

Bài 18 - Các layers quan trọng trong deep learning

02 Dec 2019 - phamdinhkhanh

Menu

- 1. Giới thiệu chung
- 2. Các layer cơ bản
 - [2.1. Time distributed](#)
 - 2.2. Batch Normalization
 - 2.3. Attention Layer
- 3. Tài liệu tham khảo

1. Giới thiệu chung

Mỗi một mô hình deep learning đều có các layers đặc trưng giúp giải quyết các tác vụ cụ thể của từng bài toán deep learning. Chẳng hạn như trong xử lý ảnh chúng ta thường sử dụng mạng CNN - convolutional neural network (<https://www.kaggle.com/phamdinhkhanh/convolutional-neural-network-p1>) để trích xuất đặc trưng trên các local regional của bức ảnh là các đường nét chính như dọc, ngang, chéo,.... Hoặc layer LSTM - long short term memory (https://phamdinhkhanh.github.io/2019/04/22/Ly_thuyet_ve_mang_LSTM.html) được sử dụng trong các mô hình dịch máy và mô hình phân loại cảm xúc văn bản (sentiment analysis). Tuy nhiên ngoài các layer trên, chúng ta sẽ còn làm quen với rất nhiều các layers khác trong các bài toán về deep learning. Việc hiểu được công dụng của từng layer cũng như trường hợp áp dụng để mang lại hiệu quả cho mô hình rất quan trọng. Chính vì thế bài viết này nhằm mục đích hệ thống lại các layers quan trọng của deep learning như một tài liệu cheat sheet (tài liệu sổ tay) để sử dụng khi cần.

Do các layer CNN và LSTM đã được trình bày ở 2 bài viết của blog nên tôi sẽ không nêu lại kiến thức của những layers này. Và tất nhiên đó là những layers rất quan trọng mà bạn đọc cần nắm vững để áp dụng trong các mô hình về xử lý ảnh và ngôn ngữ. Bài viết này chỉ hướng tới những layer khác quan trọng hơn.

2. Các layer cơ bản

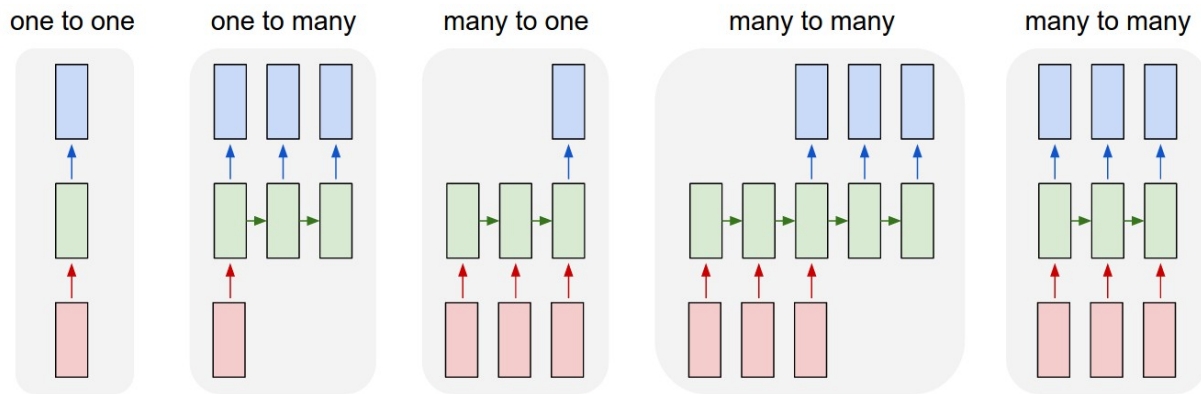
2.1. Time distributed

Dường như cái tên của layer này đã nói lên ý nghĩa của nó là xác định phân phối của dữ liệu theo thời gian.

Để hiểu rõ hơn chúng ta cùng lấy ví dụ về dự báo dạng chuỗi thông qua mạng RNN. Nhưng khoan khoan, để hiểu những gì tôi sắp viết vui lòng hiểu kĩ kiến trúc mạng RNN thông qua bài LSTM - long short term memory (https://phamdinhkhanh.github.io/2019/04/22/Ly_thuyet_ve_mang_LSTM.html).

Bạn đã đọc xong và nắm vững kiến thức về LSTM rồi chứ? Nếu chắc chắn, chúng ta hãy tiếp tục nào. Bên dưới là những dạng dự báo của RNN:

Top



Hình 1: Các dạng dự báo trong RNN

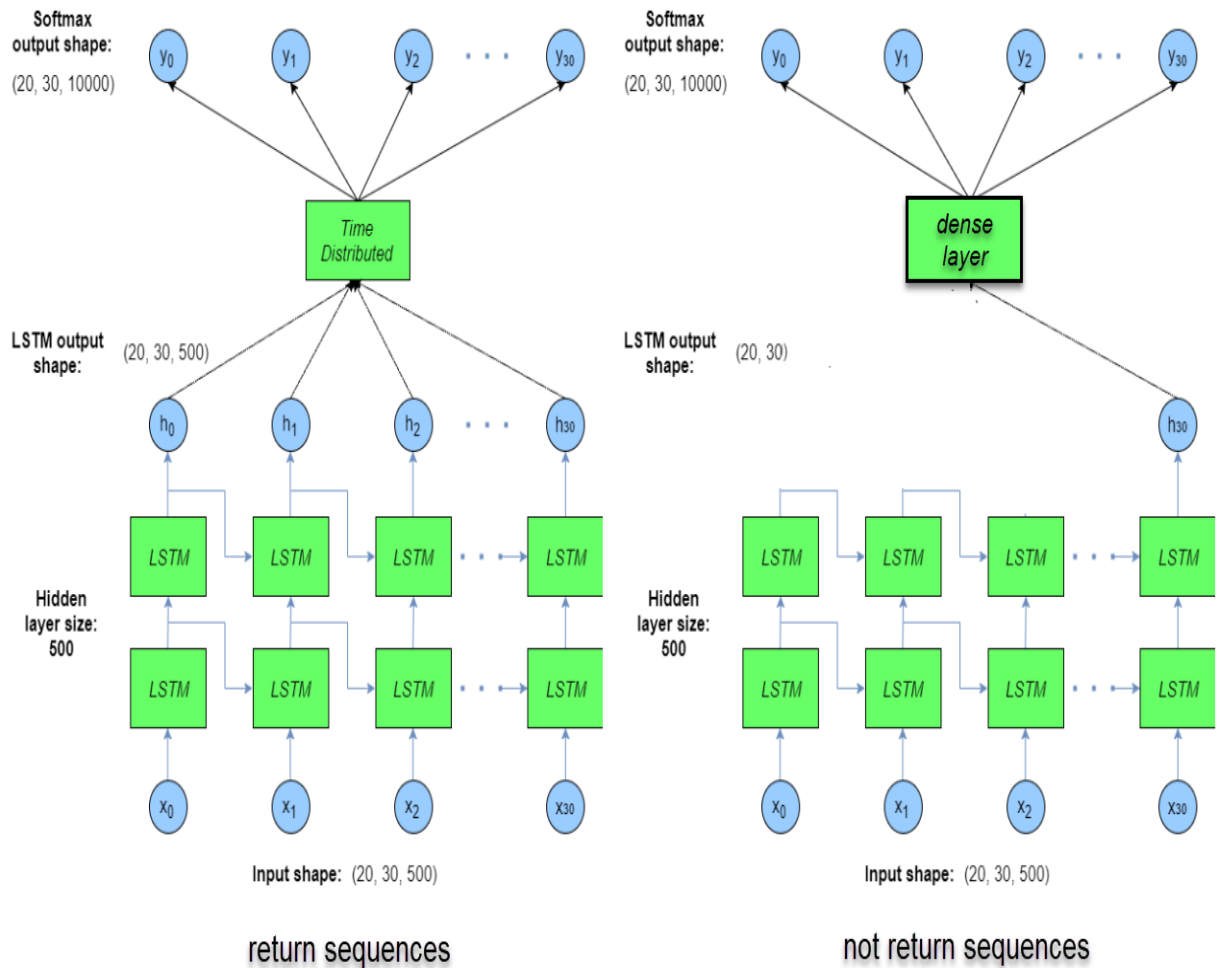
- One to one: Là trường hợp chỉ có 1 input và trả ra kết quả 1 output.
- One to many: Là trường hợp đầu vào chỉ có 1 input. Véc tơ context ở bước liền sau $t + 1$ và kết hợp với output dự báo ở bước t sẽ được sử dụng để dự báo được output ở bước $t + 1$. Cứ tiếp tục quá trình đến hết chuỗi, từ 1 input ta dự báo được một chuỗi nhiều outputs.
- Many to one: Là trường hợp ta chỉ trả ra kết quả là véc tơ output tại time step cuối cùng trong mạng LSTM (tương đương với cấu hình `return_sequence = False` trong LSTM layer của keras).
- Many to many: Đây chính là kiến trúc đặc trưng của model dịch máy. Tại mỗi bước thời gian t , input là véc tơ embedding một từ của ngôn ngữ nguồn và trả kết quả đầu ra là một véc tơ phân phối xác suất của từ output ở ngôn ngữ đích.

