# week6_final_project_image_captioning_clean

July 28, 2021

## 1 Image Captioning Final Project

In this final project you will define and train an image-to-caption model, that can produce descriptions for real world images!

Model architecture: CNN encoder and RNN decoder. (https://research.googleblog.com/2014/11/a-picture-is-worth-thousand-coherent.html)

## 2 Import stuff

```
In [1]: import sys
        sys.path.append("..")
        import grading
        import download_utils
```

```
In [2]: download_utils.link_all_keras_resources()
```

```
In [3]: import tensorflow as tf
        from tensorflow.contrib import keras
        import numpy as np
        %matplotlib inline
        import matplotlib.pyplot as plt
        L = keras.layers
        K = keras.backend
        import utils
        import time
        import zipfile
        import json
        from collections import defaultdict
        import re
        import random
        from random import choice
        import grading_utils
        import os
        from keras_utils import reset_tf_session
        import tqdm_utils
```

```
Using TensorFlow backend.
```

## 3  Fill in your Coursera token and email

To successfully submit your answers to our grader, please fill in your Coursera submission token and email

```
In [4]: grader = grading.Grader(assignment_key="NEDBg6CgEee8nQ6uE8a7OA",
                                all_parts=["19Wpv", "uJh73", "yiJkt", "rbpnH", "E2OIL", "YJR7z"]

In [34]: # token expires every 30 min
         COURSERA_TOKEN = 'ShgxVTpxTGEGIwcd'        ### YOUR TOKEN HERE
         COURSERA_EMAIL = 'knowtech94@gmail.com'      ### YOUR EMAIL HERE
```

## 4  Download data

Takes 10 hours and 20 GB. We've downloaded necessary files for you.

   Relevant links (just in case): - train images http://msvocds.blob.core.windows.net/coco2014/train2014.zip - validation images http://msvocds.blob.core.windows.net/coco2014/val2014.zip - captions for both train and validation http://msvocds.blob.core.windows.net/annotations-1-0-3/captions_train-val2014.zip

```
In [8]: # we downloaded them for you, just link them here
        download_utils.link_week_6_resources()
```

## 5  Extract image features

We will use pre-trained InceptionV3 model for CNN encoder (https://research.googleblog.com/2016/03/train-your-own-image-classifier-with.html) and extract its last hidden layer as an embedding:

```
In [9]: IMG_SIZE = 299

In [10]: # we take the last hidden layer of IncetionV3 as an image embedding
         def get_cnn_encoder():
             K.set_learning_phase(False)
             model = keras.applications.InceptionV3(include_top=False)
             preprocess_for_model = keras.applications.inception_v3.preprocess_input

             model = keras.models.Model(model.inputs, keras.layers.GlobalAveragePooling2D()(mode
             return model, preprocess_for_model
```

   Features extraction takes too much time on CPU: - Takes 16 minutes on GPU. - 25x slower (InceptionV3) on CPU and takes 7 hours. - 10x slower (MobileNet) on CPU and takes 3 hours.

   So we've done it for you with the following code:

```
# load pre-trained model
reset_tf_session()
encoder, preprocess_for_model = get_cnn_encoder()
```

```python
# extract train features
train_img_embeds, train_img_fns = utils.apply_model(
    "train2014.zip", encoder, preprocess_for_model, input_shape=(IMG_SIZE, IMG_SIZE))
utils.save_pickle(train_img_embeds, "train_img_embeds.pickle")
utils.save_pickle(train_img_fns, "train_img_fns.pickle")

# extract validation features
val_img_embeds, val_img_fns = utils.apply_model(
    "val2014.zip", encoder, preprocess_for_model, input_shape=(IMG_SIZE, IMG_SIZE))
utils.save_pickle(val_img_embeds, "val_img_embeds.pickle")
utils.save_pickle(val_img_fns, "val_img_fns.pickle")

# sample images for learners
def sample_zip(fn_in, fn_out, rate=0.01, seed=42):
    np.random.seed(seed)
    with zipfile.ZipFile(fn_in) as fin, zipfile.ZipFile(fn_out, "w") as fout:
        sampled = filter(lambda _: np.random.rand() < rate, fin.filelist)
        for zInfo in sampled:
            fout.writestr(zInfo, fin.read(zInfo))

sample_zip("train2014.zip", "train2014_sample.zip")
sample_zip("val2014.zip", "val2014_sample.zip")
```

```python
In [11]: # load prepared embeddings
         train_img_embeds = utils.read_pickle("train_img_embeds.pickle")
         train_img_fns = utils.read_pickle("train_img_fns.pickle")
         val_img_embeds = utils.read_pickle("val_img_embeds.pickle")
         val_img_fns = utils.read_pickle("val_img_fns.pickle")
         # check shapes
         print(train_img_embeds.shape, len(train_img_fns))
         print(val_img_embeds.shape, len(val_img_fns))

(82783, 2048) 82783
(40504, 2048) 40504
```

```python
In [12]: # check prepared samples of images
         list(filter(lambda x: x.endswith("_sample.zip"), os.listdir(".")))

Out[12]: ['val2014_sample.zip', 'train2014_sample.zip']
```

## 6  Extract captions for images

```python
In [13]: # extract captions from zip
         def get_captions_for_fns(fns, zip_fn, zip_json_path):
             zf = zipfile.ZipFile(zip_fn)
             j = json.loads(zf.read(zip_json_path).decode("utf8"))
             id_to_fn = {img["id"]: img["file_name"] for img in j["images"]}
```

