal']):
             slangwords[slang] = formal # Mapping kata slang dan formal dan memasukkan ke dalam dictionary se
         cara berulang
         slangword file = open('slangword/slangword.txt', "r").read() # Membuka file yang berisi kata slang d
         an kata formal dan mengkonversi kedalam string
         slangwords_text = slangword_file.split('\n') # Memisahkan kata berdasarkan baris namun kata slang da
         n kata formal masih belum terpisah. output : (['slang:formal', ...])
         #print(slangwords text)
         for slang in slangwords text:
             split slang = slang.split(":") # Memisahkan semua kata slang dan kata formal berdasarkan "titik
             slangwords[split slang[0]] = split slang[1] # Mapping semua kata slang dan kata formal ke dalam
          dictionary. Output : {'slang' : 'formal', ...}
         #print(slangwords)
         2-3. Tokenisasi dan Filtering

    Mapping semua tweet dan menyimpannya kedalam variabel list array_text

           • Menggabungkan semua teks kedalam satu teks panjang long text
           · Memisahkan teks perkata
           • Menghapus simbol, ASCII, link (https: atau www.), dan kata lainnya
In [48]: array text = list(map(lambda tweet: unidecode(tweet).lower(), tweets['text'])) # Mapping semua text
          twitter dan memasukan ke dalam list. Output : ['teks panjang ...', 'teks panjang', ...]
         long text = ' '.join(array text) # Menggabungkan semua text yang berada dalam list kedalam satu text
         tokenized = tw tokenizer().tokenize(long text) # Memisahkan kata dalam text berasarkan "spasi"
         filtered alfanumeric = [w for w in tokenized if w.isalnum()] # Filtering kata yang hanya berisi kara
         kater a-z dan 0-9 (Menghapus url, hashtag, mention)
         print(filtered alfanumeric[:100])
         ['whahaha', 'iya', 'sih', 'bener', 'ledakan', 'mercon', 'besar', 'menghancurkan', 'rumah', 'nanan
         g', 'menewaskan', 'asih', 'pembantunya', 'se', 'kalideres', 'arah', 'batu', 'ceper', 'tangerang',
         'tol', 'cawang', 'tmii', 'cibubur', 'bogor', 'ciawi', 'well', 'makin', 'kesini', 'ku', 'tahu', 'y
         ang', 'mana', 'yang', 'minta', 'tolong', 'dan', 'yang', 'manfaatin', 'doang', 'honestly', 'setela
         h', 'balik', 'ke', 'indonesia', 'hal', 'yg', 'paling', 'berasa', 'beda', 'di', 'gw', 'itu', 'adal
         ah', 'cuaca', 'astagah', 'tiap', 'jam', 'keringeta', 'kang', 'boong', 'di', 'bohongin', 'gmn', 'r
         asanya', 'dikecewain', 'tapi', 'tetep', 'cinta', 'sama', 'hehe', 'dasar', 'akuu', 'hee', 'aku',
         'jadi', 'inget', 'minggu', 'kemaren', 'ga', 'sengaja', 'abis', 'nginjek', 'kocheng', 'ya', 'alla
         h', 'kochengnya', 'biasa', 'a', 'yuhh', 'winterfell', 'sunyi', 'lah', 'donnaruma', 'ngapa', 'it
         u', 'udh', 'diganti', 'aja', 'kalo', 'ga']
         4. Stopword removal
         Menghapus kata-kata yang sering muncul dan tidak memiliki makna (yang, di, kan)
In [49]: removed stopwords = [w for w in filtered alfanumeric if w not in stopwords]
          # Filtering data dengan menghapus kata yang tidak bermakna (Stopword yang diperoleh dari file)
         print(removed stopwords[:100])
         ['whahaha', 'bener', 'ledakan', 'mercon', 'menghancurkan', 'rumah', 'nanang', 'menewaskan', 'asi
         h', 'pembantunya', 'kalideres', 'arah', 'batu', 'ceper', 'tangerang', 'tol', 'cawang', 'tmii', 'c
         ibubur', 'bogor', 'ciawi', 'kesini', 'ku', 'tolong', 'manfaatin', 'honestly', 'indonesia', 'yg',
         'berasa', 'beda', 'gw', 'cuaca', 'astagah', 'jam', 'keringeta', 'kang', 'boong', 'bohongin', 'gm
