

Project: Semantic Segmentation

Sean Pavlak

April 30 2018

1 Project

For this project I set out to better understand and hopefully implement semantic segmentation. This is a topic that is at the heart of autonomous vehicles, it also is a very interesting computer vision and deep learning problem.

Segmentation is essential for image analysis tasks. Semantic segmentation describes the process of associating each pixel of an image with a class label.

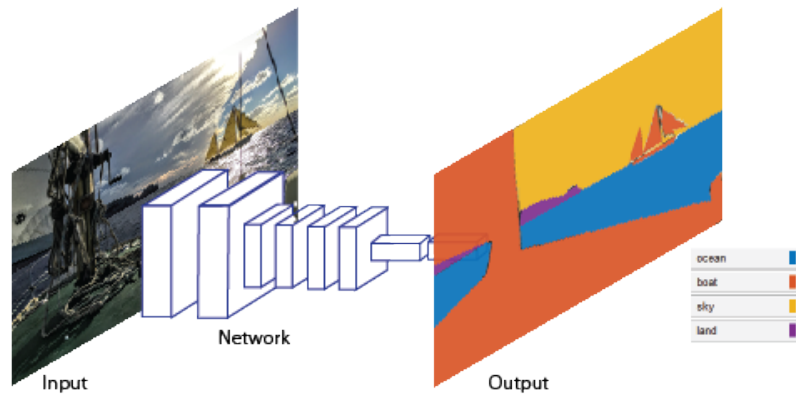


Figure 1: Demonstrating the process of semantic segmentation on an image.

As shown in figure 1, semantic segmentation begins with an image which is then processed by a trained network and outputs an color coded image. This image is then able to be deciphered and potentially yield more information than the original.

At this point in time deep learning is the chosen mechanism to perform semantic segmentation, but before deep learning took over,

people used approaches like TextonForest and Random Forest based classifiers. As with image classification, convolutional neural networks have had enormous success on segmentation problems.

One newer method of segmentation is Fully Convolutional Neural Networks. Fully convolutional indicates that the neural network is composed of convolutional layers without any fully-connected layers usually found at the end of the network. To semantic segment an image, a "fully convolutional" neural networks can take input of arbitrary size and produce correspondingly-sized output with efficient inference and learning, both learning and inference are performed whole-image-at-a-time by dense feedforward computation and back-propagation. In-network upsampling layers enable pixelwise prediction and learning in nets with subsampled pooling.

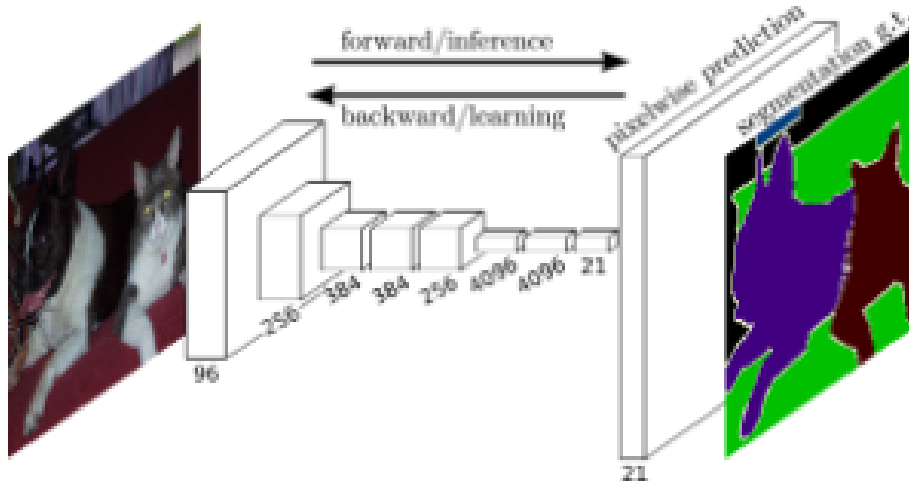


Figure 2: Demonstrating the process of fully convolutional neural network semantic segmentation on an image.

