# Lab 2 – Image Caption Generator

## Report

By

Akhil Yerrapragada (akhily@kth.se)

Gibson Chikafa (chikafa@kth.se)

Group Name – gibson_akhil

**Problem Description:**

To implement an image caption generator as described in the paper using CNN, RNN and Tensorflow.

**Data:**

The dataset used for the implementation is Flicker8k_Dataset which contains images with descriptions, training and testing sets which can be used for respective purposes.

**Implementation:**

1) A pretrained model, **InceptionV3** provided by keras is used for to extract the features from images.
2) The features extracted by Inception V3 from all the images are saved in **features.pkl** file.

```python
# Pre trained CNN

def CNN():

# CNN for feature extraction

    for name in listdir(imagesSet):
        filename = imagesSet + '/' + name
        image = load_img(filename, target_size=(299, 299))
        image = img_to_array(image)
        print(name)
        image = image.reshape(
            (1, image.shape[0], image.shape[1], image.shape[2]))
        image = preprocess_input(image)
        feature = PreTrainedCNNModel.predict(image, verbose=0)
        imgId = name.split('.')[0]
        allFeat[imgId] = feature
    saveDesc(allFeat)

def saveDesc(allFeatures):
    dump(allFeatures, open(featurespkl , 'wb'))
```

3) Now, we take all the descriptions for images from descriptions.txt file and preprocess them before we append all the descriptions for an image. The appended will look as below:

```
#Get Descriptions From Token Dataset

def loadallDesc():
    for line in open(allDescriptions, 'r').read().split('\n'):
        tokens = line.split("#")
        if len(line) < 2:
            continue
        imgId, imgDesc = tokens[0], tokens[1:]
        imgId = imgId.split('.')[0]
        imgDesc = ' '.join(imgDesc)
        if imgId not in imageDescriptions:
            imageDescriptions[imgId] = list()
        imageDescriptions[imgId].append(imgDesc)

    for key, value in imageDescriptions.items():
        for i in range(len(value)):
            eachDescription = value[i]
            eachDescription = eachDescription.split()
            eachDescription = [word.lower() for word in eachDescription]
            eachDescription = [w.translate(str.maketrans('', '', string.punctuation)) for w in eachDescription]
            eachDescription = [word for word in eachDescription if len(word)>1]
            eachDescription = [word for word in eachDescription if word.isalpha()]
            value[i] = ' '.join(eachDescription)


    for key, value in imageDescriptions.items():
        for desc in value:
            imgDescPair.append(key + ' ' + desc)

    allRel = '\n'.join(imgDescPair)
    file = open(allDesctxt , 'w')
    file.write(allRel)
    file.close()
```

4) Now that we have all the descriptions, we can begin training the data. For that, we use trainImages.txt file which contains the details of images that can be used for training.

5) Now that we know what images can be used for training, we obtain the descriptions of the training images. For all the descriptions we append a start and end sequence which will be useful later in our implementation.

```
# Preprocessing Training Data

def fromAllDesc():
    for line in open(trainImages, 'r').read().split('\n'):
        if len(line) < 1:
            continue
        imgId = line.split('.')[0]
        Images.append(imgId)

    for line in  open(allDesctxt, 'r').read().split('\n'):
        tokens = line.split()
        imgId, imgDesc = tokens[0], tokens[1:]
        if imgId in set(Images):
            if imgId not in trainingImageDescriptions:
                trainingImageDescriptions[imgId] = list()
            finalDescription = 'startseq ' + ' '.join(imgDesc) + ' endseq'
            trainingImageDescriptions[imgId].append(finalDescription)
            allTrainDescriptions.append(finalDescription)
    return trainingImageDescriptions
```

6) Similar to descriptions, we obtain the features of the image as well from **features.pkl** file.

```
def getFeatures():
    all_features = load(open(featurespkl, 'rb'))
    features = {k: all_features[k] for k in fromAllDesc()}
    return features
```

7) Now, we define a model which has two inputs, one for the input image and the other for descriptions. We combine both and obtain the output which is our prediction. The input

layer is of shape = (47, 1000) which is the shape of our input image. For the input 2 we give the length of our longest description.

```python
# Input for Image
inputs1 = layers.Input(shape=(47,1000))
fe1 = layers.Dropout(0.5)(inputs1)
fe2 = layers.Dense(256, activation='relu')(fe1)

# Input for text
inputs2 = layers.Input(shape=(getLength(),))
se1 = layers.Embedding(vocabSize, 256, mask_zero=True)(inputs2)
se2 = layers.Dropout(0.5)(se1)
se3 = layers.LSTM(256)(se2)

# Combine here
decoder1 = layers.add([fe2, se3])
decoder2 = layers.Dense(256, activation='relu')(decoder1)
outputs = layers.Dense(vocabSize, activation='softmax')(decoder2)
model = models.Model(inputs=[inputs1, inputs2], outputs=outputs)
model.compile(loss='categorical_crossentropy', optimizer='adam')
```

