

Submitted by

Harish Tata

Class ID: 25

Team 4

Objectives:

- Generate captions for your own dataset using the show and Tellmodel.
- Generate 4 captions for each image. (Beam Search k=4)
- Report your accuracy in BLEU, CIDER, METEOR and ROGUE measures.

Technologies used:

1. TensorFlow
2. Pandas

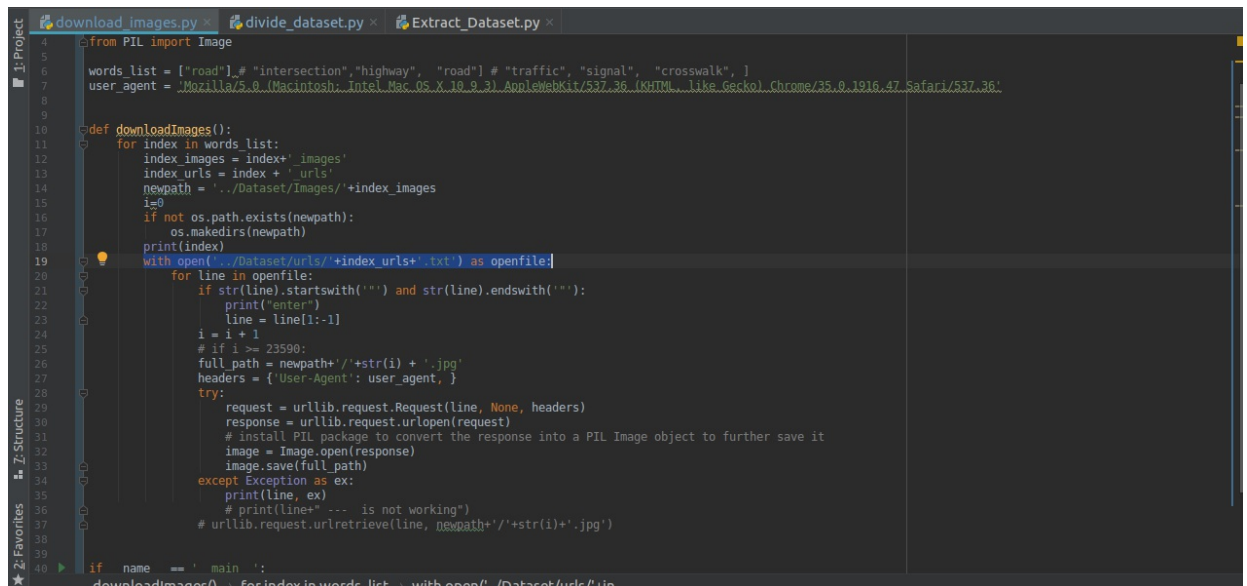
Tools Used:

Pycharm

Dataset

- Downloaded the “Google Conceptual Captions” data set.

