

## APPENDIX

Coding dibawah ini digunakan untuk memproses tf-idf

```

25. import nltk
26. nltk.download('stopwords')
27. from nltk.corpus import stopwords
28. import numpy as np
29. import matplotlib.pyplot as plt
30. import pandas as pd
31. import sklearn
32. import re
33. def tokens(text):
34.     return re.findall('[\w]+', text.lower())
35.
36. alltok = []
37. filtertrain = []
38. for t in tweet:
39.     wordstok = tokens(t)
40.     alltok.append(wordstok)
41.     filter_huruf = []
42.     for w in wordstok:
43.         if w.lower() not in Stopwords:
44.             re.sub('\W', '', w)
45.             filter_huruf.append(w.lower())
46.
47.     filtertrain.append(filter_huruf)
48. print(filtertrain)
49. from nltk.stem import PorterStemmer
50. porter = PorterStemmer()
51. stemtrain = []
52. for stem in filtertrain:
53.     wordstem = []
54.     for word in stem:
55.         wordstem.append(porter.stem(word))
56.     stemtrain.append(wordstem)
57.
58. print(stemtrain)
59. import numpy as np
60. listvectorword = [];
61. for a in stemtrain:
62.     listvectorword += a
63.
64. kataunik = np.unique(listvectorword)
65. for u in kataunik:

```

```

66.     print(u)
67.     list_jumhuruf = [];
68.
69.     for a in stemtrain:
70.         jumlah = dict.fromkeys(kataunik, 0)
71.         for huruf in a:
72.             jumlah[huruf] += 1
73.
74.         list_jumhuruf.append(jumlah)
75.
76.     print(jumlah)
77.     def computeTF(wordDict, bagOfWords):
78.         tfDiction = {}
79.         bagOfWordsCount = len(bagOfWords)
80.         for word, count in wordDict.items():
81.             tfDiction[word] = count / float(bagOfWordsCount)
82.         return tfDiction
83.
84.
85.     listtf = []
86.     i = 0
87.     for doc in stemtrain:
88.         tf = computeTF(list_jumhuruf[i], doc)
89.
90.         listtf.append(tf)
91.         i += 1
92.     def computeIDF(documents):
93.         import math
94.         N = len(documents)
95.
96.         idfDiction = dict.fromkeys(documents[0].keys(), 0)
97.         for document in documents:
98.             for word, val in document.items():
99.                 if val > 0:
100.                    idfDiction[word] += 1
101.
102.         for word, val in idfDiction.items():
103.             idfDiction[word] = math.log(N / float(val))
104.         return idfDiction
105.
106.     idf = computeIDF(list_jumhuruf)
107.     def computeTFIDF(tfBagOfWords, idf):
108.         tfidf = {}
109.         for word, val in tfBagOfWords.items():

```

```
110.         tfidf[word] = val * idf[word]
111.     return tfidf
112.
113. list_tfidf = []
114. for tf in listtf:
115.     tfidf = computeTFIDF(tf, idf)
116.     list_tfidf.append(tfidf)
117.
118. df = pd.DataFrame(list_tfidf)
```





**2.13%** PLAGIARISM  
APPROXIMATELY

## Report #11819752

Introduction Background With the development of technology today it makes it easier for us to exchange opinions or opinions through social media. This will certainly cause a problem, namely differences of opinion or opinion towards fellow social media users. Therefore this project was created to make it easier to analyze the opinions of social media users by using sentiment analysis. Sentiment Analysis itself is the process of analyzing various data in the form of opinions or views so that conclusions are generated from various existing opinions. The problem to be resolved in this project is how the opinion of Twitter social media users on the current social distancing. This problem occurs because of the Covid-19 virus which is currently endemic, one way to avoid or reduce victims of this virus, the government is implementing social distancing to reduce or prevent the spread of the Covid-19 virus. The way to solve this problem is by using the Vector Space Model algorithm and Naive Bayes algorithm. Which is then implemented into Sentiment Analysis to get maximum results. The first step of data is taken by scraping from Twitter first. Then the data will be processed by the preprocessing method before being completed using the Vector Space Model algorithm which will produce positive or negative final results. After that, it will be evaluated by calculating accuracy, precision,

REPORT #1181975217 NOV 2020, 10:00 AM

CHECKED

AUTHOR UNIKA SOEGIJAPRANATA

PAGE 1 OF 23