

Tugas NLP 4

Rahmat Kurnianto 1301188593

Hasil Kondisi awal

Bandingkan performansi dengan base model, yaitu model dengan parameter default.

```
In [18]: 1 base_model = LogisticRegression(random_state = 8)
          2 base_model.fit(features_train, labels_train)
          3 accuracy_score(labels_test, base_model.predict(features_test))
```

Out[18]: 0.9431137724550899

```
In [19]: 1 best_classifier.fit(features_train, labels_train)
          2 accuracy_score(labels_test, best_classifier.predict(features_test))
```

Out[19]: 0.9401197604790419

- ..

1. Membuat fitur

Tanpa proses normalisation

- Variable yang di ubah

```
In [12]: 1 nrows = len(df)
          2 lemmatized_text_list = []
          3
          4 for row in range(0, nrows):
          5
          6     # Create an empty list containing lemmatized words
          7     lemmatized_list = []
          8
          9     # Save the text and its words into an object
          10    #text = df.loc[row]['Content_Parsed_4']#no1
          11    text = df.loc[row]['Content']
          12    text_words = text.split(" ")
          13
          14    # Iterate through every word to lemmatize
          15    for word in text_words:
          16        lemmatized_list.append(wordnet_lemmatizer.lemmatize(word, pos="v"))
          17
          18    # Join the list
          19    lemmatized_text = " ".join(lemmatized_list)
          20
          21    # Append to the list containing the texts
          22    lemmatized_text_list.append(lemmatized_text)
```

- Hasil akurasi

```
In [18]: 1 base_model = LogisticRegression(random_state = 8)
          2 base_model.fit(features_train, labels_train)
          3 accuracy_score(labels_test, base_model.predict(features_test))
```

Out[18]: 0.9251497005988024

```
In [19]: 1 best_classifier.fit(features_train, labels_train)
          2 accuracy_score(labels_test, best_classifier.predict(features_test))
```

Out[19]: 0.9341317365269461

Tanpa proses lemmatisation

- Variabel yang diubah

```
In [17]: 1 df['Content_Parsed_6'] = df['Content_Parsed_5']
          2 #df['Content_Parsed_6'] = df['Content_Parsed_4']
          3
          4 for stop_word in stop_words:
          5
          6     regex_stopword = r"\b" + stop_word + r"\b"
          7     df['Content_Parsed_6'] = df['Content_Parsed_6'].str.replace(regex_stopword, '')
```

- Hasil Akurasi

Bandingkan performansi dengan base model, yaitu model dengan parameter default.

```
In [18]: 1 base_model = LogisticRegression(random_state = 8)
          2 base_model.fit(features_train, labels_train)
          3 accuracy_score(labels_test, base_model.predict(features_test))
```

Out[18]: 0.9251497005988024

```
In [19]: 1 best_classifier.fit(features_train, labels_train)
          2 accuracy_score(labels_test, best_classifier.predict(features_test))
```

Out[19]: 0.9341317365269461

Tanpa menghilangkan stopwords

- Variable yang di ubah

```
In [26]: 1 list_columns = ["File_Name", "Category", "Content", "Content_Parsed_6"]
          2 df = df[list_columns]
          3
          4 df = df.rename(columns={'Content_Parsed_5': 'Content_Parsed'})
          5 #df = df.rename(columns={'Content_Parsed_6': 'Content_Parsed'})
```

- Hasil Akurasi

Bandingkan performansi dengan base model, yaitu model dengan parameter default.

```
In [18]: 1 base_model = LogisticRegression(random_state = 8)
          2 base_model.fit(features_train, labels_train)
          3 accuracy_score(labels_test, base_model.predict(features_test))
```

Out[18]: 0.9461077844311377

```
In [19]: 1 best_classifier.fit(features_train, labels_train)
          2 accuracy_score(labels_test, best_classifier.predict(features_test))
```

Out[19]: 0.9461077844311377

2. tfidf dengan nilai "max_features"

- Set min 150

```
In [33]: 1 # Parameter election
          2 ngram_range = (1,2)
          3 min_df = 10
          4 max_df = 1.
          5 max_features = 150 #no 2
```

Hasil akurasi

```
In [18]: 1 base_model = LogisticRegression(random_state = 8)
          2 base_model.fit(features_train, labels_train)
          3 accuracy_score(labels_test, base_model.predict(features_test))
```

Out[18]: 0.9041916167664671

```
In [19]: 1 best_classifier.fit(features_train, labels_train)
          2 accuracy_score(labels_test, best_classifier.predict(features_test))
```

Out[19]: 0.9041916167664671

- Set maks 450

```
1 # Parameter election
2 ngram_range = (1,2)
3 min_df = 10
4 max_df = 1.
5 max_features = 450 #no 2
```

Hasil akurasi

```
In [18]: 1 base_model = LogisticRegression(random_state = 8)
          2 base_model.fit(features_train, labels_train)
          3 accuracy_score(labels_test, base_model.predict(features_test))
```

```
Out[18]: 0.9401197604790419
```

```
In [19]: 1 best_classifier.fit(features_train, labels_train)
          2 accuracy_score(labels_test, best_classifier.predict(features_test))
```

```
Out[19]: 0.9461077844311377
```

3. Jika menggunakan Bahasa Indonesia

- Pada proses Normalisation di Bahasa Inggris diganti jadi jadi proses normalisasi Bahasa Indonesia. contoh dalam bahasa Inggris untuk menyatakan jamak pada kata noun mengunakan akhiran 's menandakan kepemilikan, sedangkan pada bahasa tidak ada
- Pada proses Lemmatisation Bahasa Inggris di ubah jadi Bahasa Indonesia. contoh pada kata tempel jadi tempelan
- stopwords karena aturan penggalan kata yang tidak relevan pada bahasa Indonesia dengan Inggris berbeda, contoh euy, anjir, dll