Trong các bài toán về NLP, Khi áp dụng họ các layer RNN thì chúng ta thường có 2 lựa chọn đó là:

Lựa chọn 1: trả ra chỉ kết quả là hidden layer ở véc tơ cuối cùng hoặc

Lựa chọn 2: trả ra chuỗi các hidden véc tơ ở mỗi time step.

Theo sau đó, các layers tiếp theo là những fully connected layers có kiến trúc như một mạng MLP thông thường. Kết quả trả ra là dự báo phân phối xác suất nhận. Câu hỏi được đặt ra là ta sẽ áp dụng các fully connected layers như thế nào cho từng lựa chọn? Rõ ràng đối với lựa chọn 1 thì do đầu ra của LSTM layer là một véc tơ nên ta dễ dàng truyền qua một Dense Layer (hay còn gọi là fully connected layer) thông thường và xây dựng một chuỗi fully connected layers khá dễ dàng. Tuy nhiên đối với lựa chọn 2 làm thế nào ta có thể kết hợp một Dense Layer với một chuỗi các hidden véc tơ như output của trường hợp many to many và one to many mà các bạn thấy. Khi đó chúng ta cần một layer đặc biệt hơn, không chỉ có tác dụng như dense layer trong mạng MLP mà còn có tác dụng kết nối tới từng hidden véc tơ ở mỗi bước thời gian, đó chính là Time Distributed Layers. Để dễ hình dung hơn bạn đọc có thể xem hình so sánh bên dưới giữa Time Distributed Layer và Dense Layer.



Hình 2: Time Distributed Layer ở bên trái kết nối chung các hidden vector ở mỗi bước thời gian từ h_0 đến h_{30} tới cùng một dense layer. Trong khi ở hình bên phải, đầu vào của dense layer chính là vector cuối cùng h_{30} .

Như vậy về bản chất Time Distributed Layer không khác gì một Dense Layer thông thường. Chính vì vậy trong một issue When and How to use TimeDistributedDense (<https://github.com/keras-team/keras/issues/1029>) fchollet tác giả của keras (và rất nhiều các đầu sách về deep learning trên cả python và R) đã giải thích một cách ngắn gọn như khá khó hiểu đối với beginners:

TimeDistributedDense applies a same Dense (fully-connected) operation to every timestep of a 3D tensor.

Ngoài ra Time Distributed Layer còn được sử dụng rất nhiều trong các mô hình xử lý video. Giả định rằng đầu vào của bạn là những batch video gồm các 5 chiều: (batch_size, time, width, height, channels). Như vậy để áp dụng một mạng tích chập 2 chiều lên toàn bộ các khung hình theo thời gian chúng ta cần sử dụng Time Distributed để thu được output shape mới ở đầu ra gồm 4 chiều: (batch_size, new_width, new_height, output_channels).

Như vậy bạn đọc đã hình dung ra nguyên lý hoạt động của Time Distributed Layer rồi chứ? Để hiểu thêm cách thức sử dụng trên keras, bạn đọc có thể tham khảo tài liệu Time Distributed Layer - keras document (https://www.tensorflow.org/api_docs/python/tf/keras/layers/TimeDistributed).

Bên dưới chúng ta sẽ cùng đi xây dựng 2 lớp model phân loại mail rác sử dụng 2 phương pháp khác nhau là LSTM + Dense Layer và LSTM + Time Distributed Layer và đối chiếu kết quả thu được.

Dữ liệu được sử dụng là sms spam collection dataset (<https://www.kaggle.com/uciml/sms-spam-collection-dataset>) gồm 5574 emails tiếng anh đã được gán nhãn sẵn là các email spam (email rác)/ham (email hợp lệ). Để xây dựng model chúng ta sẽ đi qua các bước như sau:

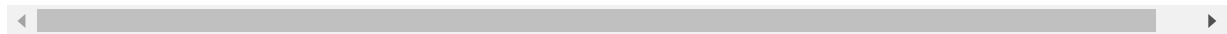
Bước 1: Xử lý dữ liệu.

Tại bước này chúng ta sẽ cần đọc và khảo sát dữ liệu để kiểm tra tính cân bằng, loại bỏ các từ stop words, dấu câu, kí tự đặc biệt và tạo từ điển để mã hóa các từ sang index.

Tất cả các công việc này được thực hiện khá dễ dàng nhờ những module có sẵn của gensim và keras .

- Đọc và khảo sát dữ liệu:

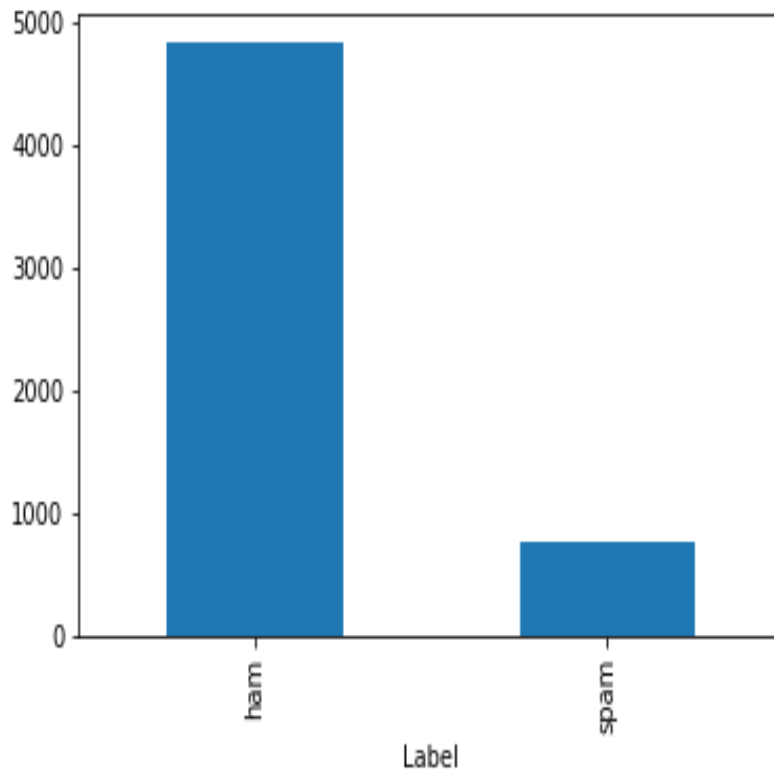
```
1 import pandas as pd
2 dataset = pd.read_csv('spam.csv', header=0, sep=',', encoding='latin-1')
3 dataset = dataset.iloc[:, :2]
4 dataset.columns = ['Label', 'Email']
5 dataset.head()
```



