


```
import numpy as np
from tensorflow.keras.layers
from tensorflow.keras.models
from tensorflow.keras.prepro
from tensorflow.keras.prepro
from sklearn.model_selection
from datasets import load_da
from keras.layers import Den
# from keras.utils.vis_utils
from tensorflow.keras.utils

# from transformers import Be
from tensorflow import keras
from sklearn.metrics import

# Load Train data
train_df = pd.read_parquet('
train_df.head(10)
```

- 0 - Scope 3: Optional s
- 1 The Group is not awa
- 2 Global climate change
- 3 Setting an investmen
- 4 Climate change the ph
- 5 Projects with potent
- 6 We emitted 13.4 million
- 7 We do not provide nor
- 8 We anticipate that
- 9 Enhancing our respon

train_df.describe()

	label
count	1000.000000
mean	0.908000
std	0.764278
min	0.000000
25%	0.000000
50%	1.000000
75%	1.250000
max	2.000000

```
train_df_data = train_df['to
train_df_label = train_df['l

# Split data
train_data, val_data, train_

# load test data
test_df = pd.read_parquet('
test_df.head(10)
```

- 0 Sustainable strategies
- 1 Verizon's environmental
- 2 In 2019, the Company
- 3 In December 2020, the
- 4 Finally, there is
- 5 Ecoefficiency Eco-effici

```
test_data = test_df['text']
test_label = test_df['label']
```

```
    Although it is interesting
print("train_label : ", len(train_label))
print("train_label : ", len(train_label))
```

```
print("val_label : ", len(val_label))
print("val_label : ", len(val_label))
```

```
print("test_data : ", len(test_data))
print("test_label : ", len(test_label))
```

```
train_label : 800
train_label : 800
val_label : 200
val_label : 200
test_data : 320
test_label : 320
```

```
# Preprocess & Tokenize
MAX_WORDS = 10000
tokenizer = Tokenizer(num_words=MAX_WORDS)
tokenizer.fit_on_texts(texts)
```

```
train_sequences = tokenizer.texts_to_sequences(train_texts)
val_sequences = tokenizer.texts_to_sequences(val_texts)
test_sequences = tokenizer.texts_to_sequences(test_texts)
```

```
train_label = np.array(train_label)
val_label = np.array(val_label)
test_label = np.array(test_label)
```

```
# Tokenize
train_data_tokenized = pad_sequences(train_sequences, maxlen=MAX_WORDS)
val_data_tokenized = pad_sequences(val_sequences, maxlen=MAX_WORDS)
test_data_tokenized = pad_sequences(test_sequences, maxlen=MAX_WORDS)
```

```
# Cast into numpy array
train_data_tokenized = np.array(train_data_tokenized)
val_data_tokenized = np.array(val_data_tokenized)
test_data_tokenized = np.array(test_data_tokenized)
```

RNN/LSTM

MODEL

```
# Define Model
# Hyper parameter sama dengan model sebelumnya
model_rnn = Sequential()
model_rnn.add(Embedding(input_dim=MAX_WORDS, output_dim=EMBEDDING_DIM))
model_rnn.add(Bidirectional(LSTM(LSTM_DIM, return_sequences=True)))
model_rnn.add(Bidirectional(LSTM(LSTM_DIM, return_sequences=True)))
model_rnn.add(Dense(1, activation='sigmoid'))
```

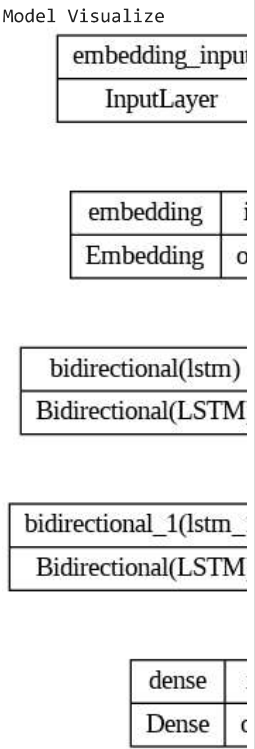
```
#compile model
model_rnn.compile(loss='binary_crossentropy', optimizer='adam')
print(model_rnn.summary())
```

```
print("\n\nModel Visualize")
plot_model(model_rnn, to_file="model_visualization.png")
```

Model: "sequential"

Layer (type)
embedding (Embedding)
bidirectional (Bidirectional)
bidirectional_1 (Bidirectional)
dense (Dense)

=====
Total params: 1420097 (1.42M)
Trainable params: 1420097 (1.42M)
Non-trainable params: 0 (0.00M)
=====
None



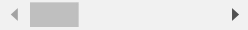
```
# Train
model_rnn.fit(train_data_token_embeddings, train_data_targets,
               validation_data=(validation_data_token_embeddings, validation_data_targets),
               epochs=10,
               callbacks=[HijackCallback()])
```

