words or phrases to others (Davidson et al., 2017). In Indonesia, abusive language commonly stems from various unfavorable conditions, including mental disorders, sexual deviance, physical disabilities, lack of modernization, breaches of etiquette, religious prohibitions, and other unfortunate circumstances (Wijana and Rohmadi, 2010). Twitter, a significant social media platform in Indonesia, frequently becomes a conduit for spreading hate speech (Alfina et al., 2017, 2018; Putri, 2018). During the analysis period, Twitter, referred to as X for anonymity purposes, retained its presence under the twitter.com domain. The training dataset comprised tweets gathered over approximately seven months from March 20, 2018, to September 10, 2018 (Ibrohim and Budi, 2019). This dataset aimed to capture tweet patterns resembling those from trending tweets following the second presidential debate on January 7, 2024. 1. Preparation 1.1. Import Library The first thing to do in the preparation stage is import all the libraries needed for this analysis such as pandas for processing dataframes, numpy for processing arrays, matplotlib, seaborn, and wordcloud for visualization, libraries for managing text such as nltk, re, string, and Sastrawi, and finally the Sklearn model for machine learning purposes. In [1]: **import** pandas **as** pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns import re import string import nltk import warnings from nltk.corpus import stopwords from nltk.tokenize import word\_tokenize from Sastrawi.Stemmer.StemmerFactory import StemmerFactory from wordcloud import WordCloud from sklearn.feature\_extraction.text import TfidfVectorizer from sklearn.model\_selection import train\_test\_split from sklearn.metrics import classification\_report, accuracy\_score, precision\_score, recall\_score, f1\_score, confusion\_matrix from sklearn.neighbors import KNeighborsClassifier from sklearn.linear\_model import LogisticRegression from sklearn.svm import SVC from sklearn.naive\_bayes import ComplementNB, BernoulliNB from sklearn.tree import DecisionTreeClassifier from sklearn.ensemble import RandomForestClassifier from sklearn.model\_selection import GridSearchCV warnings.filterwarnings("ignore") In [11]: nltk.download('stopwords') [nltk\_data] Downloading package stopwords to [nltk data] C:\Users\Rivaldi\AppData\Roaming\nltk\_data... [nltk\_data] Package stopwords is already up-to-date! True Out[11]: 1.2. Create a Function

Project Manajemen Data dan Informasi

Indonesian Hate Speech Analysis on Twitter/X After the Second Presidential Debate in 2024

This analysis seeks to investigate the nature and prevalence of hate speech in some of the tweets that trended after the second Indonesian Presidential Debate in 2024 on Twitter/X. Utilizing various machine learning models, including K-Nearest Neighbors (KNN), Logistic Regression, SVM, Naive Bayes, Decision Tree, and Random Forest, this research aims to find out whether tweets that spike in popularity after being posted are mostly hateful or not. By carefully deploying these models against a dataset of trending tweets, this research aims to determine whether there is a significant trend toward hate and hatred in

post-debate attention-grabbing discourse. These findings aim to shed light on whether trending tweets align with respectful dialogue or tend to amplify hate, thereby offering insights for readers aimed at

Hate speech constitutes direct or indirect expressions of hatred aimed at individuals or groups based on inherent characteristics, such as ethnicity, religion, disability, gender, or sexual orientation (Komnas HAM, 2015). The propagation of hate speech, especially prevalent on social media platforms, often involves the use of abusive language, which encompasses verbally or in written form conveying abusive

**Table of Contents** 

• 5. Machine Learning

• 6. Prediction Conclusions References

**Abstract** 

Introduction

best\_accuracy = 0 best\_model = None

for model\_name, model in models.items(): model.fit(X\_train\_scal, Y\_train) Y\_pred = model.predict(X\_test\_scal)

if accuracy > best\_accuracy: best\_accuracy = accuracy

best\_model = model

f1 = f1\_score(Y\_test, Y\_pred)

accuracy = accuracy\_score(Y\_test, Y\_pred)

accuracy = accuracy\_score(Y\_test, Y\_pred) precision = precision\_score(Y\_test, Y\_pred)

recall = recall\_score(Y\_test, Y\_pred)

scores['Model'].append(model\_name) scores['Accuracy'].append(accuracy) scores['Precision'].append(precision)

scores['Recall'].append(recall) scores['F1 Score'].append(f1)

best\_model\_instance = best\_model

print("No best model found")

df1 = pd.read\_csv("dataset/tweet\_hs.csv") df2 = pd.read\_csv("dataset/kamus\_alay.csv")

best\_model\_instance.fit(X\_train\_scal, Y\_train)

plot\_confusion\_matrix(Y\_test, Y\_pred\_best)

