# 用于文本的深度学习

## 自然语言处理简介

## 文本数据准备

### 文本标准化

### 文本分割（Tokenization）

### 词汇表索引（Vocabulary Indexing）

## `TextVectorization` 层

```python
import string

class Vectorizer:
    def standardize(self, text):
        text = text.lower()
        return "".join(char for char in text if char not in string.punctuation)

    def tokenize(self, text):
        return text.split()

    def make_vocabulary(self, dataset):
        self.vocabulary = {"": 0, "[UNK]": 1}
        for text in dataset:
            text = self.standardize(text)
            tokens = self.tokenize(text)
            for token in tokens:
                if token not in self.vocabulary:
                    self.vocabulary[token] = len(self.vocabulary)
        self.inverse_vocabulary = dict(
            (v, k) for k, v in self.vocabulary.items())

    def encode(self, text):
        text = self.standardize(text)
        tokens = self.tokenize(text)
        return [self.vocabulary.get(token, 1) for token in tokens]

    def decode(self, int_sequence):
        return " ".join(
            self.inverse_vocabulary.get(i, "[UNK]") for i in int_sequence)

vectorizer = Vectorizer()
dataset = [
    "I write, erase, rewrite",
    "Erase again, and then",
    "A poppy blooms.",
]
vectorizer.make_vocabulary(dataset)
```

```python
test_sentence = "I write, rewrite, and still rewrite again"
encoded_sentence = vectorizer.encode(test_sentence)
print(encoded_sentence)
```

```python
decoded_sentence = vectorizer.decode(encoded_sentence)
print(decoded_sentence)
```

```python
from tensorflow.keras.layers import TextVectorization
text_vectorization = TextVectorization(
```

```
        output_mode="int",
    )
```

```
    import re
    import string
    import tensorflow as tf

    def custom_standardization_fn(string_tensor):
        lowercase_string = tf.strings.lower(string_tensor)
        return tf.strings.regex_replace(
            lowercase_string, f"[{re.escape(string.punctuation)}]", "")

    def custom_split_fn(string_tensor):
        return tf.strings.split(string_tensor)

    text_vectorization = TextVectorization(
        output_mode="int",
        standardize=custom_standardization_fn,
        split=custom_split_fn,
    )
```

```
    dataset = [
        "I write, erase, rewrite",
        "Erase again, and then",
        "A poppy blooms.",
    ]
    text_vectorization.adapt(dataset)
```

词汇输**出**

```
    text_vectorization.get_vocabulary()
```

```
    vocabulary = text_vectorization.get_vocabulary()
    test_sentence = "I write, rewrite, and still rewrite again"
    encoded_sentence = text_vectorization(test_sentence)
    print(encoded_sentence)
```

```
    inverse_vocab = dict(enumerate(vocabulary))
    decoded_sentence = " ".join(inverse_vocab[int(i)] for i in encoded_sentence)
    print(decoded_sentence)
```

## ⌄ 表示单词集合的两种方法：集合 (Set) 和 序列 (Sequence)

## ⌄ 准备 IMDB 电影评论数据

```
    !curl -O https://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.tar.gz
    !tar -xf aclImdb_v1.tar.gz
```

```
    !rm -r aclImdb/train/unsup
```

```
    !cat aclImdb/train/pos/4077_10.txt
```

```
    import os, pathlib, shutil, random

    base_dir = pathlib.Path("aclImdb")
    val_dir = base_dir / "val"
    train_dir = base_dir / "train"
    for category in ("neg", "pos"):
        os.makedirs(val_dir / category)
        files = os.listdir(train_dir / category)
        random.Random(1337).shuffle(files)
        num_val_samples = int(0.2 * len(files))
        val_files = files[-num_val_samples:]
        for fname in val_files:
            shutil.move(train_dir / category / fname,
                        val_dir / category / fname)
```

```
from tensorflow import keras
batch_size = 32

train_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/train", batch_size=batch_size
)
val_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/val", batch_size=batch_size
)
test_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/test", batch_size=batch_size
)
```

**打印第一批的大小和 dtype**

```
for inputs, targets in train_ds:
    print("inputs.shape:", inputs.shape)
    print("inputs.dtype:", inputs.dtype)
    print("targets.shape:", targets.shape)
    print("targets.dtype:", targets.dtype)
    print("inputs[0]:", inputs[0])
    print("targets[0]:", targets[0])
    break
```

