

LAMPIRAN

Lampiran 1 Pengambilan Data

```
1 from google_play_scraper import Sort, reviews
2
3 result, continuation_token = reviews(
4 'id.go.kpu.sirekap2024',
5 lang='id',
6 country='id',
7 sort=Sort.MOST_RELEVANT,
8 count=10000000,
9 filter_score_with=None
10 )
```

Scraping Data

```
1 import pandas as pd
2 import numpy as np
3
4 df = pd.DataFrame(np.array(result), columns=['review'])
5 df = df.join(pd.DataFrame(df.pop('review').tolist()))
6
7 df.head()
```

Dataframe Pandas

```
1 df = df[['userName', 'score', 'at', 'content',
2 'thumbsUpCount']]
3 df.sort_values(by='at', ascending=False)
4 df.head()
```

Filtering dan Sorting Kolom

Lampiran 2 Analisis Data

```
1 import pandas as pd
2
3 file_csv = 'scrapped_data.csv'
4
5 df = pd.read_csv(file_csv, sep='|')
6
7 print(df)
```

Input Data File CSV

```
1 deskripsi_statistik = df.describe()
2
3 mean_rating = df['Rating'].mean()
4 median_rating = df['Rating'].median()
5 mode_rating = df['Rating'].mode()[0]
6
7 print("Statistik Deskriptif untuk Dataframe berdasarkan
8 kolom:")
9 print(deskripsi_statistik)
10
11 print("\nStatistik Deskriptif untuk Kolom 'Rating':")
12 print(f"Mean Rating: {mean_rating}")
13 print(f"Median Rating: {median_rating}")
14 print(f"Mode Rating: {mode_rating}")
```

Statistik Deskriptif

```

1 import matplotlib.pyplot as plt
2
3 plt.figure(figsize=(8, 6))
4 ratings, counts =
5 zip(*sorted(dict(df['Rating'].value_counts()).items()))
6 plt.bar(ratings, counts, alpha=0.7, edgecolor='black',
7 align='center')
8 plt.title('Distribusi Rating')
9 plt.xlabel('Rating')
10 plt.ylabel('Frekuensi')
11 plt.xticks(ratings)
12 plt.grid(axis='y', linestyle='--', alpha=0.7)
13 plt.show()
14
15 rating_counts = df['Rating'].value_counts().sort_index()
16
17 for rating, count in rating_counts.items():
18     print(f"Rating {rating}: {count} kali")

```

Distribusi Rating

```

1 import seaborn as sns
2
3 df['Panjang Komentar'] = df['Komentar'].apply(lambda x:
4 len(x.split()))

```

Relasi Rating Dengan Panjang Komentar

```

1 plt.figure(figsize=(10, 6))
2 sns.scatterplot(x='Rating', y='Panjang Komentar', data=df,
3 color='blue', alpha=0.7)
4 plt.title('Hubungan antara Rating dan Panjang Komentar')
5 plt.xlabel('Rating')
6 plt.ylabel('Panjang Komentar')
7 plt.show()

```

Scatter Plot

```

1 plt.figure(figsize=(10, 6))
2 sns.boxplot(x='Rating', y='Panjang Komentar', data=df,
3 color='blue')
4 plt.title('Distribusi Panjang Komentar untuk Setiap Rating')
5 plt.xlabel('Rating')
6 plt.ylabel('Panjang Komentar')
7 plt.show()

```

Box Plot

```

1 info_shortest_comments =
2 df.loc[df.groupby('Rating')['Panjang
3 Komentar'].idxmin()][['Rating', 'Panjang Komentar',
4 'Komentar']]
5 info_longest_comments = df.loc[df.groupby('Rating')['Panjang
6 Komentar'].idxmax()][['Rating', 'Panjang Komentar',
7 'Komentar']]
8
9 for index, row in info_shortest_comments.iterrows():
10     print(f"Rating {row['Rating']}: Komentar Paling Pendek -
11 Index: {index}, Panjang: {row['Panjang Komentar']},
12 Komentar: {row['Komentar']}")
13
14 print("\n")
15
16 for index, row in info_longest_comments.iterrows():