There are 4 datasets used at the beginning and loaded into 4 variables.

df4 = pd.read\_excel("dataset/tweet\_capres\_2024.xlsx")

- disaat semua cowok berusaha melacak perhatia... 1

41. Kadang aku berfikir, kenapa aku tetap perc... 0

RT USER: USER siapa yang telat ngasih tau elu?...

USER USER Kaum cebong kapir udah keliatan dong... 1

alay

3 USER USER AKU ITU AKU\n\nKU TAU MATAMU SIPIT T... 0

Y\_pred\_best = best\_model\_instance.predict(X\_test\_scal)

return best\_model, best\_accuracy, scores\_df, X\_test\_scal\_new

scores\_df = pd.DataFrame(scores)

if best\_model:

2. Data Collection

2.1. Hate Speech Dataset

In [13]: df1.drop("Unnamed: 0",axis=1, inplace=True)

2.2. *Alay* Dictionary Dataset

0 anakjakartaasikasik

pakcikdahtua

t3tapjokowi

Зх

pakcikmudalagi

2.3. Indonesian Stopword Dataset

df1.head()

df2.head()

df3.head()

0

1

3

0

3 4 stopword

ada

adalah

adanya

adapun

df4.head()

agak

1

2

Unnamed: 0

0

1

2

Out[13]:

In [14]:

Out[14]:

Out[15]:

In [16]

Out[16]

scores = {'Model': [], 'Accuracy': [], 'Precision': [], 'Recall': [], 'F1 Score': []}

print(f"Best Model: {best\_model} with Accuracy: {best\_accuracy}")

#url: https://www.kaggle.com/datasets/ilhamfp31/indonesian-abusive-and-hate-speech-twitter-text?select=data.csv

The first dataset consists of tweet data taken from March 20, 2018 to September 10, 2018 which have been labeled as hate speech or not.

The second dataset contains alay terms that are often used by netizens and their meanings. This dataset aims to normalize the tweet data in the first dataset.

The fourth dataset contains tweets taken on January 9, 2024 after the second presidential debate. The data taken is in the form of the 30 most popular tweets on that date.

Preprocessing was carried out for the first and fourth datasets to normalize and remove unnecessary characters which could later influence the prediction results with machine learning later. Preprocessing is

The balancing process is carried out to balance the frequency of hate speech and non-hate speech into a 50:50 ratio. The balancing process begins by looking at the number of datasets that are labeled as

plt.pie(comparation, shadow=True, explode=(0, 0.1), autopct='%1.1f%%', labels=["Non-Hate speech", "Hate speech"], colors=colors, startangle=270)

Non-Hate speech

df3 = pd.read\_csv("dataset/stopwordbahasa.csv", header=None).rename(columns={0: 'stopword'})

Tweet HS

arti

anak jakarta asyik asyik

pak cik sudah tua

pak cik muda lagi

tetap jokowi

tiga kali

The third dataset contains stopwords that will be matched with tweets in the first dataset to be removed.

2.4. Tweets After the Second Presidential Debate in 2024 Dataset

df1["Tweet"] = df1["Tweet"].fillna('').astype(str).apply(lambda x: preprocess(x)) df4["Tweet"] = df4["Tweet"].fillna('').astype(str).apply(lambda x: preprocess(x))

There are 5561 data that are indicated as hate speech. To balance the ratio, 5561 tweet data that are labeled as non-hate speech will be taken.

After that, export the preprocessed dataset for prediction purposes using machine learning.