Filtered the data set with the keywords used for our project theme



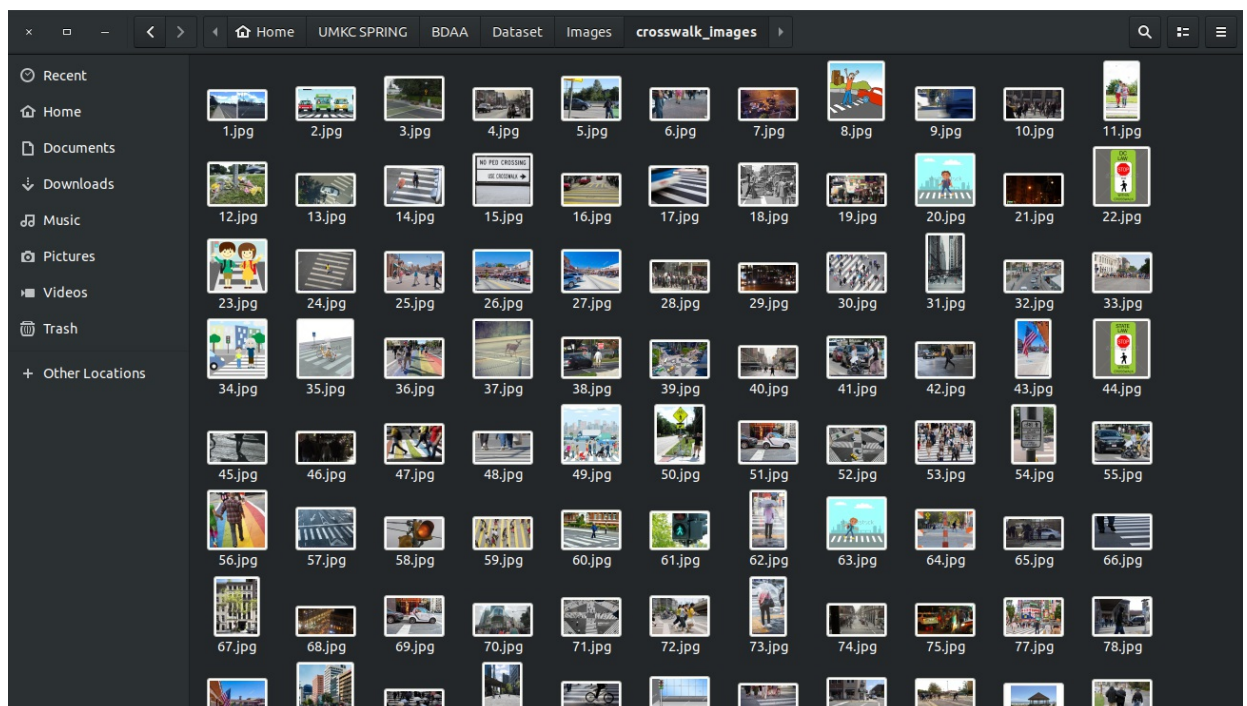
```

1 from PIL import Image
2
3 words_list = ["road", "intersection", "highway", "road"] # "traffic", "signal", "crosswalk", ]
4 user_agent = 'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_9_3) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/35.0.1916.47 Safari/537.36'
5
6
7
8
9
10 def downloadImages():
11     for index in words_list:
12         index_images = index+' images'
13         index_urls = index + ' urls'
14         newpath = '../Dataset/Images/'+index_images
15         i=0
16         if not os.path.exists(newpath):
17             os.makedirs(newpath)
18         print(index)
19         with open('../Dataset/urls/'+index_urls+'.txt') as openfile:
20             for line in openfile:
21                 if str(line).startswith('') and str(line).endswith(''):
22                     print("enter")
23                     line = line[1:-1]
24                     i = i + 1
25                     # if i >= 23590:
26                     full_path = newpath+'/'+str(i) + '.jpg'
27                     headers = {'User-Agent': user_agent, }
28                     try:
29                         request = urllib.request.Request(line, None, headers)
30                         response = urllib.request.urlopen(request)
31                         # install PIL package to convert the response into a PIL Image object to further save it
32                         image = Image.open(response)
33                         image.save(full_path)
34                     except Exception as ex:
35                         print(line, ex)
36                         # print(line+" --- is not working")
37                         # urllib.request.urlretrieve(line, newpath+'/'+str(i)+'.jpg')
38
39
40 if __name__ == '__main__':
41     downloadImages()

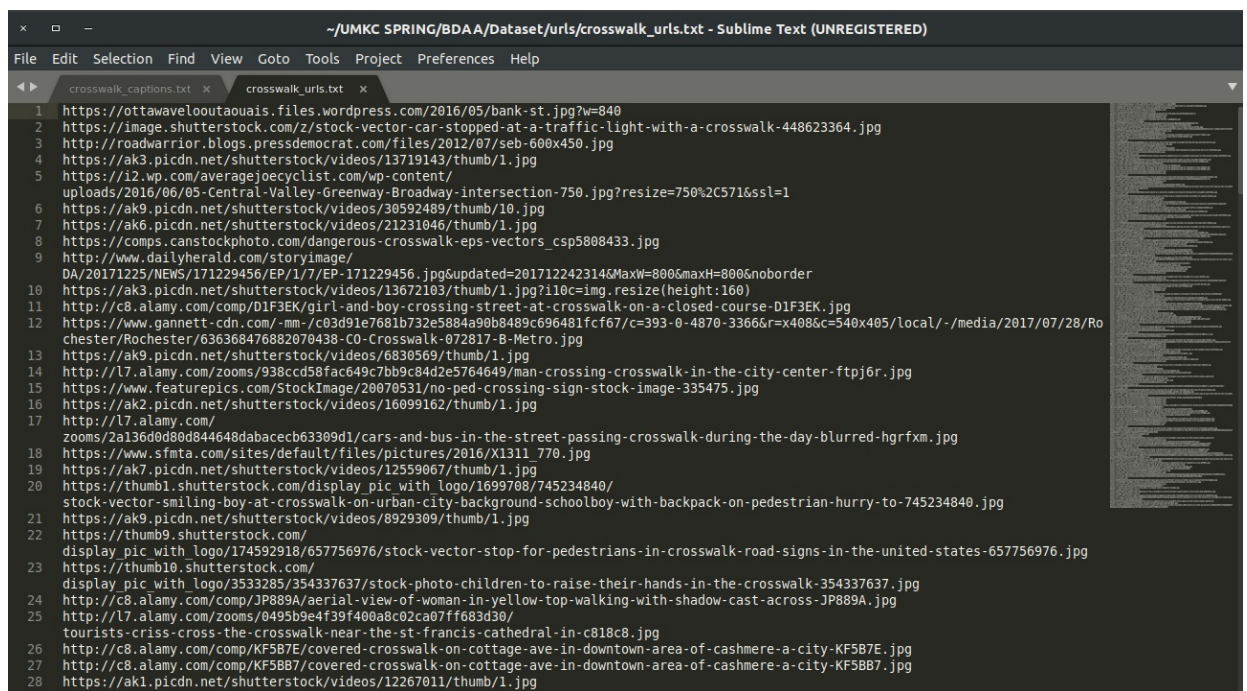
```