3

```
        fn_to_caps = defaultdict(list)
        for cap in j['annotations']:
            fn_to_caps[id_to_fn[cap['image_id']]].append(cap['caption'])
        fn_to_caps = dict(fn_to_caps)
        return list(map(lambda x: fn_to_caps[x], fns))

    train_captions = get_captions_for_fns(train_img_fns, "captions_train-val2014.zip",
                                           "annotations/captions_train2014.json")

    val_captions = get_captions_for_fns(val_img_fns, "captions_train-val2014.zip",
                                         "annotations/captions_val2014.json")

    # check shape
    print(len(train_img_fns), len(train_captions))
    print(len(val_img_fns), len(val_captions))

82783 82783
40504 40504


In [14]: # look at training example (each has 5 captions)
    def show_trainig_example(train_img_fns, train_captions, example_idx=0):
        """
        You can change example_idx and see different images
        """
        zf = zipfile.ZipFile("train2014_sample.zip")
        captions_by_file = dict(zip(train_img_fns, train_captions))
        all_files = set(train_img_fns)
        found_files = list(filter(lambda x: x.filename.rsplit("/")[-1] in all_files, zf.fil
        example = found_files[example_idx]
        img = utils.decode_image_from_buf(zf.read(example))
        plt.imshow(utils.image_center_crop(img))
        plt.title("\n".join(captions_by_file[example.filename.rsplit("/")[-1]]))
        plt.show()

    show_trainig_example(train_img_fns, train_captions, example_idx=142)
```
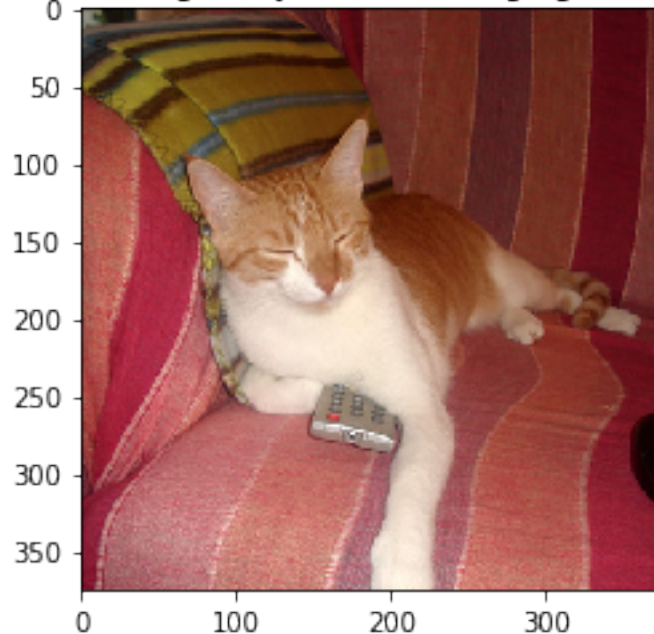
4

A cat sitting on a pink stripped couch
An orange and white cat sitting in a striped chair.
An orange and white cat sleeping on a remote
Brown and white cat sleeping on couch while lying on remote.
A cat closing its eyes while lounging on a chair.



# 7 Prepare captions for training

```
In [15]:  # preview captions data
          train_captions[:2]
```

```
Out[15]:  [['A long dirt road going through a forest.',
            'A SCENE OF WATER AND A PATH WAY',
            'A sandy path surrounded by trees leads to a beach.',
            'Ocean view through a dirt road surrounded by a forested area. ',
            'dirt path leading beneath barren trees to open plains'],
           ['A group of zebra standing next to each other.',
            'This is an image of of zebras drinking',
            'ZEBRAS AND BIRDS SHARING THE SAME WATERING HOLE',
            'Zebras that are bent over and drinking water together.',
            'a number of zebras drinking water near one another']]
```

```
In [16]:  # special tokens
          PAD = "#PAD#"
          UNK = "#UNK#"
          START = "#START#"
```

```python
END = "#END#"

# split sentence into tokens (split into lowercased words)
def split_sentence(sentence):
    return list(filter(lambda x: len(x) > 0, re.split('\W+', sentence.lower())))


def generate_vocabulary(train_captions):
    """
    Return {token: index} for all train tokens (words) that occur 5 times or more,
        `index` should be from 0 to N, where N is a number of unique tokens in the resu
    Use `split_sentence` function to split sentence into tokens.
    Also, add PAD (for batch padding), UNK (unknown, out of vocabulary),
        START (start of sentence) and END (end of sentence) tokens into the vocabulary.
    """
    #vocab = ### YOUR CODE HERE ###

    words = [words for captions in train_captions for caption in captions for words in
    words = np.reshape(words, -1)
    unique_words = np.unique(words)
    word_count = {word: 0 for word in unique_words}
    for word in words:
      word_count[word] = word_count[word] + 1

    vocab = [w for w, c in word_count.items() if c >= 5]
    vocab.append(PAD)
    vocab.append(UNK)
    vocab.append(START)
    vocab.append(END)

    return {pal:count for count, pal in enumerate(sorted(vocab))}

def caption_tokens_to_indices(captions, vocab):
    """
    `captions` argument is an array of arrays:
    [
        [
            "image1 caption1",
            "image1 caption2",
            ...
        ],
        [
            "image2 caption1",
            "image2 caption2",
            ...
        ],
        ...
    ]
    Use `split_sentence` function to split sentence into tokens.
```

```
            Replace all tokens with vocabulary indices, use UNK for unknown words (out of vocab
            Add START and END tokens to start and end of each sentence respectively.
            For the example above you should produce the following:
            [
                [
                    [vocab[START], vocab["image1"], vocab["caption1"], vocab[END]],
                    [vocab[START], vocab["image1"], vocab["caption2"], vocab[END]],
                    ...
                ],
                ...
            ]
            """
            #res = ### YOUR CODE HERE ###

            vocab_words = lambda words: [vocab[START]] + [vocab[word] if word in vocab else voc
            total_words = lambda captions: [vocab_words(split_sentence(caption.lower())) for ca
            response = [total_words(caption_list) for caption_list in captions]

            return response
```

```
In [17]: # prepare vocabulary
         vocab = generate_vocabulary(train_captions)
         vocab_inverse = {idx: w for w, idx in vocab.items()}
         print(len(vocab))
```