Tweet HS

ku tau mata sipit lihat 0

After the balancing process is carried out, export the dataset for prediction purposes.

4.1. Hate speech vs Non-Hate speech Comparison

plt.title("Hate speech vs Non-Hate Speech Comparison")

4. Exploratory Data Analysis and Visualization

Comparing the number of hate speech and non-hate speech from the dataset using a pie chart.

plt.savefig('charts/pie\_chart\_hate\_vs\_non.png', bbox\_inches='tight')

Hate speech vs Non-Hate Speech Comparison

57.8%

From the visualization provided, most of the data is non-hate speech, around 57.8%, the rest is hate speech.

Visualization uses wordcloud to see what words appear most often in tweets that are indicated as hate speech.

whs = WordCloud(width=1000, height=800, background\_color="white").generate(text)

ahok isi cebong rezim

gant hancur presiden

dukung diam

From the visualization, the words that appear most often are Jokowi (name of president), cebong, ganti presiden, and so on.

Visualization uses wordcloud to see what words appear most often in tweets that are indicated as non-hate speech.

whs = WordCloud(width=1000, height=800, background\_color="white").generate(text)

dengar masuk apakaijokowi 🎜 iyaan bodoh

plt.savefig('charts/non\_hate\_speech\_wordcloud.png', bbox\_inches='tight')

-- ekonomi

islam

From the visualization, the words that appear most often are orang, gue, Indonesia, and so on.

Tweet

whs = WordCloud(width=1000, height=800, background\_color="white").generate(text)

drama sahnangisin jejak

badandukung C

From the visualization, the words that appear most often are Prabowo and Anies (presidential candidate), Jokowi, and so on.

best\_model, best\_accuracy, scores\_df, X\_test\_scal\_new = vector\_model(data\_HS2, 'HS', (1,2))

286

1395

Hate Speech

Recall F1 Score

plt.bar(x=[pos + i \* 0.1 for pos in x], height=values, width=0.1, label=metric)

SVM

Because the best model is SVM, this model is used to predict tweets after the second presidential debate in 2024

The prediction results show that, of all tweets, 96.7% of tweets are indicated as hate speech

Hate Speech vs Not Hate Speech Comparison

96.7%

along with similar abusive language. This suggests a persistent usage of similar word patterns over the years.

(Rekayasa Sistem Dan Teknologi Informasi), 3(2), 176 - 183. https://doi.org/10.29207/resti.v3i2.935

plt.savefig('charts/pie\_chart\_debat\_comp.png', bbox\_inches='tight')

Metrics for Different Models

Complement NB Bernoulli NB

plt.pie(comparation, shadow=True, explode=(0, 0.1), autopct='%1.1f%%', labels=["Hate Speech", "Not Hate Speech"], colors=colors, startangle=180)

Hate Speech

as LSTM instead of a conventional machine learning model might enhance the ability to discern patterns within tweets, potentially elevating prediction accuracy.

• Komnas HAM. 2015. Buku Saku Penanganan Ujaran Kebencian (Hate Speech). Komisi Nasional Hak Asasi Manusia, Jakarta.

• I Dewa Putu Wijana and Muhammad Rohmadi. 2010. Sosiolinguistik: Kajian, Teori, dan Analisis. Pustaka Pelajar, Yogyakarta.

• Muhammad Okky Ibrohim and Indra Budi. 2019. Multi-label Hate Speech and Abusive Language Detection in Indonesian Twitter.

Based on the analysis of 30 popular tweets post the 2024 second presidential debate, a staggering 96.7% of these tweets were identified as hate speech. The recurrently used words in these tweets included "Prabowo," "Anies," "Jokowi," alongside various abusive terms. Interestingly, these align closely with the frequently occurring terms observed in the training dataset from 2018, featuring "Jokowi" and "Prabowo"

To enhance future analysis, augmenting the training dataset with data from both the first and second presidential debates could offer new tweet patterns. Additionally, employing an neural-network model such

• Thomas Davidson, Dana Warmsley, Michael W. Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In International AAAI Conference on Web and