```

```

10 print(f"Rating {row['Rating']}: Komentar Paling Panjang -
      Index: {index}, Panjang: {row['Panjang Komentar']},
      Komentar: {row['Komentar']}")

```

Kalimat Terpanjang dan Terpendek Setiap *Rating*

```

1  from wordcloud import WordCloud
2
3  def generate_wordcloud(text, title):
4  wordcloud = WordCloud(width=800, height=400,
5  background_color='white').generate(text)
6
7  plt.figure(figsize=(15, 10))
8  plt.imshow(wordcloud, interpolation='bilinear')
9  plt.axis('off')
10 plt.title(title)
11 plt.show()
12
13 sorted_ratings = df['Rating'].unique()[::-1]
14
15 for rating in sorted_ratings:
16 comments = ' '.join(df[df['Rating'] == rating]['Komentar'])
17 generate_wordcloud(comments, f'Word Cloud Rating {rating}')

```

Word Cloud Setiap *Rating*

Lampiran 3 Pra Proses

```

1  import pandas as pd
2
3  file_csv = 'scrapped_data.csv'
4
5  df = pd.read_csv(file_csv, sep='|')
6
7  print(df)

```

Input Data CSV

```

1  digit_count = len(str(len(df)))
2
3  df.insert(0, 'ID', df.index.map(lambda x: 'd' +
4  str(x+1).zfill(digit_count)))
5
6  print(df)

```

Memberikan *ID* Disetiap Data

```

1  import time
2
3  def lowercase(text):
4  return text.lower()
5
6  start_time = time.time()
7  df['komentar_lowercase'] = df['Komentar'].apply(lowercase)
8  end_time = time.time()
9  time_taken = end_time - start_time
10
11 print(df[['Komentar', 'komentar_lowercase']])
12 print("Waktu proses:", time taken, "detik")

```

Lower Case

```

1  import re
2

```

```

3 def remove_unnecessary_char(text):
4     text =
5         re.sub('((www\.|^s+)| (https?://[^s+)| (http?://[^s+))', '
6         ',text)
7     text = re.sub('\n',' ',text)
8     text = re.sub('\r',' ',text)
9     text = re.sub(r'\\x..',' ',text)
10    text = re.sub(' +', ' ', text)
11    text = re.sub('[^0-9a-zA-Z]+', ' ', text)
12    return text
13
14 start_time = time.time()
15 df['komentar_remove_char'] =
16     df['komentar_lowercase'].apply(remove_unnecessary_char)
17 end_time = time.time()
18 time_taken = end_time - start_time
19
20 print(df[['komentar_lowercase', 'komentar_remove_char']])
21 print("Waktu proses:", time taken, "detik")

```

Remove Unnecessary Character

```

1 alay_dict = pd.read_csv('kamus_tidak_baku.csv', sep=";",
2 encoding='latin-1', header=None)
3 alay_dict = alay_dict.rename(columns={0: 'original',
4                                     1: 'replacement'})
5
6 alay_dict_map = dict(zip(alay_dict['original'],
7 alay_dict['replacement']))

```

Input Kamus Tidak Baku

```

1 df['komentar_remove_char'] =
2     df['komentar_remove_char'].astype(str)
3 def normalize_alay(text):
4     return ' '.join([str(alay_dict_map.get(word, word)) for word
5 in text.split(' ')])
6
7 start_time = time.time()
8 df['komentar_spellchecker'] =
9     df['komentar_remove_char'].apply(normalize_alay)
10 end_time = time.time()
11 time_taken = end_time - start_time
12
13 print(df[['komentar_remove_char', 'komentar_spellchecker']])
14 print("Waktu proses:", time taken, "detik")

```

Proses Spellchecker

```

1 from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
2 factory = StemmerFactory()
3 stemmer = factory.create_stemmer()
4
5 def stemming(text):
6     return stemmer.stem(text)
7
8 start_time = time.time()
9
10 df['komentar_stemming'] =
11     df['komentar_spellchecker'].apply(stemming)
12
13 end_time = time.time()
14 time_taken = end_time - start_time