## ⌄ 将单词作为集合处理：BoW（Bag of Words）方法

## ⌄ Single words (unigrams) with binary encoding

**将数据预处理到 "TextVectorization" 层**

```
text_vectorization = TextVectorization(
    max_tokens=20000,
    output_mode="multi_hot",
)
text_only_train_ds = train_ds.map(lambda x, y: x)
text_vectorization.adapt(text_only_train_ds)

binary_1gram_train_ds = train_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
binary_1gram_val_ds = val_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
binary_1gram_test_ds = test_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
```

**检查二进制Unigram数据集的输出**

```
for inputs, targets in binary_1gram_train_ds:
    print("inputs.shape:", inputs.shape)
    print("inputs.dtype:", inputs.dtype)
    print("targets.shape:", targets.shape)
    print("targets.dtype:", targets.dtype)
    print("inputs[0]:", inputs[0])
    print("targets[0]:", targets[0])
    break
```

**模型生成实用程序**

```
from tensorflow import keras
from tensorflow.keras import layers

def get_model(max_tokens=20000, hidden_dim=16):
    inputs = keras.Input(shape=(max_tokens,))
    x = layers.Dense(hidden_dim, activation="relu")(inputs)
    x = layers.Dropout(0.5)(x)
    outputs = layers.Dense(1, activation="sigmoid")(x)
```

```
        model = keras.Model(inputs, outputs)
        model.compile(optimizer="rmsprop",
                      loss="binary_crossentropy",
                      metrics=["accuracy"])
        return model
```

训练**和**测试**二进制Unigram模型**

```
    model = get_model()
    model.summary()
    callbacks = [
        keras.callbacks.ModelCheckpoint("binary_1gram.h5",
                                        save_best_only=True)
    ]
```

```
    model.fit(binary_1gram_train_ds.cache(),
              validation_data=binary_1gram_val_ds.cache(),
              epochs=10,
              callbacks=callbacks)
```

```
    model = keras.models.load_model("binary_1gram.h5")
```

```
    print(f"테스트 정확도: {model.evaluate(binary_1gram_test_ds)[1]:.3f}")
```

⌄   二进制编码的 bigram

**바이그램을 반환하는** `TextVectorization` **층 만들기**

```
    text_vectorization = TextVectorization(
        ngrams=2,
        max_tokens=20000,
        output_mode="multi_hot",
    )
```

训练**和**测试**二进制bigram模型**

```
    text_vectorization.adapt(text_only_train_ds)
    binary_2gram_train_ds = train_ds.map(
        lambda x, y: (text_vectorization(x), y),
        num_parallel_calls=4)
    binary_2gram_val_ds = val_ds.map(
        lambda x, y: (text_vectorization(x), y),
        num_parallel_calls=4)
    binary_2gram_test_ds = test_ds.map(
        lambda x, y: (text_vectorization(x), y),
        num_parallel_calls=4)

    model = get_model()
    model.summary()
    callbacks = [
        keras.callbacks.ModelCheckpoint("binary_2gram.h5",
                                        save_best_only=True)
    ]
```

```
    model.fit(binary_2gram_train_ds.cache(),
              validation_data=binary_2gram_val_ds.cache(),
              epochs=10,
              callbacks=callbacks)
```

```
    model = keras.models.load_model("binary_2gram.h5")
```

```
    print(f"테스트 정확도: {model.evaluate(binary_2gram_test_ds)[1]:.3f}")
```

⌄   使用 TF-IDF 编码的 bigram

**使用 "TextVectorization"** 层**返回令牌**计数

```
text_vectorization = TextVectorization(
    ngrams=2,
    max_tokens=20000,
    output_mode="count"
)
```

**返回应用 TF-IDF 权重的输出的 "TextVectorization" 层**

```
text_vectorization = TextVectorization(
    ngrams=2,
    max_tokens=20000,
    output_mode="tf_idf",
)
```

**训练和测试TF-IDF bigram模型**

```
# Tensorflow 2.8.x 版本中的 TF-IDF 编码在 GPU 中执行时可能会出错。
# 虽然在Tensorflow2.9中解决了此问题，但由于Corab的Tensorflow版本在测试代码时为2.8.2，为了避免误差，使用CPU转换文本。

with tf.device("cpu"):
    text_vectorization.adapt(text_only_train_ds)

tfidf_2gram_train_ds = train_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
tfidf_2gram_val_ds = val_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
tfidf_2gram_test_ds = test_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)

model = get_model()
model.summary()
```

```
callbacks = [
    keras.callbacks.ModelCheckpoint("tfidf_2gram.h5",
                                    save_best_only=True)
]
```

```
model.fit(tfidf_2gram_train_ds.cache(),
          validation_data=tfidf_2gram_val_ds.cache(),
          epochs=10,
          callbacks=callbacks)
```

```
model = keras.models.load_model("tfidf_2gram.h5")
```

```
print(f"테스트 정확도: {model.evaluate(tfidf_2gram_test_ds)[1]:.3f}")
```

```
inputs = keras.Input(shape=(1,), dtype="string")
processed_inputs = text_vectorization(inputs)
outputs = model(processed_inputs)
inference_model = keras.Model(inputs, outputs)
```

```
import tensorflow as tf
raw_text_data = tf.convert_to_tensor([
    ["That was an excellent movie, I loved it."],
])
predictions = inference_model(raw_text_data)
print(f"긍정적인 리뷰일 확률: {float(predictions[0] * 100):.2f} 퍼센트")
```