8769

```
In [18]: # replace tokens with indices
         train_captions_indexed = caption_tokens_to_indices(train_captions, vocab)
         val_captions_indexed = caption_tokens_to_indices(val_captions, vocab)
```

Captions have different length, but we need to batch them, that's why we will add PAD tokens so that all sentences have an equal length.

We will crunch LSTM through all the tokens, but we will ignore padding tokens during loss calculation.

```
In [19]: # we will use this during training
         def batch_captions_to_matrix(batch_captions, pad_idx, max_len=None):
             """
             `batch_captions` is an array of arrays:
             [
                 [vocab[START], ..., vocab[END]],
                 [vocab[START], ..., vocab[END]],
                 ...
             ]
             Put vocabulary indexed captions into np.array of shape (len(batch_captions), column
                 where "columns" is max(map(len, batch_captions)) when max_len is None
                 and "columns" = min(max_len, max(map(len, batch_captions))) otherwise.
```

```python
            Add padding with pad_idx where necessary.
            Input example: [[1, 2, 3], [4, 5]]
            Output example: np.array([[1, 2, 3], [4, 5, pad_idx]]) if max_len=None
            Output example: np.array([[1, 2], [4, 5]]) if max_len=2
            Output example: np.array([[1, 2, 3], [4, 5, pad_idx]]) if max_len=100
            Try to use numpy, we need this function to be fast!
            """
            #matrix = ###YOUR CODE HERE###

            max_len = max_len or max(map(len, batch_captions))
            max_len = min(max_len, max(map(len, batch_captions)))
            matrix = np.empty((len(batch_captions), max_len))
            matrix.fill(pad_idx)
            for i in range(len(batch_captions)):
                line_ix = list(batch_captions[i])[:max_len]
                matrix[i,:len(line_ix)] = line_ix

            return matrix

In [20]: ## GRADED PART, DO NOT CHANGE!
         # Vocabulary creation
         grader.set_answer("19Wpv", grading_utils.test_vocab(vocab, PAD, UNK, START, END))
         # Captions indexing
         grader.set_answer("uJh73", grading_utils.test_captions_indexing(train_captions_indexed,
         # Captions batching
         grader.set_answer("yiJkt", grading_utils.test_captions_batching(batch_captions_to_matri

In [21]: # you can make submission with answers so far to check yourself at this stage
         grader.submit(COURSERA_EMAIL, COURSERA_TOKEN)

Submitted to Coursera platform. See results on assignment page!


In [22]: # make sure you use correct argument in caption_tokens_to_indices
         assert len(caption_tokens_to_indices(train_captions[:10], vocab)) == 10
         assert len(caption_tokens_to_indices(train_captions[:5], vocab)) == 5
```

# 8 Training

## 8.1 Define architecture

Since our problem is to generate image captions, RNN text generator should be conditioned on image. The idea is to use image features as an initial state for RNN instead of zeros.

Remember that you should transform image feature vector to RNN hidden state size by fully-connected layer and then pass it to RNN.

During training we will feed ground truth tokens into the lstm to get predictions of next tokens.

Notice that we don't need to feed last token (END) as input (http://cs.stanford.edu/people/karpathy/):

```
In [23]: IMG_EMBED_SIZE = train_img_embeds.shape[1]
         IMG_EMBED_BOTTLENECK = 120
         WORD_EMBED_SIZE = 100
         LSTM_UNITS = 300
         LOGIT_BOTTLENECK = 120
         pad_idx = vocab[PAD]

In [24]: # remember to reset your graph if you want to start building it from scratch!
         s = reset_tf_session()
         tf.set_random_seed(42)
```

Here we define decoder graph.

We use Keras layers where possible because we can use them in functional style with weights reuse like this:

```
dense_layer = L.Dense(42, input_shape=(None, 100) activation='relu')
a = tf.placeholder('float32', [None, 100])
b = tf.placeholder('float32', [None, 100])
dense_layer(a)  # that's how we applied dense layer!
dense_layer(b)  # and again
```

