

Lampiran 1. Riwayat Hidup

RIWAYAT HIDUP



Ida Ayu kadek Bintang Wijayanti lahir di Singaraja pada 5 Maret 2002. Penulis lahir dari pasangan suami istri, Bapak Ida Bagus Putu Adi Wijaya dan Ibu Luh Budi Wardani. Penulis berstatus Warga Negara Indonesia (WNI) dan beragama Hindu. Alamat tinggal penulis saat ini berada di Jalan Udayana No. 19 Singaraja, Kelurahan Kaliuntu, Kecamatan Buleleng, Kabupaten Buleleng, Provinsi Bali.

Penulis menyelesaikan pendidikan dasar di SD Negeri 3 Banjar Jawa dan lulus pada tahun 2014. Kemudian penulis melanjutkan di SMP Lab Singaraja dan lulus pada tahun 2017. Setelah lulus dari jenjang SMP, kemudian penulis melanjutkan ke jenjang SMA di SMAN 4 Singaraja dan lulus pada tahun 2020. Setelah penulis lulus dari jenjang SMA, penulis melanjutkan pendidikan ke Perguruan Tinggi di Universitas Pendidikan Ganesha dengan mengambil Program Studi (S1) Sistem Informasi, Jurusan Teknik Informatika, Fakultas Teknik dan Kejuruan.



Lampiran 2. Sertifikat Pendidik





Lampiran 3. Foto Bersama Pakar (Guru Bahasa Indonesia)





Lampiran 4. Surat Validasi


**KEMENTERIAN PENDIDIKAN, KEBUDAYAAN, RISET DAN TEKNOLOGI
UNIVERSITAS PENDIDIKAN GANESHA
FAKULTAS TEKNIK DAN KEJURUAN**
 Jalan Udayana, No. 11, Singaraja Bali Kode Pos 81116
 Telepon (0362) 21541 Fax. (0362) 27561
 Laman: <https://undiksha.ac.id>

**SURAT KETERANGAN VALIDASI
LABELING DATASET**

Yang bertanda tangan di bawah ini.

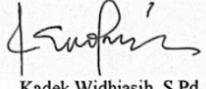
Nama	:	Kadek Widhiasih, S.Pd.
NIP	:	1966123119907022006

Menerangkan bahwa Mahasiswa Universitas Pendidikan Ganesha di bawah ini:

Nama	:	Ida Ayu Kadek Bintang Wijayanti
NIM	:	2115091083
Prodi/Jurusan	:	S1 Sistem Informasi

Memang benar bahwa dataset yang sudah dilabeli telah divalidasi pada tanggal 12 Juni 2025. Demikian surat keterangan ini dibuat dengan sebenarnya untuk dapat digunakan sebagaimana mestinya.

Singaraja, 12 Juni 2025
Ahli Pakar,


Kadek Widhiasih, S.Pd.
 NIP. 1966123119907022006



Catatan :

- UU ITE No. 11 Tahun 2008 Pasal 5 ayat 1 "Informasi Elektronik dan/atau Dokumen Elektronik dan/atau hasil cetaknya merupakan alat bukti hukum yang sah"
- Dokumen ini tertanda ditandatangani secara elektronik menggunakan sertifikat elektronik yang diterbitkan BsrE
- Surat ini dapat dibuktikan keasliannya dengan menggunakan *qr code* yang telah tersedia



**KEMENTERIAN PENDIDIKAN, KEBUDAYAAN, RISET DAN TEKNOLOGI
UNIVERSITAS PENDIDIKAN GANESHA**

FAKULTAS TEKNIK DAN KEJURUAN

Jalan Udayana, No. 11, Singaraja Bali Kode Pos 81116

Telepon (0362) 21541 Fax. (0362) 27561

Laman: <https://undiksha.ac.id>

**SURAT KETERANGAN VALIDASI
LABELING DATASET**

Yang bertanda tangan di bawah ini.

Nama : Ni Nengah Sudiantari, S.Pd.

NIP : 198207162009022001

Menerangkan bahwa Mahasiswa Universitas Pendidikan Ganesha di bawah ini:

Nama : Ida Ayu Kadek Bintang Wijayanti

NIM : 2115091083

Prodi/Jurusan : S1 Sistem Informasi

Memang benar bahwa dataset yang sudah dilabeli telah divalidasi pada tanggal 12 Juni 2025. Demikian surat keterangan ini dibuat dengan sebenarnya untuk dapat digunakan sebagaimana mestinya.