         n', 'dikecewain', 'tetep', 'cinta', 'dasar', 'akuu', 'hee', 'inget', 'minggu', 'kemaren', 'ga',
         'sengaja', 'abis', 'nginjek', 'kocheng', 'allah', 'kochengnya', 'yuhh', 'winterfell', 'sunyi', 'd
         onnaruma', 'ngapa', 'udh', 'diganti', 'ga', 'bales', 'tdr', 'ga', 'tdr', 'siang', 'tidur', 'sian
         g', 'dibilang', 'tdr', 'mulu', 'ga', 'tdr', 'siang', 'dikasianin', 'bucin', 'banget', 'gak', 'op
         o', 'opo', 'iso', 'kecuali', 'siji', 'tangi', 'esuk', '1', 'nggak', 'dikasih', 'nggak', 'ngembe
         r', 'nggak', 'cerita', 'alesannya', 'simpl', 'nggak', 'baca', 'sepotong', 'berita']
         5. Handling SlangWord
         Mengubah kata-kata yang sudah melalui tahap Stopword removal dari kata slang menjadi kata formal
In [50]: handled slangword = list(map(lambda w : slangwords[w] if w in slangwords.keys() else w, removed stop
         # Mengubah kata slang menjadi kata formal (kata slang dan kata formal yang diperoleh dari dictionar
         print(handled slangword[:100])
         ['whahaha', 'benar', 'ledakan', 'mercon', 'menghancurkan', 'rumah', 'nanang', 'menewaskan', 'asi
         h', 'pembantunya', 'kalideres', 'arah', 'batu', 'pendek', 'tangerang', 'tol', 'cawang', 'tmii',
         'cibubur', 'bogor', 'ciawi', 'kesini', 'ku', 'tolong', 'manfaatin', 'honestly', 'indonesia', 'yan
         g', 'berasa', 'beda', 'saya', 'cuaca', 'astagah', 'jam', 'keringeta', 'kang', 'bohong', 'bohongi
         n', 'bagaimana', 'dikecewain', 'tetap', 'cinta', 'dasar', 'aku', 'hee', 'ingat', 'minggu', 'kemar
         in', 'tidak', 'sengaja', 'habis', 'menginjak', 'kocheng', 'allah', 'kochengnya', 'yuhh', 'winterf
         ell', 'sunyi', 'donnaruma', 'mengapa', 'sudah', 'diganti', 'tidak', 'balas', 'tidur', 'tidak', 't
         idur', 'siang', 'tidur', 'siang', 'dibilang', 'tidur', 'mulu', 'tidak', 'tidur', 'siang', 'dikasi
         anin', 'bucin', 'banget', 'tidak', 'opo', 'opo', 'iso', 'kecuali', 'siji', 'tangi', 'esuk', '1',
         'tidak', 'dikasih', 'tidak', 'ngember', 'tidak', 'cerita', 'alesannya', 'simpl', 'tidak', 'baca',
         'sepotong', 'berita']
         6. Lemmalization
         mengubah kata-kata dalam dataset tweet menjadi kata dasar dengan menggunakan modul spacy dan sastrawi
         a.) Sastrawi
         Untuk data yang sangat besar, tidak disarankan menggunakan sastrawi karena membutuhkan running time yang sangat
         besar. Sampai laporan ini dibuat, belum bisa lemmatize kata dari dataset dengan menggunakan sastrawi
         stemmed by sastrawi = [stemmer.stem(w) for w in handled slangword]
In [51]: start time sastrawi = time.time()
         from Sastrawi.Stemmer.StemmerFactory import StemmerFactory # Import modul stemmer dari sastrawi
         # create stemmer
         factory = StemmerFactory() # Membuat factori untuk stemmer dengan memanggi objek StemmerFactory()
         stemmer = factory.create stemmer() # Membuat stemmer
         # stemmed by sastrawi = [stemmer.stem(w) for w in handled slangword] # Mengubah kata-kata menjadi ka
         ta dasar dengan menggunakan modul sastrawi
         print(stemmer.stem("Melihat"))
         end time sastrawi = time.time()
         total time sastrawi = end time sastrawi - start time sastrawi
         print("Hasil stem menggunakan modul Sastrawi adalah", total time sastrawi, "detik")
         lihat
         Hasil stem menggunakan modul Sastrawi adalah 0.4459097385406494 detik
         b.) Spacy
In [52]: start time spacy = time.time()
         from spacy.lang.id import Indonesian # Import modul spacy bahasa indonesia