Here's a figure to help you with flattening in decoder:

```
In [25]: class decoder:
             # [batch_size, IMG_EMBED_SIZE] of CNN image features
             img_embeds = tf.placeholder('float32', [None, IMG_EMBED_SIZE])
             # [batch_size, time steps] of word ids
             sentences = tf.placeholder('int32', [None, None])

             # we use bottleneck here to reduce the number of parameters
             # image embedding -> bottleneck
             img_embed_to_bottleneck = L.Dense(IMG_EMBED_BOTTLENECK,
                                         input_shape=(None, IMG_EMBED_SIZE),
                                         activation='elu')
             # image embedding bottleneck -> lstm initial state
             img_embed_bottleneck_to_h0 = L.Dense(LSTM_UNITS,
                                         input_shape=(None, IMG_EMBED_BOTTLENECK),
                                         activation='elu')
             # word -> embedding
             word_embed = L.Embedding(len(vocab), WORD_EMBED_SIZE)
             # lstm cell (from tensorflow)
             lstm = tf.nn.rnn_cell.LSTMCell(LSTM_UNITS)

             # we use bottleneck here to reduce model complexity
             # lstm output -> logits bottleneck
             token_logits_bottleneck = L.Dense(LOGIT_BOTTLENECK,
                                         input_shape=(None, LSTM_UNITS),
                                         activation="elu")
             # logits bottleneck -> logits for next token prediction
```

9

```python
token_logits = L.Dense(len(vocab),
                       input_shape=(None, LOGIT_BOTTLENECK))

# initial lstm cell state of shape (None, LSTM_UNITS),
# we need to condition it on `img_embeds` placeholder.
#c0 = h0 = ### YOUR CODE HERE ###
c0 = h0 = img_embed_bottleneck_to_h0(img_embed_to_bottleneck(img_embeds))

# embed all tokens but the last for lstm input,
# remember that L.Embedding is callable,
# use `sentences` placeholder as input.
#word_embeds = ### YOUR CODE HERE ###
word_embeds = word_embed(sentences[:, :-1])

# during training we use ground truth tokens `word_embeds` as context for next toke
# that means that we know all the inputs for our lstm and can get
# all the hidden states with one tensorflow operation (tf.nn.dynamic_rnn).
# `hidden_states` has a shape of [batch_size, time steps, LSTM_UNITS].
hidden_states, _ = tf.nn.dynamic_rnn(lstm, word_embeds,
                                     initial_state=tf.nn.rnn_cell.LSTMStateTuple(c0

# now we need to calculate token logits for all the hidden states

# first, we reshape `hidden_states` to [-1, LSTM_UNITS]
#flat_hidden_states = ### YOUR CODE HERE ###
flat_hidden_states = tf.reshape(hidden_states, [-1, LSTM_UNITS])

# then, we calculate logits for next tokens using `token_logits_bottleneck` and `to
#flat_token_logits = ### YOUR CODE HERE ###
flat_token_logits = token_logits(token_logits_bottleneck(flat_hidden_states))

# then, we flatten the ground truth token ids.
# remember, that we predict next tokens for each time step,
# use `sentences` placeholder.
#flat_ground_truth = ### YOUR CODE HERE ###
flat_ground_truth = tf.reshape(sentences[:, 1:], [-1])

# we need to know where we have real tokens (not padding) in `flat_ground_truth`,
# we don't want to propagate the loss for padded output tokens,
# fill `flat_loss_mask` with 1.0 for real tokens (not pad_idx) and 0.0 otherwise.
#flat_loss_mask = ### YOUR CODE HERE ###
flat_loss_mask = tf.not_equal(flat_ground_truth, pad_idx)

# compute cross-entropy between `flat_ground_truth` and `flat_token_logits` predict
xent = tf.nn.sparse_softmax_cross_entropy_with_logits(
    labels=flat_ground_truth,
    logits=flat_token_logits
)
```

```python
        # compute average `xent` over tokens with nonzero `flat_loss_mask`.
        # we don't want to account misclassification of PAD tokens, because that doesn't ma
        # we have PAD tokens for batching purposes only!
        #loss = ### YOUR CODE HERE ###
        loss = tf.reduce_mean(tf.boolean_mask(xent, flat_loss_mask))
```

```
In [26]:  # define optimizer operation to minimize the loss
          optimizer = tf.train.AdamOptimizer(learning_rate=0.001)
          train_step = optimizer.minimize(decoder.loss)

          # will be used to save/load network weights.
          # you need to reset your default graph and define it in the same way to be able to load
          saver = tf.train.Saver()

          # intialize all variables
          s.run(tf.global_variables_initializer())
```

```
/opt/conda/lib/python3.6/site-packages/tensorflow/python/ops/gradients_impl.py:93: UserWarning:
  "Converting sparse IndexedSlices to a dense Tensor of unknown shape. "
```

```
In [27]:  ## GRADED PART, DO NOT CHANGE!
          # Decoder shapes test
          grader.set_answer("rbpnH", grading_utils.test_decoder_shapes(decoder, IMG_EMBED_SIZE, v
          # Decoder random loss test
          grader.set_answer("E2OIL", grading_utils.test_random_decoder_loss(decoder, IMG_EMBED_SI
```

```
In [28]:  # you can make submission with answers so far to check yourself at this stage
          grader.submit(COURSERA_EMAIL, COURSERA_TOKEN)
```

```
Submitted to Coursera platform. See results on assignment page!
```

