```python
import matplotlib.pyplot as plt
import os
import re
import shutil
import string
import tensorflow as tf

from tensorflow.keras import layers
from tensorflow.keras import losses

print(tf.__version__)

url = "https://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.tar.gz"

dataset = tf.keras.utils.get_file("aclImdb_v1", url,
                                    untar=True, cache_dir='.',
                                    cache_subdir='')

dataset_dir = os.path.join(os.path.dirname(dataset), 'aclImdb')

os.listdir(dataset_dir)

train_dir = os.path.join(dataset_dir, 'train')
os.listdir(train_dir)

sample_file = os.path.join(train_dir, 'pos/1181_9.txt')
with open(sample_file) as f:
  print(f.read())

remove_dir = os.path.join(train_dir, 'unsup')
shutil.rmtree(remove_dir)

batch_size = 32
seed = 42

raw_train_ds = tf.keras.utils.text_dataset_from_directory(
    'aclImdb/train',
    batch_size=batch_size,
    validation_split=0.2,
    subset='training',
    seed=seed)

for text_batch, label_batch in raw_train_ds.take(1):
  for i in range(3):
    print("Review", text_batch.numpy()[i])
    print("Label", label_batch.numpy()[i])

print("Label 0 corresponds to", raw_train_ds.class_names[0])
print("Label 1 corresponds to", raw_train_ds.class_names[1])

raw_val_ds = tf.keras.utils.text_dataset_from_directory(
    'aclImdb/train',
    batch_size=batch_size,
    validation_split=0.2,
    subset='validation',
    seed=seed)

raw_test_ds = tf.keras.utils.text_dataset_from_directory(
    'aclImdb/test',
    batch_size=batch_size)

def custom_standardization(input_data):
  lowercase = tf.strings.lower(input_data)
  stripped_html = tf.strings.regex_replace(lowercase, '<br />', ' ')
  return tf.strings.regex_replace(stripped_html,
                                   '[%s]' % re.escape(string.punctuation),
                                   '')

max_features = 10000
sequence_length = 250

vectorize_layer = layers.TextVectorization(
    standardize=custom_standardization,
    max_tokens=max_features,
    output_mode='int',
    output_sequence_length=sequence_length)
```

```python
train_text = raw_train_ds.map(lambda x, y: x)
vectorize_layer.adapt(train_text)

def vectorize_text(text, label):
  text = tf.expand_dims(text, -1)
  return vectorize_layer(text), label

text_batch, label_batch = next(iter(raw_train_ds))
first_review, first_label = text_batch[0], label_batch[0]
print("Review", first_review)
print("Label", raw_train_ds.class_names[first_label])
print("Vectorized review", vectorize_text(first_review, first_label))

print("1287 ---> ",vectorize_layer.get_vocabulary()[1287])
print(" 313 ---> ",vectorize_layer.get_vocabulary()[313])
print('Vocabulary size: {}'.format(len(vectorize_layer.get_vocabulary())))

train_ds = raw_train_ds.map(vectorize_text)
val_ds = raw_val_ds.map(vectorize_text)
test_ds = raw_test_ds.map(vectorize_text)

AUTOTUNE = tf.data.AUTOTUNE

train_ds = train_ds.cache().prefetch(buffer_size=AUTOTUNE)
val_ds = val_ds.cache().prefetch(buffer_size=AUTOTUNE)
test_ds = test_ds.cache().prefetch(buffer_size=AUTOTUNE)

embedding_dim = 16
model = tf.keras.Sequential([
  layers.Embedding(max_features, embedding_dim),
  layers.Dropout(0.2),
  layers.GlobalAveragePooling1D(),
  layers.Dropout(0.2),
  layers.Dense(1, activation='sigmoid')])

model.summary()

model.compile(loss=losses.BinaryCrossentropy(),
              optimizer='adam',
              metrics=[tf.metrics.BinaryAccuracy(threshold=0.5)])

epochs = 10
history = model.fit(
    train_ds,
    validation_data=val_ds,
    epochs=epochs)

loss, accuracy = model.evaluate(test_ds)

print("Loss: ", loss)
print("Accuracy: ", accuracy)

history_dict = history.history
history_dict.keys()

acc = history_dict['binary_accuracy']
val_acc = history_dict['val_binary_accuracy']
loss = history_dict['loss']
val_loss = history_dict['val_loss']

epochs = range(1, len(acc) + 1)


plt.plot(epochs, loss, 'bo', label='Training loss')

plt.plot(epochs, val_loss, 'b', label='Validation loss')
plt.title('Training and validation loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()

plt.show()

plt.plot(epochs, acc, 'bo', label='Training acc')
plt.plot(epochs, val_acc, 'b', label='Validation acc')
plt.title('Training and validation accuracy')
```

```python
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.legend(loc='lower right')

plt.show()

export_model = tf.keras.Sequential([
  vectorize_layer,
  model,
  layers.Activation('sigmoid')
])

export_model.compile(
    loss=losses.BinaryCrossentropy(from_logits=False), optimizer="adam", metrics=['accuracy']
)


metrics = export_model.evaluate(raw_test_ds, return_dict=True)
print(metrics)

examples = tf.constant([
  "The movie was great!",
  "The movie was okay.",
  "The movie was terrible..."
])

export_model.predict(examples)
```

