

Lampiran 1 Source Code

Crawling Data Instagram

```
❶ pip install selenium bs4
pip install beautifulsoup4
pip install pandas
lapt-get update
lapt-get install -y chromium-browser
lapt install chromium-chromedriver
```

```
] from selenium import webdriver
from selenium.webdriver.common.keys import Keys
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.common.by import By
from selenium.webdriver.support import expected_conditions as EC
from bs4 import BeautifulSoup
import time
import pandas as pd
from pandas import ExcelWriter
import os.path
import csv
```

```
# path ke Chromedriver
chromedriver_path = '/usr/lib/chromium-browser/chromedriver'

# Set option Chrome untuk webdriver
options = webdriver.ChromeOptions()
options.add_argument('--no-sandbox')
options.add_argument('--headless')
options.add_argument('--disable-gpu')
options.add_argument('--disable-dev-shm-usage')

# Inisialisasi driver Chrome
driver = webdriver.Chrome(options=options)
```

```
url = 'https://www.instagram.com/accounts/login/'
driver.get(url)
username_input = 'username'
password_input = 'PASSWORD'
# klik di username input
username_webdriverwait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="text" value="" ]'))).send_keys(username_input)
# klik di password input
password_webdriverwait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="password" value="" ]'))).send_keys(password_input)
# press ENTER untuk login
WebDriverWait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="submit" value="Log In" ]'))).send_keys(Keys.ENTER)
```

```
from selenium import webdriver
from selenium.webdriver.common.by import By
from selenium.webdriver.support import expected_conditions as EC
from selenium.webdriver.support.ui import WebDriverWait
from selenium.common.exceptions import NoSuchElementException
from selenium.common.exceptions import TimeoutException
from selenium.webdriver.common.keys import Keys
import time
import pandas as pd
import os
import csv

# Inisialisasi driver Chrome
driver = webdriver.Chrome()

# url yang akan di akses
url = 'https://www.instagram.com/accounts/login/'
driver.get(url)

# klik di username input
username_webdriverwait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="text" value="" ]'))).send_keys('username')

# klik di password input
password_webdriverwait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="password" value="" ]'))).send_keys('PASSWORD')

# press ENTER untuk login
WebDriverWait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="submit" value="Log In" ]'))).send_keys(Keys.ENTER)

# klik di tombol login
login_webdriverwait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="submit" value="Log In" ]'))).click()

# klik di tombol follow
follow_webdriverwait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="button" value="Follow" ]'))).click()

# klik di tombol like
like_webdriverwait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="button" value="Like" ]'))).click()

# klik di tombol comment
comment_webdriverwait(driver, 10).until(EC.element_to_be_clickable((By.XPATH, '//input[type="text" value="" ]'))).send_keys('comment')
```

Preprocessing Data

```
!pip install Sastrawi
import re
import pandas as pd
from Sastrawi.Stemmer.stemmerfactory import StemmerFactory
```

```
D data_crawling = 'data.csv'
df = pd.read_csv(data_crawling)
# Menampilkan DataFrame
print(df)
```

```
] # case folding
df['comment'] = df['comment'].str.lower()
print('case folding result : \n')
print(df['comment'])
print('\n\n')
```

```
def clean_text(text):
    #cleaning
    # Remove tab, new line, and backslash
    text = re.sub(r'[\t|\n|\\]', ' ', text)
    # Remove non-ASCII characters
    text = re.sub(r'[^\x00-\x7f]+', '', text)
    # Remove mention and username, link, and hashtag
    text = re.sub(r'(@[\w]+|https?://\w+|#[\w]+)', '', text)
    # remove title
    text = re.sub(r'^(http://.*|)', '', text)
    # remove incomplete url
    text = re.sub(r'^(https?://\w+|)', '', text)
    # Remove numbers
    text = re.sub(r'^\d+', '', text)
    # remove punctuation
    text = re.sub(r'[^\w\s]', '', text)
    # Remove leading and trailing white space
    text = text.strip()
    # Remove multiple whitespace into single whitespace
    text = re.sub(r'\s+', ' ', text)
    # remove single character
    text = re.sub(r'^[\w\s]{1}$', '', text)
    return text
```

```
# Apply this origin into full data using pandas apply() method
df = df[df['comment'].notna() & df['comment'] != '']

# Cleaning text
df['cleaning'] = df['comment'].apply(lambda x: clean_text(x))

print('Cleaning text \n')
print(df['cleaning'].head())
```

```
factory = StemmerFactory()
stemmer = factory.create_stemmer()
def stem_text(text):
    return stemmer.stem(text)
df['stemming'] = df['cleaning'].apply(stem_text)
print('Stemming text:\n')
print(f'{df["stemming"].head()}\n\n')
```

```
from Sastrawi.StopwordRemover.StopwordRemoverFactory import StopwordRemoverFactory
factory = StopwordRemoverFactory()
stopword = factory.create_stop_word_remover()
```

```
def remove_stopwords(text):
    return stopword.remove(text)
df['stopwords'] = df['stemming'].apply(remove_stopwords)
print('Stopwords Text:\n')
print(f'{df["stopwords"].head()}\n\n')
```

```

# Load data slangswords from file.csv
df_slangswords = pd.read_csv("slangswords_data.csv")

# Load data stopwords from dataset
df_stopwords = df[df["id"] != "1"]

# Fungsi untuk menghapus slangswords dalam teks
def remove_slangswords(text):
    words = text.split()
    normalized_words = [slangswords[slang_word] if slang_word in slangswords else slang_word for slang_word in words]
    return " ".join(normalized_words)

# Aplikasikan fungsi pada kolom "slangswords"
df["slangswords"] = df["slangswords"].apply(remove_slangswords)

# Tampilkan hasil
df.head()