Label	Email
0ham	Go until jurong point, crazy.. Available only ...
1ham	Ok lar... Joking wif u oni...
2spam	Free entry in 2 a wkly comp to win FA Cup fina...
3ham	U dun say so early hor... U c already then say...
4ham	Nah I don't think he goes to usf, he lives aro...

Đồ thị phân phối các classes:

```
1 dataset.groupby('Label').Email.count().plot.bar()
```



- Làm sạch dữ liệu bằng cách chuẩn hóa các từ viết hoa thành viết thường, loại bỏ dấu câu, loại bỏ chữ số, tách số dính liền với từ và loại bỏ stop words,... thông qua package gensim.

```

1  import gensim
2  from gensim.parsing.preprocessing import strip_non_alphanum, strip_mu
3  import re
4  import nltk
5  from nltk.corpus import stopwords
6  nltk.download('stopwords')
7  stop_words = stopwords.words('english')
8
9  # Chuyển chữ hoa sang chữ thường
10 def lower_case(docs):
11     return [doc.lower() for doc in docs]
12
13 # Loại bỏ các dấu câu và kí tự đặc biệt
14 def remove_punc(docs):
15     return [strip_non_alphanum(doc).strip() for doc in docs]
16
17 # Tách các số và chữ liền nhau
18 def separate_num(docs):
19     return [split_alphanum(doc) for doc in docs]
20
21 # Loại bỏ những từ gồm 1 chữ cái đứng đơn lẻ
22 def remove_one_letter_word(docs):
23     return [strip_short(doc) for doc in docs]
24
25 # Loại bỏ các con số trong văn bản vì chúng không có nhiều ý nghĩa tr
26 def remove_number(docs):
27     return [strip_numeric(doc) for doc in docs]
28
29 # Thay thế nhiều khoảng spaces bằng 1 khoảng space
30 def replace_multiple_whitespaces(docs):
31     return [strip_multiple_whitespaces(doc) for doc in docs]
32
33 # Loại bỏ các stop words
34 def remove_stopwords(docs):
35     return [[word for word in doc.split() if word not in stop_words]
36
37 docs = lower_case(dataset['Email'])
38 docs = remove_punc(docs)
39 docs = separate_num(docs)
40 docs = remove_one_letter_word(docs)
41 docs = remove_number(docs)
42 docs = replace_multiple_whitespaces(docs)
43 docs = remove_stopwords(docs)
44 dataset['Content_Clean'] = docs

```

- Khởi tạo tokenizer để mã hóa các email.

Để chuyển các câu văn thành ma trận số, chúng ta cần tạo ra một từ điển mapping mỗi từ với index tương ứng của nó. module tokenizer dễ dàng giúp ta thực hiện việc này.

Top

```

1  from tensorflow.keras.preprocessing.text import Tokenizer
2  from tensorflow.keras.preprocessing import sequence
3
4  def _tokenize_matrix(docs, max_words=1000, max_len=150):
5      #max_words: Số lượng từ lớn nhất xuất hiện trong tokenizer được lấy
6      #max_len: Số lượng các từ lớn nhất trong một câu văn.
7      #docs: Tập hợp các đoạn email.
8      tok = Tokenizer(num_words=max_words)
9      tok.fit_on_texts(docs)
10     X_tok = tok.texts_to_sequences(docs)
11     X = sequence.pad_sequences(X_tok, maxlen=max_len)
12     return X_tok, X, tok
13
14     X_tok, X, tok = _tokenize_matrix(docs=dataset['Content_Clean'], max_w

```

Mã hóa nhãn spam/ham về biến one-hot.

```

1  from sklearn.preprocessing import LabelEncoder
2  le = LabelEncoder()
3  y = le.fit_transform(dataset['Label'])

```

Bước 2: Xây dựng model LSTM

Sau khi đã preprocessing dữ liệu, chúng ta thu được đầu vào là các ma trận padding X mà các dòng của nó là những câu văn có độ dài bằng nhau. Từ ma trận padding, ta cần đi qua một layer embedding để tạo véc tơ nhúng cho mỗi từ trong câu, và sau đó đi vào mạng LSTM. Chúng ta sẽ áp dụng cả 2 kiến trúc mô hình là Dense Layer và Time Distributed Layer.

- Khởi tạo model Dense Layer

Top

```

1  from tensorflow.keras.layers import Input, Embedding, LSTM, Dense, Activation
2  from tensorflow.keras.models import Model, Sequential
3  from sklearn.preprocessing import LabelEncoder
4  from tensorflow.keras.optimizers import RMSprop, Adam
5  from tensorflow.keras.utils import to_categorical
6  from tensorflow.keras.callbacks import EarlyStopping
7
8  def RNN_Dense(maxword=1000, embedding_size=100, max_len=150, n_unit_lstm=64, n_unit_dense=256):
9      inputs = Input(name='inputs', shape=[max_len])
10     layer = Embedding(maxword, embedding_size, input_length=max_len)(inputs)
11     # Embedding (input_dim: size of vocabulary,
12     # output_dim: dimension of dense embedding,
13     # input_length: length of input sequence)
14     layer = LSTM(n_unit_lstm)(layer)
15     layer = Dense(n_unit_dense, name='FC1')(layer)
16     layer = Activation('relu')(layer)
17     layer = Dropout(0.5)(layer)
18     layer = Dense(1, name='out_layer')(layer)
19     layer = Activation('sigmoid')(layer)
20     model = Model(inputs=inputs, outputs=layer)
21     return model
22
23     lstm_dense=RNN_Dense()
24     lstm_dense.summary()
25     lstm_dense.compile(loss='binary_crossentropy', optimizer=Adam(), metrics=['accuracy'])

```

```

1  Model: "model_14"
2
3  Layer (type)                Output Shape                Param #
4  =====
5  inputs (InputLayer)         [(None, 150)]               0
6
7  embedding_20 (Embedding)     (None, 150, 100)           1000000
8
9  lstm_20 (LSTM)               (None, 64)                  42240
10
11 FC1 (Dense)                  (None, 256)                 16640
12
13 activation_29 (Activation)    (None, 256)                 0
14
15 dropout_14 (Dropout)         (None, 256)                 0
16
17 out_layer (Dense)            (None, 1)                   257
18
19 activation_30 (Activation)    (None, 1)                   0
20  =====
21  Total params: 159,137
22  Trainable params: 159,137
23  Non-trainable params: 0
24

