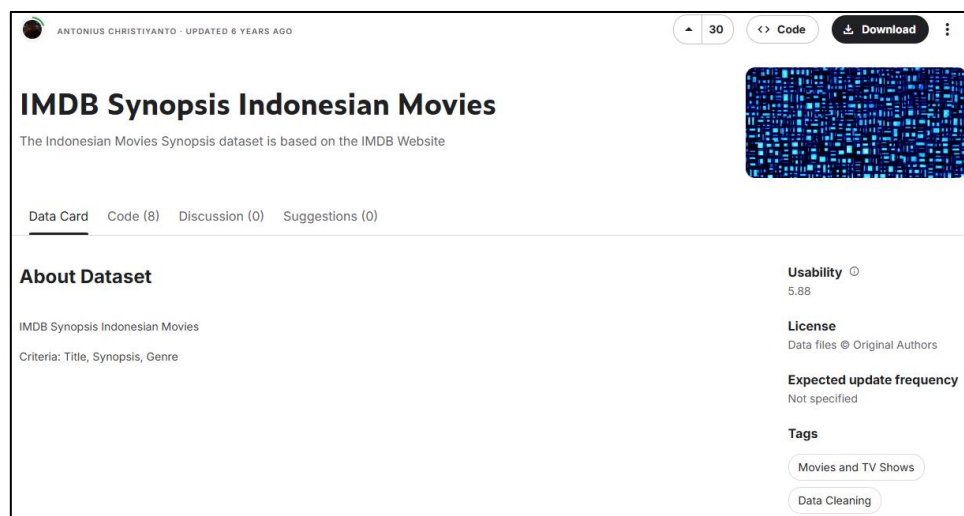


Nama Dosen : Teguh Iman Hermanto, M.Kom
Mata Kuliah : Machine Learning 2
Pembahasan : Basic NLP
Pokok Pemb : Membuat model Natural Language Processing Sederhana

DOWNLOAD DATASET

Import dataset ke dalam google colab menggunakan Alamat dataset berikut (import dataset menggunakan kaggle.json):

<https://www.kaggle.com/datasets/antoniusscs/imdb-synopsis-indonesian-movies>



IMPORT LIBRARY DAN LOAD DATASET

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from tensorflow.keras.preprocessing.text import Tokenizer
5 from tensorflow.keras.preprocessing.sequence import pad_sequences
6 import tensorflow as tf
7 from tensorflow.keras.utils import plot_model
```

```
1 df = pd.read_csv('imdb_indonesian_movies_2.csv')
```

PREPROCESSING DATA



```
1 df = df.drop(columns=['judul_film'])
```



```
1 df.head()
```

	ringkasan_sinopsis	genre
0	Raden Mas Said putra sulung Tumenggung Wilarik...	Drama
1	Soe Hok Gie adalah seorang aktivis yang hidup ...	Drama
2	Guru Bangsa Tjokroaminoto menceritakan tentang...	Drama
3	POL menceritakan kisah hidup yang luar biasa d...	Drama
4	Perjalanan pahlawan Indonesia KH Ahmad Dahlan ...	Drama



```
1 category = pd.get_dummies(df.genre)
2 df_baru = pd.concat([df, category], axis=1)
3 df_baru = df_baru.drop(columns='genre')
4 df_baru
```

	ringkasan_sinopsis	Drama	Horor	Komedi	Laga	Romantis
0	Raden Mas Said putra sulung Tumenggung Wilarik...	True	False	False	False	False
1	Soe Hok Gie adalah seorang aktivis yang hidup ...	True	False	False	False	False
2	Guru Bangsa Tjokroaminoto menceritakan tentang...	True	False	False	False	False
3	POL menceritakan kisah hidup yang luar biasa d...	True	False	False	False	False
4	Perjalanan pahlawan Indonesia KH Ahmad Dahlan ...	True	False	False	False	False
...
1000	Winter in Tokyo berpusat pada kehidupan Ishida...	False	False	False	False	True
1001	Markonah melarikan diri ke Jakarta karena akan...	False	False	False	False	True
1002	Tempat aking lebih dari 36 jam, Last Night ada...	False	False	False	False	True
1003	Proyek baru ini adalah tentang seorang lelaki ...	False	False	False	False	True
1004	Atika (Meriam Bellina) mantan penyanyi tenar, ...	False	False	False	False	True

1005 rows x 6 columns



```
1 sinopsis = df_baru['ringkasan_sinopsis'].values
2 label = df_baru[['Drama', 'Horor', 'Komedi', 'Laga', 'Romantis']].values
```



```
1 sinopsis_latih, sinopsis_test, label_latih, label_test = train_test_split(sinopsis, label, test_size=0.2)
```