## 8.2 Training loop

Evaluate train and validation metrics through training and log them. Ensure that loss decreases.

```
In [29]:  train_captions_indexed = np.array(train_captions_indexed)
          val_captions_indexed = np.array(val_captions_indexed)
```

```
In [30]:  # generate batch via random sampling of images and captions for them,
          # we use `max_len` parameter to control the length of the captions (truncating long cap
          def generate_batch(images_embeddings, indexed_captions, batch_size, max_len=None):
              """
              `images_embeddings` is a np.array of shape [number of images, IMG_EMBED_SIZE].
              `indexed_captions` holds 5 vocabulary indexed captions for each image:
              [
```

```
                [
                    [vocab[START], vocab["image1"], vocab["caption1"], vocab[END]],
                    [vocab[START], vocab["image1"], vocab["caption2"], vocab[END]],
                    ...
                ],
                ...
            ]
            Generate a random batch of size `batch_size`.
            Take random images and choose one random caption for each image.
            Remember to use `batch_captions_to_matrix` for padding and respect `max_len` parame
            Return feed dict {decoder.img_embeds: ..., decoder.sentences: ...}.
            """
            #batch_image_embeddings = ### YOUR CODE HERE ###
            batch_start_idx = np.random.randint(0, high = len(images_embeddings) - batch_size)
            idxs = np.random.permutation(range(len(images_embeddings)))[batch_start_idx: batch_
            batch_image_embeddings = np.take(images_embeddings, idxs, axis = 0)

            #batch_captions_matrix = ### YOUR CODE HERE ###
            caption_inner_idxs = np.random.choice(range(5), batch_size)
            batch_captions_matrix = [indexed_captions[m][n] for m, n in zip(idxs, caption_inner
            batch_captions_matrix = batch_captions_to_matrix(batch_captions_matrix, pad_idx, ma

            return {decoder.img_embeds: batch_image_embeddings,
                    decoder.sentences: batch_captions_matrix}

In [32]: batch_size = 64
         n_epochs = 10
         n_batches_per_epoch = 1000
         n_validation_batches = 100  # how many batches are used for validation after each epoch

In [30]: # you can load trained weights here
         # you can load "weights_{epoch}" and continue training
         # uncomment the next line if you need to load weights
         # saver.restore(s, os.path.abspath("weights"))
```

Look at the training and validation loss, they should be decreasing!

```
In [33]: # actual training loop
         MAX_LEN = 20  # truncate long captions to speed up training

         # to make training reproducible
         np.random.seed(42)
         random.seed(42)

         for epoch in range(n_epochs):

             train_loss = 0
             pbar = tqdm_utils.tqdm_notebook_failsafe(range(n_batches_per_epoch))
             counter = 0
```

```python
            for _ in pbar:
                train_loss += s.run([decoder.loss, train_step],
                                generate_batch(train_img_embeds,
                                            train_captions_indexed,
                                            batch_size,
                                            MAX_LEN))[0]
                counter += 1
                pbar.set_description("Training loss: %f" % (train_loss / counter))

            train_loss /= n_batches_per_epoch

            val_loss = 0
            for _ in range(n_validation_batches):
                val_loss += s.run(decoder.loss, generate_batch(val_img_embeds,
                                                    val_captions_indexed,
                                                    batch_size,
                                                    MAX_LEN))
            val_loss /= n_validation_batches

            print('Epoch: {}, train loss: {}, val loss: {}'.format(epoch, train_loss, val_loss)

            # save weights after finishing epoch
            saver.save(s, os.path.abspath("weights_{}".format(epoch)))

        print("Finished!")
```

A Jupyter Widget


Epoch: 0, train loss: 4.291811956882476, val loss: 3.625504195690155


A Jupyter Widget


Epoch: 1, train loss: 3.3434647524356844, val loss: 3.166829035282135


A Jupyter Widget


Epoch: 2, train loss: 3.0545015444755554, val loss: 2.9616402840614318


A Jupyter Widget

Epoch: 3, train loss: 2.903541965007782, val loss: 2.8957045578956606

A Jupyter Widget

Epoch: 4, train loss: 2.805233480453491, val loss: 2.825015242099762

A Jupyter Widget

Epoch: 5, train loss: 2.7360985076427458, val loss: 2.7711693978309633

A Jupyter Widget

Epoch: 6, train loss: 2.692197933912277, val loss: 2.7144248127937316

A Jupyter Widget

Epoch: 7, train loss: 2.6467096557617187, val loss: 2.6921755361557005

A Jupyter Widget

Epoch: 8, train loss: 2.6048688378334046, val loss: 2.674499177932739

A Jupyter Widget

Epoch: 9, train loss: 2.579082692861557, val loss: 2.6435467195510864
Finished!

```
In [35]: ## GRADED PART, DO NOT CHANGE!
         # Validation loss
         grader.set_answer("YJR7z", grading_utils.test_validation_loss(
             decoder, s, generate_batch, val_img_embeds, val_captions_indexed))
```

A Jupyter Widget

```
In [36]: # you can make submission with answers so far to check yourself at this stage
         grader.submit(COURSERA_EMAIL, COURSERA_TOKEN)
```

Submitted to Coursera platform. See results on assignment page!

```
In [37]: # check that it's learnt something, outputs accuracy of next word prediction (should be
         from sklearn.metrics import accuracy_score, log_loss

         def decode_sentence(sentence_indices):
             return " ".join(list(map(vocab_inverse.get, sentence_indices)))

         def check_after_training(n_examples):
             fd = generate_batch(train_img_embeds, train_captions_indexed, batch_size)
             logits = decoder.flat_token_logits.eval(fd)
             truth = decoder.flat_ground_truth.eval(fd)
             mask = decoder.flat_loss_mask.eval(fd).astype(bool)
             print("Loss:", decoder.loss.eval(fd))
             print("Accuracy:", accuracy_score(logits.argmax(axis=1)[mask], truth[mask]))
             for example_idx in range(n_examples):
                 print("Example", example_idx)
                 print("Predicted:", decode_sentence(logits.argmax(axis=1).reshape((batch_size,
                 print("Truth:", decode_sentence(truth.reshape((batch_size, -1))[example_idx]))
                 print("")

         check_after_training(3)
```