Singaraja, 12 Juni 2025
Ahli Pakar,

Ni Nengah Sudiantari, S.Pd.
NIP. 198207162009022001



Catatan :

- UU ITE No. 11 Tahun 2008 Pasal 5 ayat 1 "Informasi Elektronik dan/atau Dokumen Elektronik dan/atau hasil cetaknya merupakan alat bukti hukum yang sah"
- Dokumen ini tertanda ditandatangani secara elektronik menggunakan sertifikat elektronik yang diterbitkan BsrE
- Surat ini dapat dibuktikan keasliannya dengan menggunakan *qr code* yang telah tersedia



**KEMENTERIAN PENDIDIKAN, KEBUDAYAAN, RISET DAN TEKNOLOGI
UNIVERSITAS PENDIDIKAN GANESHA
FAKULTAS TEKNIK DAN KEJURUAN**
Jalan Udayana, No. 11, Singaraja Bali Kode Pos 81116
Telepon (0362) 21541 Fax. (0362) 27561
Laman: <https://undiksha.ac.id>

**SURAT KETERANGAN VALIDASI
LABELING DATASET**

Yang bertanda tangan di bawah ini.

Nama : Ida Ayu Asritia Utami, S.Pd.
NIP : 199309122023212045

Menerangkan bahwa Mahasiswa Universitas Pendidikan Ganesha di bawah ini:

Nama : Ida Ayu Kadek Bintang Wijayanti
NIM : 2115091083
Prodi/Jurusan : S1 Sistem Informasi

Memang benar bahwa dataset yang sudah dilabeli telah divalidasi pada tanggal 12 Juni 2025. Demikian surat keterangan ini dibuat dengan sebenarnya untuk dapat digunakan sebagaimana mestinya.

Singaraja, 12 Juni 2025
Ahli Pakar,

Ida Ayu Asritia Utami, S.Pd.
NIP. 199309122023212045



Catatan :

- UU ITE No. 11 Tahun 2008 Pasal 5 ayat 1 "Informasi Elektronik dan/atau Dokumen Elektronik dan/atau hasil cetaknya merupakan alat bukti hukum yang sah"
- Dokumen ini tertanda ditandatangani secara elektronik menggunakan sertifikat elektronik yang diterbitkan BsrE
- Surat ini dapat dibuktikan keasliananya dengan menggunakan *qr code* yang telah tersedia

Lampiran 5. Kode Program Proses *Crawling* Data

```

● ● ●

1 # 1. Set Twitter Auth Token
2 twitter_auth_token = 'd9fbcef32ccb379c459749ceec558e071145dbf' # 1. Change this auth token with your own
3
4 # 2. Install pandas package
5 !pip install pandas # 2. Install pandas for data processing
6
7 # 3. Install Node.js (required for tweet-harvest)
8 !sudo apt-get update # 3. Update package list
9 !sudo apt-get install -y ca-certificates curl gnupg # 4. Install certificates, curl, gnupg
10 !sudo mkdir -p /etc/apt/keyrings # 5. Create keyrings directory
11 !curl -fsSL https://deb.nodesource.com/gpgkey/nodesource-repo.gpg.key | \
12     sudo gpg --dearmor -o /etc/apt/keyrings/nodesource.gpg # 6. Download and save Node.js GPG key
13 !NODE_MAJOR=20 && echo \
14 "deb [signed-by=/etc/apt/keyrings/nodesource.gpg] https://deb.nodesource.com/node_$NODE_MAJOR.x nodistro main" | \
15     sudo tee /etc/apt/sources.list.d/nodesource.list # 7. Add Node.js repository
16 !sudo apt-get update # 8. Update package list again
17 !sudo apt-get install nodejs -y # 9. Install Node.js
18 !node -v # 10. Check Node.js version
19
20 # 4. Crawl Data from Twitter using tweet-harvest
21 filename = 'ometv.csv' # 11. Output filename
22 search_keyword = 'ometv lang:id since:2020-01-01 until:2025-4-30' # 12. Search keyword with filters
23 limit = 6000 # 13. Limit data to crawl
24
25 # 5. Run tweet-harvest command
26 !npx -y tweet-harvest@2.6.1 \
27 -o "{filename}" \ # 14. Output file path
28 -s "{search_keyword}" \ # 15. Search keyword for crawling
29 --tab "LATEST" -l {limit} --token {twitter_auth_token} # 16. Additional options: latest tab, limit, and token
30
31 # 6. Import pandas for reading data
32 import pandas as pd # 17. Import pandas library
33
34 # 7. Set file path for reading the CSV
35 filename = '/content/tweets-data/ome.csv' # 18. File path (adjust if different in your drive)
36
37 # 8. Assign file path to variable
38 file_path = filename # 19. Assign path
39
40 # 9. Read CSV into pandas DataFrame
41 df = pd.read_csv(file_path, delimiter=",", encoding='utf-8', on_bad_lines='skip') # 20. Read CSV, skip bad lines
42
43 # 10. Display the DataFrame
44 display(df) # 21. Display crawled data for checking