Output

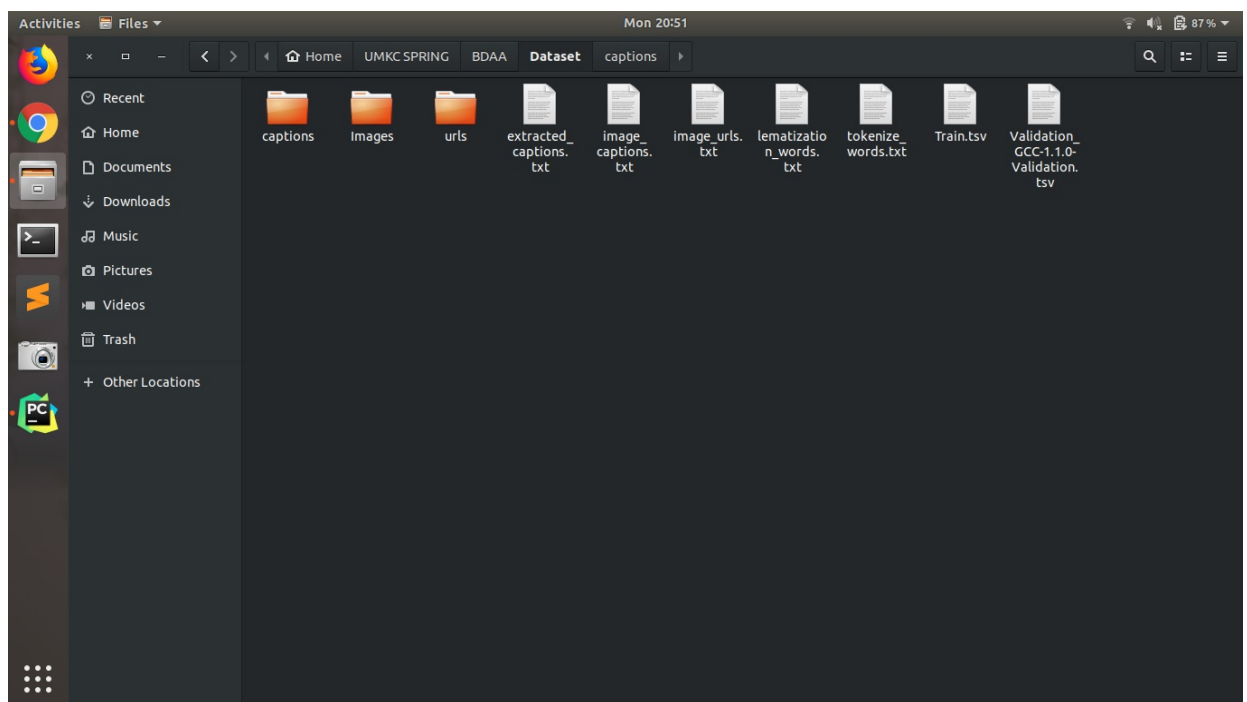
- This is output of downloaded images data set