2.17.1
Downloading data from https://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.tar.gz
84125825/84125825 ━━━━━━━━━━━━━━━━━━━━ **6s** 0us/step
Rachel Griffiths writes and directs this award winning short film. A heartwarming story about coping with grief and cherishing the memor
Found 25000 files belonging to 2 classes.
Using 20000 files for training.
Review b'"Pandemonium" is a horror movie spoof that comes off more stupid than funny. Believe me when I tell you, I love comedies. Espec
Label 0
Review b"David Mamet is a very interesting and a very un-equal director. His first movie 'House of Games' was the one I liked best, and
Label 0
Review b'Great documentary about the lives of NY firefighters during the worst terrorist attack of all time.. That reason alone is why t
Label 1
Label 0 corresponds to neg
Label 1 corresponds to pos
Found 25000 files belonging to 2 classes.
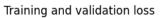Using 5000 files for validation.
Found 25000 files belonging to 2 classes.
Review tf.Tensor(b'Silent Night, Deadly Night 5 is the very last of the series, and like part 4, it\'s unrelated to the first three exce
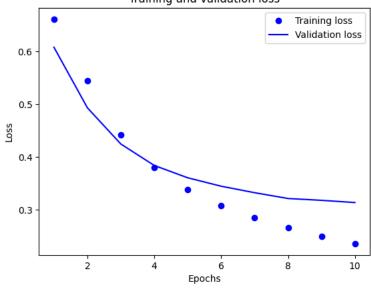Label neg
Vectorized review (<tf.Tensor: shape=(1, 250), dtype=int64, numpy=
array([[1287,  313, 2380,  313,  661,    7,    2,   52,  229,    5,    2,
         200,    3,   38,  170,  669,   29, 5492,    6,    2,   83,  297,
         549,   32,  410,    3,    2,  186,   12,   29,    4,    1,  191,
         510,  549,    6,    2, 8229,  212,   46,  576,  175,  168,   20,
           1, 5361,  290,    4,    1,  761,  969,    1,    3,   24,  935,
        2271,  393,    7,    1, 1675,    4, 3747,  250,  148,    4,  112,
         436,  761, 3529,  548,    4, 3633,   31,    2, 1331,   28, 2096,
           3, 2912,    9,    6,  163,    4, 1006,   20,    2,    1,   15,
          85,   53,  147,    9,  292,   89,  959, 2314,  984,   27,  762,
           6,  959,    9,  564,   18,    7, 2140,   32,   24, 1254,   36,
           1,   85,    3, 3298,   85,    6, 1410,    3, 1936,    2, 3408,
         301,  965,    7,    4,  112,  740, 1977,   12,    1, 2014, 2772,
           3,    4,  428,    3, 5177,    6,  512, 1254,    1,  278,   27,
         139,   25,  308,    1,  579,    5,  259, 3529,    7,   92, 8981,
          32,    2, 3842,  230,   27,  289,    9,   35,    2, 5712,   18,
          27,  144, 2166,   56,    6,   26,   46,  466, 2014,   27,   40,
        2745,  657,  212,    4, 1376, 3002, 7080,  183,   36,  180,   52,
         920,    8,    2, 4028,   12,  969,    1,  158,   71,   53,   67,
          85, 2754,    4,  734,   51,    1, 1611,  294,   85,    6,    2,
        1164,    6,  163,    4, 3408,   15,   85,    6,  717,   85,   44,
           5,   24, 7158,    3,   48,  604,    7,   11,  225,  384,   73,
          65,   21,  242,   18,   27,  120,  295,    6,   26,  667,  129,
        4028,  948,    6,   67,   48,  158,   93,    1]])>, <tf.Tensor: shape=(), dtype=int32, numpy=0>)
1287 ---> silent
 313 ---> night
Vocabulary size: 10000
**Model: "sequential"**