```

Lampiran 6. Kode Program Proses *Cleaning Data*

```

● ● ●

1 # 1. Install library yang dibutuhkan
2 !pip install emoji # 1. Install library emoji
3 !pip install langdetect # 2. Install library langdetect
4
5 # 2. Import library
6 import pandas as pd # 3. Import pandas
7 import re # 4. Import regex
8 import emoji # 5. Import emoji
9 from langdetect import detect, DetectorFactory # 6. Import language detection
10
11 DetectorFactory.seed = 0 # 7. Agar deteksi bahasa konsisten
12
13 # 3. BACA DATA
14 df = pd.read_csv('data_berlabel_awal.csv', encoding='latin1') # 8. Baca file data
15
16 print("==> DATA AWAL ==") # 9. Print judul data awal
17 print(df.head()) # 10. Tampilkan 5 data pertama
18 print(df.info()) # 11. Informasi data
19
20 # 4. HAPUS DATA MISSING & DUPLIKAT
21 df = df.dropna() # 12. Hapus data missing/null
22 df = df.drop_duplicates() # 13. Hapus data duplikat
23
24 # 5. HAPUS TWEET MENGANDUNG LINK
25 df = df[~df['tweet'].str.contains(r'http\S+|www\.\S+', na=False)] # 14. Hapus tweet mengandung link
26
27 # 6. HAPUS TWEET MENGANDUNG KATA IKLAN
28 pattern_iklan=r'RESULT|SHIO|JACKPOT|Para Pemenang|SYAIR|ANGKA' # 15. Pattern kata iklan
29 df = df[~df['tweet'].str.contains(pattern_iklan, case=False, na=False)] # 16. Hapus tweet iklan
30
31 # 7. HAPUS TWEET NON-BAHASA INDONESIA
32 def is_indonesian(text): # 17. Fungsi cek bahasa Indonesia
33     try:
34         return detect(text) == 'id'
35     except:
36         return False
37
38 df = df[df['tweet'].apply(is_indonesian)] # 18. Filter hanya bahasa Indonesia
39
40 # 8. HAPUS TWEET DENGAN SIMBOL BERULANG (!!! , ???, ***)
41 def has_repeated_symbols(text): # 19. Fungsi cek simbol berulang
42     return bool(re.search(r'\?!\*\*/\.,\.{3,}', text))
43
44 df = df[~df['tweet'].apply(has_repeated_symbols)] # 20. Hapus tweet dengan simbol berulang
45
46 # 9. CLEANING TEXT (MENTION, HASHTAG, EMOJI, ANGKA, TANDA BACA, SPASI GANDA)
47 def bersihkan_teks(text): # 21. Fungsi cleaning text
48     text = re.sub(r'@\w+', '', text) # 22. Hapus mention
49     text = re.sub(r'http\S+|www\.\S+', '', text) # 23. Hapus URL (jaga-jaga)
50     text = re.sub(r'#\w+', '', text) # 24. Hapus hashtag
51     text = emoji.replace_emoji(text, replace='') # 25. Hapus emoji
52     text = re.sub(r'\d+', '', text) # 26. Hapus angka
53     text = re.sub(r'[\^@\s]', '', text) # 27. Hapus tanda baca
54     text = re.sub(r'\s+', ' ', text) # 28. Spasi ganda jadi satu
55     return text.strip() # 29. Hapus spasi awal/akhir
56
57 df['cleaning_data'] = df['tweet'].apply(bersihkan_teks) # 30. Terapkan cleaning ke kolom tweet
58
59 # 10. TAMPILKAN HASIL
60 print("\n==> DATA BERSIH ==") # 31. Judul data bersih
61 display(df[['tweet', 'cleaning_data']].head(10)) # 32. Tampilkan 10 data hasil cleaning
62
63 # 11. SIMPAN HASIL
64 df.to_csv('final_cleaning_data.csv', index=False) # 33. Simpan data hasil cleaning ke file
65 print("\n✓ Data cleaning selesai dan disimpan ke 'final_cleaning_data.csv'") # 34. Print selesai

```

Lampiran 7. Kode Program Proses Case Folding

```
● ● ●

1 # 1. Import library pandas
2 import pandas as pd
3
4 #PROSES CASE FOLDING#
5
6 # 3. Membaca dataset final_cleaning_data.csv
7 df = pd.read_csv('final_cleaning_data.csv')
8
9 # 4. Fungsi case folding: ubah huruf menjadi huruf kecil
10 def case_folding(text):
11     return str(text).lower()
12
13 # 5. Terapkan fungsi case folding ke kolom cleaning_data
14 df['case_folding'] = df['cleaning_data'].apply(case_folding)
15
16 # 6. Simpan hasilnya ke file baru (opsional)
17 df.to_csv('data_casefolding.csv', index=False)
18
19 # 7. Tampilkan hasil sebelum dan sesudah case folding (10 contoh)
20 from IPython.display import display
21 print("\n==> TABEL HASIL CASEFOLDING (SEBELUM & SESUDAH) ==>") display(df[['cleaning_data', 'case_folding']].head(10))
```