Loss: 2.74344
Accuracy: 0.444600280505
Example 0
Predicted: a woman eating and a bowl of a of a pizza #END# food #END# #END# #END# #END# #END#
Truth: a kid drink from a glass in front of a food of plate #END# #PAD# #PAD# #PAD# #PAD#

Example 1
Predicted: a toilets of a on a store on a middle #END# #END# #END# #END# #END# #END# #END# #END#
Truth: some bottles of cleaner in a spot in the wall #END# #PAD# #PAD# #PAD# #PAD# #PAD# #PAD# #

Example 2
Predicted: a bananas and on a outdoor plate box #END# a table plate #END# #END# #END# #END# #END
Truth: sliced bananas sit on an open faced sandwich on a paper plate #END# #PAD# #PAD# #PAD# #PA

```
In [38]: # save graph weights to file!
         saver.save(s, os.path.abspath("weights"))

Out[38]: '/home/jovyan/work/week6/weights'
```

# 9   Applying model

Here we construct a graph for our final model.
    It will work as follows: - take an image as an input and embed it - condition lstm on that
embedding - predict the next token given a START input token - use predicted token as an input
at next time step - iterate until you predict an END token

```
In [39]: class final_model:
             # CNN encoder
             encoder, preprocess_for_model = get_cnn_encoder()
             saver.restore(s, os.path.abspath("weights"))  # keras applications corrupt our grap

             # containers for current lstm state
             lstm_c = tf.Variable(tf.zeros([1, LSTM_UNITS]), name="cell")
             lstm_h = tf.Variable(tf.zeros([1, LSTM_UNITS]), name="hidden")

             # input images
             input_images = tf.placeholder('float32', [1, IMG_SIZE, IMG_SIZE, 3], name='images')

             # get image embeddings
             img_embeds = encoder(input_images)

             # initialize lstm state conditioned on image
             init_c = init_h = decoder.img_embed_bottleneck_to_h0(decoder.img_embed_to_bottlenec
             init_lstm = tf.assign(lstm_c, init_c), tf.assign(lstm_h, init_h)

             # current word index
             current_word = tf.placeholder('int32', [1], name='current_input')

             # embedding for current word
             word_embed = decoder.word_embed(current_word)

             # apply lstm cell, get new lstm states
             new_c, new_h = decoder.lstm(word_embed, tf.nn.rnn_cell.LSTMStateTuple(lstm_c, lstm_

             # compute logits for next token
             new_logits = decoder.token_logits(decoder.token_logits_bottleneck(new_h))
             # compute probabilities for next token
             new_probs = tf.nn.softmax(new_logits)

             # `one_step` outputs probabilities of next token and updates lstm hidden state
             one_step = new_probs, tf.assign(lstm_c, new_c), tf.assign(lstm_h, new_h)
```

16

```
INFO:tensorflow:Restoring parameters from /home/jovyan/work/week6/weights


In [40]: # look at how temperature works for probability distributions
         # for high temperature we have more uniform distribution
         _ = np.array([0.5, 0.4, 0.1])
         for t in [0.01, 0.1, 1, 10, 100]:
             print(" ".join(map(str, _**(1/t) / np.sum(_**(1/t)))), "with temperature", t)

0.999999999796 2.03703597592e-10 1.26765059997e-70 with temperature 0.01
0.903037043325 0.0969628642039 9.24709932365e-08 with temperature 0.1
0.5 0.4 0.1 with temperature 1
0.353447726392 0.345648113606 0.300904160002 with temperature 10
0.335367280481 0.334619764349 0.33001295517 with temperature 100


In [41]: # this is an actual prediction loop
         def generate_caption(image, t=1, sample=False, max_len=20):
             """
             Generate caption for given image.
             if `sample` is True, we will sample next token from predicted probability distribut
             `t` is a temperature during that sampling,
                 higher `t` causes more uniform-like distribution = more chaos.
             """
             # condition lstm on the image
             s.run(final_model.init_lstm,
                   {final_model.input_images: [image]})

             # current caption
             # start with only START token
             caption = [vocab[START]]

             for _ in range(max_len):
                 next_word_probs = s.run(final_model.one_step,
                                         {final_model.current_word: [caption[-1]]})[0]
                 next_word_probs = next_word_probs.ravel()

                 # apply temperature
                 next_word_probs = next_word_probs**(1/t) / np.sum(next_word_probs**(1/t))

                 if sample:
                     next_word = np.random.choice(range(len(vocab)), p=next_word_probs)
                 else:
                     next_word = np.argmax(next_word_probs)

                 caption.append(next_word)
                 if next_word == vocab[END]:
                     break
```

```
            return list(map(vocab_inverse.get, caption))

In [42]:  # look at validation prediction example
          def apply_model_to_image_raw_bytes(raw):
              img = utils.decode_image_from_buf(raw)
              fig = plt.figure(figsize=(7, 7))
              plt.grid('off')
              plt.axis('off')
              plt.imshow(img)
              img = utils.crop_and_preprocess(img, (IMG_SIZE, IMG_SIZE), final_model.preprocess_f
              print(' '.join(generate_caption(img)[1:-1]))
              plt.show()


          def show_valid_example(val_img_fns, example_idx=0):
              zf = zipfile.ZipFile("val2014_sample.zip")
              all_files = set(val_img_fns)
              found_files = list(filter(lambda x: x.filename.rsplit("/")[-1] in all_files, zf.fil
              example = found_files[example_idx]
              apply_model_to_image_raw_bytes(zf.read(example))


          show_valid_example(val_img_fns, example_idx=100)

a baseball player holding a bat on a field
```

```
In [43]: # sample more images from validation
         for idx in np.random.choice(range(len(zipfile.ZipFile("val2014_sample.zip").filelist) -
             show_valid_example(val_img_fns, example_idx=idx)
             time.sleep(1)
```

