

In [2]:

```
1 from tensorflow.keras import Input, Model
2 from tensorflow.keras.layers import Embedding, Dense, Conv1D, GlobalMaxPooling1D, Concatenate,
3
4 class TextCNN(object):
5     def __init__(self, maxlen, max_features, embedding_dims,
6                 class_num=14,
7                 last_activation='softmax'):
8         self.maxlen = maxlen # 文本词的最大长度
9         self.max_features = max_features #最大特征数
10        self.embedding_dims = embedding_dims #词向量的大小
11        self.class_num = class_num #类别数
12        self.last_activation = last_activation #最近的激活函数
13
14    def get_model(self):
15        input = Input((self.maxlen,))
16        embedding = Embedding(self.max_features, self.embedding_dims, input_length=self.maxlen)
17        convs = []
18        for kernel_size in [3, 4, 5]:
19            c = Conv1D(128, kernel_size, activation='relu')(embedding)
20            c = GlobalMaxPooling1D()(c)
21            convs.append(c)
22        x = Concatenate()(convs)
23
24        output = Dense(self.class_num, activation=self.last_activation)(x)
25        model = Model(inputs=input, outputs=output)
26        return model
```

In [3]:

```
1 from tensorflow.keras.preprocessing import sequence
2 import random
3 from sklearn.model_selection import train_test_split
4 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint
5 from tensorflow.keras.utils import to_categorical
6 from utils import *
7 from keras.utils.np_utils import *
8
9 # 路径等配置
10 data_dir = './data/train_set.csv'
11 vocab_file = "./vocab/vocab.txt"
12 # vocab_size = 40000
13 vocab_size = 5000
14 # 神经网络配置
15 # max_features = 40001
16 max_features = 5001
17 maxlen = 100
18 batch_size = 64
19 embedding_dims = 50
20 epochs = 20
21
22 print('数据预处理与加载数据...')
23 # 如果不存在词汇表, 重建
24 if not os.path.exists(vocab_file):
25     build_vocab(data_dir, vocab_file, vocab_size)
26
```

Using TensorFlow backend.

数据预处理与加载数据...

In [4]:

```
1 os.path.exists(vocab_file)
```

Out[4]:

True

In [5]:

```
1 # 获得 词汇/类别 与id映射字典
2 categories, cat_to_id = read_category() # cat_to_id:字典
3 ''' ['科技','股票','体育','娱乐','时政','社会','教育','财经','家居','游戏','房产','时尚','彩票']
4 words, word_to_id = read_vocab(vocab_file) # 词汇整成字典格式
5
6 # 全部数据
7 x, y = read_csv_file(data_dir)
8 data = list(zip(x,y))
9 del x,y
10 # 乱序
11 random.shuffle(data)
12 # 切分训练集和测试集
13 train_data, test_data = train_test_split(data)
14 # 对文本的词id和类别id进行编码
15 x_train = encode_sentences([content[0] for content in train_data], word_to_id)
16
17 y_train = to_categorical([content[1] for content in train_data])
```

In [6]:

```
1 x_train.shape, x_test.shape
```

```
-----  
-  
AttributeError                                Traceback (most recent call last)  
<ipython-input-6-f73e7a7cdf2c> in <module>  
----> 1 x_train.shape, x_test.shape
```

AttributeError: 'list' object has no attribute 'shape'

In [7]:

```
1 x_test = encode_sentences([content[0] for content in test_data], word_to_id)  
2 y_test = to_categorical([content[1] for content in test_data])  
3 print('对序列做padding, 保证是 samples*timestep 的维度')  
4 x_train = sequence.pad_sequences(x_train, maxlen=maxlen)  
5 x_test = sequence.pad_sequences(x_test, maxlen=maxlen)  
6 print('x_train shape:', x_train.shape)  
7 print('x_test shape:', x_test.shape)
```

对序列做padding, 保证是 samples*timestep 的维度

x_train shape: (15000, 100)

x_test shape: (5000, 100)

In [8]:

```
1 print('x_train shape:', y_train.shape)  
2 print('x_test shape:', y_test.shape)
```

x_train shape: (15000, 14)

x_test shape: (5000, 14)

In [9]:

```
1 y_test[0]
```

Out[9]:

```
array([0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.],  
      dtype=float32)
```