         nlp = Indonesian() # memanggi objek Indonesian() pada modul spacy
         def stem spacy(text): # Fungsi untuk mengubah kata-kata menjadi kata dasar
             for txt in nlp(text):
                 t = txt.lemma
             return t
         stemmed by spacy = [stem spacy(w) for w in handled_slangword] # Mengubah kata-kata menjadi kata das
         ar dengan menggunakan modul spacy
         end_time_spacy = time.time()
         total time spacy = end time spacy - start time spacy
         print("Hasil stem menggunakan modul Spacy adalah", total_time_spacy, "detik\n") # Hasil running tera
         khir adalah 67.62633967399597 detik
         print(stemmed_by_spacy[:100])
         Hasil stem menggunakan modul Spacy adalah 43.73394703865051 detik
         ['whahaha', 'benar', 'ledak', 'mercon', 'hancur', 'rumah', 'nanang', 'tewas', 'asih', 'bantu', 'k
         alideres', 'arah', 'batu', 'pendek', 'tangerang', 'tol', 'cawang', 'tmii', 'cibubur', 'bogor', 'c
         iawi', 'kesini', 'ku', 'tolong', 'manfaatin', 'honestly', 'indonesia', 'yang', 'rasa', 'beda', 's
         aya', 'cuaca', 'astagah', 'jam', 'keringeta', 'kang', 'bohong', 'bohongin', 'bagaimana', 'dikecew
         ain', 'tetap', 'cinta', 'dasar', 'aku', 'hee', 'ingat', 'minggu', 'kemarin', 'tidak', 'sengaja',
         'habis', 'injak', 'kocheng', 'allah', 'kochengnya', 'yuhh', 'winterfell', 'sunyi', 'donnaruma',
         'apa', 'sudah', 'diganti', 'tidak', 'balas', 'tidur', 'tidak', 'tidur', 'siang', 'tidur', 'sian
         g', 'dibilang', 'tidur', 'mulu', 'tidak', 'tidur', 'siang', 'dikasianin', 'bucin', 'banget', 'tid
         ak', 'opo', 'opo', 'iso', 'kecuali', 'siji', 'tangi', 'esuk', '1', 'tidak', 'dikasih', 'tidak',
         'ngember', 'tidak', 'cerita', 'alesannya', 'simpl', 'tidak', 'baca', 'sepotong', 'berita']
         Menghitung frekuensi semua kata yang muncul
         Kata yang belum melalui tahap 2-6 disimpan ke dalam variabel count_words_before
         Kata yang sudah melalui tahap 2-6 disimpan ke dalam variabel count_words_after
In [53]: from collections import Counter
         Counter1 = Counter(tokenized) # Menghitung frekuensi muncul semua kata
         count_words_before = Counter1.items() # Menghitung kata
         #print(count words before)
         Counter2 = Counter(stemmed_by_spacy) # Menghitung frekuensi muncul semua kata
         count words after = Counter2.items() # Menghitung kata (Tidak termasuk kata yang tidak bermakna)
         #print(count words after)
         Fungsi untuk membuat wordcloud
         Fungsi tersebut untuk membuat wordcloud dan file yang berisi kata beserta frekuensinya
         Dengan parameter:

    count_words, yaitu untuk menginput dictionary count_words_before dan count_words_after

           • create_file, yaitu boolean untuk membuat file atau tidak
           • after, yaitu boolean untuk menentukan dia sudah melalui proses 2-6 atau belum
In [60]: import matplotlib.pyplot as plt
         from wordcloud import WordCloud
         def create wordcloud(count words, create file = False, after = True):
             text = ""
             data_tweet = {}
             for element in count_words:
                 text += str(element[0]) + " " + str(element[1]) + "\n"
                 data_tweet[element[0]] = element[1]
                 nama file = "sesudah"
             else:
                 nama file = "sebelum"
             if create_file:
                 file = open("data/" + str(time.time()) + "_wordcloud_" + nama_file + "_filtering.txt", "w").
