

LSTM_without_preprocessing

```
In [2]: ## LSTM without preprocessing
import pandas as pd
import numpy as np
import torch
from torch import cuda
import random
import os
import re
import string
import torch
from torch import nn
from collections import Counter
from sklearn.model_selection import train_test_split
from sklearn.metrics import f1_score, classification_report
from sklearn.utils.class_weight import compute_class_weight
from matplotlib import pyplot as plt
print(torch.__version__)

1.10.2
```

```
In [4]: # data import
#Training data
train_data = pd.read_csv(r"C:\Users\Sunny\Downloads\train.csv")
#test data import
test_data = pd.read_csv(r"C:\Users\Sunny\Downloads\test.csv")
```

```
In [6]: #Print out the columns name and remove the header
print(train_data.columns)
train_data.head(10)
```

Out[6]:

	id	keyword	location	text	target
0	1	NaN	NaN	Our Deeds are the Reason of this #earthquake M...	1
1	4	NaN	NaN	Forest fire near La Ronge Sask. Canada	1
2	5	NaN	NaN	All residents asked to 'shelter in place' are ...	1
3	6	NaN	NaN	13,000 people receive #wildfires evacuation or...	1
4	7	NaN	NaN	Just got sent this photo from Ruby #Alaska as ...	1
5	8	NaN	NaN	#RockyFire Update => California Hwy. 20 closed...	1
6	10	NaN	NaN	#flood #disaster Heavy rain causes flash flood...	1
7	13	NaN	NaN	I'm on top of the hill and I can see a fire in...	1
8	14	NaN	NaN	There's an emergency evacuation happening now ...	1
9	15	NaN	NaN	I'm afraid that the tornado is coming to our a...	1

```
In [7]: #create Target
train_label = train_data.target.to_numpy()
y_train = train_label.reshape(-1)
train_data.drop(['id', 'keyword', 'location', 'target'], axis=1, inplace=True)
test_data.drop(['id', 'keyword', 'location'], axis=1, inplace=True)
train_data = train_data.to_numpy()
test_data = test_data.to_numpy()
x_train = train_data.reshape(-1)
x_test = test_data.reshape(-1)
print(x_train)
print(y_train)

['Our Deeds are the Reason of this #earthquake May ALLAH Forgive us all'
'Forest fire near La Ronge Sask. Canada'
"All residents asked to 'shelter in place' are being notified by officers. No other evacuation or shelter in
place orders are expected"
... 'M1.94 [01:04 UTC]?5km S of Volcano Hawaii. http://t.co/zDtoy8EbJ'
'Police investigating after an e-bike collided with a car in Little Portugal. E-bike rider suffered serious n
on-life threatening injuries.'
'The Latest: More Homes Razed by Northern California Wildfire - ABC News http://t.co/YmY4rSkQ3d']
[1 1 1 ... 1 1 1]
```

```
In [8]: #Model constant
max_features = 20000
embedding_dim = 128
max_len = 500
from keras.layers import*
vectorize_layer = TextVectorization(
    max_tokens=max_features,
    output_mode='int',
    output_sequence_length=max_len)
vectorize_layer.adapt(train_data)
import tensorflow as tf
def vectorize_text(text, label):
    text = tf.expand_dims(text, -1)
    return vectorize_layer(text), label
```

```
In [12]: # Start to build bidirection LSTM model
from keras.metrics import AUC
from keras import optimizers, losses
def build_model():

    text_input = tf.keras.Input(shape=(1,), dtype=tf.string, name='text')
    x = vectorize_layer(text_input)
    x = Embedding(max_features + 1, embedding_dim)(x)
    x = Dropout(0.5)(x)

    x1 = Bidirectional(LSTM(units=128, return_sequences=True))(x)
    x2 = Bidirectional(LSTM(units=64, return_sequences=True))(x1)

    z2 = Bidirectional(GRU(units=128, return_sequences=True))(x)
    z3 = Bidirectional(GRU(units=64, return_sequences=True))(z2)
    c = Concatenate(axis=2)([x2, z3])

    x3 = Bidirectional(LSTM(units=128, return_sequences=True))(c)
    # Conv1D + global max pooling

    x = GlobalMaxPooling1D()(x3)

    # We add a vanilla hidden layer:
    x = Dense(128, activation="relu")(x)
    x = Dropout(0.5)(x)

    # We project onto a single unit output layer, and squash it with a sigmoid:
    predictions = Dense(1, activation="sigmoid", name="predictions")(x)

    model = tf.keras.Model(text_input, predictions)

    # Compile the model with binary crossentropy loss and an adam optimizer.

    model.compile(loss=losses.binary_crossentropy, optimizer=optimizers.rmsprop_v2.RMSprop(learning_rate=0.01)
    #model.compile(loss=losses.binary_crossentropy, optimizer=optimizers.rmsprop_v2.RMSprop(learning_rate=0.01)
    return model
```

```
In [13]: model = build_model()
model.summary()
```