In [10]:

```
1 # import tensorflow as tf
2 # tf.config.gpu.set_per_process_memory_growth(enabled=True)
3
4
5 print('构建模型...')
6 model = TextCNN(maxlen, max_features, embedding_dims).get_model()
7 model.compile('adam', 'categorical_crossentropy', metrics=['accuracy'])
8
9 print('训练...')
10 # 设定callbacks回调函数
11 my_callbacks = [
12     ModelCheckpoint('./cnn_model.h5', verbose=1),
13     EarlyStopping(monitor='val_accuracy', patience=2, mode='max')
14 ]
15
16 # fit拟合数据
17 history = model.fit(x_train, y_train,
18                     batch_size=batch_size,
19                     epochs=epochs,
20                     callbacks=my_callbacks,
21                     validation_data=(x_test, y_test))
22
23 #print('对测试集预测...')
24 result = model.predict(x_test)
```

curacy: 0.8259 - val_loss: 0.6235 - val_accuracy: 0.8046

Epoch 4/8

14976/15000 [=====>.] - ETA: 0s - loss: 0.4272 - accuracy: 0.8750

Epoch 00004: saving model to ./cnn_model.h5

15000/15000 [=====] - 3s 196us/sample - loss: 0.4271 - accuracy: 0.8751 - val_loss: 0.5793 - val_accuracy: 0.8220

Epoch 5/8

14784/15000 [=====>.] - ETA: 0s - loss: 0.3157 - accuracy: 0.9123

Epoch 00005: saving model to ./cnn_model.h5

15000/15000 [=====] - 3s 188us/sample - loss: 0.3157 - accuracy: 0.9124 - val_loss: 0.5647 - val_accuracy: 0.8280

Epoch 6/8

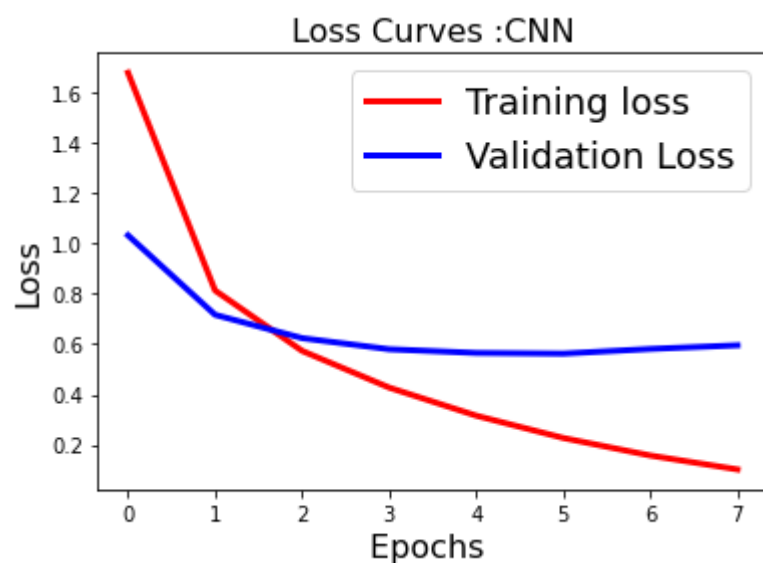
14784/15000 [=====>.] - ETA: 0s - loss: 0.2271 - accuracy: 0.9426

Epoch 00006: saving model to ./cnn_model.h5

15000/15000 [=====] - 3s 191us/sample - loss: 0.2271 - accuracy: 0.9425 - val_loss: 0.5620 - val_accuracy: 0.8300

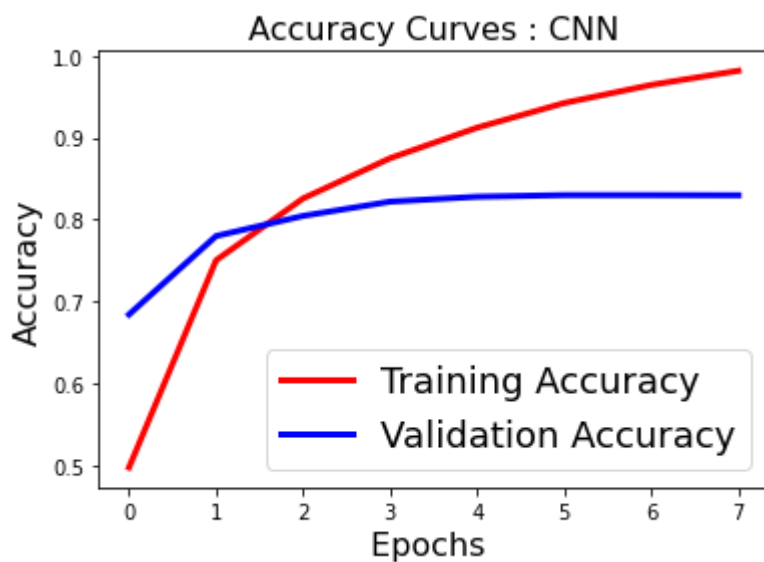
In [11]:

```
1 import matplotlib.pyplot as plt
2 plt.switch_backend('agg')
3 %matplotlib inline
4
5 fig1 = plt.figure()
6 plt.plot(history.history['loss'], 'r', linewidth=3.0)
7 plt.plot(history.history['val_loss'], 'b', linewidth=3.0)
8 plt.legend(['Training loss', 'Validation Loss'], fontsize=18)
9 plt.xlabel('Epochs ', fontsize=16)
10 plt.ylabel('Loss', fontsize=16)
11 plt.title('Loss Curves :CNN', fontsize=16)
12 fig1.savefig('loss_cnn.png')
13 plt.show()
```