```

- Huấn luyện model dense

Top

```

1  lstm_dense.fit(X, y, batch_size=128,
2                  epochs=10,
3                  validation_split=0.2,
4                  callbacks=[EarlyStopping(monitor='val_loss', min_delta=0

```

```

1  Train on 4457 samples, validate on 1115 samples
2  Epoch 1/10
3  4457/4457 [=====] - 13s 3ms/sample - loss: 0
4  Epoch 2/10
5  4457/4457 [=====] - 12s 3ms/sample - loss: 0
6  ...
7  Epoch 9/10
8  4457/4457 [=====] - 12s 3ms/sample - loss: 0
9  Epoch 10/10
10 4457/4457 [=====] - 12s 3ms/sample - loss: 0

```

- Khởi tạo model TimeDistributed Layer

```

1  def RNN_TimeDis(maxword=1000, embedding_size=100, max_len=150, n_unit
2      inputs = Input(name='inputs', shape=[max_len])
3      layer = Embedding(maxword, embedding_size, input_length=max_len)(
4      # Embedding (input_dim: size of vocabulary,
5      # output_dim: dimension of dense embedding,
6      # input_length: length of input sequence)
7      layer = LSTM(n_unit_lstm, return_sequences=True)(layer)
8      layer = TimeDistributed(Dense(n_unit_dense))(layer)
9      layer = Flatten()(layer)
10     layer = Activation('relu')(layer)
11     layer = Dropout(0.5)(layer)
12     layer = Dense(1, name='out_layer')(layer)
13     layer = Activation('sigmoid')(layer)
14     model = Model(inputs=inputs, outputs=layer)
15     return model
16
17     lstm_timedis=RNN_TimeDis()
18     lstm_timedis.summary()
19     lstm_timedis.compile(loss='binary_crossentropy', optimizer=Adam(), me

```

Top


```

1 Model: "model_15"
2
3 Layer (type)                Output Shape                Param #
4 =====
5 inputs (InputLayer)         [(None, 150)]              0
6
7 embedding_22 (Embedding)     (None, 150, 100)          1000000
8
9 lstm_22 (LSTM)               (None, 150, 64)           42240
10
11 time_distributed_3 (TimeDist (None, 150, 64)          4160
12
13 flatten_1 (Flatten)         (None, 9600)              0
14
15 activation_31 (Activation)   (None, 9600)              0
16
17 dropout_15 (Dropout)        (None, 9600)              0
18
19 out_layer (Dense)           (None, 1)                 9601
20
21 activation_32 (Activation)   (None, 1)                 0
22 =====
23 Total params: 156,001
24 Trainable params: 156,001
25 Non-trainable params: 0
26

```

- Huấn luyện model Time Distributed Layer

```

1 lstm_timedis.fit(X, y, batch_size=128,
2                  epochs=10,
3                  validation_split=0.2,
4                  callbacks=[EarlyStopping(monitor='val_loss', min_delta=0

```

```

1 Train on 4457 samples, validate on 1115 samples
2 Epoch 1/10
3 4457/4457 [=====] - 13s 3ms/sample - loss: 0
4 Epoch 2/10
5 4457/4457 [=====] - 13s 3ms/sample - loss: 0
6 Epoch 3/10
7 4457/4457 [=====] - 13s 3ms/sample - loss: 0
8 Epoch 4/10
9 4457/4457 [=====] - 13s 3ms/sample - loss: 0
10 Epoch 5/10
11 4457/4457 [=====] - 13s 3ms/sample - loss: 0
12 Epoch 6/10
13 4457/4457 [=====] - 13s 3ms/sample - loss: 0

```

Kết quả cho thấy model áp dụng Time Distributed và áp dụng Dense đều khá tốt. Mức độ chính xác thu được trên tập validation đều trên 95%.

Ngoài ra thì Time Distributed còn được áp dụng trong phân loại nội dung video rất hiệu quả. Có thời gian tôi sẽ hướng dẫn các bạn.

Top

2.2. Batch Normalization

Nếu bạn đã làm quen với các kiến trúc mô hình trong CNN như Alexnet, VGG16, VGG19 thì batch normalization là một layer được sử dụng khá nhiều. Layer này thường áp dụng ngay sau Convolutional layer và thường ở những vị trí đầu tiên của mô hình để đạt hiệu quả cao nhất. Mục đích chính của batch normalization đó là chuẩn hóa dữ liệu ở các layer theo batch về phân phối chuẩn để quá trình gradient descent hội tụ nhanh hơn.

Giả sử chúng ta có một mini-batch với giá trị của từng quan sát trong batch như sau:

$\mathcal{B} = x_1, x_2, \dots, x_m$ Khi đó Batch Normalization sẽ được xác định thông qua phép chuẩn hóa trung bình và phương sai của các phần tử trong batch theo công thức bên dưới:

$$\mu = \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2$$

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

giá trị \hat{x}_i chính là kết quả sau chuẩn hóa.

Để hiểu hơn về cách áp dụng và hiệu quả của batch normalization, chúng ta sẽ thử nghiệm xây dựng mô hình LeNet có sử dụng batch normalization và không sử dụng batch normalization và so sánh kết quả mô hình sau huấn luyện. Bộ dữ liệu mà chúng ta sử dụng là mnist gồm các bức ảnh kích thước 28 x 28 của 10 chữ số từ 0 đến 9. Thật may mắn là dữ liệu này có thể load trực tiếp từ keras, rất tiện lợi.

Bước 1: Load dữ liệu từ keras.

```

1      # Load dữ liệu.
2      from tensorflow.examples.tutorials.mnist import input_data
3      import tensorflow as tf
4      from tensorflow import keras
5      print(tf.__version__)
6
7      # mnist = keras.datasets.mnist
8      # (X_train, y_train), (X_test, y_test) = mnist.load_data()
9
10     mnist = input_data.read_data_sets("MNIST_data/", reshape=False)
11     X_train, y_train = mnist.train.images, mnist.train.labels
12     X_val, y_val = mnist.validation.images, mnist.validation.labels
13     X_test, y_test = mnist.test.images, mnist.test.labels
14     print('X_train shape: {}'.format(X_train.shape))
15     print('y_train shape: {}'.format(y_train.shape))
16     print('X_test shape: {}'.format(X_test.shape))
17     print('y_test shape: {}'.format(y_test.shape))
18     print('X_val shape: {}'.format(X_val.shape))
19     print('y_val shape: {}'.format(y_val.shape))

1      X_train shape: (55000, 28, 28, 1)
2      y_train shape: (55000,)
3      X_test shape: (10000, 28, 28, 1)
4      y_test shape: (10000,)
5      X_val shape: (5000, 28, 28, 1)
6      y_val shape: (5000,)
```