```

Translate Data

```

1.1.1. Install translate
!pip install translate

1.1.2. Load data from pd
from translate import Translator

# Load data from file.csv
df = pd.read_csv("slangswords_data.csv")

def translate_to_english(text):
    if isinstance(text, str):
        translator = Translator(to_lang="en", from_lang="id")
        translation = translator.translate(text)
        return translation
    else:
        return None

df["english_translation"] = df["slangswords"].apply(translate_to_english)
df["english_translation"].to_csv("translate_data.csv", index=False)

1.1.3. Load data from pd
# Load data from file.csv
df = pd.read_csv("slangswords_data.csv")

# Create a new column for english translation
df["english_translation"] = df["slangswords"].apply(translate_to_english)

# Save the data to a new CSV file
df.to_csv("translate_data.csv", index=False)

1.1.4. Filter data
df_filtered = df[df["english_translation"] != None]

# Save the filtered data to a new CSV file
df_filtered.to_csv("filtered_translate.csv", index=False)

```

Labelling Data

```

# Read CSV
df = pd.read_csv("filtered_translate.csv")
# Print "slangswords" column
print(df["English_Translation"])

!pip install vaderSentiment
from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer

# Create an analyzer
analyzer = SentimentIntensityAnalyzer()

# Print the score of each comment
df["Compound_Score"] = df["English_Translation"].apply(lambda x: analyzer.polarity_scores(x))
print(df)

df["Compound_Score"] = [x["compound"] for x in df["Compound_Score"].values]

```

```
hasilAnalisis = []
for score in df['compound_score']:
    sentimen = {}
    if score < 0:
        sentimen['sentiments'] = 'negatif'
        hasilAnalisis.append(sentimen)
    elif score > 0:
        sentimen['sentiments'] = 'positif'
        hasilAnalisis.append(sentimen)
    else:
        # Jika score tepat 0
        sentimen['sentiments'] = 'netral'
        hasilAnalisis.append(sentimen)

# membuat df baru dari hasilAnalisis
result_df = pd.DataFrame(hasilAnalisis)

# menggabungkan dataframe asli dengan dataframe hasil
df = pd.concat([df, result_df], axis=1)

df.head()
```

Konversi sentiment kedalam data index

```
column_to_delete = 'compound_score'
df = df.drop(column_to_delete, axis=1)

labels = df.Sentiments.unique()

labels_dict = {}
for index, labels in enumerate(labels):
    labels_dict[labels] = index

df['labels'] = df.Sentiments.replace(labels_dict)
print(df)
```

Split Data

```
from sklearn.model_selection import train_test_split

x_train, x_val, y_train, y_val = train_test_split(df.index.values,
                                                df.labels.values,
                                                test_size=0.10,
                                                random_state=2020,
                                                stratify=df.labels.values)

df['data_type'] = ['not_set'] * df.shape[0]
df.loc[x_train, 'data_type'] = 'train'
df.loc[x_val, 'data_type'] = 'val'
df.groupby(['Sentiments', 'labels', 'data_type']).count()
```

