

Laporan Eksperimen Latihan 02 KLASIFIKASI		Kelas : D Ketua : Dimas Anggota : 1. fazar 2. rifaldi 3. elsa 4. dendi 5. hilmi	
Makalah Jurnal yg dijadikan referensi Judul: ANALISA SENTIMEN TERHADAP PROGRAM KARTU PRAKERJA PADA INSTAGRAM MENGGUNAKAN METODE NAÏVE BAYES CLASSIFIER Author: MOHAMMAD RIZAL RIZWAN Tahun publikasi: 2022 Link unduh: http://repository.uin-suska.ac.id/58016/1/File%20Lengkap%20Sampai%20Lampiran%20Kecuali%20Hasil%20Penelitian.pdf			
Makalah jurnal referensi boleh lebih dari 1.			
Dataset yang digunakan	Nama dataset: DatasetIG.csv Link unduh: https://raw.githubusercontent.com/dimasw09/BelajarKoding/main/DatasetIG.csv		
	No	Nama Atribut	Type Data
	1. 2. 3. 4. dst	id sentimen komentar	int string string
Algoritma yang dipelajari	Support Vector Machine		
Penjelasan tentang algoritma <i>Berikan penjelasan mengenai algoritma yang dipelajari, baik konsep, karakteristik maupun langkah kerjanya.</i>	Support Vector Machine (SVM) adalah algoritma pembelajaran mesin yang digunakan untuk klasifikasi dan regresi. Tujuannya adalah mencari hiperplane terbaik yang memisahkan dua kelas data dengan margin terbesar. SVM menggunakan support vectors, yaitu titik-titik data penting di sekitar margin, untuk melakukan klasifikasi. Algoritma ini juga dapat menggunakan kernel untuk memproyeksikan data ke dalam ruang fitur yang lebih tinggi. SVM memiliki kelebihan seperti kemampuan menangani data dengan fitur yang banyak, ketahanan terhadap overfitting, dan keakuratan yang tinggi.		
Persiapan Awal Data <i>Berikan penjelasan mengenai kegiatan data preparation yang dilakukan (optional/ jika ada).</i>	Pengumpulan Data: Dilakukan scraping data komentar dari postingan akun Instagram Kartu Prakerja menggunakan teknik tertentu, seperti menggunakan API Instagram atau library web scraping. Data komentar tersebut kemudian disimpan dalam format yang sesuai, seperti file CSV atau spreadsheet. Pembersihan Data: Tahap ini melibatkan pembersihan data yang mungkin mengandung karakter khusus, tautan, emoji, atau teks yang tidak relevan. Data yang tidak diperlukan juga dapat dihapus, seperti kolom waktu atau nama pengguna yang tidak relevan dengan analisis yang akan dilakukan. Normalisasi Teks: Dilakukan normalisasi teks untuk mengubah semua huruf menjadi lowercase, menghapus tanda baca, dan menghilangkan spasi berlebih. Hal ini akan membantu dalam konsistensi analisis kata-kata. Penghapusan Stop Words: Stop words adalah kata-kata umum yang tidak memberikan makna khusus dalam analisis teks, seperti "dan", "di", atau "dari". Menghapus stop words dapat membantu fokus pada kata-kata penting yang lebih relevan dalam analisis. Tokenisasi: Data komentar perlu dipecah menjadi token-token terpisah, yaitu unit-unit terkecil seperti kata atau frasa. Proses		