```

```

13
14 print(df[['komentar_spellchecker', 'komentar_stemming']])
15 print("Waktu proses:", time taken, "detik")

```

Stemming Sastrawi

```

1 import numpy as np
2
3 df['komentar_stemming'] = df['komentar_stemming'].replace('',
np.nan)
4
5 baris_nan = df[df['komentar_stemming'].isnull()]
6 jumlah_nan = df['komentar_stemming'].isnull().sum()
7 print(f"Jumlah NaN pada kolom 'komentar_stemming':
{jumlah_nan}")
8 print(baris_nan)
9 #print("ID dengan 'komentar_stemming' NaN:",
baris_nan['ID'].tolist())

```

Cek Data Kosong

```

1 df_cleaned = df.dropna(subset=['komentar_stemming'])
2 df_cleaned

```

Hapus Data Kosong

Lampiran 4 WordNet

```

1 import pandas as pd
2
3 file_csv = 'output_praproses.csv'
4 3
5 df = pd.read_csv(file_csv, sep='|')
6
7 print(df)

```

Input Data

```

1 from deep_translator import GoogleTranslator
2 import time
3
4 def translate_komentar(text):
5 text_translated = GoogleTranslator(source='id',
target='en').translate(text)
6 return text_translated
7
8 from tqdm import tqdm
9
10 tqdm.pandas()
11
12 start_time = time.time()
13 df['komentar_stemming'] = df['komentar_stemming'].fillna('')
14
15 df['Translated'] =
df['komentar_stemming'].progress_apply(translate_komentar)
16 end_time = time.time()
17 time_taken = end_time - start_time
18
19 print(df[['komentar_stemming', 'Translated']])
20 print("Waktu proses:", time taken, "detik")

```

Fungsi Translated

```

1 from textblob import TextBlob
2
3 def scoring_sentiment(text):

```

```

4 analysis = TextBlob(text)
5 sentiment_score = analysis.sentiment.polarity
6 # print("Sentiment polarity:", sentiment_score)
7 score_wordnet = 0
8
9 if sentiment_score > 0.5:
10     score_wordnet = 5
11 elif 0.2 <= sentiment_score <= 0.5:
12     score_wordnet = 4
13 elif -0.2 < sentiment_score < 0.2:
14     score_wordnet = 3
15 elif -0.5 <= sentiment_score <= -0.2:
16     score_wordnet = 2
17 else:
18     score_wordnet = 1
19
20 return score_wordnet, sentiment_score

```

Hitung Score

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.metrics import confusion_matrix, f1_score
5 from matplotlib.ticker import FixedLocator
6
7 cm = confusion_matrix(df['Rating'], df['Score_Wordnet'])
8
9 macro_f1 = f1_score(df['Rating'], df['Score_Wordnet'],
10 average='macro')
11 print(f"Macro F1-Score: {macro_f1}")
12
13 fig, ax = plt.subplots(figsize=(10, 8), dpi=100)
14 cax = ax.matshow(cm, cmap=plt.cm.Blues)
15 fig.colorbar(cax)
16
17 tick_marks = np.arange(len(np.unique(df['Rating'])))
18 ax.set_xticks(tick_marks)
19 ax.set_yticks(tick_marks)
20 ax.set_xticklabels(list(range(1, 6)))
21 ax.set_yticklabels(list(range(1, 6)))
22
23 plt.xlabel('Predicted')
24 plt.ylabel('Actual')
25 plt.title('Confusion Matrix WordNet')
26
27 for i in range(len(cm)):
28     for j in range(len(cm[i])):
29         ax.text(j, i, str(cm[i][j]), va='center', ha='center',
30                 color='black')
31
32 plt.show()

```

Confusion Matrix F1-Score

Lampiran 5 Klasifikasi

```

1 import pandas as pd
2
3 file_csv = 'output_praproses.csv'
4
5 df = pd.read_csv(file_csv, sep='|')

```