## 将单词作为序列处理：序列模型方法

**Data download**

```
!rm -r aclImdb
!curl -O https://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.tar.gz
!tar -xf aclImdb_v1.tar.gz
!rm -r aclImdb/train/unsup
```

**Prepare data**

```
import os, pathlib, shutil, random
from tensorflow import keras
batch_size = 32
base_dir = pathlib.Path("aclImdb")
val_dir = base_dir / "val"
train_dir = base_dir / "train"

for category in ("neg", "pos"):
    os.makedirs(val_dir / category)
    files = os.listdir(train_dir / category)
    random.Random(1337).shuffle(files)
    num_val_samples = int(0.2 * len(files))
    val_files = files[-num_val_samples:]
    for fname in val_files:
        shutil.move(train_dir / category / fname,
                    val_dir / category / fname)

train_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/train", batch_size=batch_size
)
val_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/val", batch_size=batch_size
)
test_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/test", batch_size=batch_size
)
text_only_train_ds = train_ds.map(lambda x, y: x)
```

准备**整数序列数据集**

```
from tensorflow.keras import layers

max_length = 600
max_tokens = 20000
text_vectorization = layers.TextVectorization(
    max_tokens=max_tokens,
    output_mode="int",
    output_sequence_length=max_length,
)
text_vectorization.adapt(text_only_train_ds)

int_train_ds = train_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
int_val_ds = val_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
int_test_ds = test_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
```

用 **One Hot** 编码**的矢量序列**创**建序列模型**

⌄ 由于错误，省略了该代码。

理解词嵌入

⌄ 使用 Embedding 层学习词嵌入

创建 `Embedding` 层

```
        embedding_layer = layers.Embedding(input_dim=max_tokens, output_dim=256)
```

**使用从底层开始**训练的 `Embedding` 层**的模型**

```
        inputs = keras.Input(shape=(None,), dtype="int64")
        embedded = layers.Embedding(input_dim=max_tokens, output_dim=256)(inputs)
        x = layers.Bidirectional(layers.LSTM(32))(embedded)
        x = layers.Dropout(0.5)(x)
        outputs = layers.Dense(1, activation="sigmoid")(x)
        model = keras.Model(inputs, outputs)
        model.compile(optimizer="rmsprop",
                      loss="binary_crossentropy",
                      metrics=["accuracy"])
        model.summary()

        callbacks = [
            keras.callbacks.ModelCheckpoint("embeddings_bidir_lstm.keras",
                                            save_best_only=True)
        ]
```

```
        #model.fit(int_train_ds, validation_data=int_val_ds, epochs=10, callbacks=callbacks)
        import gdown
        url = "https://drive.google.com/uc?id=1D3thztDz4Ke_8KDMasD36avyNWGXdNFc"
        gdown.download(url, output="embeddings_bidir_lstm.keras")
```

```
        model = keras.models.load_model("embeddings_bidir_lstm.keras")
        print(f"테스트 정확도: {model.evaluate(int_test_ds)[1]:.3f}")
```

⌄  理解填充（padding）与掩码（masking）

**启用 masking（掩码） 功能**

```
        inputs = keras.Input(shape=(None,), dtype="int64")
        embedded = layers.Embedding(
            input_dim=max_tokens, output_dim=256, mask_zero=True)(inputs)
        x = layers.Bidirectional(layers.LSTM(32))(embedded)
        x = layers.Dropout(0.5)(x)
        outputs = layers.Dense(1, activation="sigmoid")(x)
        model = keras.Model(inputs, outputs)
        model.compile(optimizer="rmsprop",
                      loss="binary_crossentropy",
                      metrics=["accuracy"])
        model.summary()

        callbacks = [
            keras.callbacks.ModelCheckpoint("embeddings_bidir_lstm_with_masking.keras",
                                            save_best_only=True)
        ]
```

```
        #model.fit(int_train_ds, validation_data=int_val_ds, epochs=10, callbacks=callbacks)
        import gdown
        url = "https://drive.google.com/uc?id=1ickUOrTlpNcsGT3GlmOBVRpYBsoQrZZr"
        gdown.download(url, output="embeddings_bidir_lstm_with_masking.keras")
```

```
        model = keras.models.load_model("embeddings_bidir_lstm_with_masking.keras")
        print(f"테스트 정확도: {model.evaluate(int_test_ds)[1]:.3f}")
```