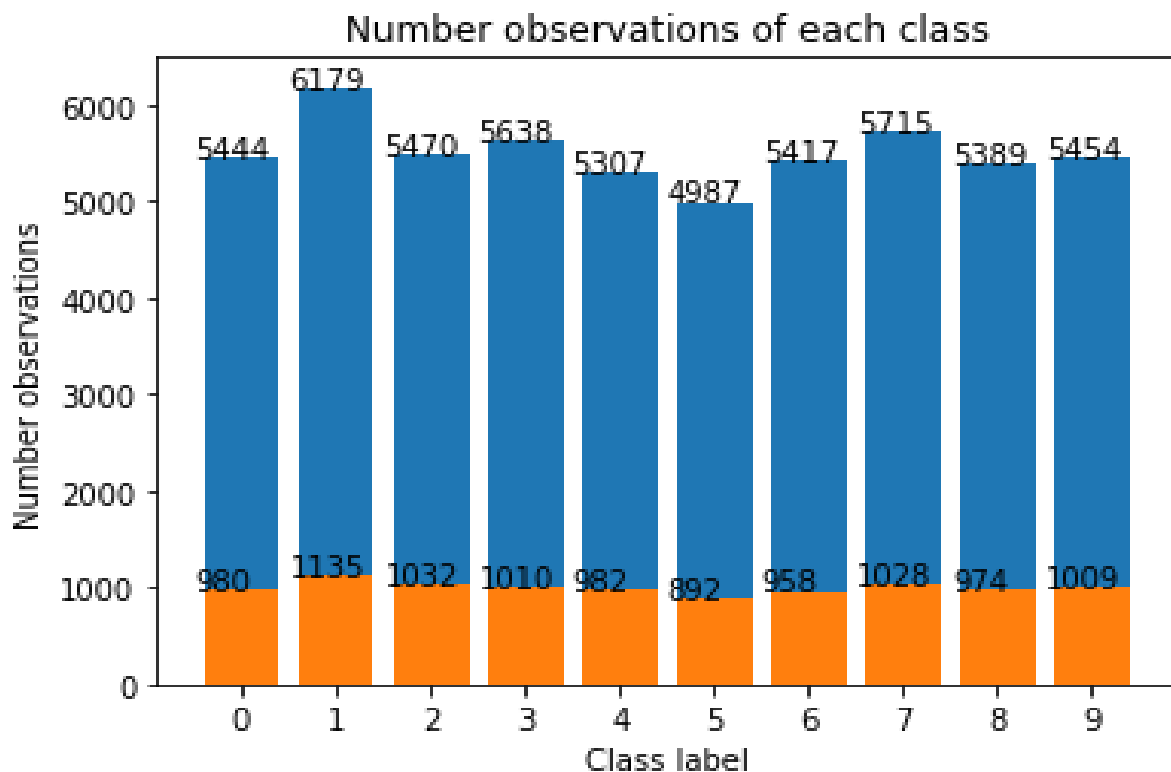
Top

Bước tiếp theo tuy đơn giản nhưng vô cùng quan trọng. Rất nhiều beginner thường bỏ qua vì chưa có kinh nghiệm. Đó là kiểm tra phân phối số quan sát trên các nhóm.

```

1  import matplotlib.pyplot as plt
2  import numpy as np
3
4  # Thống kê số lượng ở mỗi classes
5  def _plot_bar(x, title = 'Number observations of each class'):
6      x, y = np.unique(x, return_counts=True)
7      x_lab = [str(lab) for lab in x]
8      plt.bar(x_lab, y)
9      plt.xlabel('Class label')
10     plt.ylabel('Number observations')
11     plt.title(title)
12
13     for i in range(len(x)): # your number of bars
14         plt.text(x = x[i]-0.5, #takes your x values as horizontal position
15                 y = y[i]+1, #takes your y values as vertical positioning argument
16                 s = y[i], # the labels you want to add to the data
17                 size = 10)
18
19     _plot_bar(y_train)
20     _plot_bar(y_test)

```



Do input của mạng LeNet có kích thước 32×32 nên ta cần padding thêm các chiều để có kích thước là 32×32 . Thực hiện dễ dàng như sau:

Top

```

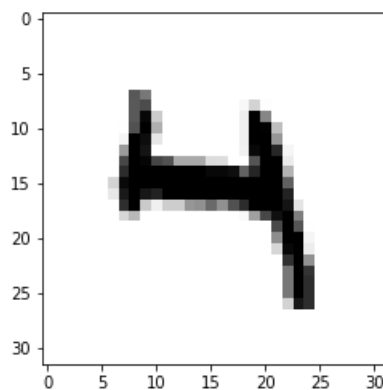
1 X_train = np.pad(X_train, ((0, 0), (2, 2), (2, 2), (0, 0)), 'constant')
2 X_val = np.pad(X_val, ((0, 0), (2, 2), (2, 2), (0, 0)), 'constant')
3 X_test = np.pad(X_test, ((0, 0), (2, 2), (2, 2), (0, 0)), 'constant')
4
5 print('X_train shape: {}'.format(X_train.shape))
6 print('X_test shape: {}'.format(X_test.shape))
7 print('X_val shape: {}'.format(X_val.shape))
8
9 # Để chắc chắn đã padding đúng thì hiển thị ra 1 hình ảnh số để kiểm
10 plt.imshow(np.squeeze(X_train[np.random.randint(55000)], axis=2), cma

```

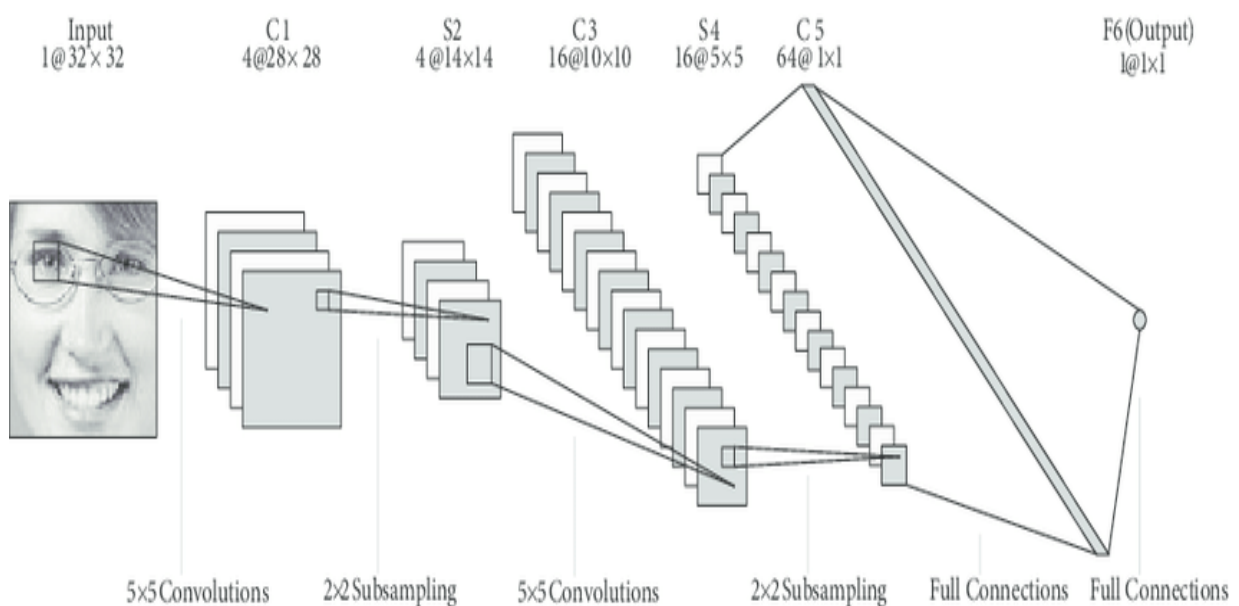
```

1 X_train shape: (55000, 32, 32, 1)
2 X_test shape: (10000, 32, 32, 1)
3 X_val shape: (5000, 32, 32, 1)

```



LeNet là mạng CNN đơn giản được tạo ra vào 1998 bởi Yan LeCun khi ông áp dụng thêm các layers tích chập (Convolutional) kết hợp với MaxPooling. Trong đó layer tích chập giúp nhận diện các đường nét dọc, ngang, chéo,... của vật thể còn MaxPooling nhằm giảm chiều dữ liệu mà không thay đổi các đặc trưng của ảnh giúp nhận diện vật thể. Cụ thể về kiến trúc từng layers các bạn xem hình vẽ bên dưới.