```
6
7 print(df)
```

Input Data

```
1 X = df['komentar_stemming']
2 y = df['Rating']
```

Menentukan Input dan Kelas

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 X.fillna('', inplace=True)
4 vectorizer = TfidfVectorizer()
5 features = vectorizer.fit_transform(X)
```

Ekstraksi Fitur TF-IDF

```
1 from sklearn.model_selection import KFold
2 from sklearn.model_selection import cross_val_predict,
  StratifiedKFold
3 from sklearn.metrics import accuracy_score, precision_score,
  recall_score, f1_score, confusion_matrix
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 def cross_validation(model, _X, _y, _cv):
8     kf = KFold(n_splits=_cv, shuffle=True)
9     fold_count = 1
10    results = {}
11    fold_predictions = []
12    overall_cm = np.zeros((len(np.unique(_y)),
13                          len(np.unique(_y))), dtype=int)
14
15    for train_index, test_index in kf.split(_X):
16        X_train, X_test = _X[train_index], _X[test_index]
17        y_train, y_test = _y[train_index], _y[test_index]
18        model.fit(X_train, y_train)
19        y_pred = model.predict(X_test)
20        fold_predictions.append(y_pred)
21        cm = confusion_matrix(y_test, y_pred)
22
23        results[f"Fold {fold_count}"] = {
24            "Confusion Matrix": cm,
25            "Precision": precision_score(y_test, y_pred,
26                                       average='macro'),
27            "Recall": recall_score(y_test, y_pred, average='macro'),
28            "F1 Score": f1_score(y_test, y_pred, average='macro')
29        }
30        overall_cm += cm
31        fold_count += 1
32
33    for key, value in results.items():
34        print(key + ":")
35        print("Confusion Matrix:")
36        print(value["Confusion Matrix"])
37        print("Precision:", value["Precision"])
38        print("Recall:", value["Recall"])
39        print("F1 Score:", value["F1 Score"])
40        print()
41    fig, ax = plt.subplots(figsize=(10, 8), dpi=100)
```

```

42 im = ax.imshow(value["Confusion Matrix"],
43 interpolation='nearest', cmap=plt.cm.Blues)
44
45 class_names = np.unique(_y)
46 ax.set(xticks=np.arange(len(class_names)),
47 yticks=np.arange(len(class_names)),
48 xticklabels=class_names, yticklabels=class_names,
49 title=f'Confusion Matrix - {key}',
50 ylabel='True label',
51 xlabel='Predicted label')
52
53 thresh = value["Confusion Matrix"].max() / 2.
54 for i in range(value["Confusion Matrix"].shape[0]):
55 for j in range(value["Confusion Matrix"].shape[1]):
56 ax.text(j, i, format(value["Confusion Matrix"][i, j], 'd'),
57 ha="center", va="center",
58 color="white" if value["Confusion Matrix"][i, j] > thresh
59 else "black")
60 plt.show()
61
62 print("Overall Confusion Matrix:")
63 print(overall_cm)
64 print()
65
66 overall_precision = np.diag(overall_cm).sum() /
67 overall_cm.sum(axis=0).sum()
68 overall_recall = np.diag(overall_cm).sum() /
69 overall_cm.sum(axis=1).sum()
70 overall_f1_score = 2 * (overall_precision * overall_recall) /
71 (overall_precision + overall_recall)
72
73 #print("Overall Evaluation:")
74 #print("Precision:", overall_precision)
75 #print("Recall:", overall_recall)
76 #print("F1 Score:", overall_f1_score)
77
78 y_pred = cross_val_predict(model, _X, _y, cv=_cv)
79
80 cm_overall = confusion_matrix(_y, y_pred)
81 #print("Confusion Matrix:")
82 #print(cm_overall)
83
84 fig, ax = plt.subplots(figsize=(10, 8), dpi=100)
85 cax = ax.matshow(cm_overall, cmap=plt.cm.Blues)
86 fig.colorbar(cax)
87
88 tick_marks = np.arange(len(np.unique(_y)))
89 ax.set_xticks(tick_marks)
90 ax.set_yticks(tick_marks)
91 ax.set_xticklabels(list(range(1, 6)))
92 ax.set_yticklabels(list(range(1, 6)))
93
94 plt.xlabel('Predicted')
95 plt.ylabel('Actual')
96 plt.title('Confusion Matrix Overall')
97
98 thresh = overall_cm.max()/2.
99
100 for i in range(cm_overall.shape[0]):