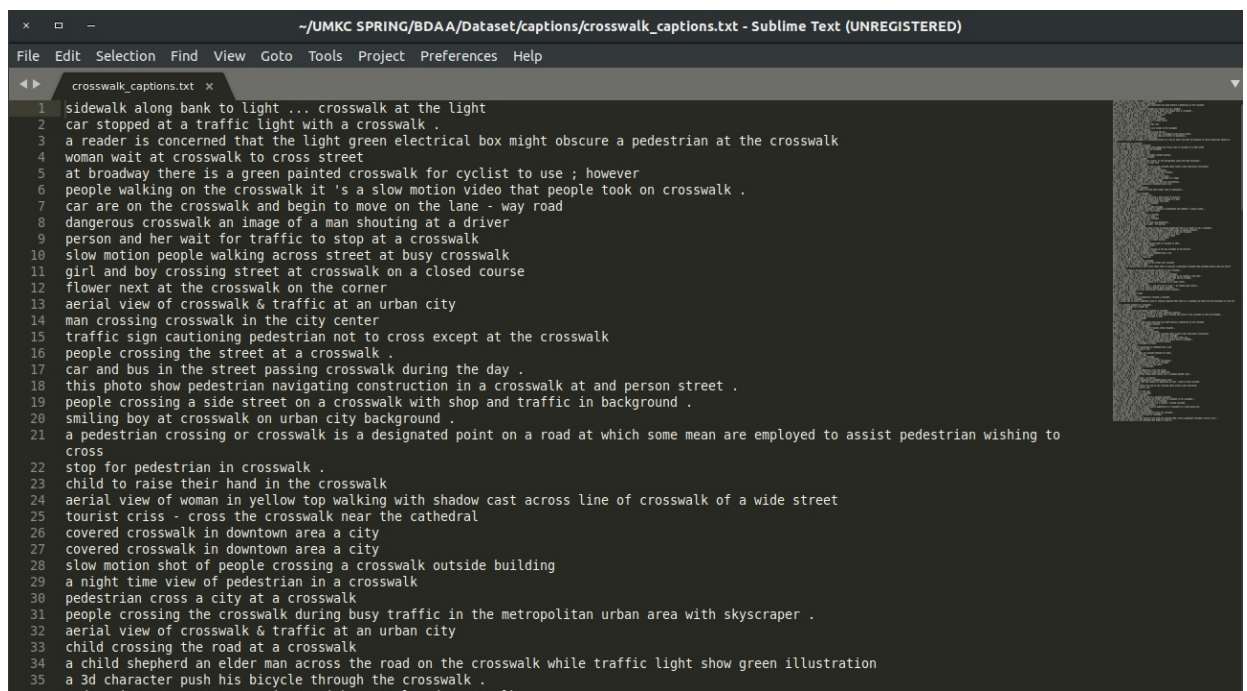
- This is output i.e urls text file of Crosswalk keyword



- This is the data set folder consists of images, url text file and caption text file of filtered keywords



- These are the captions of Crosswalk keyword



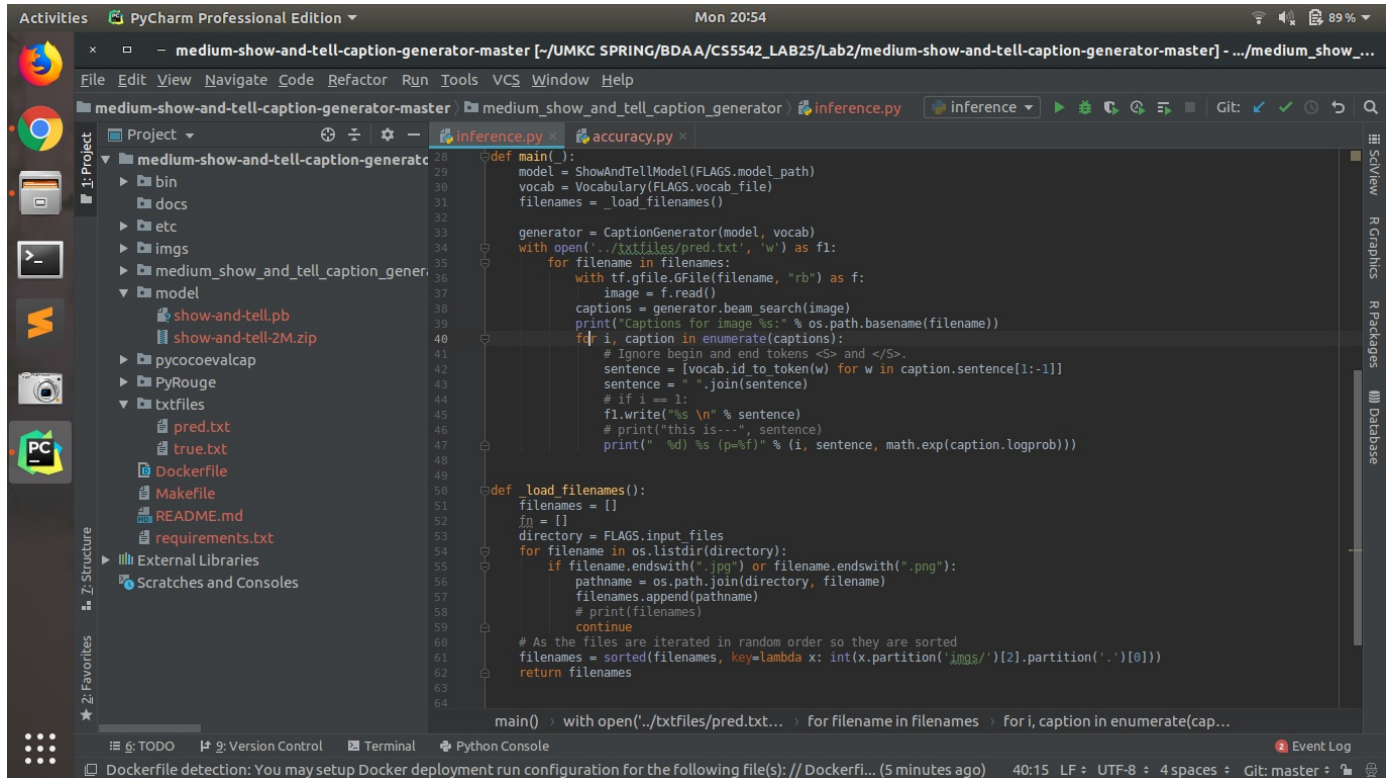
1. Generate captions for your own dataset using the show and Tellmodel.