Epoch 1/10
25/25 [=====]
Epoch 2/10
25/25 [=====]
Epoch 3/10
25/25 [=====]
Epoch 4/10
25/25 [=====]
Epoch 5/10
25/25 [=====]
Epoch 6/10
25/25 [=====]
Epoch 7/10
25/25 [=====]
Epoch 8/10
25/25 [=====]
Epoch 9/10
25/25 [=====]
Epoch 10/10
25/25 [=====]
<keras.src.callbacks.HijackCallback at 0x79e925fa7160>

```
# Evaluate
```

```
loss, acc = model_rnn.evaluate
print("loss: ", loss)
print("accuracy: ", acc)
```

```
10/10 [=====]
loss: -1.1683006286621
accuracy: 0.6531249886
```



```
# Prediction
```

```
prediction = model_rnn.predict
```

```
for text, prediction, groundtruth in zip(
    sentiment = "positive" if prediction == 1 else "negative"
    groundtruth = "positive" if groundtruth == 1 else "negative"
    print(f"Text: {text} \n Predicted Sentiment: {prediction} \n Groundtruth: {groundtruth}")
```

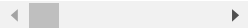
```
1/1 [=====]
Text: sustainable strategies for a better world
Predicted Sentiment: positive
Groundtruth: negative
```

```
Text: environmental health is a priority
Predicted Sentiment: positive
Groundtruth: negative
```

```
Text: in 2019 the company reported a record profit
Predicted Sentiment: positive
Groundtruth: negative
```

```
Text: which would normally be considered a positive sign
Predicted Sentiment: positive
Groundtruth: negative
```

```
Text: finally there is a chance for a better future
Predicted Sentiment: positive
Groundtruth: negative
```



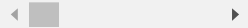
Word2Vec Embedding

```
from gensim.models import Word2Vec
```

```
word2vec_model = Word2Vec(sentences)
```

```
word2vec_model.save("word2vec_model.pkl")
```

```
WARNING:gensim.models.word2vec: Saving model to file
```



```
embedding_matrix = np.zeros((MAX_WORDS, EMBEDDING_DIM))
for word, i in tokenizer.word_index.items():
    if i < MAX_WORDS:
        if word in word2vec_model.wv:
            embedding_matrix[i] = word2vec_model.wv[word]
```

```
# define model
```

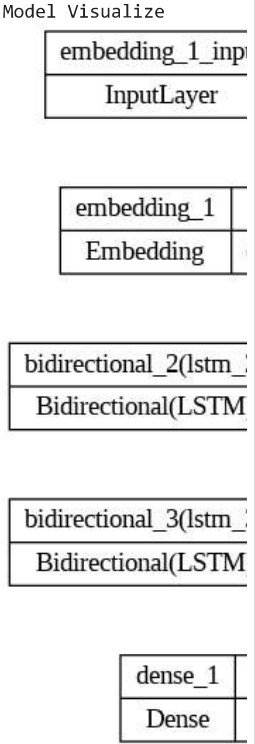
```
# Hyper parameter sama dengan model sebelumnya
word2vec_model = Sequential()
word2vec_model.add(Embedding(MAX_WORDS, EMBEDDING_DIM))
word2vec_model.add(Bidirectional([LSTM(EMBEDDING_DIM), LSTM(EMBEDDING_DIM)]))
word2vec_model.add(Bidirectional([LSTM(EMBEDDING_DIM), LSTM(EMBEDDING_DIM)]))
word2vec_model.add(Dense(1, activation='sigmoid'))
```

```
#compile model
word2vec_model.compile(loss=
print(word2vec_model.summary
print("\n\nModel Visualize")
plot_model(word2vec_model, to
```

Model: "sequential_1"

Layer (type)
embedding_1 (Embedding)
bidirectional_2 (Bidirectional)
bidirectional_3 (Bidirectional)
dense_1 (Dense)

=====
Total params: 1420097 (
Trainable params: 14200
Non-trainable params: 0
=====
None



```
# Train
word2vec_model.fit(train_data

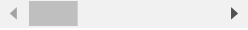
Epoch 1/10
25/25 [=====
Epoch 2/10
25/25 [=====
Epoch 3/10
25/25 [=====
Epoch 4/10
25/25 [=====
Epoch 5/10
25/25 [=====
Epoch 6/10
25/25 [=====
Epoch 7/10
25/25 [=====
Epoch 8/10
25/25 [=====
Epoch 9/10
25/25 [=====
Epoch 10/10
25/25 [=====
```

```
<keras.src.callbacks.Hi
at 0x79e925284910>
```

```
# Evaluate
```

```
loss, acc = word2vec_model.e
print("loss: ", loss)
print("accuracy: ", acc)
```

```
10/10 [=====]
loss: -2.031978607177
accuracy: 0.6781250238
```



```
# Prediction
```

```
prediction = word2vec_model.p
```

```
for text, prediction, ground
    sentiment = "positive" if
    groundtruth = "positive"
    print(f"Text: {text} \n I
```