         write(text) # membuat file wordcloud
             wordcloud = WordCloud(background_color = 'white', max_words=100, contour_width=2)
             wordcloud.generate_from_frequencies(frequencies=data_tweet)
             plt.figure(figsize=(20,10))
             plt.title("Wordloud hasil dari dataset " + nama_file + " melalui proses 2-6\n", fontsize=40)
             plt.imshow(wordcloud)
             plt.axis("off")
             plt.tight_layout(pad=0)
             plt.show()
         7. Menampilkan semua kata unik sebelum dilakukan proses poin
         2-6
In [61]: create_wordcloud(count_words_before, True, False)
                Wordloud hasil dari dataset sebelum melalui proses 2-6
         8. Menampilkan semua kata unik sesudah dilakukan proses poin
         2-6
In [62]: create_wordcloud(count_words_after, True)
                Wordloud hasil dari dataset sesudah melalui proses 2-6
                                                   • minggubelum ■ alhamdulillah karena • kemarin jadi
                     menang
                                    Olagi
            sakit
In [63]: data user = pd.DataFrame()
         data user['username'] = tweets['user screen name']
         data user['nama'] = tweets['user name']
         data user['total like'] = tweets['user favourites count']
         data user['quoted total retweet'] = tweets['quoted status retweet count']
         9. Menampilkan Top 10 user yang paling banyak di retweet
In [64]: data user sorted by retweet = data user.sort values(by=['quoted total retweet'], ascending=False)
         data user droped duplicates retweet = data user sorted by retweet.drop duplicates(subset='username',
         keep='first')
         top ten retweet = data user droped duplicates retweet.head(10)
         top ten retweet
Out[64]:
                                                |total_like|quoted_total_retweet
                    username
                                                 891
                                                          411554.0
          24264
                RanjitaTessa
                              supergirl:3
          14511
                voear7
                                                          374673.0
                              Voeller Ari7onang
          20025
                glossy
                        baby
                                                 16745
                                                          362080.0
                              زحر
          20001
                              GRACE
                                                 2472
                grciael
                                                          361667.0
                                                 217
          5292
                BudiNugroho29 | Mendoan Addict.
                                                          279108.0
          16113
                HeliHuriah
                                                 158
                                                          278130.0
                              thkim id
          16049
                park nurull
                              Pipit N Fhitriyah
                                                 365
                                                          275919.0
          68924
                salsa zahra24
                              'call me by your name'
                                                1359
                                                          236910.0
          28622
                sexy4yennie
                                                 15057
                                                          232927.0
          26862 | KimHyunSun3
                                                 9424
                                                          202560.0
                              김휸순□□
         Membuat grafik 10 user retweet terbanyak
In [65]: plot = top ten retweet.plot(kind='bar', x='username', y='quoted total retweet')
         plot.set title('Top 10 user dengan retweet terbanyak', fontsize=15)
         plot.set xlabel('Username', weight='bold', labelpad=15)
         plot.set ylabel('Jumlah retweet', weight='bold', labelpad=15)
         plot.tick params(axis='x', pad=5)
                     Top 10 user dengan retweet terbanyak
                                           quoted total retweet
             400000
             350000
             300000
          umlah retweet
             250000
             200000
             150000
             100000
              50000
                        voear7
                             baby
                                 grciael
                                             park_nurull
                                     BudiNugroho29
                                                  salsa_zahra24
                    RanjitaTessa
                                         HeliHuriah
                                     Username
         10. Menampilkan Top 10 user dengan likes paling banyak
         data user sorted by like = data user.sort values(by=['total like'], ascending=False)
         data user droped duplicates = data user sorted by like.drop duplicates(subset='username', keep='firs
         top_ten_like = data_user_droped_duplicates.head(10)
         top ten like
Out[66]:
                                                            quoted total retweet
                     username
                                             nama
                                                   total_like
          55156
                susi079
                                                   468890
                               Francisca Susi
                                                            NaN
          24340
                riskifebriyan94
                                                  214262
                               Riski Febriyan Isaputra
                                                            NaN
          58606
                weatherarena
                               weatherarena
                                                   188499
                                                            NaN
          33678
                trackoftear
                               amirah
                                                   124330
                                                            NaN
```

TUGAS III

PENGANTAR DATA MINING

Nama Anggota:

Membuat grafik 10 user like terbanyak

In [67]: plot = top\_ten\_like.plot(kind='bar', x='username', y='total\_like')
 plot.set\_title('Top 10 user dengan likes terbanyak', fontsize=15)
 plot.set\_xlabel('Username', weight='bold', labelpad=15)
 plot.set\_ylabel('Jumlah like', weight='bold', labelpad=15)

118008

114416

113514

104290

103947

98978

NaN

NaN

NaN

NaN

NaN

NaN

plot.set\_ylabel('Jumlah like', weight='bold', labelpad=15)
plot.tick\_params(axis='x', pad=5)

nabilanabilan2

jessiej99142542

MuhammadDarry

okta\_nindya\_22

yuliarti22

ewin8923

nabilanabilan

Irenlim□

ewin winarti

Yuliarti22Hardjono

Mohammad Darry

fadillah oktanindya

69579

53630

22767

70342

34477

52658