In [12]:

```
1 fig2=plt.figure()
2 plt.plot(history.history['accuracy'],'r',linewidth=3.0)
3 plt.plot(history.history['val_accuracy'],'b',linewidth=3.0)
4 plt.legend(['Training Accuracy', 'Validation Accuracy'],fontsize=18)
5 plt.xlabel('Epochs ',fontsize=16)
6 plt.ylabel('Accuracy',fontsize=16)
7 plt.title('Accuracy Curves : CNN',fontsize=16)
8 fig2.savefig('../accuracy_cnn.png')
9 plt.show()
```



In []:

1

In [25]:

```
1 from tensorflow.keras.utils import plot_model
2 model.summary()
3 plot_model(model, show_shapes=True, show_layer_names=True)
```

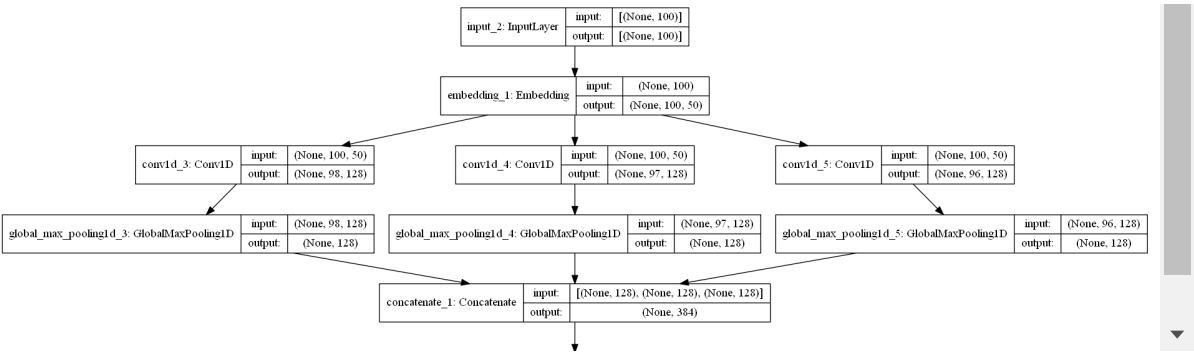
Model: "model_1"

Layer (type)	Output Shape	Param #	Connected to
=====			
input_2 (InputLayer)	[(None, 100)]	0	
=====			
embedding_1 (Embedding)	(None, 100, 50)	250050	input_2[0][0]
=====			
conv1d_3 (Conv1D)	(None, 98, 128)	19328	embedding_1[0][0]
=====			
conv1d_4 (Conv1D)	(None, 97, 128)	25728	embedding_1[0][0]
=====			
conv1d_5 (Conv1D)	(None, 96, 128)	32128	embedding_1[0][0]
=====			
global_max_pooling1d_3 (GlobalM	(None, 128)	0	conv1d_3[0][0]
=====			
global_max_pooling1d_4 (GlobalM	(None, 128)	0	conv1d_4[0][0]
=====			
global_max_pooling1d_5 (GlobalM	(None, 128)	0	conv1d_5[0][0]
=====			
concatenate_1 (Concatenate)	(None, 384)	0	global_max_pooling1 d_3[0][0] global_max_pooling1 d_4[0][0] global_max_pooling1 d_5[0][0]
=====			
dense_1 (Dense)	(None, 14)	5390	concatenate_1[0][0]
=====			
Total params: 332,624			
Trainable params: 332,624			
Non-trainable params: 0			



Out[25]:





In []:

1	
---	--