BERT Tokenizer dan Encode Data

```
!pip install transformers
```

```
from transformers import BertTokenizer
from torch.utils.data import TensorDataset

tokenizer = BertTokenizer.from_pretrained('bert-base-multilingual-uncased',
                                         do_lower_case=True)
```

```
# Encode sentences dan indexes ke dalam list untuk setiap kalimat
tokenized_sentences = [tokenizer.tokenize(sentence) for sentence in sentences]
token_ids_list = [tokenizer.convert_tokens_to_ids(tokenized_sentence) for tokenized_sentence in tokenized_sentences]

# Mendefinisikan dataset sederhana dari sekumpulan data
max_sentence_length = max(len(token_ids) for token_ids in token_ids_list)

# Menyiapkan embeddings dan embeddings()
encoder_embeddings = torch.nn.LstmEmbedder(max_sentence_length)
encoder_embeddings.requires_grad_()
encoder_embeddings.eval()

from torch.nn import LstmEmbedder, LstmCell
print(encoder_embeddings)
print(encoder_embeddings.eval())
```

```
#encode training data
encoder_data_loader = torch.utils.data.DataLoader(
    dataset=train_loader_loader(torchvision.transforms.Compose(
        [transforms.Resize(224),
        transforms.CenterCrop(224),
        transforms.ToTensor(),
        transforms.Normalize(mean=[0.485, 0.456, 0.605],
                             std=[0.229, 0.224, 0.225])])),
    batch_size=batch_size,
    shuffle=True,
    num_workers=4,
    pin_memory=True)

#encode validation data
encoder_data_loader_val = torch.utils.data.DataLoader(
    dataset=val_loader_loader(torchvision.transforms.Compose(
        [transforms.Resize(224),
        transforms.CenterCrop(224),
        transforms.ToTensor(),
        transforms.Normalize(mean=[0.485, 0.456, 0.605],
                             std=[0.229, 0.224, 0.225])])),
    batch_size=batch_size_val,
    shuffle=False,
    num_workers=4,
    pin_memory=True)

# extract input_ids, attention_mask, and labels from training data
input_ids_train = encoder_data_loader.dataset.input_ids
attention_mask_train = encoder_data_loader.dataset.attention_mask
labels_train = torch.tensor([label for label in labels_loader_loader])

# extract input_ids, attention_mask, and labels from validation data
input_ids_val = encoder_data_loader_val.dataset.input_ids
attention_mask_val = encoder_data_loader_val.dataset.attention_mask
labels_val = torch.tensor([label for label in labels_loader_loader_val])
```

```
# create Pytorch DataLoader for training and validation data
dataset_train = ConcatDataset(input_ids_train, attention_mask_train, labels_train)
dataset_val = ConcatDataset(input_ids_val, attention_mask_val, labels_val)
```

BERT Pretrained Model

```
from transformers import BertForSequenceClassification
model = BertForSequenceClassification.from_pretrained('bert-base-multilingual-uncased',
                                                    num_labels=len(labels),
                                                    output_attentions=True,
                                                    output_hidden_states=True)
```

```
params = list(model.named_parameters())

print(f'No. model has {len(params)} different sized parameters')

print(f'--- embedding layer ---')
for p in params[0:4]:
    print(f'vocab {p[0]}, {p[1]}, {p[2]}, {p[3]}')

print(f'--- First Transformer ---')
for p in params[4:10]:
    print(f'layer {p[0]}, {p[1]}, {p[2]}, {p[3]}, {p[4]}, {p[5]}')

print(f'--- Output layer ---')
for p in params[11:13]:
    print(f'vocab {p[0]}, {p[1]}, {p[2]}')
```

```
from torch.utils.data import DataLoader, RandomSampler, SequentialSampler
```

```
# get train and validation loaders
train_loader = DataLoader(
    dataset=dataset_train,
    sampler=RandomSampler(dataset_train),
    batch_size=batch_size,
    num_workers=4,
    pin_memory=True)

val_loader = DataLoader(
    dataset=dataset_val,
    sampler=SequentialSampler(dataset_val),
    batch_size=batch_size_val,
    num_workers=4,
    pin_memory=True)
```

```
from torch.utils.data import DataLoader, RandomSampler, SequentialSampler
from torch.utils.data import ConcatDataset, RandomSampler, SequentialSampler

dataset_train = ConcatDataset(input_ids_train, attention_mask_train, labels_train)
dataset_val = ConcatDataset(input_ids_val, attention_mask_val, labels_val)

train_loader = DataLoader(
    dataset=dataset_train,
    sampler=RandomSampler(dataset_train),
    batch_size=batch_size,
    num_workers=4,
    pin_memory=True)
```

```

import numpy as np
from sklearn.metrics import f1_score, accuracy_score, precision_score, recall_score

def main():
    # Load data
    data = load_data()
    X_train, X_test, y_train, y_test = train_test_split(data, test_size=0.2, random_state=42)

    # Train model
    model = LogisticRegression()
    model.fit(X_train, y_train)

    # Evaluate model
    y_pred = model.predict(X_test)
    f1 = f1_score(y_test, y_pred)
    acc = accuracy_score(y_test, y_pred)
    prec = precision_score(y_test, y_pred)
    rec = recall_score(y_test, y_pred)

    print(f'F1 Score: {f1}')
    print(f'Accuracy: {acc}')
    print(f'Precision: {prec}')
    print(f'Recall: {rec}')

def load_data():
    # Example data loading logic
    data = {}
    # ... (data loading logic) ...
    return data

if __name__ == '__main__':
    main()
    
```

```

# Example code for tokenization and model training
import re
import nltk
from nltk.tokenize import word_tokenize
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression

# Tokenize text
def tokenize(text):
    return word_tokenize(text)

# Vectorize text
def vectorize(texts):
    vectorizer = TfidfVectorizer()
    return vectorizer.fit_transform(texts)

# Train model
def train_model(texts, labels):
    X = vectorize(texts)
    model = LogisticRegression()
    model.fit(X, labels)
    return model

# Evaluate model
def evaluate(model, texts, labels):
    X = vectorize(texts)
    y_pred = model.predict(X)
    return accuracy_score(y_pred, labels)

# Main function
def main():
    # Example data
    texts = ["This is a test sentence.", "Another test sentence."]
    labels = [0, 1]

    # Train model
    model = train_model(texts, labels)

    # Evaluate model
    accuracy = evaluate(model, texts, labels)
    print(f'Accuracy: {accuracy}')

if __name__ == '__main__':
    main()
    
```