Top

Hình 3: Kiến trúc mạng LeNet. Trong đó input layer nhận đầu vào là những ảnh kích thước 32×32 . Tiếp theo là 2 lượt (Convolutional 2D + Maxpooling) giúp giảm chiều dữ liệu từ 32×32 xuống còn 5×5 . Flatten kết quả thành véc tơ và chuyển sang một mạng fully connected thông thường ta sẽ thu được giá trị dự báo xác suất ở output.

Chúng ta khởi tạo model LeNet không có BatchNormalization.

Top

```
1  from tensorflow.keras.layers import Flatten, Dense, Input, Activation
2  from tensorflow.keras.optimizers import Adam, SGD
3  from tensorflow.keras.models import Sequential, Model
4
5  def _Lenet_No_BatchNorm():
6      inp = Input(shape=(32, 32, 1))
7      conv1 = Conv2D(
8          filters=4,
9          kernel_size=5,
10         padding='valid',
11         strides=1,
12         activation='relu')(inp)
13
14     maxpool1 = MaxPooling2D(
15         pool_size=2,
16         strides=2
17     )(conv1)
18
19     conv2 = Conv2D(
20         filters=16,
21         kernel_size=5,
22         padding='valid',
23         strides=1,
24         activation='relu')(maxpool1)
25
26     maxpool2 = MaxPooling2D(
27         pool_size=2,
28         strides=2
29     )(conv2)
30
31     conv3 = Conv2D(
32         filters=64,
33         kernel_size=5,
34         padding='valid',
35         strides=1,
36         activation='relu')(maxpool2)
37
38     flatten = Flatten()(conv3)
39     dense1 = Dense(units=64, activation='relu')(flatten)
40     dense2 = Dense(units=32, activation='relu')(dense1)
41     dense3 = Dense(10, activation='sigmoid')(dense2)
42     output = Activation('softmax')(dense3)
43     model = Model(inputs = inp, outputs=output)
44     model.summary()
45     return model
46
47     lenet_no_batchnorm = _Lenet_No_BatchNorm()
```

[Top](#)

1 Model: "model_17"

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

Layer (type)	Output Shape	Param #
input_2 (InputLayer)	[(None, 32, 32, 1)]	0
conv2d_3 (Conv2D)	(None, 28, 28, 4)	104
max_pooling2d_2 (MaxPooling2	(None, 14, 14, 4)	0
conv2d_4 (Conv2D)	(None, 10, 10, 16)	1616
max_pooling2d_3 (MaxPooling2	(None, 5, 5, 16)	0
conv2d_5 (Conv2D)	(None, 1, 1, 64)	25664
flatten_3 (Flatten)	(None, 64)	0
dense_7 (Dense)	(None, 64)	4160
dense_8 (Dense)	(None, 32)	2080
dense_9 (Dense)	(None, 10)	330
activation_34 (Activation)	(None, 10)	0
Total params: 33,954		
Trainable params: 33,954		
Non-trainable params: 0		

Huấn luyện model LeNet khi không có Batch Normalization

```

1  from tensorflow.keras.callbacks import EarlyStopping
2  # Compile model
3  optimizer = Adam(learning_rate=0.005, beta_1=0.9, beta_2=0.999, amsgrad=True)
4  lenet_no_batchnorm.compile(loss='sparse_categorical_crossentropy', optimizer=optimizer)
5
6  # Khởi tạo model checkpoint
7  earlyStopping = EarlyStopping(monitor='val_loss', min_delta=0.0001, patience=10)
8  # Huấn luyện model
9  lenet_no_batchnorm.fit(X_train, y_train,
10                        batch_size=256,
11                        # validation_data=[X_val, y_val],
12                        validation_split=0.2,
13                        epochs=10,
14                        shuffle=True,
15                        callbacks=[earlyStopping])

```

Top

```
1    Train on 44000 samples, validate on 11000 samples
2    Epoch 1/10
3    44000/44000 [=====] - 21s 479us/sample - los:
4    Epoch 2/10
5    44000/44000 [=====] - 20s 462us/sample - los:
6    ...
7    Epoch 6/10
8    44000/44000 [=====] - 20s 463us/sample - los:
9    Epoch 7/10
10   44000/44000 [=====] - 20s 465us/sample - los:
```



Khởi tạo model LeNet khi có BatchNormalization

Top


```
1  from tensorflow.keras.layers import Flatten, Dense, Input, Activation
2  from tensorflow.keras.optimizers import Adam, SGD
3  from tensorflow.keras.models import Sequential, Model
4
5  def _Lenet_BatchNorm():
6      inp = Input(shape=(32, 32, 1))
7      conv1 = Conv2D(
8          filters=4,
9          kernel_size=5,
10         padding='valid',
11         strides=1,
12         activation='relu')(inp)
13
14     batch_norm1 = BatchNormalization()(conv1)
15
16     maxpool1 = MaxPooling2D(
17         pool_size=2,
18         strides=2
19     )(batch_norm1)
20
21     conv2 = Conv2D(
22         filters=16,
23         kernel_size=5,
24         padding='valid',
25         strides=1,
26         activation='relu')(maxpool1)
27
28     batch_norm2 = BatchNormalization()(conv2)
29
30     maxpool2 = MaxPooling2D(
31         pool_size=2,
32         strides=2
33     )(batch_norm2)
34
35     conv3 = Conv2D(
36         filters=64,
37         kernel_size=5,
38         padding='valid',
39         strides=1,
40         activation='relu')(maxpool2)
41
42     flatten = Flatten()(conv3)
43     dense1 = Dense(units=64, activation='relu')(flatten)
44     dense2 = Dense(units=32, activation='relu')(dense1)
45     dense3 = Dense(10, activation='sigmoid')(dense2)
46     output = Activation('softmax')(dense3)
47     model = Model(inputs = inp, outputs=output)
48     model.summary()
49     return model
50
51 lenet_batchnorm = _Lenet_BatchNorm()
```

[Top](#)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

Model: "model_22"

Layer (type)	Output Shape	Param #
=====		
input_9 (InputLayer)	[(None, 32, 32, 1)]	0

conv2d_20 (Conv2D)	(None, 28, 28, 4)	104

batch_normalization_6 (Batch Normalization)	(None, 28, 28, 4)	16

max_pooling2d_14 (MaxPooling2D)	(None, 14, 14, 4)	0

conv2d_21 (Conv2D)	(None, 10, 10, 16)	1616

batch_normalization_7 (Batch Normalization)	(None, 10, 10, 16)	64

max_pooling2d_15 (MaxPooling2D)	(None, 5, 5, 16)	0

conv2d_22 (Conv2D)	(None, 1, 1, 64)	25664

flatten_8 (Flatten)	(None, 64)	0

dense_22 (Dense)	(None, 64)	4160

dense_23 (Dense)	(None, 32)	2080

dense_24 (Dense)	(None, 10)	330

activation_39 (Activation)	(None, 10)	0
=====		
Total params: 34,034		
Trainable params: 33,994		
Non-trainable params: 40		