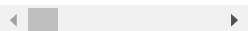
```
1/1 [=====]
Text: sustainable strat
Predicted Sentiment: p
Groundtruth: negative
```

```
Text: environmental hea
Predicted Sentiment: p
Groundtruth: negative
```

```
Text: in 2019 the compa
Predicted Sentiment: p
Groundtruth: negative
```

```
Text: which would norma
Predicted Sentiment: r
Groundtruth: negative
```

```
Text: finally there is
Predicted Sentiment: r
Groundtruth: negative
```



Attention Based Model

```
import tensorflow as tf
import torch
from transformers import Auto
from transformers import Tra
```

```
from sklearn.model_selection
from datasets import load_da
```

```
from transformers import TFA
import transformers
from tensorflow.keras.layers
from tensorflow.keras.models
from tensorflow.keras.optimi
```

```
import random
```

```
from datasets import load_da
from datasets import Dataset
dataset_train_pd = pd.read_p
dataset_train = Dataset.from
```

```
# dataset_val_pd = pd.read_parquet('dataset_val.parquet')
# dataset_val = Dataset.from_pandas(dataset_val_pd)
dataset_test_pd = pd.read_parquet('dataset_test.parquet')
dataset_test = Dataset.from_pandas(dataset_test_pd)
```

```
0    - Scope 3: Optional sources
1    The Group is not aware of
2    Global climate change
3    Setting an investment
4    Climate change the ph
5    Projects with potent
6    We emitted 13.4 million
7    We do not provide nor
8    We anticipate that
9    Enhancing our respon
```

```
dataset_test_pd
```

```
0    Sustainable stra
1    Verizon's environm
2    In 2019, the Comp
3    In December 2020, tl
4    Finally, there
...
315    Indirect emission
316    All data in this
317    Outcome: The bank
318    In 2020, Banco do F
319    Climate change is pr
320 rows x 2 columns
```

```
import datasets
```

```
labels_train = np.array(data:
# labels_val = np.array(data:
labels_test = np.array(data:
```

```
my_dataset = datasets.Dataset
```

```
my_dataset
```

```
DatasetDict({
  train: Dataset({
    features:
    ['text', 'label'],
    num_rows: 10
  })
  test: Dataset({
    features:
    ['text', 'label'],
    num_rows: 320
  })
})
```



```
tokenizer = AutoTokenizer.from_pretrained('gpt2')

def tokenize_function(example):
    return tokenizer(example['text'])

tokenized_datasets = my_data_loader.map(tokenize_function, batched=True)

Map: 100%
Map: 100%

model = AutoModelForSequenceClassification.from_pretrained('gpt2', num_labels=10)

Some weights of the model were not initialized from the random state.
- This IS expected if you are initializing a model that was previously frozen and then doing uninitialized weights.
- This IS NOT expected to happen.
Some weights of BertForSequenceClassification were not initialized from the random state.
You should probably TRAIN the model on a down-stream task to be able to use it for sequence classification.

import numpy as np
import evaluate

metric = evaluate.load("accuracy")

def compute_metrics(eval_preds):
    logits, labels = eval_preds
    predictions = np.argmax(logits, axis=-1)
    return metric.compute(predictions=predictions, references=labels)

from transformers import GradientDescentWrapper, GradientDescentWrapperTrainer
# from transformers import AutoModelForSequenceClassification
training_args = TrainingArguments(output_dir='./results', num_train_epochs=10, eval_steps=100, save_steps=100, logging_dir='./logs', logging_level='error', logging_first_step=True, logging_steps=10, save_total_limit=1, report_to='none', disable_tqdm=True)

trainer = GradientDescentWrapperTrainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_datasets['train'],
    eval_dataset=tokenized_datasets['validation'],
    compute_metrics=compute_metrics,
)

trainer.train()

# Eval
result = trainer.evaluate()

result

# make prediction
result = trainer.predict(tokenized_datasets['validation'])

for i in range(10):
    print(f"Text: {tokenized_datasets['validation'][i]['text']}")
    print(f"label: {tokenized_datasets['validation'][i]['label']}")
    print(f"Prediction: {np.argmax(result[i], axis=-1)}")
```

Report

Performance

- LSTM MODEL

Training Accuracy :
0.71

Test Accuracy : 0.70

- Word2Vec Embedding

Training Accuracy :
0.70

Test Accuracy : 0.64

- Attention - Based
Model

Training Accuracy :
0.90

Test Accuracy :

Reference :

- <https://huggingface.co/distilbert-base-uncased>
- https://www.tensorflow.org/text/tutorials/classify_text_with_bert

Double-click (or enter) to edit

Change runtime type

1m 44s completed at 11:21 PM