Lampiran 8. Kode Program Proses *Stopword Removal*

```

● ● ●

1 # 1. Import library yang dibutuhkan
2 import pandas as pd
3 import re
4 import nltk
5 from nltk.corpus import stopwords
6 from IPython.display import display
7
8 # 2. Download stopwords bahasa Indonesia dari NLTK
9 nltk.download('stopwords')
10
11 # 3. Baca dataset hasil case folding
12 df = pd.read_csv('data_casetolding.csv')
13
14 # 4. Ambil stopwords bahasa Indonesia dari NLTK dan tambahkan kata tambahan
15 additional_stopwords = {
16     'aku', 'gue', 'gwa', 'eh', 'ig', 'anjir', 'anjirrrr', 'sih', 'njir', 'bjir',
17     'wkwk', 'wkkwkwk', 'wkkwkwk', 'nya', 'yang', 'yg', 'yuk', 'yukk', 'yokk', 'yok', 'ya',
18     'trs', 'tp', 'gw', 'pas', 'ig', 'tau', 'gitu', 'gtu', 'gt', 'lu', 'lo', 'pa', 'gak', 'aja',
19     'kek', 'ni', 'si', 'biar', 'tu', 'tuh', 'dah', 'najis', 'plis', 'pd', 'dri', 'dn', 'on', 'sm',
20     'bgt', 'bgnt', 'bsa', 'udh', 'kalo', 'klo', 'mw', 'mo', 'prnh', 'sih', 'aduh', 'pada', 'juga', 'karena',
21     'atau', 'tapi', 'jadi', 'aku', 'gue', 'gw', 'gwa', 'lu', 'org', 'hufit', 'trus', 'trs', 'jadi'
22 }
23 stop_words = set(stopwords.words('indonesian')).union(additional_stopwords)
24
25 # 5. Fungsi untuk hapus stopword tanpa tokenisasi eksplisit
26 def remove_stopwords(text):
27     text = str(text).lower()
28     text = re.sub(r'^\w\s', '', text)
29     words = text.split()
30     filtered = [word for word in words if word not in stop_words]
31     return ' '.join(filtered)
32
33 # 6. Terapkan fungsi ke kolom 'case_folding'
34 df['stopword_removal'] = df['case_folding'].apply(remove_stopwords)
35
36 # 7. Simpan hasil ke file baru
37 df.to_csv('data_stopword_removal.csv', index=False)
38
39 # 8. Tampilkan hasil sebelum dan sesudah stopword removal
40 print("✅ Stopword removal selesai!")
41 display(df[['case_folding', 'stopword_removal']].head(10))

```



Lampiran 9. Kode Program Proses *Tokenizing*

```
1 # 1. Import library yang dibutuhkan
2 import re
3 import pandas as pd
4
5 # 2. Baca dataset hasil stopword removal
6 df = pd.read_csv('data_stopword_removal.csv')
7
8 # 3. Isi kolom stopword_removal dengan string kosong jika ada NaN
9 df['stopword_removal'] = df['stopword_removal'].fillna('')
10
11 # 4. Fungsi tokenize (contoh sederhana)
12 def tokenize(text):
13     if not isinstance(text, str):
14         return []
15     return text.split()
16
17 # 5. Terapkan fungsi tokenize ke kolom 'stopword_removal'
18 df['tokens'] = df['stopword_removal'].apply(tokenize)
19
20 # 6. Simpan hasil ke file baru
21 df.to_csv('hasil_tokenizing.csv', index=False)
22
23 # 7. Tampilkan hasil sebelum dan sesudah tokenize
24 print("✅ Tokenizing selesai!")
25 display(df[['stopword_removal', 'tokens']].head(10))
```



Lampiran 10. Kode Program Proses *Normalization*

```
1 # 1. Import library yang dibutuhkan
2 import pandas as pd
3 import json
4 import ast
5
6 # 2. Upload slang-indo.json dari lokal
7 from google.colab import files
8 uploaded = files.upload()
9
10 # 3. Baca dataset hasil tokenizing
11 df = pd.read_csv('hasil_tokenizing.csv')
12 df['tokens'] = df['tokens'].apply(ast.literal_eval)
13
14 # 4. Baca kamus slang dari file JSON
15 with open('slang-indo.json', 'r', encoding='utf-8') as f:
16     slang_dict = json.load(f)
17
18 # 5. Fungsi normalisasi slang
19 def normalize_tokens(tokens):
20     return [slang_dict.get(token.lower(), token) for token in tokens]
21
22 # 6. Terapkan normalisasi ke kolom 'tokens'
23 df['normalized'] = df['tokens'].apply(normalize_tokens)
24
25 # 7. Simpan hasil ke file baru
26 df.to_csv('hasil_normalization.csv', index=False)
27
28 # 8. Tampilkan hasil sebelum dan sesudah normalisasi
29 print("✅ Normalized selesai!")
30 display(df[['tokens', 'normalized']].head(10))
```



