

APPENDIX

CRAWLING CODING

```
1. import tweepy
2. import csv
3. from tweepy import OAuthHandler
4. import string
5.
6. consumer_key = 'Gq6DQpjsxV4tMo3BH2zZmmD51T'
7. consumer_secret = 'Bg7XsPSmnlWyKeZLtZrsjW34WdULkGcUWHxwACIwa1J4exDlsK'
8. access_token = '1456889597883912228-CePiSfi4SZWXVGPOkDRLa9n8WDxKJc'
9. access_secret = 'yTNTjtituna3Xa2pEAaatbhgmt8GRbcW5nh8GiyM68WP5'
10.
11. auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
12. auth.set_access_token(access_token, access_secret)
13. api = tweepy.API(auth,wait_on_rate_limit=True)
14.
15. csvFile = open('/content/sample_data/hasilcrawling.csv', 'a', encoding=
    'utf-8')
16. csvWriter = csv.writer(csvFile)
17.
18. from google.colab import drive
19. drive.mount('/content/drive')
20.
21. for tweet in tweepy.Cursor(api.search,q="vaksin, covid", lang="id", sin
    ce="2021-11-01").items(10000):
22. print (tweet.created_at,tweet.author.screen_name,tweet.text)
23. csvWriter.writerow([tweet.created_at, tweet.author.screen_name, tweet.t
    ext])
```

PRE-PROCESSING

```
1. !pip install Sastrawi
2. !pip install nltk
3.
4. import pandas as pd
5. import numpy as np
6. import matplotlib.pyplot as plt
7. from sklearn.feature_extraction.text import CountVectorizer
8.
9. #download the data
10. import nltk
11. nltk.download('stopwords')
12.
13. #using the stopwords
14. From nltk.corpus import stopwords
15.
16. #initialize the stopwords
```

```

17. Stoplist = stopwords.words('indonesian')
18.
19. import string
20. import re
21. from Sastrawi.StopWordRemover.SopWordRemoverFactory import
    StopWordRemoverFactory
22. import csv
23.
24. From google.colab import drive
25. Drive.mount('/content/drive')
26.
27. df = pd.read_csv('/content/drive/Mydrive/coba/datacrawling.csv',
    encoding = 'unicode_escape')
28.
29. df.head()
30.
31. Print('Dataset size:',df.shape)
32. Print('Columns are:',df.columns)
33.
34. #case folding
35. df['TWEET'] = df['tweet'].str.lower()
36. df.head()
37.
38. def remove_punct(text):
39.     #text = "".join([char for char in text if char not in
    string.punctuation])
40.     text = re.sub(r'^[a-zA-Z0-9]', ' ', str(text))
41.     text = re.sub(r'\b\w{1,2}\b', '', text) #menghilangkan 2 kata
42.     text = re.sub(r'\s\s+', ' ', text)
43.     text = re.sub(r"\d+", "", text)
44.     return text
45. df['tweet_clean'] = df['TWEET'].apply(lambda x: remove_punct(x))
46. df.head()
47.
48. #tokenization
49. def tokenization(text):
50.     text = re.split('\W+', text)
51.     return text
52.
53. df['TOKENIZATION'] = df['tweet_clean'].apply(lambda x:
    tokenization(x.lower()))
54. df.head()
55.
56. #filtering stopwords removal
57. Stopword = nltk.corpus.stopwords.words('indonesian')
58.
59. def remove_stopwords(text):
60.     text = [word for word in text if word not in stopword]
61.     return text
62.
63. df['STOP_REMOVAL'] = df['TOKENIZATION'].apply(lambda x:
    remove_stopwords(x))

```

```

64. df.head()
65.
66. Stop_removal = df[['STOP_REMOVAL']]
67.
68. def fit_stopwords(text):
69.     text = np.array(text)
70.     text = ' '.join(text)
71.     return text
72.
73. df['STOP_REMOVAL'] = df['STOP_REMOVAL'].apply(lambda x:
fit_stopwords(x))
74. df.head()
75.
76. df.drop('tweet',axis=1)
77.
78. df.to_csv('/content/drive/MyDrive/dataskripsi/hasilpreprocessing.csv',
sep= ',' , encoding='utf-8')

```

CODING SVM

```

1. from flask import Flask, render_template, url_for
2. import numpy as np
3. import pandas as pd
4. import csv
5. import matplotlib.pyplot as plt
6. from sklearn import model_selection
7. from sklearn.model_selection import train_test_split
8. from sklearn.feature_extraction.text import TfidfVectorizer
9. from sklearn import svm
10. from sklearn.metrics import accuracy_score
11.
12. # Packages for visuals
13. import matplotlib.pyplot as plt
14. import seaborn as sns; sns.set(font_scale=1.2)
15.
16. #labelling
17. import pandas as pd
18. from textblob import TextBlob
19.
20. from google.colab import drive
21. drive.mount('/content/drive')

```

```

22.
23. df = pd.read_csv('/content/drive/MyDrive/dataskripsi/
    hasilpreprocessing.csv', encoding = 'unicode_escape')
24. df['label'] = ''
25. for i,x in df.tweet_akhir.iteritems():
26.     label = TextBlob(x)
27.     df['label'][i] = label.sentiment.polarity
28.     print("Index: ",i, "label", label.sentiment.polarity)
29.
30. def polarity_to_label(x):
31.     if(x >= -1 and x < 0):
32.         return 'negatif'
33.     if(x == 0):
34.         return 'neutral'
35.     if(x > 0 and x <=1):
36.         return 'positif'
37. df.label = df.label.apply(polarity_to_label)
38. df
39.
40. df.to_csv('/content/drive/MyDrive/dataskripsi/
    hasillabelling.csv', sep=',', encoding='utf-8')
41.
42. # split data
43. # Split into train and test data
44. Train_X, test_X, train_Y, test_Y=model_selection.train_test_split
    (df[' tweet_akhir '], df['label'], test_size = 0.1, random_state = 0)
45. # state = 0 here there is no randomization on the split data which
    means the order is still the same
46.
47. df_train = pd.DataFrame()
48. df_train[' tweet_akhir '] = train_X
49. df_train['label'] = train_Y
50.
51. df_test = pd.DataFrame()
52. df_test[' tweet_akhir '] = test_X
53. df_test['label'] = test_Y