A general semantic segmentation architecture can be broadly thought of as an encoder network followed by a decoder network. The encoder is usually a pre-trained classification network like VGG/ResNet followed by a decoder network. The decoder network/mechanism is mostly where these architectures differ. The task of the decoder is to semantically project the discriminating features (lower resolution) learned by the encoder onto the pixel space (higher resolution) to get a dense classification.

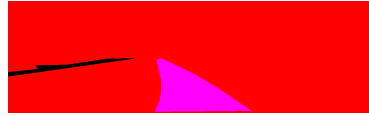
A convolutional encoder-decoder neural network is composed of a stack of encoders followed by a corresponding decoder stack which

feeds into a soft-max classification layer. The decoders help map low resolution feature maps at the output of the encoder stack to full input image size feature maps.

The dataset I am using is the KITTI road data. KITTI is a project of Karlsruhe Institute of Technology and Toyota Technological Institute at Chicago. The dataset includes three different categories of road; urban unmarked, urban marked, and urban multiple marked lanes. This data is also broken up into preset training and testing sets.



(a) Urban marked road training image from the KITTI dataset



(b) Colorized urban marked road training image from the KITTI dataset

Figure 3: (a) and (b) are the training image pair used to train the CNN model.

Using the images provided by KITTI and a CNN designed for semantic segmentation, I set out to train a model on this dataset.

1.1 Code Description

In order to perform semantic segmentation I defined a set of function what would load in my training data, setup a tensorflow model, and then use this defined model to generate a semantic segmentation image on the testing images provided by the dataset.

```

1 def generate_batch_function(data_folder, image_shape):
2     def get_batches_function(batch_size):
3         image_paths = glob(os.path.join(data_folder, 'image_2',
4             ↳ '*.png'))
5         label_paths = {
6             re.sub(r'_(lane|road)_', '_',
7             ↳ os.path.basename(path)): path
8             for path in glob(os.path.join(data_folder,
9             ↳ 'gt_image_2', '*_road_*.png'))}
10        background_color = np.array([255, 0, 0])
11        random.shuffle(image_paths)

```

```

10     for i in range(0, len(image_paths), batch_size):
11         images = []
12         gt_images = []
13
14         for image_file in image_paths[i:i+batch_size]:
15             gt_image_file =
16                 ↳ label_paths[os.path.basename(image_file)]
17
18             image =
19                 ↳ scipy.misc.imresize(scipy.misc.imread(image_file),
20                 ↳ image_shape)
21             gt_image =
22                 ↳ scipy.misc.imresize(scipy.misc.imread(gt_image_file),
23                 ↳ image_shape)
24
25             gt_bg = np.all(gt_image == background_color,
26                 ↳ axis=2)
27             gt_bg = gt_bg.reshape(*gt_bg.shape, 1)
28             gt_image = np.concatenate((gt_bg,
29                 ↳ np.invert(gt_bg)), axis=2)
30
31             images.append(image)
32             gt_images.append(gt_image)
33
34         yield np.array(images), np.array(gt_images)
35     return get_batches_function

```

The first functions `generate_batch_function` and `get_batches_function` simply generate two arrays containing the street images and their colored pairs, similar to the pair shown in figure 3.

```

1  def generate_test_output(sess, logits, keep_prob, image_pl,
2     ↳ data_folder, image_shape):
3     for image_file in glob(os.path.join(data_folder, 'image_2',
4     ↳ '*.png')):
5         image =
6             ↳ scipy.misc.imresize(scipy.misc.imread(image_file),
7             ↳ image_shape)
8
9         im_softmax = sess.run([tf.nn.softmax(logits)],
10            ↳ {keep_prob: 1.0, image_pl: [image]})
11         im_softmax = im_softmax[0][:, 1].reshape(image_shape[0],
12            ↳ image_shape[1])
13
14

```

```

8         segmentation = (im_softmax > 0.5).reshape(image_shape[0],
9             ↪ image_shape[1], 1)
10
11         mask = np.dot(segmentation, np.array([[0, 255, 0, 127]]))
12         mask = scipy.misc.toimage(mask, mode='RGBA')
13
14         street_im = scipy.misc.toimage(image)
15         street_im.paste(mask, box=None, mask=mask)
16
17         yield os.path.basename(image_file), np.array(street_im)

```

The next function `generate_test_output` is of a similar nature. Rather than loading in training data, this function loads in the testing data. It also produces an array of images.

