



UTS SEMESTER GANJIL TAHUN AKADEMIK 2022/2023

Mata Kuliah : Pembelajaran Mesin / Machine Learning
Dosen : Adevian Fairuz Pratama,S.ST.,M.Eng
Kelas : TI 3E, TI 3F

Deteksi Emosi Pengguna Twitter

Deteksi emosi merupakan salah satu permasalahan yang dihadapi pada **Natural Language Processing** (NLP). Alasannya diantaranya adalah kurangnya dataset berlabel untuk mengklasifikasi emosi berdasarkan data twitter. Selain itu, sifat dari data twitter yang dapat memiliki banyak label emosi (**multi-class**). Manusia memiliki berbagai emosi dan sulit untuk mengumpulkan data yang cukup untuk setiap emosi. Oleh karena itu, masalah ketidakseimbangan kelas akan muncul (**class imbalance**). Pada Ujian Tengah Semester (UTS) kali ini, Anda telah disediakan dataset teks twitter yang sudah memiliki label untuk beberapa kelas emosi. Tugas utama Anda adalah membuat model yang mumpuni untuk kebutuhan klasifikasi emosi berdasarkan teks.

Informasi Data

Dataset yang akan digunakan adalah ***tweet_emotion.csv**. Berikut merupakan informasi tentang dataset yang dapat membantu Anda.

- Total data: 40000 data
- Label emosi: anger, boredom, empty, enthusiasm, fun, happiness, hate, love, neutral, relief, sadness, surprise, worry
- Jumlah data untuk setiap label tidak sama (**class imbalance**)
- Terdapat 3 kolom = 'tweet_id', 'sentiment', 'content'

Penilaian UTS

UTS akan dinilai berdasarkan 4 proses yang akan Anda lakukan, yaitu pra pengolahan data, ekstraksi fitur, pembuatan model machine learning, dan evaluasi.

Pra Pengolahan Data

Perhatian

Sebelum Anda melakukan sesuatu terhadap data Anda, pastikan data yang Anda miliki sudah "baik", bebas dari data yang hilang, menggunakan tipe data yang sesuai, dan sebagainya.

Data tweeter yang ada didapatkan merupakan sebuah data mentah, maka beberapa hal dapat Anda lakukan (namun tidak terbatas pada) yaitu,

1. Case Folding
2. Tokenizing
3. Filtering
4. Stemming

CATATAN: PADA DATA TWITTER TERDAPAT MENTION (@something) YANG ANDA HARUS TANGANI SEBELUM MASUK KE TAHAP EKSTRAKSI FITUR

Ekstraksi Fitur

Anda dapat menggunakan beberapa metode, diantaranya

1. Bag of Words (Count / TF-IDF)
2. N-gram
3. dan sebagainya

Pembuatan Model

Anda dibebaskan dalam memilih algoritma klasifikasi. Anda dapat menggunakan algoritma yang telah diajarkan didalam kelas atau yang lain, namun dengan catatan. Berdasarkan asas akuntabilitas pada pengembangan model machine learning, Anda harus dapat menjelaskan bagaimana model Anda dapat menghasilkan nilai tertentu.

Evaluasi

Pada proses evaluasi, minimal Anda harus menggunakan metric akurasi. Akan tetapi Anda juga dapat menambahkan metric lain seperti Recall, Precision, F1-Score, detail Confusion Metric, ataupun Area Under Curve (AUC).

Lembar Pengerjaan

Lembar pengerjaan dimulai dari cell dibawah ini

```
import numpy as np
import pandas as pd

df = pd.read_csv('data/tweet_emotions.csv')
df.head()
```

	tweet_id	sentiment	content
0	1956967341	empty	@tiffanylue i know i was listenin to bad habi...
1	1956967666	sadness	Layin n bed with a headache ughhhh...waitin o...
2	1956967696	sadness	Funeral ceremony...gloomy friday...
3	1956967789	enthusiasm	wants to hang out with friends SOON!
4	1956968416	neutral	@dannycastillo We want to trade with someone w...