```
1 tokenizer = Tokenizer(num_words=5000, oov_token='x')
2 tokenizer.fit_on_texts(sinopsis_latih)
3 tokenizer.fit_on_texts(sinopsis_test)
4
5 sekuens_latih = tokenizer.texts_to_sequences(sinopsis_latih)
6 sekuens_test = tokenizer.texts_to_sequences(sinopsis_test)
7
8 padded_latih = pad_sequences(sekuens_latih)
9 padded_test = pad_sequences(sekuens_test)
```

MODELING



```

1 input_length = padded_latih.shape[1]
2
3 model = tf.keras.Sequential([
4     tf.keras.layers.Embedding(input_dim=5000, output_dim=16),
5     tf.keras.layers.LSTM(64),
6     tf.keras.layers.Dense(128, activation='relu'),
7     tf.keras.layers.Dense(64, activation='relu'),
8     tf.keras.layers.Dense(5, activation='softmax')
9 ])

```



```

1 model.build(input_shape=(None, input_length))

```



```

1 model.compile(loss='categorical_crossentropy',
2               optimizer='adam',
3               metrics=['accuracy'])

```



```

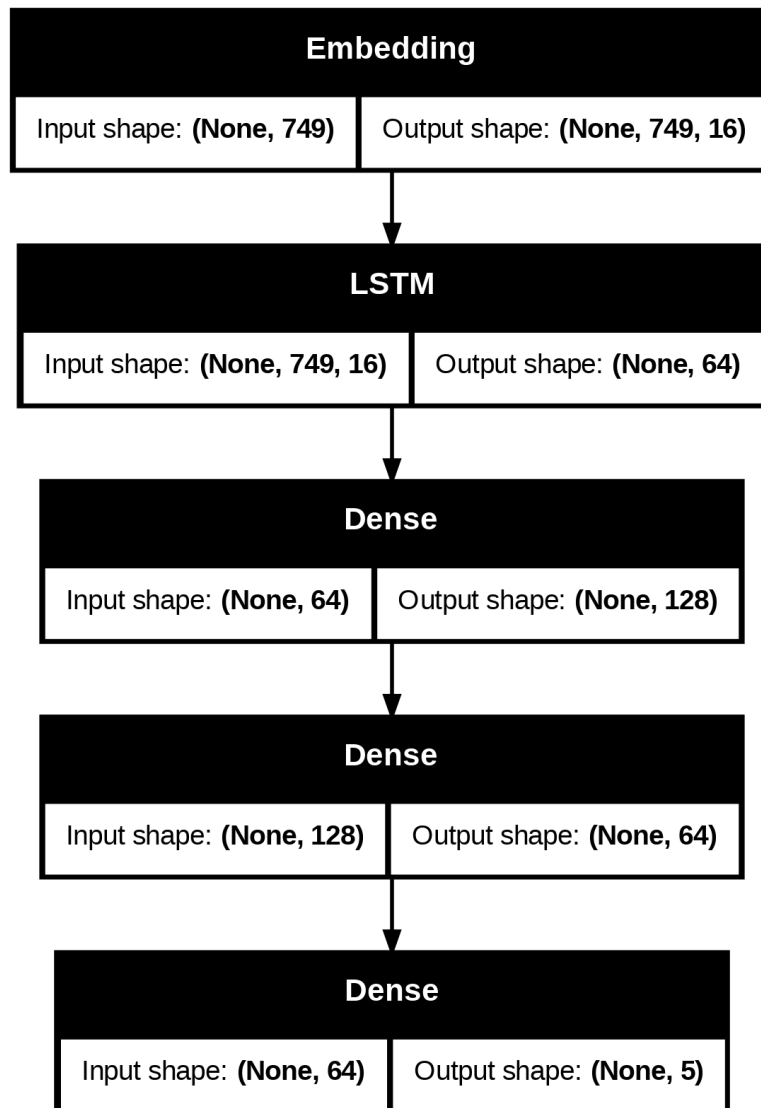
1 model.summary()

```

Layer (type)	Output Shape	Param #
embedding_6 (Embedding)	(None, 749, 16)	80,000
lstm_6 (LSTM)	(None, 64)	20,736
dense_17 (Dense)	(None, 128)	8,320
dense_18 (Dense)	(None, 64)	8,256
dense_19 (Dense)	(None, 5)	325



```
1 plot_model(model, show_shapes = True)
```



```
1 history = model.fit(padded_latih, label_latih, epochs=15,  
2                     validation_data=(padded_test, label_test))
```

SIMULASI MODEL

```
1  hasil = ["Kisah cinta antara dua yuppies urban yang sinis yang harus saling berhadapan sebagai pengacara di ruang sidang"]
2  hasil = tokenizer.texts_to_sequences(hasil)
3  hasil = pad_sequences(hasil)
4  #print(model.predict(hasil))
5  hasil = np.argmax(model.predict(hasil), axis=1)
6  if hasil == [0]:
7      print('Drama')
8  if hasil == [1]:
9      print('Horor')
10 if hasil == [2]:
11     print('Komedi')
12 if hasil == [3]:
13     print('Laga')
14 if hasil == [4]:
15     print('Romantis')
```