	<p>tokenisasi ini mempermudah analisis selanjutnya, seperti menghitung frekuensi kata atau melakukan analisis sentimen.</p> <p>Lematisasi/Stemming: Dalam tahap ini, kata-kata dalam data komentar dapat diubah ke bentuk dasar atau kata dasar, menggunakan teknik lemmatisasi atau stemming. Tujuannya adalah untuk mengurangi variasi kata yang memiliki akar yang sama, sehingga dapat meningkatkan efisiensi dan akurasi analisis.</p> <p>Analisis Sentimen (Opsional): Jika tujuan dari analisis data komentar adalah untuk mengukur sentimen atau pendapat pengguna terhadap Kartu Prakerja, maka dapat dilakukan analisis sentimen. Ini melibatkan pengklasifikasian komentar menjadi positif, negatif, atau netral berdasarkan makna dan konteksnya.</p> <p>Visualisasi/Data Exploration: Setelah data telah dipersiapkan, tahap ini melibatkan visualisasi data, seperti membuat grafik atau word cloud, untuk membantu memahami pola dan tren dalam komentar pengguna.</p>
<p>Langkah pengerjaan secara manual</p> <p><i>Berikan penjelasan langkah per langkah bagaimana algoritma yang dipilih melakukan proses clustering pada dataset</i></p>	<p>Persiapan Data</p> <ul style="list-style-type: none"> • Mengumpulkan dataset yang akan digunakan untuk pelatihan dan pengujian model SVM. • Membagi dataset menjadi dua subset, yaitu data latih (training data) dan data uji (testing data). Data latih digunakan untuk melatih model SVM, sedangkan data uji digunakan untuk menguji kinerja model yang telah dilatih. <p>Pemilihan Kernel</p> <ul style="list-style-type: none"> • Memilih jenis kernel yang sesuai untuk model SVM. Kernel linier, kernel polinomial, dan kernel RBF (radial basis function) adalah beberapa kernel yang umum digunakan. • Kernel ini akan membantu dalam mentransformasikan data ke dalam ruang fitur yang lebih tinggi, jika diperlukan, untuk meningkatkan pemisahan antara kelas-kelas data. <p>Pelatihan Model</p> <ul style="list-style-type: none"> • Melatih model SVM dengan menggunakan data latih. • Selama proses pelatihan, SVM akan mencari hiperplane terbaik yang memisahkan kedua kelas data dengan margin terbesar. • Jika menggunakan kernel, SVM akan menghitung dan mengoptimalkan parameter kernel yang sesuai selama pelatihan. <p>Validasi dan Evaluasi</p> <ul style="list-style-type: none"> • Menggunakan data uji untuk menguji kinerja model SVM yang telah dilatih. • Model SVM akan melakukan prediksi kelas untuk setiap data uji berdasarkan hiperplane yang telah ditentukan selama pelatihan. • Menghitung metrik evaluasi seperti akurasi, presisi, recall, dan F1-score untuk mengukur kinerja model SVM. <p>Penyetelan Parameter (Opsional)</p> <ul style="list-style-type: none"> • Jika hasil evaluasi tidak memuaskan, dapat dilakukan penyetelan parameter SVM untuk meningkatkan kinerja model. • Misalnya, parameter C dapat disesuaikan untuk mengatur trade-off antara margin yang lebih besar dan jumlah kesalahan klasifikasi. <p>Prediksi</p> <ul style="list-style-type: none"> • Setelah model SVM dilatih dan dievaluasi, model tersebut dapat digunakan untuk melakukan prediksi pada data baru yang tidak terlihat sebelumnya. • SVM akan mengklasifikasikan data baru berdasarkan hiperplane yang telah ditentukan selama pelatihan.

Isi dengan screenshot model, parameter yang digunakan serta penjelasan mengenai bagaimana penerapan algoritma pada alat data mining yang dipilih

```
[ ] dataUrl = 'https://raw.githubusercontent.com/dimasw09/BelajarKoding/main/DatasetIG.csv'
```

```
[ ] # membersihkan karakter data agar data dapat diolah
import re
import emoji
import contractions

def cleaningtext(text):
    text = re.sub('RT|u', '', text) # replace RT tag
    text = re.sub('@[a-zA-Z0-9_]+', '', text) # replace @_username
    text = emoji.demojize(text) # replace emoji with text
    text = re.sub('http|https://|/|5+', '', text) # replace URL
    text = re.sub('#+', '', text) # replace #something
    text = text.lower() # lower case kalimat
    text = re.sub('(.)(1+)', r'\1', text) # replace word repetition
    text = re.sub(r'[\u0017-\u001f]', '', text)
    text = re.sub(r'[\&]+', '', text)
    text = re.sub(r'[\]\[\#\&:<0>]+=[\|,_,", " ', text)
    text = re.sub(r'["a-zA-Z]', ' ', text)
    text = contractions.fix(text) # replace contractions

    return text
```