Code Snippet

Steps

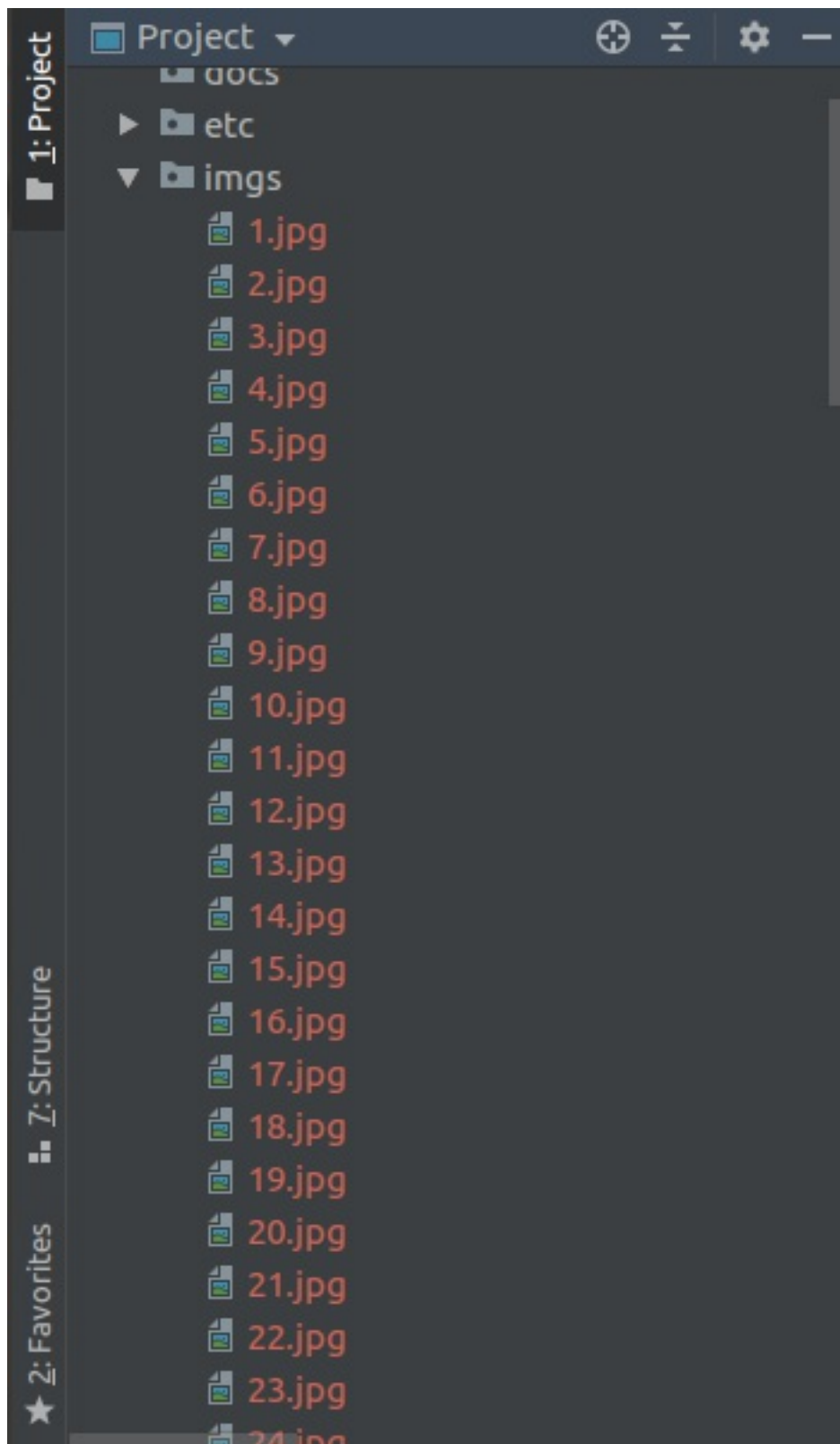
1. First, Download the images from image URL's and make ready the dataset
2. Change the path names which are initialized for the flags
3. Now iterate through the images directory and pass all images to the main method(code is in `_load_filenames`)
4. In the `caption_generator.py` change the beam size to 4 i.e. for generating 4 captions.
5. Now give the image to the model it will generate the 4 captions