⌄  启用预训练单词嵌入

```
        # !wget http://nlp.stanford.edu/data/glove.6B.zip
        # !unzip -q glove.6B.zip

        import gdown
        url = "https://drive.google.com/uc?id=1i30Js4x3YPqNV9vLFZHmqsNEqSaqmdPV"
        gdown.download(url, output="glove.6B.100d.txt")
```

**解析 GloVe 单词嵌入文件**

```
import numpy as np
path_to_glove_file = "glove.6B.100d.txt"

embeddings_index = {}
with open(path_to_glove_file) as f:
    for line in f:
        word, coefs = line.split(maxsplit=1)
        coefs = np.fromstring(coefs, "f", sep=" ")
        embeddings_index[word] = coefs

print(f"단어 벡터 개수: {len(embeddings_index)}")
```

**准备嵌入 GloVe 单词矩阵**

```
embedding_dim = 100

vocabulary = text_vectorization.get_vocabulary()
word_index = dict(zip(vocabulary, range(len(vocabulary))))

embedding_matrix = np.zeros((max_tokens, embedding_dim))
for word, i in word_index.items():
    if i < max_tokens:
        embedding_vector = embeddings_index.get(word)
    if embedding_vector is not None:
        embedding_matrix[i] = embedding_vector
```

```
embedding_layer = layers.Embedding(
    max_tokens,
    embedding_dim,
    embeddings_initializer=keras.initializers.Constant(embedding_matrix),
    trainable=False,
    mask_zero=True,
)
```

**使用预训练嵌入的模型**

```
inputs = keras.Input(shape=(None,), dtype="int64")
embedded = embedding_layer(inputs)
x = layers.Bidirectional(layers.LSTM(32))(embedded)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs, outputs)
model.compile(optimizer="rmsprop",
              loss="binary_crossentropy",
              metrics=["accuracy"])
model.summary()

callbacks = [
    keras.callbacks.ModelCheckpoint("glove_embeddings_sequence_model.keras",
                                    save_best_only=True)
]
```

```
#model.fit(int_train_ds, validation_data=int_val_ds, epochs=10, callbacks=callbacks)

import gdown
url = "https://drive.google.com/uc?id=1kJB4fTZpPv4fZkQdNIFwoGUBL2gKPU0i"
gdown.download(url, output="glove_embeddings_sequence_model.keras")
```

```
model = keras.models.load_model("glove_embeddings_sequence_model.keras")
print(f"테스트 정확도: {model.evaluate(int_test_ds)[1]:.3f}")
```

## ⌄ Transformer 架构

理解自注意力（Self-Attention）

多头注意力（Multi-Head Attention）

**Load data**

```
!rm -r aclImdb
!curl -O https://ai.stanford.edu/~amaas/data/sentiment/aclImdb_v1.tar.gz
!tar -xf aclImdb_v1.tar.gz
!rm -r aclImdb/train/unsup
```

**Prepare data**

```
import os, pathlib, shutil, random
from tensorflow import keras
batch_size = 32
base_dir = pathlib.Path("aclImdb")
val_dir = base_dir / "val"
train_dir = base_dir / "train"
for category in ("neg", "pos"):
    os.makedirs(val_dir / category)
    files = os.listdir(train_dir / category)
    random.Random(1337).shuffle(files)
    num_val_samples = int(0.2 * len(files))
    val_files = files[-num_val_samples:]
    for fname in val_files:
        shutil.move(train_dir / category / fname,
                    val_dir / category / fname)

train_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/train", batch_size=batch_size
)
val_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/val", batch_size=batch_size
)
test_ds = keras.utils.text_dataset_from_directory(
    "aclImdb/test", batch_size=batch_size
)
text_only_train_ds = train_ds.map(lambda x, y: x)
```

**Data vectorization**

```
from tensorflow.keras import layers

max_length = 600
max_tokens = 20000
text_vectorization = layers.TextVectorization(
    max_tokens=max_tokens,
    output_mode="int",
    output_sequence_length=max_length,
)
text_vectorization.adapt(text_only_train_ds)

int_train_ds = train_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
int_val_ds = val_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
int_test_ds = test_ds.map(
    lambda x, y: (text_vectorization(x), y),
    num_parallel_calls=4)
```