• Ika Alfina, Rio Mulia, Mohamad Ivan Fanany, and Yudo Ekananta. 2017. Hate speech detection in the indonesian language: A dataset and preliminary study. In International Conference on Advanced

• Tumasjan, A., Andranik, S., Sprenger, T., Sandner, P., & Welpe, I. (2010). Predicting Elections with Twitter: What 140 Characters Reveal about Political Sentiment. Word: Journal of the International

• Kurniawan, S., Gata, W., Puspitawati, D. A., -, N., Tabrani, M., & Novel, K. (2019). Perbandingan Metode Klasifikasi Analisis Sentimen Tokoh Politik Pada Komentar Media Berita Online. Jurnal RESTI

• Omran, E., Al Tararwah, E., & Al Qundus, J. (2023). A comparative analysis of machine learning algorithms for hate speech detection in social media. Online Journal of Communication and Media

• Tansa Trisna Astono Putri. 2018. Analisis dan deteksi hate speech pada sosial twitter berbahasa indonesia. Master's thesis, Faculty of Computer Science, Universitas Indonesia.

• Ansari, M. Z., Aziz, M. B., Siddiqui, M. O., Mehra, H., & Singh, K. P. (2020). Analysis of Political Sentiment Orientations on Twitter. Procedia Computer Science, 167, 1821-1828.

Models

Decision Tree Random Forest

Accuracy Precision Recall

F1 Score

Using a balanced dataset, training is carried out using the model training function. The training results show that the best model is SVM with an accuracy of 83% with prediction results represented using a

1200

- 1000

800

600

400

4.4. The Most Frequently Used Words for Tweets After the Second Presidential Debate in 2024

Visualization uses wordcloud to see what words appear most often in popular tweets after the second presidential debate.

<sup>hina</sup> pakai

4.2. The Most Frequently Used Words for Hate Speech

plt.savefig('charts/hate\_speech\_wordcloud.png', bbox\_inches='tight')

Komunis

4.3. The Most Frequently Used Words for Non-hate Speech

text = ' '.join(nhs['Tweet'].fillna('').astype(str))

hat baca nih

politik

budayaallah

In [18]: df\_debat = pd.read\_csv("dataset/cleaned\_tweet\_capres\_2024.csv") df\_debat.drop("Unnamed: 0", axis=1, inplace=True)

saudara saudara singgung singgung tanah tanah ...

allhamdulilah jabar 80 prabowo baliho doang ma...

istirahat news sumber data anies debat calon p... orang pimpin rendah pijak orang puji bullyan I...

text = ' '.join(df\_debat['Tweet'].fillna('').astype(str))

selamat

plt.savefig('charts/debat\_wordcloud.png', bbox\_inches='tight')

2 ahahahahahahaha teman kelas gue pro prabowo...

plt.imshow(whs, interpolation="bilinear")

ajar

kepa⊥a

text = ' '.join(hs['Tweet'].fillna('').astype(str))

plt.imshow(whs, interpolation="bilinear")

anjing inpimpin cina negara presiden ganti Sina negara korupsi benci ulama

plt.imshow(whs, interpolation="bilinear")

pilih dasar

"Saudara-saudara ada pula yang nyinggung-nying... Allhamdulilah Jabar 80% Prabowo tapi balihonya...

BREAKING NEWS: Sumber Data Yang Disampaikan An...

Seorang pemimpin tidak akan pernah merendahkan...

AHAHAHAHAHAHAHA TEMEN KELAS GW YANG PRO PRA..

carried out with previously declared text processing functions.

df4.to\_csv("dataset/cleaned\_tweet\_capres\_2024.csv")

df = pd.read\_csv("dataset/cleaned\_tweets.csv") df.drop("Unnamed: 0", axis=1, inplace=True)

**0** cowok usaha lacak perhati gue lantas remeh per...

telat tau edan sarap gue gaul cigax jifla cal ...

41 kadang pikir percaya tuhan jatuh kali kali ...