- Code which iterates the image dataset



```
28 def main():
29     model = ShowAndTellModel(FLAGS.model_path)
30     vocab = Vocabulary(FLAGS.vocab_file)
31     filenames = _load_filenames()
32
33     generator = CaptionGenerator(model, vocab)
34     with open('../txtfiles/pred.txt', 'w') as f1:
35         for filename in filenames:
36             with tf.gfile.GFile(filename, "rb") as f:
37                 image = f.read()
38                 captions = generator.beam_search(image)
39                 print("Captions for image %s:" % os.path.basename(filename))
40                 for i, caption in enumerate(captions):
41                     # Ignore begin and end tokens <S> and </S>.
42                     sentence = [vocab.id_to_token(w) for w in caption.sentence[1:-1]]
43                     sentence = " ".join(sentence)
44                     # if i == 1:
45                     f1.write("%s \n" % sentence)
46                     # print("this is---", sentence)
47                     print('%d) %s (p=%f)' % (i, sentence, math.exp(caption.logprob)))
48
49
50 def _load_filenames():
51     filenames = []
52     fn = []
53     directory = FLAGS.input_files
54     for filename in os.listdir(directory):
55         if filename.endswith('.jpg') or filename.endswith('.png'):
56             pathname = os.path.join(directory, filename)
57             filenames.append(pathname)
58             # print(filenames)
59             continue
60     # As the files are iterated in random order so they are sorted
61     filenames = sorted(filenames, key=lambda x: int(x.partition('imgs/')[2].partition('.')[0]))
62     return filenames
63
64
main() > with open('../txtfiles/pred.txt... > for filename in filenames > for i, caption in enumerate(cap...
```

- Image dataset is as follows



2. Generate 4 captions for each image. (Beam Search k=4)

Code Snippet

- Code for generating captions

```
12 from medium_show_and_tell_caption_generator.model import ShowAndTellModel
13 from medium_show_and_tell_caption_generator.vocabulary import Vocabulary
14
15 FLAGS = tf.flags.FLAGS
16
17 tf.flags.DEFINE_string("model_path", "../model/show-and-tell.pb", "Model graph def path")
18 tf.flags.DEFINE_string("vocab_file", "../etc/word_counts.txt", "Text file containing the vocabulary.")
19 tf.flags.DEFINE_string("input_files", "../images/",
20                       "File pattern or comma-separated list of file patterns "
21                       "of image files.")
22
23 logging.basicConfig(level=logging.INFO)
24
25 logger = logging.getLogger(__name__)
26
27
28 def main(_):
29     model = ShowAndTellModel(FLAGS.model_path)
30     vocab = Vocabulary(FLAGS.vocab_file)
31     filenames = _load_filenames()
32
33     generator = CaptionGenerator(model, vocab)
34     with open('../txtfiles/pred.txt', 'w') as f1:
35         for filename in filenames:
36             with tf.gfile.GFile(filename, "rb") as f:
37                 image = f.read()
38                 captions = generator.beam_search(image)
39                 print("Captions for image '%s':" % os.path.basename(filename))
40                 for i, caption in enumerate(captions):
41                     # Ignore begin and end tokens <S> and </S>.
42                     sentence = [vocab.id_to_token(w) for w in caption.sentence[1:-1]]
43                     sentence = " ".join(sentence)
44                     # if i == 1:
45                     f1.write("%s \n" % sentence)
46                     # print("this is---", sentence)
47                     print(" %d) %s (p=%f)" % (i, sentence, math.exp(caption.logprob)))
48
49 _load_filenames()
```

- Set Beam size as 4

```
113
114
115 def __init__(self,
116             model,
117             vocab,
118             beam_size=4,
119             max_caption_length=20,
120             length_normalization_factor=0.0):
121
122     self.vocab = vocab
123     self.model = model
124
```