After the images are in memory, I defined a function, `load_vgg`, which then loads in the TensorFlow VGG architecture. This architecture is pre-trained to handle semantic segmentation. However, as shown in the results it still needs to be extensively trained in order to handle the task at hand.

```

1 def load_vgg(sess, vgg_path):
2     model = tf.saved_model.loader.load(sess, ['vgg16'], vgg_path)
3
4     graph = tf.get_default_graph()
5     image_input = graph.get_tensor_by_name('image_input:0')
6     keep_prob = graph.get_tensor_by_name('keep_prob:0')
7     layer3 = graph.get_tensor_by_name('layer3_out:0')
8     layer4 = graph.get_tensor_by_name('layer4_out:0')
9     layer7 = graph.get_tensor_by_name('layer7_out:0')
10
11     return image_input, keep_prob, layer3, layer4, layer7

```

This pre-trained model generates layers for the CNN that will train on the training data.

Once the pre-trained model is set and those layers are defined another function named `layers` takes those pre-trained layers and convolves them in order to generate a decoder layer. This layer will be used to semantically project the discriminating features learned by the encoding layers onto the pixel space.

```

1  def layers(vgg_layer3_out, vgg_layer4_out, vgg_layer7_out,
    ↪ num_classes = NUMBER_OF_CLASSES):
2      layer3x = tf.layers.conv2d(inputs = vgg_layer3_out,
3                                  filters = NUMBER_OF_CLASSES,
4                                  kernel_size = (1, 1),
5                                  strides = (1, 1),
6                                  name = 'layer3conv1x1')
7
8      layer4x = tf.layers.conv2d(inputs = vgg_layer4_out,
9                                  filters = NUMBER_OF_CLASSES,
10                                 kernel_size = (1, 1),
11                                 strides = (1, 1),
12                                 name = 'layer4conv1x1')
13
14     layer7x = tf.layers.conv2d(inputs = vgg_layer7_out,
15                                 filters = NUMBER_OF_CLASSES,
16                                 kernel_size = (1, 1),
17                                 strides = (1, 1),
18                                 name = 'layer7conv1x1')
19
20     decoderlayer1 = tf.layers.conv2d_transpose(inputs = layer7x,
21                                                  filters =
22                                                  ↪ NUMBER_OF_CLASSES,
23                                                  kernel_size = (4,
24                                                  ↪ 4),
25                                                  strides = (2, 2),
26                                                  padding = 'same',
27                                                  name =
28                                                  ↪ 'decoderlayer1')
29
30     decoderlayer2 = tf.add(decoderlayer1, layer4x, name =
31     ↪ 'decoderlayer2')
32
33     decoderlayer3 = tf.layers.conv2d_transpose(inputs =
34     ↪ decoderlayer2,
35                                                  filters =
36                                                  ↪ NUMBER_OF_CLASSES,
37                                                  kernel_size = (4,
38                                                  ↪ 4),
39                                                  strides = (2, 2),
40                                                  padding = 'same',
41                                                  name =
42                                                  ↪ 'decoderlayer3')
43
44     decoderlayer4 = tf.add(decoderlayer3, layer3x, name =
45     ↪ 'decoderlayer4')

```



```

2     start = time.time()
3     print('epoch: 0', '/', EPOCHS, 'training loss: N/A')
4
5     for epoch in range(EPOCHS):
6         losses, i = [], 0
7         for images, labels in get_batches_function(BATCH_SIZE):
8             i += 1
9             feed = {input_image: images,
10                    correct_label: labels,
11                    keep_prob: DROPOUT,
12                    learning_rate: LEARNING_RATE}
13
14             _, partial_loss = sess.run([train_optimizer,
15                                         ↪ cross_entropy_loss], feed_dict = feed)
16
17             print('-> iteration: ', i, '/', NUMBER_OF_IMAGES,
18                   ↪ 'partial loss: ', partial_loss)
19             losses.append(partial_loss)
20
21         training_loss = sum(losses) / len(losses)
22
23         end = time.time()
24         training_time = end - start
25
26         print('epoch: ', epoch + 1, '/', EPOCHS, 'training loss:
27               ↪ ', training_loss)
28     print('training time: ', training_time)