 $df_HS_0 = df[df['HS'] == 0].iloc[:5561]$ 

df\_HS\_bal.to\_csv("dataset/HS.csv")

comparation = df["HS"].value\_counts()

42.2%

colors = ["#6BBCD1", "#f69c9c"]

kaum cebong kafir lihat dongok dungu haha 1

 $df_HS_bal = pd.concat([df_HS_1, df_HS_0], axis=0)$ 

df1.to\_csv("dataset/cleaned\_tweets.csv")

3. Data Wrangling

3.1. Preprocessing Data

3.2. Balancing Data

df1 = df.drop(['Tweet'], axis=1)

hate speech.

df1.sum()

df.head()

1

2

3

4

In [63]:  $df_HS_1 = df[df['HS'] == 1]$ 

df\_HS\_bal.shape

(11122, 2)

plt.show()

Hate speech

In [41]: hs = df.loc[df["HS"]==1]

plt.axis("off")

iya agama<sub>turun</sub> islam bangsat

prabowo kerja bilang ad

nam dungu

In [42]: nhs = df.loc[df["HS"]==0]

plt.axis("off")

ubah bagus

cinta

conton lama

df\_debat.head()

plt.axis("off")

babe hektar

confusion matrix.

Not Hate Speech

Speech

scores\_df

1 Logistic Regression

Complement NB

Bernoulli NB

**Decision Tree** 

0

2

4

True labels

In [20]:

Out[20]:

kalah<sup>nyata</sup>

5. Machine Learning

data\_HS2 = pd.read\_csv('dataset/HS.csv')

1378

278

Not Hate Speech

**Model Accuracy Precision** 

0.823494

0.830986

0.817501

Visualize metrics using bar graphs for each model

metrics = viz\_scores.columns.tolist() for i, metric in enumerate(metrics):

x = range(len(viz\_scores))

viz\_scores.set\_index('Model', inplace=True)

values = viz\_scores[metric].values

plt.title('Metrics for Different Models')

plt.xticks([i + 0.15 for i in x], viz\_scores.index) plt.legend(bbox\_to\_anchor=(1.02, 1), loc='upper left')

plt.savefig('charts/metrics\_bar\_chart.png', bbox\_inches='tight')

Logistic Regression

Y\_pred = best\_model.predict(X\_test\_scal\_new)

comparation = df4["HS"].value\_counts()

plt.title("Hate Speech vs Not Hate Speech Comparison")

3.3%

colors = ["#fc8a77", "#a2d7d8"]

Metrics table of all trained models.

KNN

SVM

viz\_scores = scores\_df.copy()

plt.figure(figsize=(10, 4))

plt.ylabel('Values')

plt.tight\_layout()

plt.show()

1.0

0.8

0.6

0.4

0.2

0.0

In [23]: df4["HS"] = Y\_pred

plt.show()

Not Hate Speech

Conclusions

References

Social Media (ICWSM), pages 512-515.

https://doi.org/10.1016/j.procs.2020.03.201

Linguistic Association, 10.

Computer Science and Information Systems (ICACSIS), pages 233–238.

Technologies, 13(4), e202348. https://doi.org/10.30935/ojcmt/13603

In [45]:

KNN

6. Prediction

Values

Best Model: SVC() with Accuracy: 0.8309859154929577

Confusion Matrix

Predicted labels

0.772251 0.751793 0.814704 0.781985

Random Forest 0.820497 0.815882 0.829050 0.822413

0.818449 0.832636 0.825481

0.829863 0.833831 0.831843

0.782979 0.879857 0.828596

plt.show()

swedan air \_ \_ \_

Out[18]:

1

3

cina

nama

kepal

plt.show()

bangun

cobarakyatgoblok

partai komunis

plt.show()

bikin

Out[42]:

Out[6]:

Out[63]:

In [40]:

5561

dtype: int64

**Tweet** 

• 4. Exploratory Data Analysis and Visualization

encouraging more constructive and responsible online conversations around political events on Twitter/X.