Output

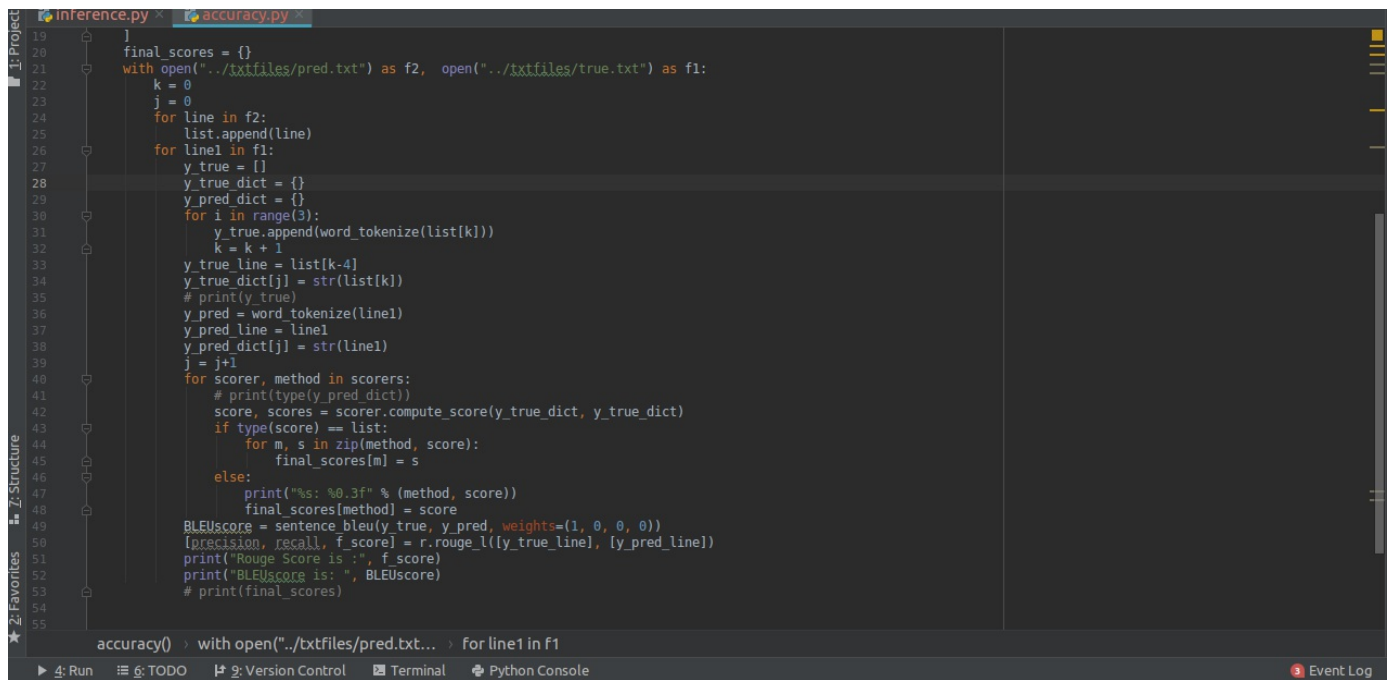
```
Run: Inference x
Captions for image 1.jpg:
0) a stop sign on the side of the road . (p=0.000107)
1) a stop sign on the side of a road . (p=0.000085)
2) a stop sign on the side of the road (p=0.000043)
3) a stop sign on the side of a road (p=0.000024)
Captions for image 2.jpg:
0) a double decker bus driving down a street . (p=0.001092)
1) a double decker bus is driving down the street . (p=0.000428)
2) a double decker bus driving down the street . (p=0.000418)
3) a double decker bus is driving down the road . (p=0.000233)
Captions for image 3.jpg:
0) a city street with a traffic light and street signs . (p=0.000012)
1) a street with a lot of cars and a traffic light (p=0.000008)
2) a city street with a traffic light and street sign . (p=0.000007)
3) a city street with a traffic light and street signs (p=0.000006)
Captions for image 4.jpg:
0) a group of people walking down a street . (p=0.001588)
1) a group of people walking down a city street . (p=0.000826)
2) a group of people crossing a street in the rain . (p=0.000423)
3) a group of people crossing a street in the rain (p=0.000081)
Captions for image 5.jpg:
0) a man riding a bike down a street . (p=0.000540)
1) a man is riding a bike down the street (p=0.000237)
2) a man is riding a bike down the street . (p=0.000188)
3) a man riding a bike down a street next to a sign . (p=0.000112)
Captions for image 6.jpg:
0) a group of people walking down a street . (p=0.004239)
1) a group of people walking down a sidewalk . (p=0.001803)
```