Model: "model"			
Layer (type)	Output Shape	Param #	Connected to
=====			
text (InputLayer)	[(None, 1)]	0	[]
text_vectorization (TextVectorization)	(None, 50)	0	['text[0][0]']
embedding (Embedding)	(None, 50, 128)	2560128	['text_vectorization[0][0]']
dropout (Dropout)	(None, 50, 128)	0	['embedding[0][0]']
bidirectional (Bidirectional)	(None, 50, 256)	263168	['dropout[0][0]']
bidirectional_2 (Bidirectional)	(None, 50, 256)	198144	['dropout[0][0]']
bidirectional_1 (Bidirectional)	(None, 50, 128)	164352	['bidirectional[0][0]']
bidirectional_3 (Bidirectional)	(None, 50, 128)	123648	['bidirectional_2[0][0]']
concatenate (Concatenate)	(None, 50, 256)	0	['bidirectional_1[0][0]', 'bidirectional_3[0][0]']
bidirectional_4 (Bidirectional)	(None, 50, 256)	394240	['concatenate[0][0]']
global_max_pooling1d (GlobalMaxPooling1D)	(None, 256)	0	['bidirectional_4[0][0]']
dense (Dense)	(None, 128)	32896	['global_max_pooling1d[0][0]']
dropout_1 (Dropout)	(None, 128)	0	['dense[0][0]']
predictions (Dense)	(None, 1)	129	['dropout_1[0][0]']
=====			
Total params: 3,736,705			
Trainable params: 3,736,705			
Non-trainable params: 0			
=====			

```
In [14]: #Split the training set into 2 part, train and valid
from sklearn.model_selection import train_test_split
x_train, x_valid, y_train, y_valid = train_test_split(x_train, y_train, test_size=0.2, random_state=1)
```

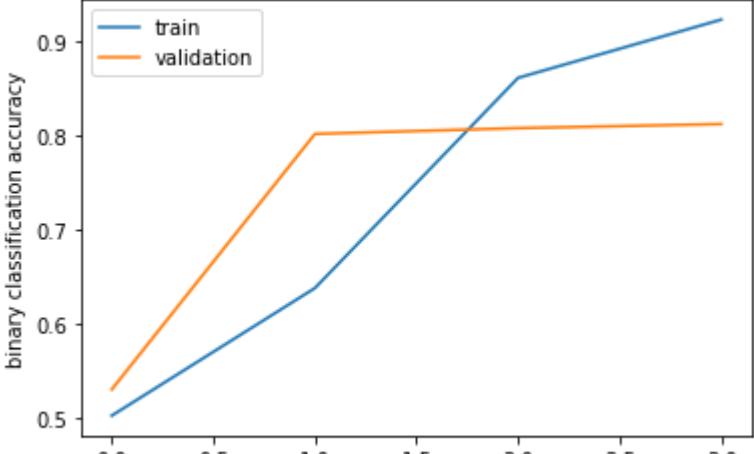
```
In [15]: from keras.callbacks import *

lr = ReduceLROnPlateau(monitor="val_auc", mode='max', factor=0.7, patience=4, verbose=False)
es = EarlyStopping(monitor='val_auc',mode='max', patience=10, verbose=False,restore_best_weights=True)

history = model.fit(x_train, y_train, validation_data=(x_valid, y_valid), epochs=4, batch_size=256,
                    callbacks=[es,lr])

Epoch 1/4
24/24 [=====] - 114s 4s/step - loss: 1.1843 - auc: 0.5026 - val_loss: 0.6870 - val_auc: 0.5305 - lr: 0.0100
Epoch 2/4
24/24 [=====] - 76s 3s/step - loss: 0.6593 - auc: 0.6383 - val_loss: 0.5318 - val_auc: 0.8024 - lr: 0.0100
Epoch 3/4
24/24 [=====] - 75s 3s/step - loss: 0.4565 - auc: 0.8620 - val_loss: 0.5612 - val_auc: 0.8085 - lr: 0.0100
Epoch 4/4
24/24 [=====] - 79s 3s/step - loss: 0.3396 - auc: 0.9240 - val_loss: 0.6868 - val_auc: 0.8129 - lr: 0.0100
```

```
In [17]: #Plot the performance
def plot_hist(hist, metric='auc'):
    plt.plot(hist.history[metric])
    plt.plot(hist.history['val_ ' + metric])
    plt.title(f"model performance")
    plt.ylabel("binary classification accuracy")
    plt.xlabel("epoch")
    plt.legend(["train", "validation"], loc="upper left")
    plt.show()
    return
plot_hist(hist=history)
```



#LSTM with preprocessing

[illegible]

[illegible]

```

auc: 0.7209 - lr: 0.0100
Epoch 4/4
24/24 [=====] - 14s 570ms/step - loss: 0.4733 - auc: 0.8393 - val_loss: 0.5404 - val_
auc: 0.7919 - lr: 0.0100

In [40]:
#visualization of the result
def plot_hist(hist, metrics='auc'):
    plt.plot(hist.history[metric])
    plt.plot(hist.history["val_" + metric])
    plt.title("Model performance")
    plt.ylabel("Binary classification accuracy")
    plt.xlabel("epoch")
    plt.legend(["train", "validation"], loc="upper left")
    plt.show()
    return
plot_hist(hist=history)

```

epoch	train	validation
0.0	0.51	0.51
1.0	0.52	0.51
2.0	0.72	0.71
3.0	0.84	0.79