```

Then I defined the `train_nn` function. Which iterates through each epoch and each training image pair in order to train the CNN over many separate cycles. This function also keeps track of how long each set of epochs took in order to train.

```

1 def save_inference_samples(runs_dir, data_dir, sess, image_shape,
2     ↪ logits, keep_prob, input_image):
3     output_dir = os.path.join(runs_dir, 'epochs
4     ↪ {0}'.format(EPOCHS))
5
6     if os.path.exists(output_dir):
7         shutil.rmtree(output_dir)
8     os.makedirs(output_dir)
9
10    print('Training Finished. Saving test images to:
11          ↪ {}'.format(output_dir))

```



```

9     image_outputs = generate_test_output(sess, logits, keep_prob,
    ↪     input_image, os.path.join(data_dir, 'data_road/testing'),
    ↪     image_shape)
10
11     for name, image in image_outputs:
12         scipy.misc.imsave(os.path.join(output_dir, name), image)

```

Lastly, I defined a function that saves the semantic segmentation image for each test case. This allowed me to empirically determine how successful each set of epochs was.

Once all the functions were defined, I defined the paths and the learning parameters. Then training and testing was as simple as loading in the data and tensorflow model, training, and testing.

```

1 DATA_DIRECTORY = './data'
2 RUNS_DIRECTORY = './runs'
3 TRAINING_DATA_DIRECTORY = './data/data_road/training'
4 NUMBER_OF_IMAGES =
    ↪ len(glob('./data/data_road/training/calib/*.png'))
5 VGG_PATH = './data/vgg'
6
7 NUMBER_OF_CLASSES = 2
8 IMAGE_SHAPE = (160, 576)
9
10 EPOCHS = 20
11 BATCH_SIZE = 1
12
13 LEARNING_RATE = 0.0001
14 DROPOUT = 0.75
15
16 correct_label = tf.placeholder(tf.float32, [None, IMAGE_SHAPE[0],
    ↪     IMAGE_SHAPE[1], NUMBER_OF_CLASSES])
17 learning_rate = tf.placeholder(tf.float32)
18 keep_prob = tf.placeholder(tf.float32)
19
20 get_batches_function =
    ↪ generate_batch_function(TRAINING_DATA_DIRECTORY, IMAGE_SHAPE)
21
22 with tf.Session() as session:
23     image_input, keep_prob, layer3, layer4, layer7 =
    ↪ load_vgg(session, VGG_PATH)
24     model_output = layers(layer3, layer4, layer7,
    ↪     NUMBER_OF_CLASSES)

```

```

25 logits, train_optimizer, cross_entropy_loss =
26     ↪ optimize(model_output, correct_label, learning_rate,
27     ↪ NUMBER_OF_CLASSES)
28
29 session.run([tf.global_variables_initializer(),
30     ↪ tf.local_variables_initializer()])
31 train_nn(session, EPOCHS, BATCH_SIZE, get_batches_function,
    ↪ train_optimizer, cross_entropy_loss, image_input,
    ↪ correct_label, keep_prob, learning_rate)
32
33 save_inference_samples(RUNS_DIRECTORY, DATA_DIRECTORY,
    ↪ session, IMAGE_SHAPE, logits, keep_prob, image_input)

```

The very last step was calling `save_inference_samples` so that the test images that underwent semantic segmentation could be saved into a separate directory and reviewed.

1.2 Results

The final results of this project were more impressive than I originally expected. Semantic segmentation is significantly more accurate at defining objects in a view than most of the edge detectors we explored in this course.

The approach I took in order to evaluate this model was to generate three separate models, which differed only in their epochs for training. The first only went through a single epoch, and took a few minutes to train. The next was over 10 epochs, and took over an hour to train. The last went through 20 epochs, and took multiple hours to train.

If I were training on GPU rather than CPU and I wasn't running this code on a local machine I may have been able to train in significantly shorter times.