Tokenize

```
import nltk
nltk.download('punkt')
from nltk.tokenize import word_tokenize
def tokenizing(text):
    data = word_tokenize(text)
    return data
df['tokenizing'] = df['Cleaning'].apply(tokenizing)
df
```

Normalisasi Kata

Normalisasi Kata

	Id	Sentimen	Komentar	Case Folding	Cleaning	Tokenizing	Formalisasi
	0	Positif	Yah semoga dan semoga. Baru pertama kali nyoba...	yah semoga dan semoga. baru pertama kali nyoba...	yah semoga dan semoga baru pertama kali nyoba...	[yah, semoga, dan, semoga, baru, pertama, kali, nyoba, ...]	[yah, semoga, dan, semoga, baru, pertama, kali, nyoba, ...]
	1	Positif	BISMILLAAHH JUMAT BERKAH COBA BISA...	bismillaah jumat berkah coba kash red nea...	bismillaah jumat berkah coba kash red nea...	[bismillaah, jumat, berkah, coba, kash, red, nea, ...]	[bismillaah, jumat, berkah, coba, kash, red, nea, ...]
	2	Positif	MININ DOAKAN AGAR BISA LOLOS YA DI GELOMBANG 5...	minin doakan agar bisa lolos ya di gelombang 5...	minin doakan agar bisa lolos ya di gelombang 5...	[minin, doakan, agar, bisa, lolos, ya, di, gelombang, 5, ...]	[minin, doakan, agar, bisa, lolos, ya, di, gelombang, 5, ...]
	3	Netral	Kio mau beli petalannya bagusnya dimana yah...	Kio mau beli petalannya bagusnya dimana yah...	Kio mau beli petalannya bagusnya dimana yah...	[kio, mau, beli, petalannya, bagusnya, diman, ...]	[kialau, mau, beli, petalannya, bagusnya, diman, ...]
	4	Netral	Min, yg lagi kuliah bisa ikut prakteja?	min, yg lagi kuliah bisa ikut prakteja?	min yg lagi kuliah bisa ikut prakteja	[min, yg, lagi, kuliah, bisa, ikut, prakteja]	[min, yang, lagi, kuliah, bisa, ikut, prakteja]

453	454	Negatif	Syarat Petalihan nya kok kebanyakan harus merti...	syarat petalihan nya kok kebanyakan harus merti...	syarat petalihan nya kok kebanyakan harus merti...	[syarat, petalihan, nya, kok, kebanyakan, haru...	[syarat, petalihan, nya, kok, kebanyakan, haru...
454	455	Negatif	Kapok iutan ngga pernah bisa lolos	kapok iutan ngga pernah bisa lolos	kapok iutan ngga pernah bisa lolos	[kapok, iutan, ngga, pernah, bisa, lolos]	[kapok, iutan, tngga, pernah, bisa, lolos]
455	456	Negatif	Untung gw dah lolos tau lah 🤔🤔🤔 600K 4X...	untung gw dah lolos tau lah 🤔🤔🤔 600K 4X...	untung gw dah lolos tau lah face wit tears...	[untung, gw, dah, lolos, tau, lah, face, wit, tears, ...]	[untung, saya, sudah, lolos, tauh, lah, face, wit, tears, ...]
456	457	Negatif	Gak pernah lolos males	gak pernah lolos males	gak pernah lolos males	[gak, pernah, lolos, males]	[gak, pernah, lolos, males]
457	458	Negatif	Gk enak di awal2 prakteja 600b sebanyak 4x	gk enak di awal2 prakteja 600b sebanyak 4x	gk enak di awal prakteja rd sebanyak x	[gk, enak, di, awal, prakteja, rd, sebanyak, x]	[gk, enak, di, awal, prakteja, rd, sebanyak, x]

```

4 from nltk.corpus import stopwords
nltk.download("stopwords")

defn stopword = stopwords.words("indonesian")
defn stopword = pd.read_csv("https://raw.githubusercontent.com/dimas09/BelajarNLP/main/stopword_list_vaja.csv")
defn stopword = set(stopword.stopword)

defn stopwordTest(words):
    return [word for word in words if word not in defn stopword]

df["stopword"] = df["formalisasi"].apply(stopwordTest)
df