a person on a ski slope with a ski lift



a woman holding a hot dog in front of a sandwich

a man is holding a cell phone in his hand

a large piece of cake with a large number of toppings

a bathroom with a sink and a sink

a bathroom with a sink and a sink

a man is throwing a frisbee in a park

a baseball player is getting ready to swing at a pitch

a piece of cake with a fork and a fork

a man is walking down a street with a surfboard

You can download any image from the Internet and appply your model to it!

```
In [49]: download_utils.download_file(
            "https://h7f7z2r7.stackpathcdn.com/sites/default/files/images/articles/newyorkmain_
            "newyorkmain_0.jpg"
         )

HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443): Max retries exceeded with url:


Traceback (most recent call last):
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 595, in urlopen
    self._prepare_proxy(conn)
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 816, in _prepare
    conn.connect()
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connection.py", line 294, in connect
    self._tunnel()
  File "/opt/conda/lib/python3.6/http/client.py", line 919, in _tunnel
    message.strip()))
OSError: Tunnel connection failed: 403 Forbidden

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "/opt/conda/lib/python3.6/site-packages/requests/adapters.py", line 440, in send
    timeout=timeout
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 639, in urlopen
    _stacktrace=sys.exc_info()[2])
  File "/opt/conda/lib/python3.6/site-packages/urllib3/util/retry.py", line 388, in increment
    raise MaxRetryError(_pool, url, error or ResponseError(cause))
urllib3.exceptions.MaxRetryError: HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "../download_utils.py", line 21, in f_retry
    return f(*args, **kwargs)
  File "../download_utils.py", line 39, in download_file
    r = requests.get(url, stream=True)
  File "/opt/conda/lib/python3.6/site-packages/requests/api.py", line 72, in get
    return request('get', url, params=params, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/api.py", line 58, in request
    return session.request(method=method, url=url, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/sessions.py", line 508, in request
    resp = self.send(prep, **send_kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/sessions.py", line 618, in send
```

```
    r = adapter.send(request, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/adapters.py", line 502, in send
    raise ProxyError(e, request=request)
requests.exceptions.ProxyError: HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443):


HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443): Max retries exceeded with url:


Traceback (most recent call last):
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 595, in urlopen
    self._prepare_proxy(conn)
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 816, in _prepare
    conn.connect()
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connection.py", line 294, in connect
    self._tunnel()
  File "/opt/conda/lib/python3.6/http/client.py", line 919, in _tunnel
    message.strip()))
OSError: Tunnel connection failed: 403 Forbidden

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "/opt/conda/lib/python3.6/site-packages/requests/adapters.py", line 440, in send
    timeout=timeout
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 639, in urlopen
    _stacktrace=sys.exc_info()[2])
  File "/opt/conda/lib/python3.6/site-packages/urllib3/util/retry.py", line 388, in increment
    raise MaxRetryError(_pool, url, error or ResponseError(cause))
urllib3.exceptions.MaxRetryError: HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "../download_utils.py", line 21, in f_retry
    return f(*args, **kwargs)
  File "../download_utils.py", line 39, in download_file
    r = requests.get(url, stream=True)
  File "/opt/conda/lib/python3.6/site-packages/requests/api.py", line 72, in get
    return request('get', url, params=params, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/api.py", line 58, in request
    return session.request(method=method, url=url, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/sessions.py", line 508, in request
    resp = self.send(prep, **send_kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/sessions.py", line 618, in send
    r = adapter.send(request, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/adapters.py", line 502, in send
    raise ProxyError(e, request=request)
```

```
requests.exceptions.ProxyError: HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443):


HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443): Max retries exceeded with url:


Traceback (most recent call last):
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 595, in urlopen
    self._prepare_proxy(conn)
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 816, in _prepare
    conn.connect()
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connection.py", line 294, in connect
    self._tunnel()
  File "/opt/conda/lib/python3.6/http/client.py", line 919, in _tunnel
    message.strip()))
OSError: Tunnel connection failed: 403 Forbidden

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "/opt/conda/lib/python3.6/site-packages/requests/adapters.py", line 440, in send
    timeout=timeout
  File "/opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py", line 639, in urlopen
    _stacktrace=sys.exc_info()[2])
  File "/opt/conda/lib/python3.6/site-packages/urllib3/util/retry.py", line 388, in increment
    raise MaxRetryError(_pool, url, error or ResponseError(cause))
urllib3.exceptions.MaxRetryError: HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "../download_utils.py", line 21, in f_retry
    return f(*args, **kwargs)
  File "../download_utils.py", line 39, in download_file
    r = requests.get(url, stream=True)
  File "/opt/conda/lib/python3.6/site-packages/requests/api.py", line 72, in get
    return request('get', url, params=params, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/api.py", line 58, in request