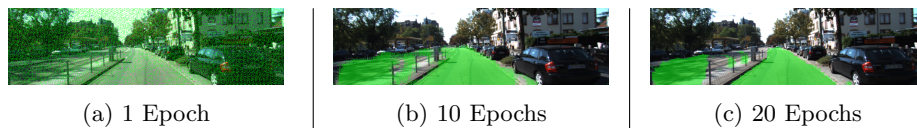


Figure 4: (a),(b), and (c) are the generated semantic segmentation images on the urban unmarked road images.

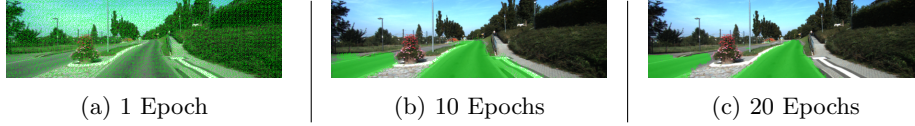


Figure 5: (a),(b), and (c) are the generated semantic segmentation images on the urban unmarked road images.

As shown in these images of urban unmarked road, it wasn't until the 10 epoch trial that some semantic segmentation began to take hold.

However, between the 10 and 20 epochs there does not appear to be very much improvement.

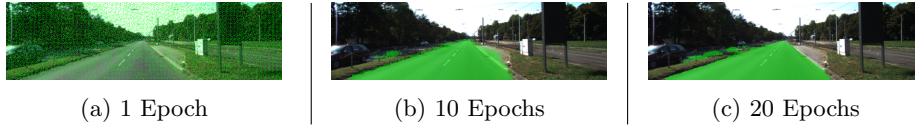


Figure 6: (a),(b), and (c) are the generated semantic segmentation images on the urban marked road images.

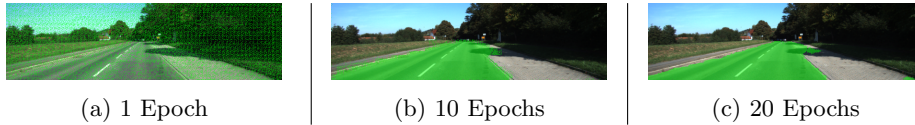


Figure 7: (a),(b), and (c) are the generated semantic segmentation images on the urban marked road images.

When performed on the urban marked road images both the 1 epoch and 10 epochs cases fell short. Neither was able to accurately perform semantic segmentation on the lanes without extending past the bounds of the road.

The 20 epoch case was able to differentiate between road and non-road area well. This testing set still appears to give this model some difficulties shown in image c of figure 7. This image has a rough edge on the right side of the road along with errors caused by shadows on the road.

Potentially running a 30 epoch case may solve this, or tweaking the model may help.

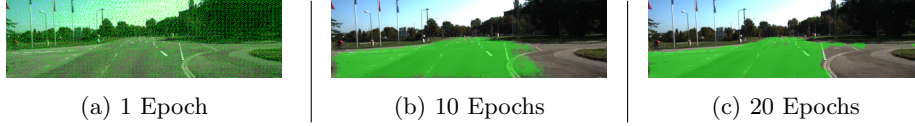


Figure 8: (a),(b), and (c) are the generated semantic segmentation images on the urban multiple marked lanes road images.

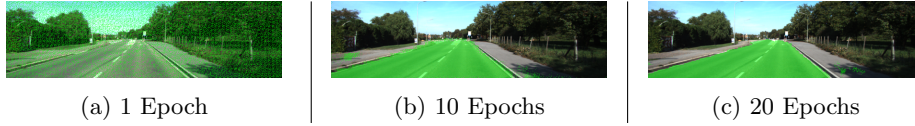


Figure 9: (a),(b), and (c) are the generated semantic segmentation images on the urban multiple marked lanes road images.

Again semantic segmentation on this test set of urban multiple marked lanes images appears to only hold true for the 20 epochs case. This set does appear to be better suited for my model. Except when approaching turning lanes, but this may be corrected with more epochs or examples of turning lanes presented in the training data.

Overall all three cases for the 10 or 20 epochs performed very well. The only changes I would make would be to extend to a 30 epochs case and maybe define a function to compare the generated image with a hand labeled image to calculate an accuracy value for these images.

Unless the KITTI dataset has these labeled test images somewhere this may not be feasible as there are 300 images that would need labeled. Which wouldn't be time efficient. Regardless, without labels this project's accuracy can only be estimated.