继承`Layer`层实现的变压器编码器

```
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers

class TransformerEncoder(layers.Layer):
    def __init__(self, embed_dim, dense_dim, num_heads, **kwargs):
```

```python
        super().__init__(**kwargs)
        self.embed_dim = embed_dim
        self.dense_dim = dense_dim
        self.num_heads = num_heads
        self.attention = layers.MultiHeadAttention(
            num_heads=num_heads, key_dim=embed_dim)
        self.dense_proj = keras.Sequential(
            [layers.Dense(dense_dim, activation="relu"),
             layers.Dense(embed_dim),]
        )
        self.layernorm_1 = layers.LayerNormalization()
        self.layernorm_2 = layers.LayerNormalization()

    def call(self, inputs, mask=None):
        if mask is not None:
            mask = mask[:, tf.newaxis, :]
        attention_output = self.attention(
            inputs, inputs, attention_mask=mask)
        proj_input = self.layernorm_1(inputs + attention_output)
        proj_output = self.dense_proj(proj_input)
        return self.layernorm_2(proj_input + proj_output)

    def get_config(self):
        config = super().get_config()
        config.update({
            "embed_dim": self.embed_dim,
            "num_heads": self.num_heads,
            "dense_dim": self.dense_dim,
        })
        return config
```

**使用**变压**器**编码**器对文本**进**行分**类

```python
vocab_size = 20000
embed_dim = 256
num_heads = 2
dense_dim = 32

inputs = keras.Input(shape=(None,), dtype="int64")
x = layers.Embedding(vocab_size, embed_dim)(inputs)
x = TransformerEncoder(embed_dim, dense_dim, num_heads)(x)
x = layers.GlobalMaxPooling1D()(x)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs, outputs)
model.compile(optimizer="rmsprop",
              loss="binary_crossentropy",
              metrics=["accuracy"])
model.summary()
```

**基于**变压**器**编码**器的模型**训练**和评**估

```python
callbacks = [
    keras.callbacks.ModelCheckpoint("transformer_encoder.keras",
                                    save_best_only=True)
]
```

```python
#model.fit(int_train_ds, validation_data=int_val_ds, epochs=20, callbacks=callbacks)
```

```python
import gdown
url = "https://drive.google.com/uc?id=1q7Zyh-pX7m6WPQSN4y4irAOPkYwpQbpA"
gdown.download(url, output="transformer_encoder.keras")
```

```python
model = keras.models.load_model(
    "transformer_encoder.keras",
    custom_objects={"TransformerEncoder": TransformerEncoder})

print(f"테스트 정확도: {model.evaluate(int_test_ds)[1]:.3f}")
```

## 使用位置编码位置信息注入位置信息

**通**过亚类实现**位置嵌入**

```python
class PositionalEmbedding(layers.Layer):
    def __init__(self, sequence_length, input_dim, output_dim, **kwargs):
        super().__init__(**kwargs)
        self.token_embeddings = layers.Embedding(
            input_dim=input_dim, output_dim=output_dim)
        self.position_embeddings = layers.Embedding(
            input_dim=sequence_length, output_dim=output_dim)
        self.sequence_length = sequence_length
        self.input_dim = input_dim
        self.output_dim = output_dim

    def call(self, inputs):
        length = tf.shape(inputs)[-1]
        positions = tf.range(start=0, limit=length, delta=1)
        embedded_tokens = self.token_embeddings(inputs)
        embedded_positions = self.position_embeddings(positions)
        return embedded_tokens + embedded_positions

    def compute_mask(self, inputs, mask=None):
        #return tf.math.not_equal(inputs, 0)
        return self.token_embeddings.compute_mask(inputs)

    def get_config(self):
        config = super().get_config()
        config.update({
            "output_dim": self.output_dim,
            "sequence_length": self.sequence_length,
            "input_dim": self.input_dim,
        })
        return config
```