```



```

97 for j in range(cm_overall.shape[1]):
98 ax.text(j, i, format(overall_cm[i][j], 'd'), va='center',
    ha='center', color='black' if overall_cm[i,j] > thresh else
    "black")
99
100 plt.show()
101
102 precision_macro = precision_score(y, y_pred,
    average='macro')
103 print(f"\nPrecision Macro: {precision_macro}")
104
105 recall_macro = recall_score(y, y_pred, average='macro')
106 print(f"Recall Macro: {recall_macro}")
107
108 f1_score_macro = f1_score(y, y_pred, average='macro')
109 print(f"F1-Score Macro: {f1_score_macro}")
110
111
112 return results, fold_predictions

```

Cross Validation 10 Fold

```

1
2 kfold = KFold(n_splits=10, shuffle=False)
3 for fold, (train_index, test_index) in
    enumerate(kfold.split(X, y), 1):
4 if fold == 10:
5
6 X_train, X_test = X[train_index], X[test_index]
7 y_train, y_test = y[train_index], y[test_index]
8
9 print(f'Fold {fold}:')
10 print(f' - Train data: {len(X_train)} samples')
11 print(f' - Test data: {len(X_test)} samples')
12 print(f' - Train Index: {train_index}')
13 print(f' - Test Index: {test_index}')

```

Menampilkan Isi Data Setiap Pembagian Fold

```

1 from sklearn.neighbors import KNeighborsClassifier
2 # from sklearn.neighbors import NearestNeighbors
3 knn_model = KNeighborsClassifier()

```

Klasifikasi K-NN

```

1 import time
2
3 start_time = time.time()
4 cv_results, fold_predictions = cross_validation(knn_model,
    features, y, _cv=10)
5
6 end_time = time.time()
7 time_taken = end_time - start_time
8 print("Waktu proses:", time taken, "detik")


```

Menjalankanss Cross Validation

Lampiran 6 Kartu Kendali Bimbingan

KARTU KENDALI BIMBINGAN LAPORAN KARYA ILMIAH

Nama : Emyzar Hafliida Tanjung
 NIM : 2011102441240
 Nama Dosen Pembimbing : Naufal Azmi Verdikha, S.Kom., M.Eng.
 Judul Penelitian : Perbandingan Analisis *Wordnet* dan *K-Nearest Neighbor* Pada Ulasan Aplikasi Sirekap 2024.

No	Tanggal	Uraian Pembimbingan	Paraf Dosen
1	27 Januari 2024	penentuan algoritma yang akan digunakan dan pembuatan judul besar.	
2	29 Januari 2024	arahan proses codingan awal dan pembuatan kamus tidak baku	
3	31 Januari 2024	1. koreksi pembuatan kamus tidak baku 2. penjelasan codingan dan mencari alasan mengapa harus melakukan klasifikasi sentimen.	
4	1 Februari 2024	arahan point penting yang harus ditulis diproposal.	
5	18 Februari 2024	1. koreksi penulisan latar belakang 2. pembuatan tahap alur penelitian (metodologi)	
6	19 Februari 2024	koreksi bab 1 dan bab 2 proposal	
7	21 Februari 2024	bimbingan koreksi revisi proposal	
8	11 Maret 2024	arahan membuat kerangka codingan	
9	25 Juni 2024	koreksi revisi proposal	
10	13 Juni 2024	revisi bab 3	
11	14 Juni 2024	1. koreksi penyesuaian penulisan bab 3 2. koreksi revisian bab 3	
12	22 Juni 2024	koreksi revisian bab 3	