Lampiran 11. Kode Program Proses *Stemming*

```
● ● ●  
1 # 1. Install dan import library yang dibutuhkan  
2 !pip install PySastrawi  
3 import pandas as pd  
4 import ast  
5 from Sastrawi.Stemmer.StemmerFactory import StemmerFactory  
6  
7 # 2. Buat stemmer  
8 factory = StemmerFactory()  
9 stemmer = factory.create_stemmer()  
10  
11 # 3. Baca file hasil normalisasi  
12 df = pd.read_csv('hasil_normalization.csv')  
13 df['normalized'] = df['normalized'].apply(ast.literal_eval)  
14  
15 # 4. Fungsi stemming untuk setiap token  
16 def stem_tokens(tokens):  
17     return [stemmer.stem(token) for token in tokens]  
18  
19 # 5. Terapkan stemming ke kolom 'normalized'  
20 df['stemmed'] = df['normalized'].apply(stem_tokens)  
21  
22 # 6. Simpan hasil ke file baru  
23 df.to_csv('hasil_stemming.csv', index=False)  
24  
25 # 7. Tampilkan hasil sebelum dan sesudah stemming  
26 print("✓ Stemming selesai!")  
27 display(df[['normalized', 'stemmed']].head(10))
```



Lampiran 12. Kode Program Proses TF-IDF

```

● ● ●

1 # 1. Import library yang dibutuhkan
2 import pandas as pd
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.preprocessing import LabelEncoder
5
6 # 2. Baca dataset yang sudah pre-processing (kolom: 'stemmed', 'label')
7 df = pd.read_csv('hasil_stemming.csv')
8
9 # 3. TF-IDF dari kolom 'stemmed'
10 corpus = df['stemmed'].astype(str).tolist()
11 vectorizer = TfidfVectorizer()
12 tfidf_matrix = vectorizer.fit_transform(corpus)
13
14 # 4. Ubah hasil TF-IDF menjadi DataFrame
15 tfidf_df = pd.DataFrame(tfidf_matrix.toarray(), columns=vectorizer.get_feature_names_out())
16
17 # 5. Label Encoding dari kolom 'label'
18 le = LabelEncoder()
19 df['label_encoded'] = le.fit_transform(df['label'])
20
21 # 6. Print distribusi label asli
22 print("\nDistribusi label teks:")
23 print(df['label'].value_counts())
24
25 # 7. Gabungkan TF-IDF (fitur) + label_encoded (target)
26 final_df = tfidf_df.copy()
27 final_df['label'] = df['label_encoded']
28
29 # 8. Tampilkan 5 data pertama dari dataset final
30 print(final_df.head())
31
32 # 9. Simpan ke CSV jika perlu
33 final_df.to_csv('dataset_tfidf_label_baru.csv', index=False)
34
35 # 10. Menampilkan 10 kata dengan rata-rata TF-IDF tertinggi
36 mean_tfidf = tfidf_df.mean(axis=0)
37 top_tfidf = mean_tfidf.sort_values(ascending=False)
38 print("\n10 kata dengan rata-rata TF-IDF tertinggi:")
39 print(top_tfidf.head(10))
40
41 # 11. (Opsional) Lihat mapping label aslinya
42 label_mapping = dict(zip(le.classes_, le.transform(le.classes_)))
43 print("Mapping label ke angka:", label_mapping)

```

Lampiran 13. Kode Program *Naive Bayes Fold-5 + Confusion Matrix*

```

● ● ●

1 # 1. Import library yang dibutuhkan
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import StratifiedKFold
5 from sklearn.naive_bayes import MultinomialNB
6 from sklearn.metrics import classification_report, accuracy_score, confusion_matrix, ConfusionMatrixDisplay
7 import matplotlib.pyplot as plt
8
9 === NAIVE BAYES TANPA SMOTE ===
10
11 # 2. Load data
12 df = pd.read_csv('dataset_tfidf_label_baru.csv')
13 X = df.drop(columns=['label'])
14 y = df['label']
15
16 # 3. Tampilkan distribusi label di awal
17 print("== Distribusi Label Awal ==")
18 print(y.value_counts())
19
20 # 4. Setup Stratified K-Fold
21 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
22
23 # 5. List untuk menyimpan metrik setiap fold
24 macro_precisions, macro_recalls, macro_f1s, accuracies = [], [], [], []
25
26 # 6. Simpan semua label dan prediksi dari semua fold untuk report gabungan
27 all_y_true = []
28 all_y_pred = []
29
30 print("== Naive Bayes Tanpa SMOTE + Confusion Matrix ==")
31
32 # 7. Loop K-Fold cross validation
33 for fold, (train_idx, test_idx) in enumerate(skf.split(X, y), 1):
34     X_train, y_train = X.iloc[train_idx], y.iloc[train_idx]
35     X_test, y_test = X.iloc[test_idx], y.iloc[test_idx]
36     model = MultinomialNB()
37     model.fit(X_train, y_train)
38     y_pred = model.predict(X_test)
39
40     # 8. Hitung metrik per fold (macro avg & akurasi)
41     report = classification_report(y_test, y_pred, output_dict=True)
42     acc = accuracy_score(y_test, y_pred)
43
44     macro_precisions.append(report['macro avg']['precision'])
45     macro_recalls.append(report['macro avg']['recall'])
46     macro_f1s.append(report['macro avg']['f1-score'])
47     accuracies.append(acc)
48
49     # 9. Simpan label dan prediksi untuk report gabungan
50     all_y_true.extend(y_test)
51     all_y_pred.extend(y_pred)
52
53     # 10. === Classification Report Gabungan ===
54     print("\n== Classification Report Gabungan (Tanpa SMOTE) ==")
55     print(classification_report(all_y_true, all_y_pred, digits=2))
56
57     # 11. === Rata-rata metrik ===
58     print("\n== Rata-rata Hasil Evaluasi (Tanpa SMOTE) ===")
59     print(f"Precision: {np.mean(macro_precisions):.4f}")
60     print(f"Recall: {np.mean(macro_recalls):.4f}")
61     print(f"F1-Score: {np.mean(macro_f1s):.4f}")
62     print(f"Accuracy: {np.mean(accuracies):.4f}")
63
64     # 12. === Confusion Matrix Gabungan ===
65     cm = confusion_matrix(all_y_true, all_y_pred)
66     disp = ConfusionMatrixDisplay(confusion_matrix=cm)
67     disp.plot(cmap='Blues')
68     plt.title("Confusion Matrix (Tanpa SMOTE - Gabungan Semua Fold)") plt.show()