## 文本分类变压器

**合并**变压**器编码器和位置embedding**

```python
vocab_size = 20000
sequence_length = 600
embed_dim = 256
num_heads = 2
dense_dim = 32

inputs = keras.Input(shape=(None,), dtype="int64")
x = PositionalEmbedding(sequence_length, vocab_size, embed_dim)(inputs)
x = TransformerEncoder(embed_dim, dense_dim, num_heads)(x)
x = layers.GlobalMaxPooling1D()(x)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs, outputs)
model.compile(optimizer="rmsprop",
              loss="binary_crossentropy",
              metrics=["accuracy"])
model.summary()


#from google.colab import files
#files.download("transformer_encoder.keras")
```

```python
callbacks = [
    keras.callbacks.ModelCheckpoint("full_transformer_encoder.keras",
                                    save_best_only=True)
]
```

```python
#model.fit(int_train_ds, validation_data=int_val_ds, epochs=20, callbacks=callbacks)
```

```
import gdown
url = "https://drive.google.com/uc?id=1iBU6PSFYK9uxwHo9_VjqhE_C9_5sFDxW"
gdown.download(url, output="full_transformer_encoder.keras")
```

```
model = keras.models.load_model(
    "full_transformer_encoder.keras",
    custom_objects={"TransformerEncoder": TransformerEncoder,
                    "PositionalEmbedding": PositionalEmbedding})
```

```
print(f"테스트 정확도: {model.evaluate(int_test_ds)[1]:.3f}")
```

## 超越文本分类：序列到序列学习

```
!wget http://storage.googleapis.com/download.tensorflow.org/data/spa-eng.zip
!unzip -q spa-eng.zip
```

```
text_file = "spa-eng/spa.txt"
with open(text_file) as f:
    lines = f.read().split("\n")[:-1]
text_pairs = []
for line in lines:
    english, spanish = line.split("\t")
    spanish = "[start] " + spanish + " [end]"
    text_pairs.append((english, spanish))
```

```
import random
print(random.choice(text_pairs))
```

```
import random
random.shuffle(text_pairs)
num_val_samples = int(0.15 * len(text_pairs))
num_train_samples = len(text_pairs) - 2 * num_val_samples
train_pairs = text_pairs[:num_train_samples]
val_pairs = text_pairs[num_train_samples:num_train_samples + num_val_samples]
test_pairs = text_pairs[num_train_samples + num_val_samples:]
```

** 英文字母和西班牙语文本配对**

```
import tensorflow as tf
from tensorflow.keras import layers
import string
import re

strip_chars = string.punctuation + "¿"
strip_chars = strip_chars.replace("[", "")
strip_chars = strip_chars.replace("]", "")

def custom_standardization(input_string):
    lowercase = tf.strings.lower(input_string)
    return tf.strings.regex_replace(
        lowercase, f"[{re.escape(strip_chars)}]", "")

vocab_size = 15000
sequence_length = 20

source_vectorization = layers.TextVectorization(
    max_tokens=vocab_size,
    output_mode="int",
    output_sequence_length=sequence_length,
)
target_vectorization = layers.TextVectorization(
    max_tokens=vocab_size,
    output_mode="int",
    output_sequence_length=sequence_length + 1,
    standardize=custom_standardization,
)
train_english_texts = [pair[0] for pair in train_pairs]
train_spanish_texts = [pair[1] for pair in train_pairs]
source_vectorization.adapt(train_english_texts)
```

```
        target_vectorization.adapt(train_spanish_texts)
```

为**翻译工作准**备数**据集**

```
    batch_size = 64

    def format_dataset(eng, spa):
        eng = source_vectorization(eng)
        spa = target_vectorization(spa)
        return ({
            "english": eng,
            "spanish": spa[:, :-1],
        }, spa[:, 1:])

    def make_dataset(pairs):
        eng_texts, spa_texts = zip(*pairs)
        eng_texts = list(eng_texts)
        spa_texts = list(spa_texts)
        dataset = tf.data.Dataset.from_tensor_slices((eng_texts, spa_texts))
        dataset = dataset.batch(batch_size)
        dataset = dataset.map(format_dataset, num_parallel_calls=4)
        return dataset.shuffle(2048).prefetch(16).cache()

    train_ds = make_dataset(train_pairs)
    val_ds = make_dataset(val_pairs)
```

```
    for inputs, targets in train_ds.take(1):
        print(f"inputs['english'].shape: {inputs['english'].shape}")
        print(f"inputs['spanish'].shape: {inputs['spanish'].shape}")
        print(f"targets.shape: {targets.shape}")
```

## ⌄ 使用 RNN 的序列到序列模型

**基**础编码**器**

```
    from tensorflow import keras
    from tensorflow.keras import layers

    embed_dim = 256
    latent_dim = 1024

    source = keras.Input(shape=(None,), dtype="int64", name="english")
    x = layers.Embedding(vocab_size, embed_dim, mask_zero=True)(source)
    encoded_source = layers.Bidirectional(
        layers.GRU(latent_dim), merge_mode="sum")(x)
```