8)  Now that we have our model defined and training data ready, we can proceed with training.
9)  For training the model, we give the training descriptions, features and all the vocabulary available in the training descriptions and the length of longest description.
10) First, we convert all the descriptions that we are having to sequence using **text_to_sequences** provided by tokenizer which is also used to obtain vocabulary of our descriptions. We train the model by giving the features and input sequence as an input and the output sequence. Both will be given as a tuple for the model.

```python
def txtToseq(tokenizer, maxLen, value, feat):
    X1, X2, y = list(), list(), list()
    for desc in value:
        seq = tokenizer.texts_to_sequences([desc])[0]
        for i in range(1, len(seq)):
            inSeq, outSeq = seq[:i], seq[i]
            inSeq = pad_sequences([inSeq], maxlen=maxLen)[0]
            outSeq = to_categorical([outSeq], num_classes=vocabSize)[0]

            X1.append(feat)
            X2.append(inSeq)
            y.append(outSeq)
    return array(X1), array(X2), array(y)


def trainingData(descriptions, trainFeat, tokenizer, maxLen):
    while 1:
        for key, value in descriptions.items():
            feat = trainFeat[key][0]
            recvFeatVect, inSeqVect, outSeqVect = txtToseq(tokenizer, maxLen, value, feat)
            yield [recvFeatVect, inSeqVect], outSeqVect
```

**Training:**

1)  First, we extract the features using the **CNN** function which has our pre trained model.

```python
# Step 1: Uncomment below to initiate pre training and extract features from all images

"""

CNN()

"""
```
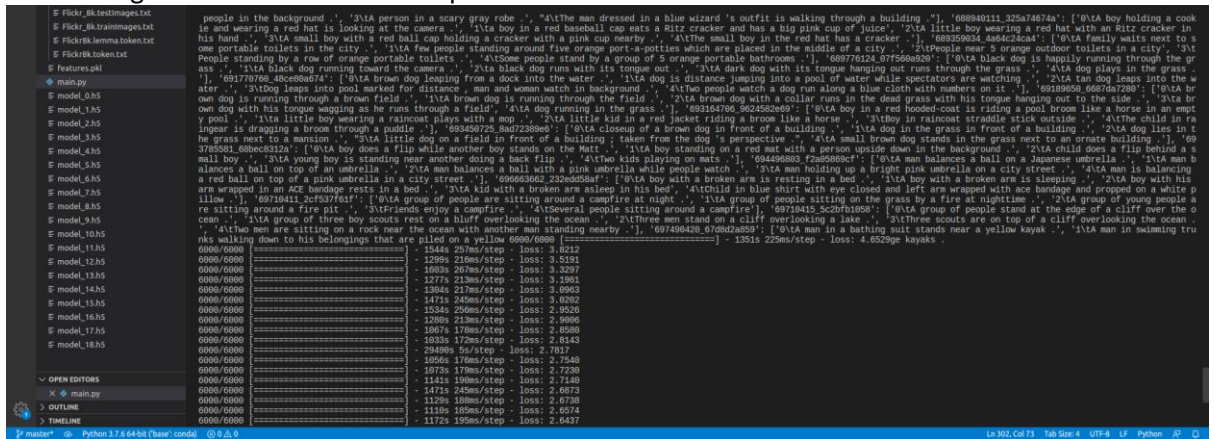
2)  Now we load all the training descriptions and their features, available vocabulary.

```
# Step 2: load all training descriptions, training features and available vocabulary

loadallDesc()
print('Loaded training descriptions from all desc ', len(fromAllDesc()))
tokenizer = getVocab()
vocabSize = len(tokenizer.word_index) + 1
print('Total Vocabulary from training descriptions ', vocabSize)
print('Longest sentence for RNN input layer ',  getLength())
trainFeat = getFeatures()
```
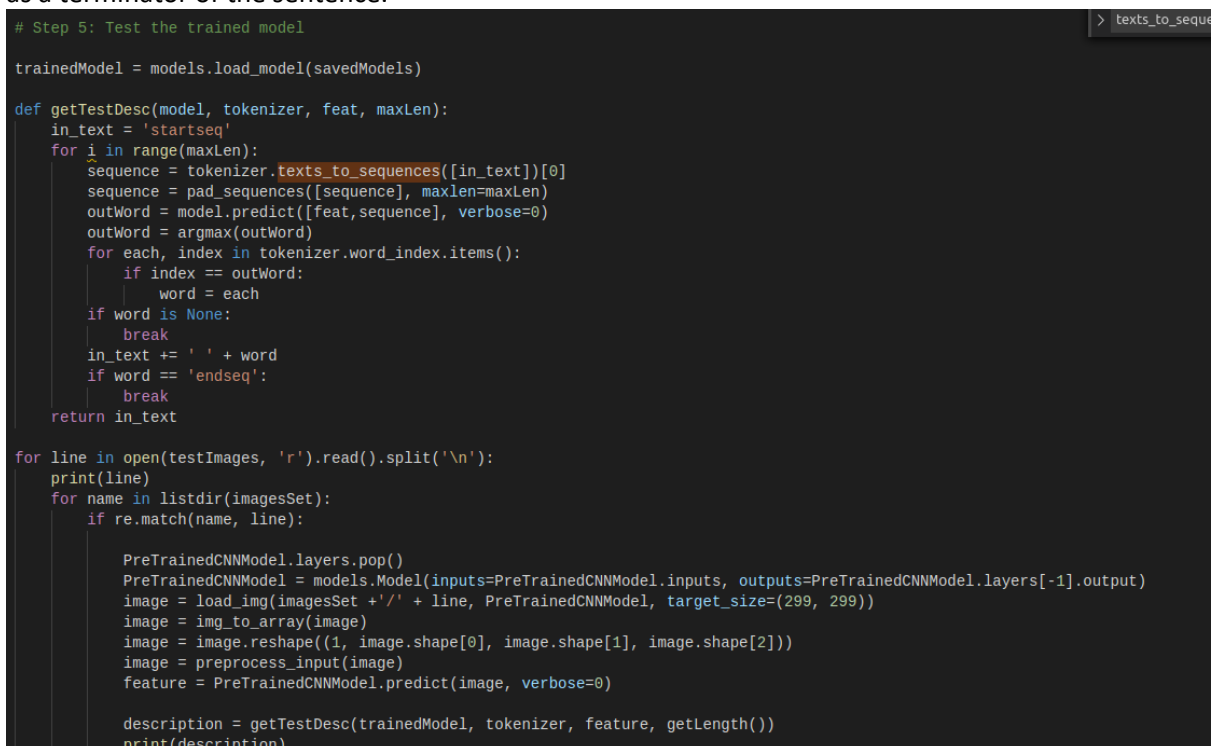
3) We train the model using the above information for 18 epochs each for 6000 steps considering that we have 6000 descriptions.