```

Lampiran 14. Kode Program *Naive Bayes Fold-10 + Confusion Matrix*

```

● ● ●

1 # 1. Import library yang dibutuhkan
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import StratifiedKFold
5 from sklearn.naive_bayes import MultinomialNB
6 from sklearn.metrics import classification_report, accuracy_score, confusion_matrix, ConfusionMatrixDisplay
7 import matplotlib.pyplot as plt
8
9 # === NAIVE BAYES TANPA SMOTE ===
10
11 # 2. Load data
12 df = pd.read_csv('dataset_tfidf_label_baru.csv')
13 X = df.drop(columns=['label'])
14 y = df['label']
15
16 # 3. Tampilkan distribusi label di awal
17 print("== Distribusi Label Awal ==")
18 print(y.value_counts())
19
20 # 4. Setup Stratified K-Fold
21 skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
22
23 # 5. List untuk menyimpan metrik setiap fold
24 macro_precisions, macro_recalls, macro_f1s, accuracies = [], [], [], []
25
26 # 6. Simpan semua label dan prediksi dari semua fold untuk report gabungan
27 all_y_true = []
28 all_y_pred = []
29 print("== Naive Bayes Tanpa SMOTE + Confusion Matrix ==")
30
31 # 7. Loop K-Fold cross validation
32 for fold, (train_idx, test_idx) in enumerate(skf.split(X, y), 1):
33     X_train, y_train = X.iloc[train_idx], y.iloc[train_idx]
34     X_test, y_test = X.iloc[test_idx], y.iloc[test_idx]
35     model = MultinomialNB()
36     model.fit(X_train, y_train)
37     y_pred = model.predict(X_test)
38
39 # 8. Hitung metrik per fold (macro avg & akurasi)
40     report = classification_report(y_test, y_pred, output_dict=True)
41     acc = accuracy_score(y_test, y_pred)
42
43     macro_precisions.append(report['macro avg']['precision'])
44     macro_recalls.append(report['macro avg']['recall'])
45     macro_f1s.append(report['macro avg']['f1-score'])
46     accuracies.append(acc)
47
48 # 9. Simpan label dan prediksi untuk report gabungan
49     all_y_true.extend(y_test)
50     all_y_pred.extend(y_pred)
51
52 # 10. === Classification Report Gabungan ===
53 print("\n== Classification Report Gabungan (Tanpa SMOTE) ==")
54 print(classification_report(all_y_true, all_y_pred, digits=2))
55
56 # 11. === Rata-rata metrik ===
57 print("\n== Rata-rata Hasil Evaluasi (Tanpa SMOTE) ==")
58 print(f"Precision: {np.mean(macro_precisions):.4f}")
59 print(f"Recall: {np.mean(macro_recalls):.4f}")
60 print(f"F1-Score: {np.mean(macro_f1s):.4f}")
61 print(f"Accuracy: {np.mean(accuracies):.4f}")
62
63 # 12. === Confusion Matrix Gabungan ===
64 cm = confusion_matrix(all_y_true, all_y_pred)
65 disp = ConfusionMatrixDisplay(confusion_matrix=cm)
66 disp.plot(cmap='Blues') # 40. Plot confusion matrix dengan cmap blue
67 plt.title("Confusion Matrix (Tanpa SMOTE - Gabungan Semua Fold)") plt.show()