Huấn luyện model với Batch Normalization

```

1  from tensorflow.keras.callbacks import EarlyStopping
2  # Compile model
3  optimizer = Adam(learning_rate=0.005, beta_1=0.9, beta_2=0.999, amsgrad=True)
4  lenet_batchnorm.compile(loss='sparse_categorical_crossentropy', optimizer=optimizer)
5
6  # Khởi tạo model checkpoint
7  earlyStopping = EarlyStopping(monitor='val_loss', min_delta=0.0001, patience=10)
8  # Huấn luyện model
9  lenet_batchnorm.fit(X_train, y_train,
10                     batch_size=256,
11                     # validation_data=[X_val, y_val],
12                     validation_split=0.2,
13                     epochs=10,
14                     shuffle=True,
15                     callbacks=[earlyStopping])

```

```

1    Train on 44000 samples, validate on 11000 samples
2    Epoch 1/10
3    44000/44000 [=====] - 30s 688us/sample - los:
4    Epoch 2/10
5    44000/44000 [=====] - 29s 659us/sample - los:
6    ...
7    Epoch 6/10
8    44000/44000 [=====] - 29s 654us/sample - los:
9    Epoch 7/10
10   44000/44000 [=====] - 29s 655us/sample - los:

```

Ta thấy rằng model khi áp dụng Batch Normalization có tốc độ hội tụ nhanh hơn và độ chính xác cũng đồng thời cao hơn chút so với không áp dụng Batch Normalization. Điều đó chứng tỏ cơ chế chuẩn hóa dữ liệu sau mỗi layer đã giúp cho các tham số hội tụ nhanh hơn. Một thay đổi nhỏ nhưng hiệu quả mang lại thật bất ngờ.

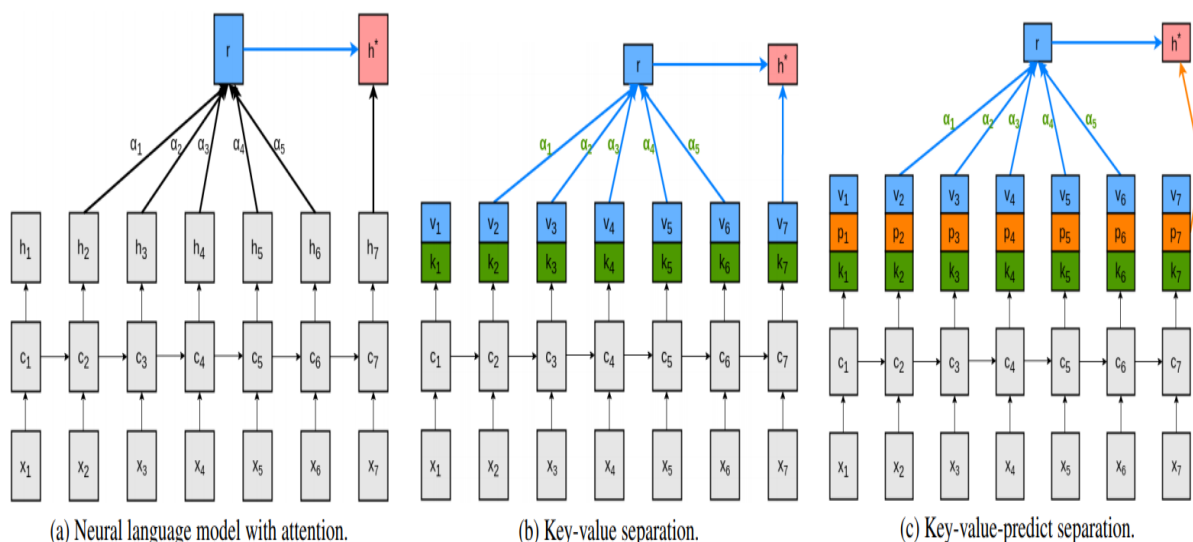
2.3. Attention Layer

Trong những năm gần đây, độ chính xác của các mô hình dịch máy được cải thiện một cách đáng kể, thành tựu đó đạt được chính là nhờ một layer đặc biệt có tác dụng phân bố lại trọng số attention weight của các một từ input lên từ output sao cho càng ở vị trí gần thì trọng số càng cao.

Cụ thể về thuật toán attention tôi sẽ không trình bày chi tiết ở đây bởi tôi đã giới thiệu ở một bài trước đó Bài 4: Attention is all you need

(<https://phamdinhhkhanh.github.io/2019/06/18/AttentionLayer.html>).

Như đã biết về cơ chế của thuật toán RNN, chúng ta yêu cầu output véc tơ tại mỗi time step liên tục lưu trữ thông tin để dự báo từ tiếp theo, tính toán attention và encode nội dung tương ứng với các step trong tương lai. Việc sử dụng output quá tải như vậy có thể ảnh hưởng đến tốc độ và hiệu năng của mô hình. Một kiến trúc key-value(-predict) attention đã được giới thiệu bởi Daniluk, 2017 (<https://arxiv.org/abs/1702.04521.pdf>) để giải quyết vấn đề trên. Kiến trúc này cho phép tách riêng các phần tính attention và dự báo output. Cụ thể kiến trúc này như bên dưới:



Top

Hình 4: Sơ đồ cơ chế key-value attention. Chúng ta chia đôi hidden output véc tơ thành các cặp key-value véc tơ để giảm tải gánh nặng cho quá trình tính toán output.

Các cặp key-value được chia ra từ output véc tơ thành 2 phần như hình (b). Công dụng của các cặp key này như sau:

- phần keys chiếm một nửa hidden output được sử dụng để tính toán attention.
- phần value chiếm một nửa hidden output được sử dụng để giải mã phân phối của từ tiếp theo và biểu diễn của véc tơ context.

Các input đầu vào của mô hình là x và ta gọi đó là query. Trong một số trường hợp, ta lựa chọn giá trị của query và value là trùng nhau. Từ giá trị keys và values chúng ta tính được ra giá trị dự báo cho các bước tiếp theo và giá trị này chính là predict. Như vậy đầu ra cuối cùng của mô hình trả ra 3 véc tơ tại mỗi bước thời gian gồm: các véc tơ v được dùng để giải mã phân phối của từ tiếp theo, các véc tơ k được sử dụng như là key cho tính toán các attention véc tơ, véc tơ p như là value cuối cùng sau khi áp dụng cơ chế attention như chúng ta nhìn thấy trong hình 4(c).

Hiện nay kiến trúc key-value attention layer cho kết quả và tốc độ tính toán tốt hơn so với các kiến trúc attention trước đó.

Trong tensorflow version 2.0.0, attention layer được khởi tạo chính là dựa trên kiến trúc key-value attention layer. Các bạn có thể tham khảo tại: Attention Layer - Tensorflow (https://www.tensorflow.org/api_docs/python/tf/keras/layers/Attention). Ngoài ra layer này còn cho phép sử dụng thêm các layer mask như query_mask và value_mask để loại bỏ attention tại một số vị trí nhất định trong query hoặc value.

Do mức độ tập trung phụ thuộc vào vị trí của từ trong câu nên attention layer được sử dụng chủ yếu ở những mô hình Seq2Seq. Ngoài ra trong các mô hình phân loại văn bản sử dụng LSTM thì attention layer cũng hoạt động khá hiệu quả.