**基于 GRU 的解码器和端到端模型**

```
    past_target = keras.Input(shape=(None,), dtype="int64", name="spanish")
    x = layers.Embedding(vocab_size, embed_dim, mask_zero=True)(past_target)
    decoder_gru = layers.GRU(latent_dim, return_sequences=True)
    x = decoder_gru(x, initial_state=encoded_source)
    x = layers.Dropout(0.5)(x)
    target_next_step = layers.Dense(vocab_size, activation="softmax")(x)
    seq2seq_rnn = keras.Model([source, past_target], target_next_step)
```

训练**基于 RNN 的序列到序列模型**

```
    seq2seq_rnn.compile(
        optimizer="rmsprop",
        loss="sparse_categorical_crossentropy",
        metrics=["accuracy"])
```

```
    # seq2seq_rnn.fit(train_ds, epochs=15, validation_data=val_ds)
```

```
import gdown
url = "https://drive.google.com/uc?id=1YdVB7_1PIOhodKkFKRX7IhT8dm9841Sr"
gdown.download(url, output="seq2seq_rnn_full.keras")
```

```
seq2seq_rnn.load_weights("seq2seq_rnn_full.keras")
```

**用 RNN 编码器和解码器翻译新句子**

```
import numpy as np
spa_vocab = target_vectorization.get_vocabulary()
spa_index_lookup = dict(zip(range(len(spa_vocab)), spa_vocab))
max_decoded_sentence_length = 20

def decode_sequence(input_sentence):
    tokenized_input_sentence = source_vectorization([input_sentence])
    decoded_sentence = "[start]"
    for i in range(max_decoded_sentence_length):
        tokenized_target_sentence = target_vectorization([decoded_sentence])
        next_token_predictions = seq2seq_rnn.predict(
            [tokenized_input_sentence, tokenized_target_sentence])
        sampled_token_index = np.argmax(next_token_predictions[0, i, :])
        sampled_token = spa_index_lookup[sampled_token_index]
        decoded_sentence += " " + sampled_token
        if sampled_token == "[end]":
            break
    return decoded_sentence

test_eng_texts = [pair[0] for pair in test_pairs]
for _ in range(20):
    input_sentence = random.choice(test_eng_texts)
    print("-")
    print(input_sentence)
    print(decode_sequence(input_sentence))
```

⌄  使用 Transformer 的序列到序列模型

⌄  TransformerDecoder

`TransformerDecoder` **Class**

```
class TransformerDecoder(layers.Layer):
    def __init__(self, embed_dim, dense_dim, num_heads, **kwargs):
        super().__init__(**kwargs)
        self.embed_dim = embed_dim
        self.dense_dim = dense_dim
        self.num_heads = num_heads
        self.attention_1 = layers.MultiHeadAttention(
            num_heads=num_heads, key_dim=embed_dim)
        self.attention_2 = layers.MultiHeadAttention(
            num_heads=num_heads, key_dim=embed_dim)
        self.dense_proj = keras.Sequential(
            [layers.Dense(dense_dim, activation="relu"),
             layers.Dense(embed_dim),]
        )
        self.layernorm_1 = layers.LayerNormalization()
        self.layernorm_2 = layers.LayerNormalization()
        self.layernorm_3 = layers.LayerNormalization()
        self.supports_masking = True

    def get_config(self):
        config = super().get_config()
        config.update({
            "embed_dim": self.embed_dim,
            "num_heads": self.num_heads,
            "dense_dim": self.dense_dim,
        })
        return config

    def call(self, inputs, encoder_outputs, mask=None):
        attention_output_1 = self.attention_1(
```

```
            query=inputs,
            value=inputs,
            key=inputs,
            use_causal_mask=True)
        attention_output_1 = self.layernorm_1(inputs + attention_output_1)
        attention_output_2 = self.attention_2(
            query=attention_output_1,
            value=encoder_outputs,
            key=encoder_outputs
        )
        attention_output_2 = self.layernorm_2(
            attention_output_1 + attention_output_2)
        proj_output = self.dense_proj(attention_output_2)
        return self.layernorm_3(attention_output_2 + proj_output)
```