    return session.request(method=method, url=url, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/sessions.py", line 508, in request
    resp = self.send(prep, **send_kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/sessions.py", line 618, in send
    r = adapter.send(request, **kwargs)
  File "/opt/conda/lib/python3.6/site-packages/requests/adapters.py", line 502, in send
    raise ProxyError(e, request=request)
requests.exceptions.ProxyError: HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443):
```

```
    -------------------------------------------------------------------------

    OSError                                       Traceback (most recent call last)

    /opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py in urlopen(self, method
    594               if is_new_proxy_conn:
--> 595                   self._prepare_proxy(conn)
    596


    /opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py in _prepare_proxy(self,
    815
--> 816           conn.connect()
    817


    /opt/conda/lib/python3.6/site-packages/urllib3/connection.py in connect(self)
    293               # self._tunnel_host below.
--> 294               self._tunnel()
    295               # Mark this connection as not reusable


    /opt/conda/lib/python3.6/http/client.py in _tunnel(self)
    918               raise OSError("Tunnel connection failed: %d %s" % (code,
--> 919                                                     message.strip()))
    920           while True:


    OSError: Tunnel connection failed: 403 Forbidden


During handling of the above exception, another exception occurred:


    MaxRetryError                                 Traceback (most recent call last)

    /opt/conda/lib/python3.6/site-packages/requests/adapters.py in send(self, request, strea
    439                   retries=self.max_retries,
--> 440                   timeout=timeout
    441               )


    /opt/conda/lib/python3.6/site-packages/urllib3/connectionpool.py in urlopen(self, method
    638               retries = retries.increment(method, url, error=e, _pool=self,
--> 639                                           _stacktrace=sys.exc_info()[2])
    640               retries.sleep()
```

```
    /opt/conda/lib/python3.6/site-packages/urllib3/util/retry.py in increment(self, method,
    387         if new_retry.is_exhausted():
--> 388             raise MaxRetryError(_pool, url, error or ResponseError(cause))
    389


    MaxRetryError: HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443): Max retr


During handling of the above exception, another exception occurred:


    ProxyError                                Traceback (most recent call last)

    <ipython-input-49-3d23f68a2091> in <module>()
      1 download_utils.download_file(
      2     "https://h7f7z2r7.stackpathcdn.com/sites/default/files/images/articles/newyorkma
----> 3     "newyorkmain_0.jpg"
      4 )


    ~/work/download_utils.py in f_retry(*args, **kwargs)
     28                 mtries -= 1
     29                 mdelay *= backoff
---> 30             return f(*args, **kwargs)
     31
     32     return f_retry  # true decorator


    ~/work/download_utils.py in download_file(url, file_path)
     37 @retry(Exception)
     38 def download_file(url, file_path):
---> 39     r = requests.get(url, stream=True)
     40     total_size = int(r.headers.get('content-length'))
     41     bar = tqdm_utils.tqdm_notebook_failsafe(total=total_size, unit='B', unit_scale=T


    /opt/conda/lib/python3.6/site-packages/requests/api.py in get(url, params, **kwargs)
     70
     71     kwargs.setdefault('allow_redirects', True)
---> 72     return request('get', url, params=params, **kwargs)
     73
     74


    /opt/conda/lib/python3.6/site-packages/requests/api.py in request(method, url, **kwargs)
     56         # cases, and look like a memory leak in others.
     57         with sessions.Session() as session:
```

```
--> 58            return session.request(method=method, url=url, **kwargs)
    59
    60


/opt/conda/lib/python3.6/site-packages/requests/sessions.py in request(self, method, url
    506          }
    507          send_kwargs.update(settings)
--> 508          resp = self.send(prep, **send_kwargs)
    509
    510          return resp


/opt/conda/lib/python3.6/site-packages/requests/sessions.py in send(self, request, **kwa
    616
    617          # Send the request
--> 618          r = adapter.send(request, **kwargs)
    619
    620          # Total elapsed time of the request (approximately)


/opt/conda/lib/python3.6/site-packages/requests/adapters.py in send(self, request, strea
    500
    501              if isinstance(e.reason, _ProxyError):
--> 502                  raise ProxyError(e, request=request)
    503
    504              if isinstance(e.reason, _SSLError):


ProxyError: HTTPSConnectionPool(host='h7f7z2r7.stackpathcdn.com', port=443): Max retries


In [47]: apply_model_to_image_raw_bytes(open("portal-cake-10.jpg", "rb").read())


    ---------------------------------------------------------------------------

    FileNotFoundError                         Traceback (most recent call last)

    <ipython-input-47-1e36111809c8> in <module>()
 ----> 1 apply_model_to_image_raw_bytes(open("http://www.bijouxandbits.com/wp-content/uploads


    FileNotFoundError: [Errno 2] No such file or directory: 'http://www.bijouxandbits.com/wp
```

Now it's time to find 10 examples where your model works good and 10 examples where it fails!

You can use images from validation set as follows:

```
show_valid_example(val_img_fns, example_idx=...)
```

You can use images from the Internet as follows:

```
! wget ...
apply_model_to_image_raw_bytes(open("...", "rb").read())
```

If you use these functions, the output will be embedded into your notebook and will be visible during peer review!

When you're done, download your noteboook using "File" -> "Download as" -> "Notebook" and prepare that file for peer review!

In [ ]: *### YOUR EXAMPLES HERE ###*

That's it!

Congratulations, you've trained your image captioning model and now can produce captions for any picture from the Internet!