 Abstract Introduction • 1. Preparation • 2. Data Collection • 3. Data Wrangling

After importing all the necessary libraries, the next stage is to create functions that will help with analysis, here are some of these functions. 1.2.1. Text processing function There are several functions created to process text, such as changing text to lowercase, removing unnecessary characters such as url, rt, user, emoji, etc. Then there are also functions for removing nonalphanumeric text, normalizing alay words, removing stop words and excess spaces, and stemming. In [49]: factory = StemmerFactory() stemmer = factory.create\_stemmer() def lowercase(text): return text.lower() def remove\_unnecessary\_char(text):  $text = re.sub('\n', ' ', text)$ text = re.sub('rt',' ',text) text = re.sub('user',' ',text) text = re.sub('url',' ', text) text = re.sub('((www\.[^\s]+)|(https?://[^\s]+))(http?://[^\s]+))',' ',text) text = re.sub(' +', ' ', text) return text def remove\_unicode(text): text =  $re.sub(r'\bx[a-fA-F0-9]{2}\b', '', text)$ text =  $re.sub(r'\bx([a-fA-F0-9]{2})', '', text)$ return text def remove\_nonaplhanumeric(text): text =  $re.sub('[^0-9a-zA-Z]+', '', text)$ **return** text alay\_dict\_map = dict(zip(df2["alay"], df2["arti"])) def normalize\_alay(text): return ' '.join([alay\_dict\_map[word] if word in alay\_dict\_map else word for word in text.split(' ')]) def remove\_stopword(text): text = ' '.join(['' if word in df3.stopword.values else word for word in text.split(' ')])
text = re.sub(' +', ' ', text) text = text.strip() return text def stemming(text): return stemmer.stem(text) def remove\_extra\_spaces(text): text =  $re.sub(r'\s+', '', text)$ text = text.strip() return text In [ ]: def preprocess(text): text = lowercase(text)

text = remove\_nonaplhanumeric(text) text = remove\_unnecessary\_char(text) text = normalize\_alay(text) text = stemming(text) text = remove\_stopword(text) text = remove\_unicode(text) text = remove\_extra\_spaces(text) return text 1.2.2. Confusion matrix generator function This function was created to see the prediction results from machine learning to match whether the prediction results are accurate or not for evaluation. In [44]: def plot\_confusion\_matrix(y, y\_predict): cm = confusion\_matrix(y, y\_predict) ax = plt.subplot() sns.heatmap(cm, annot=True, ax=ax, fmt='d') ax.set\_xlabel('Predicted labels') ax.set\_ylabel('True labels') ax.set\_title('Confusion Matrix') ax.xaxis.set\_ticklabels(['Not Hate Speech', 'Hate Speech']) ax.yaxis.set\_ticklabels(['Not Hate Speech', 'Hate Speech']) plt.savefig('charts/confusion\_matrix.png', bbox\_inches='tight') plt.show() 1.2.3. Model trainer function This function was created to carry out training on data using several models at once such as KNN, Logistic Regression, SVM, Complement Naive Bayes, Bernoulli Naive Bayes, Decision Tree, and Random Forest. After training, the function will return the module with the highest accuracy results while generating the prediction results in the form of a confusion matrix. Apart from that, the function will also return the evaluation results of all models in the form of a dataframe. In [4]: def vector\_model(data, category, ngram): X = data['Tweet'].fillna(' ') Y = data[category]

X\_train, X\_test, Y\_train, Y\_test = train\_test\_split(X, Y, test\_size=0.3, random\_state=42) vector = TfidfVectorizer(ngram\_range=ngram, stop\_words=df3['stopword'].tolist()) X\_train\_scal = vector.fit\_transform(X\_train) X\_test\_scal = vector.transform(X\_test) X\_test\_scal\_new = vector.transform(df\_debat["Tweet"]) models = { 'KNN': KNeighborsClassifier(n\_neighbors=5), 'Logistic Regression': LogisticRegression(), 'SVM': SVC(kernel='rbf'), 'Complement NB': ComplementNB(), 'Bernoulli NB': BernoulliNB(), 'Decision Tree': DecisionTreeClassifier(criterion='entropy', min\_samples\_split=2, random\_state=42), 'Random Forest': RandomForestClassifier(n\_estimators=105, min\_samples\_split=2, random\_state=42) }