```

Lampiran 15. Kode Program *Naive Bayes* SMOTE Fold-5 + *Confusion Matrix*

```

● ● ●

1 # 1. Import library yang dibutuhkan
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import StratifiedKFold
5 from sklearn.naive_bayes import MultinomialNB
6 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay, classification_report
7 import matplotlib.pyplot as plt
8 from imblearn.over_sampling import SMOTE
9
10 # 2. Load data
11 df = pd.read_csv('dataset_tfidf_label.csv')
12 X = df.drop(columns=['label'])
13 y = df['label']
14
15 # 3. Bersihkan target NaN
16 not_nan_mask = y.notna()
17 X = X[not_nan_mask]
18 y = y[not_nan_mask]
19
20 # 4. Target per label SESUDAH SMOTE
21 target_count = 2752
22
23 # 5. Pakai 1 kali SMOTE (split train test)
24 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
25 train_idx, test_idx = next(iter(skf.split(X, y)))
26 X_train, y_train = X.iloc[train_idx], y.iloc[train_idx]
27 X_test, y_test = X.iloc[test_idx], y.iloc[test_idx] # 18. Split test
28
29 # 6. SMOTE dengan target count per label
30 sampling_strategy = {label: target_count for label in y_train.unique()}
31 smote = SMOTE(sampling_strategy=sampling_strategy, random_state=42)
32 X_train_res, y_train_res = smote.fit_resample(X_train, y_train)
33
34 # 7. Print distribusi label setelah SMOTE
35 print("\n== Distribusi Label DENGAN SMOTE ==")
36 print(y_train_res.value_counts())
37
38 # 8. Train model dengan data latih hasil SMOTE
39 model = MultinomialNB()
40 model.fit(X_train_res, y_train_res)
41
42 # 9. Prediksi di data latih hasil SMOTE
43 y_train_pred = model.predict(X_train_res)
44
45 # 10. Prediksi di data uji (yang asli, tanpa SMOTE)
46 y_test_pred = model.predict(X_test)
47
48 # 11. === Classification Report pada train set ===
49 print("\n== Classification Report - Train (DENGAN SMOTE) ===")
50 print(classification_report(y_train_res, y_train_pred, digits=2))
51
52 # 12. === Rata-rata metrik ===
53 print("\n== Rata-rata Hasil Evaluasi (DENGAN SMOTE) ===")
54 print(f"Precision: {np.mean(macro_precisions):.4f}")
55 print(f"Recall: {np.mean(macro_recalls):.4f}")
56 print(f"F1-Score: {np.mean(macro_f1s):.4f}")
57 print(f"Accuracy: {np.mean(accuracies):.4f}")
58
59 # 13. === Confusion Matrix pada data latih hasil SMOTE ===
60 labels = sorted(y_train_res.unique())
61 cm_train = confusion_matrix(y_train_res, y_train_pred, labels=labels)
62 disp = ConfusionMatrixDisplay(confusion_matrix=cm_train, display_labels=labels)
63 disp.plot(cmap='Blues')
64 plt.title("Confusion Matrix - Train (DENGAN SMOTE)")
65 plt.show()
66
67 # 14. Print total per label
68 print(f"Total per label seharusnya: {target_count} x {len(labels)} = {target_count * len(labels)}")
69 print(f"Total di confusion matrix (train): {cm_train.sum()}")