3. Report your accuracy in BLEU, CIDER ROUGE measures.

Steps

1. When you run [inference.py](#) it will generate the pred.txt file which has all the predicted for captions for the dataset
2. Create a txt file which has true dataset captions
3. Iterate both files and make a list of 4 captions and one true caption give this data to sentence_bleu to get BLEU SCORE
4. Send the single caption of pred and true to the rouge you will get presion, recall and f_score(code is taken from git hub url - <https://github.com/pcyin/PyRouge>)
5. Make the dict of pred caption and true captions and set to comput_score to get the CIDEr score
6. (code is taken from github url - <https://github.com/salaniz/pycocoEvalcap>)

Code Snippet



```
19 }
20 final_scores = {}
21 with open("../txtfiles/pred.txt") as f2, open("../txtfiles/true.txt") as f1:
22     k = 0
23     j = 0
24     for line in f2:
25         list.append(line)
26         for line1 in f1:
27             y_true = []
28             y_true_dict = {}
29             y_pred_dict = {}
30             for i in range(3):
31                 y_true.append(word_tokenize(list[k]))
32                 k = k + 1
33             y_true_line = list[k-4]
34             y_true_dict[j] = str(list[k])
35             # print(y_true)
36             y_pred = word_tokenize(line1)
37             y_pred_line = line1
38             y_pred_dict[j] = str(line1)
39             j = j+1
40             for scorer, method in scorers:
41                 # print(type(y_pred_dict))
42                 score, scores = scorer.compute_score(y_true_dict, y_true_dict)
43                 if type(score) == list:
44                     for m, s in zip(method, score):
45                         final_scores[m] = s
46                 else:
47                     print("%s: %.3f" % (method, score))
48                     final_scores[method] = score
49             BLEUScore = sentence_bleu(y_true, y_pred, weights=(1, 0, 0, 0))
50             [precision, recall, f_score] = r.rouge_l([y_true_line], [y_pred_line])
51             print("Rouge Score is :", f_score)
52             print("BLEUScore is: ", BLEUScore)
53             # print(final_scores)
54
55 accuracy() > with open("../txtfiles/pred.txt... > for line1 in f1
```

BLEU

- BLEU (bilingual evaluation understudy)
- It is an algorithm for evaluating the quality of text which has been machine-translated from one natural language to another
- BLEU's output is always a number between 0 and 1.
- This value indicates how similar the candidate text is to the reference texts, with values closer to 1 representing more similar texts.

ROUGE

- ROUGE, or Recall-Oriented Understudy for Gisting Evaluation
- It is a set of metrics and a software package used for evaluating automatic summarization and machine translation software in natural language processing.
- The following five evaluation metrics are available.

- ROUGE-N :Overlap of N-grams between the system and reference summaries.
- ROUGE-L: Longest Common Subsequence (LCS)
- ROUGE-W: Weighted LCS-based statistics that favors consecutive LCSes
- ROUGE-S: Skip-bigram based co-occurrence statistics. Skip-bigram is any pair of words in their sentence order.
- ROUGE-SU: Skip-bigram plus unigram-based co-occurrence statistics.

Output

Here in this accuracy score ROUGE-L BLEU and CIDEr are as follows

```
Run: accuracy x
CIDEr: 0.000
Precision is :0.56
Recall is :0.27450980392156865
F Score is :0.3684220084833848
BLEUScore is: 0.29411764705882354
CIDEr: 0.000
Precision is :0.38181818181818183
Recall is :0.525
F Score is :0.44210621440443754
BLEUScore is: 0.09306272250443651
CIDEr: 0.000
Precision is :0.5094339622641509
Recall is :0.35526315789473684
F Score is :0.41860560275224407
BLEUScore is: 0.07142857142857141
CIDEr: 0.000
Precision is :0.5882352941176471
Recall is :0.3333333333333333
F Score is :0.4255328687188823
BLEUScore is: 0.23529411764705885
CIDEr: 0.000
Precision is :0.6
Recall is :0.36363636363636365
F Score is :0.4528311416874381
BLEUScore is: 0.13333333333333333
```

Conclusion

- According to BLEU the accuracy of the model is low.
- BLEU value ranges from 0 to 1. 1 being most accurate and 0 being not accurate at all.
- Here we are getting an average of 0.3 accuracy which is not

good enough

- According to ROUGE precision of image is around 0.45
- CIDEr score is very less near to 0.01

References

- <https://machinelearningmastery.com/calculate-bleu-score-for-text-python/>
- <https://github.com/salaniz/pycocoevalcap>
- <https://gist.github.com/kracwarlock/c979b10433fe4ac9fb97>
- <https://github.com/pcyin/PyRouge>