## 机器翻译变压器

`PositionalEmbedding` 层

```
class PositionalEmbedding(layers.Layer):
    def __init__(self, sequence_length, input_dim, output_dim, **kwargs):
        super().__init__(**kwargs)
        self.token_embeddings = layers.Embedding(
            input_dim=input_dim, output_dim=output_dim)
        self.position_embeddings = layers.Embedding(
            input_dim=sequence_length, output_dim=output_dim)
        self.sequence_length = sequence_length
        self.input_dim = input_dim
        self.output_dim = output_dim

    def call(self, inputs):
        length = tf.shape(inputs)[-1]
        positions = tf.range(start=0, limit=length, delta=1)
        embedded_tokens = self.token_embeddings(inputs)
        embedded_positions = self.position_embeddings(positions)
        return embedded_tokens + embedded_positions

    def compute_mask(self, inputs, mask=None):
        # return tf.math.not_equal(inputs, 0)
        return self.token_embeddings.compute_mask(inputs)

    def get_config(self):
        config = super(PositionalEmbedding, self).get_config()
        config.update({
            "output_dim": self.output_dim,
            "sequence_length": self.sequence_length,
            "input_dim": self.input_dim,
        })
        return config
```

### 端到端变压器

```
class TransformerEncoder(layers.Layer):
    def __init__(self, embed_dim, dense_dim, num_heads, **kwargs):
        super().__init__(**kwargs)
        self.embed_dim = embed_dim
        self.dense_dim = dense_dim
        self.num_heads = num_heads
        self.attention = layers.MultiHeadAttention(
            num_heads=num_heads, key_dim=embed_dim)
        self.dense_proj = keras.Sequential(
            [layers.Dense(dense_dim, activation="relu"),
             layers.Dense(embed_dim),]
        )
        self.layernorm_1 = layers.LayerNormalization()
        self.layernorm_2 = layers.LayerNormalization()

    def call(self, inputs, mask=None):
        if mask is not None:
            mask = mask[:, tf.newaxis, :]
        attention_output = self.attention(
```

```
                inputs, inputs, attention_mask=mask)
        proj_input = self.layernorm_1(inputs + attention_output)
        proj_output = self.dense_proj(proj_input)
        return self.layernorm_2(proj_input + proj_output)

    def get_config(self):
        config = super().get_config()
        config.update({
            "embed_dim": self.embed_dim,
            "num_heads": self.num_heads,
            "dense_dim": self.dense_dim,
        })
        return config
```

```
embed_dim = 256
dense_dim = 2048
num_heads = 8

encoder_inputs = keras.Input(shape=(None,), dtype="int64", name="english")
x = PositionalEmbedding(sequence_length, vocab_size, embed_dim)(encoder_inputs)
encoder_outputs = TransformerEncoder(embed_dim, dense_dim, num_heads)(x)

decoder_inputs = keras.Input(shape=(None,), dtype="int64", name="spanish")
x = PositionalEmbedding(sequence_length, vocab_size, embed_dim)(decoder_inputs)
x = TransformerDecoder(embed_dim, dense_dim, num_heads)(x, encoder_outputs)
x = layers.Dropout(0.5)(x)
decoder_outputs = layers.Dense(vocab_size, activation="softmax")(x)
transformer = keras.Model([encoder_inputs, decoder_inputs], decoder_outputs)
```

训练**序列到序列**变压**器**

```
transformer.compile(
    optimizer="rmsprop",
    loss="sparse_categorical_crossentropy",
    metrics=["accuracy"])
```

```
# transformer.fit(train_ds, epochs=30, validation_data=val_ds)
```

```
# transformer.save("seq2seq_transformer_full.keras")

# from google.colab import files
# files.download("seq2seq_transformer_full.keras")
```

```
import gdown
url = "https://drive.google.com/uc?id=1gWZWOu244A5LOWI87wZqjWPnm5Y_XJ_q"
gdown.download(url, output="seq2seq_transformer_full.keras")
```

```
transformer.load_weights("seq2seq_transformer_full.keras")
```

使用变形金刚**模型翻译新句子**

```
import numpy as np
spa_vocab = target_vectorization.get_vocabulary()
spa_index_lookup = dict(zip(range(len(spa_vocab)), spa_vocab))
max_decoded_sentence_length = 20

def decode_sequence(input_sentence):
    tokenized_input_sentence = source_vectorization([input_sentence])
    decoded_sentence = "[start]"
    for i in range(max_decoded_sentence_length):
        tokenized_target_sentence = target_vectorization(
            [decoded_sentence])[:, :-1]

        predictions = transformer(
            [tokenized_input_sentence, tokenized_target_sentence])

        sampled_token_index = np.argmax(predictions[0, i, :])
        sampled_token = spa_index_lookup[sampled_token_index]
        decoded_sentence += " " + sampled_token
        if sampled_token == "[end]":
```

```
            break
        return decoded_sentence

test_eng_texts = [pair[0] for pair in test_pairs]
for _ in range(20):
    input_sentence = random.choice(test_eng_texts)
    print("-")
    print(input_sentence)
    print(decode_sequence(input_sentence))
```