```

[illegible]

```

from Sastrawi.Stemmer.StemmerFactory import StemmerFactory
import swifter

factory = StemmerFactory()
stemmer = factory.create_stemmer()

def stemmed_wrapper(term):
    return stemmer.stem(term)

term_dict = {}

for document in df['Stopword']:
    for term in document:
        if term not in term_dict:
            term_dict[term] = ''

for term in term_dict:
    term_dict[term] = stemmed_wrapper(term)
    print(term, ":", term_dict[term])

def stemmingText(document):
    return [term_dict[term] for term in document]

df['Stemming'] = df['Stopword'].swifter.apply(stemmingText)

```

```

yah : yah
semoga : moga
dan : dan
baru : baru
pertama : pertama
kali : kali
coba : coba
soalnya : soal
bismillaahh : bismillaahh
jumat : jumat
berkah : berkah
kasih : kasih
redaksi : redaksi
heart : heart
buat : buat
sobat : sobat
prakerja : prakerja
yang : yang
mau : mau
rejekinya : rejekinya
di : di
gelombang : gelombang

```

Uji Coba

```

from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
from sklearn import model_selection
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder

X = df['Cleaning']
Y = df['Sentimen']

x_train, x_test, y_train, y_test = model_selection.train_test_split(X, Y, test_size=0.1)

```

```

[19] vectorizer = TfidfVectorizer()
x_train = vectorizer.fit_transform(x_train)
x_test = vectorizer.transform(x_test)
Encoder = LabelEncoder()
y_train = Encoder.fit_transform(y_train)
y_test = Encoder.fit_transform(y_test)

```

Modeling

```

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.ensemble import AdaBoostClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.preprocessing import MinMaxScaler
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import cross_val_score