Tokenize

```

import re
import numpy as np
import pandas as pd
# nltk
from nltk.stem import WordNetLemmatizer
# sklearn
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, classification_report

```

```

import pandas as pd
import nltk
from nltk.tokenize import word_tokenize

df = pd.read_csv('labelling.csv')

# Unduh model tokenisasi punkt
nltk.download('punkt')

column_name = 'English translation'

# tokenisasi
df[column_name] = df[column_name].apply(lambda x: word_tokenize(x))

# Tampilkan hasil tokenisasi
print(df[column_name])

```

TF-IDF

```

# Membuat TFIDFVectorizer dengan parameter
vectorizer = TfidfVectorizer(max_df=(1/3), max_features=50000)
# Membuat vectorizer
vectorizer.fit(df['English translation'])
# 100% of features words, dan vectorizer.get_features_names()

# Membuat dataframe untuk training data
df_train = pd.DataFrame()
df_train['English translation'] = x_train
df_train['sentiments'] = y_train

# Membuat dataframe untuk test data
df_test = pd.DataFrame()
df_test['English translation'] = x_test
df_test['sentiments'] = y_test

# Untuk melihat transformasi data menggunakan TFIDF vectorizer
> x_train_tfidf = vectorizer.transform(df_train['English translation'])
> x_test_tfidf = vectorizer.transform(df_test['English translation'])

# Untuk melihat hasil transformasi TF-IDF ke dalam DataFrame
df_tfidf = pd.DataFrame(x_train_tfidf.toarray(), columns=vectorizer.get_feature_names_out())

# Kita teras dengan weight rata-rata TF-IDF
count_tfidf = df_tfidf.count(axis=0)
top_terms = average_tfidf.sort_values(ascending=False)

print(count_tfidf)
print(top_terms.head(10))

```

Metode SVM

```

from sklearn.metrics import accuracy_score
from sklearn.svm import LinearSVC
X = np.array([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
y = LinearSVC(C)
svm = LinearSVC(C)
svm.fit(X_train, y_train)
print("Akurasi SVM Kernel Linear = %d%%" % (accuracy_score(y_val, svm.predict(x_val_test))))

```

```

from sklearn.svm import SVC
model = SVC(kernel='linear', C=1)
model.fit(x_train, y_train)

```

```

# Proses pengujian SVM 2 (1-10)
from sklearn.metrics import accuracy_score

# Melakukakan prediksi setelah pada data validasi
prediksi_val = model.predict(x_val_test)

test_predictor = pd.DataFrame()
test_predictor["Treal"] = x_val
test_predictor["Predictions"] = prediksi_val

# Menghitung akurasi prediksi
SVM_accuracy_TFIDF = accuracy_score(prediksi_val, y_val)*100
SVM_accuracy_TFIDF = round(SVM_accuracy_TFIDF, 1)

```

```

SVM Accuracy
SVM accuracy: 100.00

```

Confusion Matrix

```

# Accuracy, Precision, Recall, F1-score
from sklearn.metrics import classification_report
y_pred = svm.predict(x_val_test)
print(classification_report(y_val, y_pred, labels))

```

```

import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import confusion_matrix

# Menghitung confusion matrix
cm = confusion_matrix(y_val, y_pred, labels=labels)

# Menampilkan confusion matrix sebagai heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['0', '1', '2'], yticklabels=['0', '1', '2'])
plt.title('Confusion Matrix')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()

```