```

```
54.
55. df_train
56.
57. df_test
58.
59. df_train.to_csv(r"/content/drive/MyDrive/dataskripsi/df_train.csv")
60. df_test.to_csv(r"/content/drive/MyDrive/dataskripsi/df_test.csv")
61.
62. #TF-IDF to analyze the relationship between a phrase or sentence and a
    set of documents
63.
64. from sklearn.feature_extraction.text import TfidfVectorizer
65.
66. tfidf_vect = TfidfVectorizer(max_features = 5000)
67. tfidf_vect.fit(df[' tweet_akhir '])
68. train_X_tfidf = tfidf_vect.transform(df_train['STOP_REMOVAL'])
69. test_X_tfidf = tfidf_vect.transform(df_test['STOP_REMOVAL'])
70.
71. tfidf_vect
72.
73. print(train_X_tfidf)
74.
75. print(test_X_tfidf)
76.
77. print(train_X_tfidf.shape)
78. print(test_X_tfidf.shape)
79.
80. #the syntax below is used to see the learned vocabulary of the corpus
81. print(tfidf_vect.vocabulary_)
82.
83. # training process
84. from sklearn.svm import SVC
85.
86. model = SVC(kernel='linear')
87. model.fit(train_X_tfidf,train_Y)
88.
```

```

89. # testing process
90. from sklearn.metrics import accuracy_score
91.
92. predictions_SVM = model.predict(test_X_tfidf)
93. test_prediction = pd.DataFrame()
94. test_prediction[' tweet_akhir '] = test_X
95. test_prediction['label'] = predictions_SVM
96. SVM_accuracy = accuracy_score(predictions_SVM, test_Y)*100
97. SVM_accuracy = round(SVM_accuracy,1)
98.
99. test_prediction
100.
101. test_prediction.to_csv(r"/content/drive/MyDrive/dataskripsi/test_predi
    ction.csv")
102.
103. SVM_accuracy
104.
105. #Accuracy, Precision, Recall, f1-score
106.
107. from sklearn.metrics import classification_report
108.
109. print ("\nLaporan klasifikasi:")
110. print (classification_report(test_Y, predictions_SVM))
111.
112. #Bar plot for commentary analysis of covid vaccine
113. labels = ['negatif', 'neutral', 'positif']
114. Category1 = [80, 9227, 674]
115. plt.bar(labels, Category1, tick_label=Labels, width=0.5, color['navy',
    'c', 'blue'])
116. plt.xlabel('Kelas Sentimen')
117. plt.ylabel('Data')
118. plt.title('Diagram Bar Data Analisis Sentimen')
119. plt.savefig(r"/content/drive/MyDrive/dataskripsi/bar_data.png")
120. plt.show()
121.
122. #Bar plot untuk train set

```

```
123.labels = ['negatif', 'neutral', 'positif']
124.Category1 = [72, 8318, 592]
125.Plt.bar(labels, Category2, tick_label=Labels, width=0.5, color['navy',
    'c', 'blue'])
126.plt.xlabel('Kelas Sentimen')
127.plt.ylabel('Data')
128.plt.title('Diagram Bar Data Latih')
129.plt.savefig(r"/content/drive/MyDrive/dataskripsi/bar_datalatih.png")
130.plt.show()
131.
132.#Bar plot untuk test set
133.labels = ['negatif', 'neutral', 'positif']
134.Category1 = [8, 909, 82]
135.Plt.bar(labels, Category3, tick_label=Labels, width=0.5, color['navy',
    'c', 'blue'])
136.plt.xlabel('Kelas Sentimen')
137.plt.ylabel('Data')
138.plt.title('Diagram Bar Data Uji')
139.plt.savefig(r"/content/drive/MyDrive/dataskripsi/bar_datauji.png")
140.plt.show()
141.
142.#Bar plot untuk klasifikasi dengan SVM
143.labels = ['negatif', 'neutral', 'positif']
144.Category1 = [5, 921, 73]
145.Plt.bar(labels, Category4, tick_label=Labels, width=0.5, color['navy',
    'c', 'blue'])
146.plt.xlabel('Kelas Sentimen')
147.plt.ylabel('Data')
148.plt.title('Diagram Bar pada Klasifikasi dengan Support Vector
    Machine')
149.plt.savefig(r"/content/drive/MyDrive/dataskripsi/bar_svm.png")
150.plt.show()
```



0.78% PLAGIARISM
APPROXIMATELY

Report #14353461

INTRODUCTION Background The covid vaccine has been implemented in early 2021. But until now there are still many people who have not been vaccinated. some people refuse to be vaccinated because they think that covid doesn't exist, are worried about side effects after vaccines, and don't believe in certain vaccines. Nowadays people are free to express their opinions and opinions on social media. Twitter is one of the social media that is often used to express opinions. Currently, many Twitter users provide opinions or comments about vaccines in Indonesia. But the tweets that appear are still arranged randomly, making it difficult for readers to find out more negative comments or appear on the topic of this covid vaccine. Based on the explanation above, a sentiment analysis research was conducted on Twitter to classify Twitter user comments about the Covid vaccine in Indonesia. To conduct this research, firstly, a data crawling process will be carried out whose function is to retrieve tweet data, after that the data will be pre-processed, then