```

Support Vector Machine Split 2:8

	<pre># SVM Model svc = SVC() svc.fit(x_train, y_train) pred5 = svc.predict(x_test) print(classification_report(y_test, pred5))</pre> <table><thead><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr></thead><tbody><tr><td>0</td><td>0.70</td><td>0.91</td><td>0.79</td><td>46</td></tr><tr><td>1</td><td>0.80</td><td>0.46</td><td>0.59</td><td>26</td></tr><tr><td>2</td><td>0.82</td><td>0.70</td><td>0.76</td><td>20</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.74</td><td>92</td></tr><tr><td>macro avg</td><td>0.77</td><td>0.69</td><td>0.71</td><td>92</td></tr><tr><td>weighted avg</td><td>0.76</td><td>0.74</td><td>0.73</td><td>92</td></tr></tbody></table> <pre>[34] acc_score = cross_val_score(svc, x_train, y_train, cv=5, scoring='accuracy') pre_score = cross_val_score(svc, x_train, y_train, cv=5, scoring='precision_macro') rec_score = cross_val_score(svc, x_train, y_train, cv=5, scoring='recall_macro') f_score = cross_val_score(svc, x_train, y_train, cv=5, scoring='f1_macro') print('Hasil Accuracy : %s' % (acc_score)) print('Hasil Rata - Rata Accuracy : %s' % acc_score.mean()) print('Hasil Precision : %s' % (pre_score)) print('Hasil Rata - Rata Precision : %s' % pre_score.mean()) print('Hasil Recall : %s' % (rec_score)) print('Hasil Rata - Rata Recall : %s' % rec_score.mean()) print('Hasil F-Measure : %s' % (f_score)) print('Hasil Rata - Rata F-Measure : %s' % f_score.mean())</pre> <pre>Hasil Accuracy : [0.7027027 0.68493151 0.68493151 0.68493151 0.7260274] Hasil Rata - Rata Accuracy : 0.6967049241021843 Hasil Precision : [0.76619133 0.8308642 0.8452381 0.81908832 0.81687675] Hasil Rata - Rata Precision : 0.8156517375144826 Hasil Recall : [0.64513557 0.62626263 0.62169312 0.62247121 0.67251462] Hasil Rata - Rata Recall : 0.6376154289466983 Hasil F-Measure : [0.65877671 0.65304483 0.60601852 0.63470729 0.68545056] Hasil Rata - Rata F-Measure : 0.6475995815456327</pre>		precision	recall	f1-score	support	0	0.70	0.91	0.79	46	1	0.80	0.46	0.59	26	2	0.82	0.70	0.76	20	accuracy			0.74	92	macro avg	0.77	0.69	0.71	92	weighted avg	0.76	0.74	0.73	92
	precision	recall	f1-score	support																																
0	0.70	0.91	0.79	46																																
1	0.80	0.46	0.59	26																																
2	0.82	0.70	0.76	20																																
accuracy			0.74	92																																
macro avg	0.77	0.69	0.71	92																																
weighted avg	0.76	0.74	0.73	92																																
Hasil Eksperimen <i>Berikan penjelasan mengenai hasil yang didapatkan: pengetahuan/ pola dapat berupa pohon keputusan atau classification rules.</i>	<p>Dengan menggunakan metode split 2:8, hasil yang diperoleh menunjukkan keandalan yang tinggi. Akurasi rata-ratanya adalah 0.6967049241021843, mengindikasikan tingkat keakuratan yang sangat baik.</p> <p>Selain itu, presisi rata-rata mencapai 0.8156517375144826, menunjukkan konsistensi dan ketepatan model.</p> <p>Recall rata-ratanya adalah 0.6376154289466983, menunjukkan kemampuan model dalam mengenali dan mengingat informasi yang relevan.</p> <p>Selain itu, F-Measure rata-ratanya adalah 0.6475995815456327, mencerminkan keseimbangan yang baik antara presisi dan recall.</p> <p>Dengan demikian, evaluasi menunjukkan bahwa model ini dapat diandalkan dan memberikan hasil yang meyakinkan.</p>																																			
Metode Evaluasi <i>Berikan penjelasan mengenai bagaimana model klasifikasi dan hasilnya meliputi aspek akurasi, precision, recall & AUC</i>	<p>Akurasi: Akurasi mengukur seberapa baik model dapat mengklasifikasikan data dengan benar. Hasil akurasi rata-rata sebesar 0.6967049241021843 menunjukkan bahwa model memiliki tingkat keakuratan yang sangat baik dalam memprediksi kelas yang benar.</p> <p>Presisi: Presisi menggambarkan kemampuan model dalam memberikan hasil positif yang benar dari semua prediksi yang positif. Hasil presisi rata-rata sebesar 0.8156517375144826 menunjukkan bahwa model konsisten dan akurat dalam mengidentifikasi dan mengklasifikasikan data dengan benar.</p> <p>Recall: Recall mengukur kemampuan model dalam mengenali dan mengingat informasi yang relevan. Hasil recall rata-rata sebesar 0.6376154289466983 menunjukkan bahwa model memiliki kemampuan yang baik dalam mengenali data yang relevan dan menghindari kesalahan dalam mengklasifikasikan data yang sebenarnya positif.</p> <p>F-Measure: F-Measure adalah ukuran yang menggabungkan presisi dan recall untuk memberikan gambaran keseluruhan tentang kinerja model. Hasil F-Measure rata-rata sebesar 0.6475995815456327 menunjukkan bahwa model memiliki keseimbangan yang baik antara presisi dan recall.</p>																																			

<p>Daftar Pustaka</p> <p><i>Daftar buku, jurnal atau link referensi yang membantu proses eksperimen</i></p>	<ol style="list-style-type: none"> 1. https://journal.lppmunindra.ac.id/index.php/Faktor_Exacta/article/view/7964 2. https://openlibrarypublications.telkomuniversity.ac.id/index.php/engineering/article/view/10704 3. http://ejournal.poltektegal.ac.id/index.php/informatika/article/view/2870 4. http://download.garuda.kemdikbud.go.id/article.php?article=638347&val=10384&title=Analisis%20Sentimen%20Cyberbullying%20pada%20Komentar%20Instagram%20dengan%20Metode%20Klasifikasi%20Support%20Vector%20Machine 5. https://scholar.archive.org/work/32lxzlfclncjdfd5kbwtso/gaem/access/wayback/http://jurnal.fikom.umi.ac.id/index.php/ILKOM/article/download/597/pdf
--	--