Sau đây ta sẽ áp dụng Attention Layer đằng sau LSTM model trong tác vụ phân loại mail spam và so sánh hiệu quả so với áp dụng các layer như TimeDistributed Layer và Dense Layer.

```
1  from tensorflow.keras.layers import Layer, InputSpec
2  from tensorflow.keras import initializers
3
4  class AttLayer(Layer):
5      def __init__(self, **kwargs):
6          self.init = initializers.RandomNormal()
7          #self.input_spec = [InputSpec(ndim=3)]
8          super(AttLayer, self).__init__(**kwargs)
9
10     def build(self, input_shape):
11         assert len(input_shape)==3
12         #self.W = self.init((input_shape[-1],1))
13         self.W = self.init((input_shape[-1],))
14         #self.input_spec = [InputSpec(shape=input_shape)]
15         self.trainable_weights = [self.W]
16         super(AttLayer, self).build(input_shape) # be sure you call
17
18     def call(self, x, mask=None):
19         eij = K.tanh(K.dot(x, self.W))
20
21         ai = K.exp(eij)
22         weights = ai/K.sum(ai, axis=1).dimshuffle(0,'x')
23
24         weighted_input = x*weights.dimshuffle(0,1,'x')
25         return weighted_input.sum(axis=1)
26
27     def get_output_shape_for(self, input_shape):
28         return (input_shape[0], input_shape[-1])
```



Top

```
1  from tensorflow.keras.layers import Input, Embedding, LSTM, Dense, Activation
2  from tensorflow.keras.models import Model, Sequential
3  from sklearn.preprocessing import LabelEncoder
4  from tensorflow.keras.optimizers import RMSprop, Adam
5  from tensorflow.keras.utils import to_categorical
6  from tensorflow.keras.callbacks import EarlyStopping
7
8  def RNN_AttLayer(maxword=1000, embedding_size=100, max_len=150, n_units_lstm=128, n_units_dense=128, n_units_attention=128, n_units_dropout=128, n_units_output=1):
9      inputs = Input(name='inputs', shape=[max_len])
10     layer = Embedding(maxword, embedding_size, input_length=max_len)(inputs)
11     # Embedding (input_dim: size of vocabulary,
12     # output_dim: dimension of dense embedding,
13     # input_length: length of input sequence)
14     layer = LSTM(n_units_lstm, return_sequences=False)(layer)
15     layer = Attention()([layer, layer])
16     layer = Dense(n_units_dense, name='FC1')(layer)
17     layer = Activation('relu')(layer)
18     layer = Dropout(0.5)(layer)
19     layer = Dense(1, name='out_layer')(layer)
20     layer = Activation('sigmoid')(layer)
21     model = Model(inputs=inputs, outputs=layer)
22     return model
23
24 lstm_attlayer=RNN_AttLayer()
25 lstm_attlayer.summary()
26 lstm_attlayer.compile(loss='binary_crossentropy', optimizer=Adam(), metrics=['accuracy'])
```



1 Model: "model_24"

2

3 Layer (type)	Output Shape	Param #	Connections
4 =====	=====	=====	=====
5 inputs (InputLayer)	[(None, 150)]	0	
6			
7 embedding_28 (Embedding)	(None, 150, 100)	1000000	inputs
8			
9 lstm_28 (LSTM)	(None, 64)	42240	embedding_28
10			
11 attention_2 (Attention)	(None, 64)	0	lstm_28
12			lstm_28
13			
14 FC1 (Dense)	(None, 64)	4160	attention_2
15			
16 activation_42 (Activation)	(None, 64)	0	FC1
17			
18 dropout_17 (Dropout)	(None, 64)	0	activation_42
19			
20 out_layer (Dense)	(None, 1)	65	dropout_17
21			
22 activation_43 (Activation)	(None, 1)	0	out_layer
23			
24	Total params: 146,465		
25	Trainable params: 146,465		
26	Non-trainable params: 0		
27			

```

1 lstm_dense.fit(X, y, batch_size=128,
2               epochs=10,
3               validation_split=0.2,
4               callbacks=[EarlyStopping(monitor='val_loss', min_delta=0

```

```

1 Train on 4457 samples, validate on 1115 samples
2 Epoch 1/10
3 4457/4457 [=====] - 12s 3ms/sample - loss: 0
4 Epoch 2/10
5 4457/4457 [=====] - 12s 3ms/sample - loss: 0
6 Epoch 3/10
7 4457/4457 [=====] - 12s 3ms/sample - loss: 0
8 Epoch 4/10
9 4457/4457 [=====] - 12s 3ms/sample - loss: 0
10 Epoch 5/10
11 4457/4457 [=====] - 12s 3ms/sample - loss: 0

```

Như vậy ta thấy rằng khi áp dụng attention layer vào mô hình phân loại email thì độ chính xác của mô hình cao hơn. Ngay từ những epoch đầu tiên mô hình đã đạt được độ chính xác dường như hoàn hảo.

Điều này cho thấy attention layer rất mạnh trong việc nắm bắt các liên kết dài của chuỗi các từ trong câu và khắc phục được nhược điểm về sự phụ thuộc dài hạn kém của các mô hình Recurrent Neural Network.

Như vậy bài này tôi đã giới thiệu đến các bạn một số layers nổi bật của Deep Learning bao gồm Time Distributed, Batch Normalization và Attention Layer bên cạnh các layers chuyên biệt cho các mô hình xử lý ảnh (CNN) và ngôn ngữ tự nhiên (LSTM, GRU, RNN, BiDirectional RNN) mà các bạn chắc chắn phải nắm vững.

Qua các ví dụ minh họa chúng ta cũng so sánh được một cách tương đối hiệu quả của việc áp dụng các layers này và trường hợp nào, lớp mô hình hoặc bài toán nào thì nên áp dụng.

Ngoài ra còn rất nhiều các layers quan trọng khác của deep learning mà tôi sẽ viết ở những bài sau nữa nếu có dịp. Hi vọng rằng các bạn có thể áp dụng tốt những layers trên vào các bài toán của mình để nâng cao hiệu quả cho mô hình.

Và cuối cùng không thể thiếu trong các bài viết của Khanh Blog là Tài liệu tham khảo như một sự tôn trọng và tri ân gửi tới các tác giả mà tôi đã học hỏi.

3. Tài liệu tham khảo

1. Attention in NLP - Kate Loginova (<https://medium.com/@joealato/attention-in-nlp-734c6fa9d983>)
2. How to Use the TimeDistributed Layer in Keras - machinelearningmastery (<https://machinelearningmastery.com/timedistributed-layer-for-long-short-term-memory-networks-in-python/>)
3. How to work with Time Distributed data in a neural network - Patrice Ferlet (<https://medium.com/smileinnovation/how-to-work-with-time-distributed-data-in-a-neural-network-b8b39aa4ce00>)
4. Batch Normalization Tensorflow Keras Example - Cory Maklin (<https://towardsdatascience.com/backpropagation-and-batch-normalization-in-feedforward-neural-networks-explained-901fd6e5393e>)
5. Time Distributed Layer - tensorflow 2.0 (https://www.tensorflow.org/api_docs/python/tf/keras/layers/TimeDistributed)
6. Batch Normalization - tensorflow 2.0 (https://www.tensorflow.org/api_docs/python/tf/keras/layers/BatchNormalization)
7. Attention Layer - tensorflow 2.0 (https://www.tensorflow.org/api_docs/python/tf/keras/layers/Attention)

Top