```

Lampiran 16. Kode Program *Naive Bayes* SMOTE Fold-10 + *Confusion Matrix*

```

● ● ●

1 # 1. Import library yang dibutuhkan
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import StratifiedKFold
5 from sklearn.naive_bayes import MultinomialNB
6 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay, classification_report
7 import matplotlib.pyplot as plt
8 from imblearn.over_sampling import SMOTE
9
10 # 2. Load data
11 df = pd.read_csv('dataset_tfidf_label.csv')
12 X = df.drop(columns=['label'])
13 y = df['label']
14
15 # 3. Bersihkan target NaN
16 not_nan_mask = y.notna()
17 X = X[not_nan_mask]
18 y = y[not_nan_mask]
19
20 # 4. Target per label SESUDAH SMOTE
21 target_count = 2752
22
23 # 5. Pakai 1 kali SMOTE (split train test)
24 skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
25 train_idx, test_idx = next(iter(skf.split(X, y)))
26 X_train, y_train = X.iloc[train_idx], y.iloc[train_idx]
27 X_test, y_test = X.iloc[test_idx], y.iloc[test_idx]
28
29 # 6. SMOTE dengan target count per label
30 sampling_strategy = {label: target_count for label in y_train.unique()}
31 smote = SMOTE(sampling_strategy=sampling_strategy, random_state=42)
32 X_train_res, y_train_res = smote.fit_resample(X_train, y_train)
33
34 # 7. Print distribusi label setelah SMOTE
35 print("\n== Distribusi Label DENGAN SMOTE ==")
36 print(y_train_res.value_counts())
37
38 # 8. Train model dengan data latih hasil SMOTE
39 model = MultinomialNB()
40 model.fit(X_train_res, y_train_res)
41
42 # 9. Prediksi di data latih hasil SMOTE
43 y_train_pred = model.predict(X_train_res)
44
45 # 10. Prediksi di data uji (yang asli, tanpa SMOTE)
46 y_test_pred = model.predict(X_test)
47
48 # 11. === Classification Report pada train set ===
49 print("\n== Classification Report - Train (DENGAN SMOTE) ==")
50 print(classification_report(y_train_res, y_train_pred, digits=2))
51
52 # 12. === Rata-rata metrik ===
53 # *Note*: macro_precisions, macro_recalls, macro_f1s, accuracies belum diinisialisasi di script ini.
54 print("\n== Rata-rata Hasil Evaluasi (DENGAN SMOTE) ==")
55 print(f"Precision: {np.mean(macro_precisions):.4f}")
56 print(f"Recall: {np.mean(macro_recalls):.4f}")
57 print(f"F1-Score: {np.mean(macro_f1s):.4f}")
58 print(f"Accuracy: {np.mean(accuracies):.4f}")
59
60 # 13. === Confusion Matrix pada data latih hasil SMOTE ===
61 labels = sorted(y_train_res.unique())
62 cm_train = confusion_matrix(y_train_res, y_train_pred, labels=labels)
63 disp = ConfusionMatrixDisplay(confusion_matrix=cm_train, display_labels=labels)
64 disp.plot(cmap='Blues')
65 plt.title("Confusion Matrix - Train (DENGAN SMOTE)")
66 plt.show()
67
68 # 14. Print total per label dan total confusion matrix
69 print(f"Total per label seharusnya: {target_count} x {len(labels)} = {target_count * len(labels)}")
70 print(f"Total di confusion matrix (train): {cm_train.sum()}")

```

Lampiran 17. Kode Program Word Cloud

```

● ● ●

1 # 1. Import library
2 import pandas as pd
3 import re
4 from wordcloud import WordCloud
5 import matplotlib.pyplot as plt
6
7 # 2. Baca data
8 df = pd.read_csv('hasil_stemming.csv')
9
10 # === POSITIF ===
11 # 3. Filter POSITIF
12 positif_df = df[df['label'] == 'positif']['stemmed'].astype(str)
13
14 # 4. Gabung dan bersihkan
15 positif_teks = " ".join(positif_df) jadi 1 string
16 positif_teks_bersih = re.sub(r"[^a-zA-Z\s]", " ", positif_teks)
17
18 # 5. Buat WordCloud POSITIF
19 wc = WordCloud(
20     width=800,
21     height=400,
22     background_color='white',
23     stopwords=None,
24     max_words=100,
25     collocations=False
26 ).generate(positif_teks_bersih)
27
28 # 6. Tampilkan WordCloud POSITIF
29 plt.figure(figsize=(12, 6))
30 plt.imshow(wc, interpolation='bilinear')
31 plt.axis('off')
32 plt.title('Word Cloud - Label Positif', fontsize=20)
33 plt.show()
34
35 # === NETRAL ===
36 # 7. Filter NETRAL
37 netral_df = df[df['label'] == 'netral']['stemmed'].astype(str)
38
39 # 8. Gabung dan bersihkan
40 netral_teks = " ".join(netral_df)
41 netral_teks_bersih = re.sub(r"[^a-zA-Z\s]", " ", netral_teks)
42
43 # 9. Buat WordCloud NETRAL
44 wc = WordCloud(
45     width=800,
46     height=400,
47     background_color='white',
48     stopwords=None,
49     max_words=100,
50     collocations=False
51 ).generate(netral_teks_bersih)
52
53 # 10. Tampilkan WordCloud NETRAL
54 plt.figure(figsize=(12, 6))
55 plt.imshow(wc, interpolation='bilinear')
56 plt.axis('off')
57 plt.title('Word Cloud - Label Netral', fontsize=20)
58 plt.show()
59
60 # === NEGATIF ===
61 # 11. Filter NEGATIF
62 negatif_df = df[df['label'] == 'negatif']['stemmed'].astype(str)
63
64 # 12. Gabung dan bersihkan
65 negatif_teks = " ".join(negatif_df)
66 negatif_teks_bersih = re.sub(r"[^a-zA-Z\s]", " ", negatif_teks)
67
68 # 13. Buat WordCloud NEGATIF
69 wc = WordCloud(
70     width=800,
71     height=400,
72     background_color='white',
73     stopwords=None,
74     max_words=100,
75     collocations=False
76 ).generate(negatif_teks_bersih)
77
78 # 14. Tampilkan WordCloud NEGATIF
79 plt.figure(figsize=(12, 6))
80 plt.imshow(wc, interpolation='bilinear')
81 plt.axis('off')
82 plt.title('Word Cloud - Label Negatif', fontsize=20)
83 plt.show()

```