Dosen Pembimbing


 Naufal Azmi Verdikha, S.Kom., M.Eng.
 NIDN.1114048801

Mengetahui

Ketua Program Studi



 Emyzar Hafliida Tanjung, S.Kom., M.TI
 NIDN-1118019203

SKRIPSI EMYZAR HAFILDA TANJUNG

by Teknik Informatika Universitas Muhammadiyah Kalimantan Timur



Submission date: 25-Jul-2024 02:23PM (UTC+0800)

Submission ID: 2422166292

File name: SKRIPSI_EMYZAR_HAFILDA_TANJUNG.docx (1.95M)

Word count: 5788

Character count: 36142

SKRIPSI EMYZAR HAFILDA TANJUNG



ORIGINALITY REPORT

18%

SIMILARITY INDEX

17%

INTERNET SOURCES

7%

PUBLICATIONS

6%

STUDENT PAPERS

PRIMARY SOURCES

1	repository.ub.ac.id Internet Source	2%
2	e-journal.hamzanwadi.ac.id Internet Source	1%
3	media.neliti.com Internet Source	1%
4	core.ac.uk Internet Source	1%
5	docplayer.info Internet Source	1%
6	dspace.umkt.ac.id Internet Source	1%
7	scholar.unand.ac.id Internet Source	1%
8	digilib.polban.ac.id Internet Source	1%
9	repo.iain-tulungagung.ac.id Internet Source	<1%

Lampiran 8 Surat Ijin Penelitian



UMKT
Program Studi
Teknik Informatika
Fakultas Sains dan Teknologi

Telp. 0541-748511 Fax. 0541-766832

Website <http://informatika.umkt.ac.id>

email: informatika@umkt.ac.id



بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

Nomor : 055-001/KET/FST.1/A/2024

Lampiran : -

Perihal : **Keterangan Pengambilan Data Sekunder**

Assalamu'alaikum Warrahmatullahi Wabarrakatuh

Puji Syukur kepada Allah Subhanahu wa ta'ala yang senantiasa melimpahkan Rahmat-Nya kepada kita sekalian. Amin.

Dengan surat ini, kami menerangkan bahwa mahasiswa berikut:

No	Nama	NIM
1	Emyzar Hafliida Tanjung	2011102441240
2	Lilis Sagita	2011102441198
3	Sri Ramadani	2011102441177
4	Muhammad Fariz Ijlal Rafi	2011102441124
5	Nurlita	2011102441070

Melakukan penelitian dengan pengambilan data sekunder di Google Playstore data yang diambil yaitu data ulasan aplikasi "Sirekap" dari rating 1-5.

Demikian hal ini disampaikan, atas kerjasamanya kami ucapkan terima kasih.

Wassalamu'alaikum Warrahmatullahi Wabarrakatuh

Samarinda, 18 Dzulhijjah 1445 H
25 Juni 2024 M



Arbansyah, S.Kom., M.TI

NIDN. 1118019203

Kampus 1 : Jl. Ir. H. Juanda, No.15, Samarinda
Kampus 2 : Jl. Pelita, Pesona Mahakam, Samarinda

RIWAYAT HIDUP



Penulis bernama lengkap Emyzar Hafliida Tanjung dilahirkan di Kisaran 23 Agustus 2001, dan merupakan anak ketiga dari 4 bersaudara dari pasangan Ghazali Tanjung, SE dan Paridah. Mengawali pendidikan formal di Pendidikan Sekolah Dasar di MIN Tanah Paser (2007-2013) dan melanjutkan Pendidikan Sekolah Menengah Pertama di Mts Bina Islam Tanah Paser (2013-2016). Penulis melanjutkan jenjang pendidikan formal Sekolah Menengah Atas di SMKN 1 Tanah Paser (2016-2019). Penulis masuk di Fakultas Sains dan Teknologi Universitas Muhammadiyah Kalimantan Timur pada tahun 2020.

Untuk menyelesaikan studi di Fakultas Sains dan Teknologi Jurusan Teknik Informatika UMKT, penulis melakukan penelitian dengan judul “**Perbandingan Analisis Wordnet Dan K-Nearest Neighbor Pada Ulasan Aplikasi Sirekap 2024**” sebagai salah satu syarat untuk memperoleh gelar Sarjana Komputer.