**Test:**

1) We load one of our trained models using **load_model.**
2) We also load our pre trained CNN for feature extraction, i.e. InceptionV3.
3) We load one of the images available in our test set and give it as an input for our pre trained CNN.
4) The pre trained CNN will get the features of the image which will act as an input to or model.
5) We also give our sequence initiator along with the features and the sequence ender will act as a terminator of the sentence.

```
# Step 5: Test the trained model

trainedModel = models.load_model(savedModels)

def getTestDesc(model, tokenizer, feat, maxLen):
    in_text = 'startseq'
    for i in range(maxLen):
        sequence = tokenizer.texts_to_sequences([in_text])[0]
        sequence = pad_sequences([sequence], maxlen=maxLen)
        outWord = model.predict([feat,sequence], verbose=0)
        outWord = argmax(outWord)
        for each, index in tokenizer.word_index.items():
            if index == outWord:
                word = each
        if word is None:
            break
        in_text += ' ' + word
        if word == 'endseq':
            break
    return in_text

for line in open(testImages, 'r').read().split('\n'):
    print(line)
    for name in listdir(imagesSet):
        if re.match(name, line):

            PreTrainedCNNModel.layers.pop()
            PreTrainedCNNModel = models.Model(inputs=PreTrainedCNNModel.inputs, outputs=PreTrainedCNNModel.layers[-1].output)
            image = load_img(imagesSet +'/' + line, PreTrainedCNNModel, target_size=(299, 299))
            image = img_to_array(image)
            image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))
            image = preprocess_input(image)
            feature = PreTrainedCNNModel.predict(image, verbose=0)

            description = getTestDesc(trainedModel, tokenizer, feature, getLength())
            print(description)
```

6) We iterate until longest possible description available from our training and when we receive sequence terminator or when we complete iterations, we print the predicted output word so far.

Input Image:



Output Description:

```
Loaded training descriptions from all desc  6000
Total Vocabulary from training descriptions  7579
Longest sentence for RNN input layer  34
3385593926_d3e9c21170.jpg
WARNING:tensorflow:Model was constructed with shape (None, 47, 1000) for input Tensor("input_2_1:0", shape=(None, 47, 1000), dtype=float32), but it was called on an input with incompatible shape e, 1000).
startseq two dogs are playing in the grass endseq
```

**References:**

Show and Tell: A Neural Image Caption Generator
https://arxiv.org/abs/1411.4555

Geeks for geeks – Image Caption Generator
https://www.geeksforgeeks.org/image-caption-generator-using-deep-learning-on-flickr8k-dataset/