| Layer (type) | Output Shape | Param # |
|---|---|---|
| embedding (Embedding) | ? | 0 (unbuilt) |
| dropout (Dropout) | ? | 0 (unbuilt) |
| global_average_pooling1d (GlobalAveragePooling1D) | ? | 0 (unbuilt) |
| dropout_1 (Dropout) | ? | 0 (unbuilt) |
| dense (Dense) | ? | 0 (unbuilt) |

 **Total params:** 0 (0.00 B)
 **Trainable params:** 0 (0.00 B)
 **Non-trainable params:** 0 (0.00 B)
Epoch 1/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **7s** 10ms/step - binary_accuracy: 0.5813 - loss: 0.6805 - val_binary_accuracy: 0.7398 - val_loss: 0.6076
Epoch 2/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **8s** 6ms/step - binary_accuracy: 0.7619 - loss: 0.5742 - val_binary_accuracy: 0.8152 - val_loss: 0.4931
Epoch 3/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **5s** 6ms/step - binary_accuracy: 0.8287 - loss: 0.4617 - val_binary_accuracy: 0.8308 - val_loss: 0.4247
Epoch 4/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **6s** 7ms/step - binary_accuracy: 0.8529 - loss: 0.3907 - val_binary_accuracy: 0.8432 - val_loss: 0.3841
Epoch 5/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **6s** 9ms/step - binary_accuracy: 0.8691 - loss: 0.3457 - val_binary_accuracy: 0.8460 - val_loss: 0.3608
Epoch 6/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **5s** 9ms/step - binary_accuracy: 0.8828 - loss: 0.3133 - val_binary_accuracy: 0.8500 - val_loss: 0.3449
Epoch 7/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **9s** 6ms/step - binary_accuracy: 0.8911 - loss: 0.2885 - val_binary_accuracy: 0.8550 - val_loss: 0.3326
Epoch 8/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **6s** 8ms/step - binary_accuracy: 0.9007 - loss: 0.2678 - val_binary_accuracy: 0.8590 - val_loss: 0.3217
Epoch 9/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **10s** 7ms/step - binary_accuracy: 0.9074 - loss: 0.2522 - val_binary_accuracy: 0.8582 - val_loss: 0.3182
Epoch 10/10
625/625 ━━━━━━━━━━━━━━━━━━━━ **5s** 7ms/step - binary_accuracy: 0.9137 - loss: 0.2378 - val_binary_accuracy: 0.8606 - val_loss: 0.3140
782/782 ━━━━━━━━━━━━━━━━━━━━ **4s** 4ms/step - binary accuracy: 0.8541 - loss: 0.3345

Loss: 0.33228880167007446
Accuracy: 0.854200005531311

## Training and validation loss



## Training and validation accuracy