

THE NEURAL PERSPECTIVE

DEEP LEARNING SIMPLIFIED

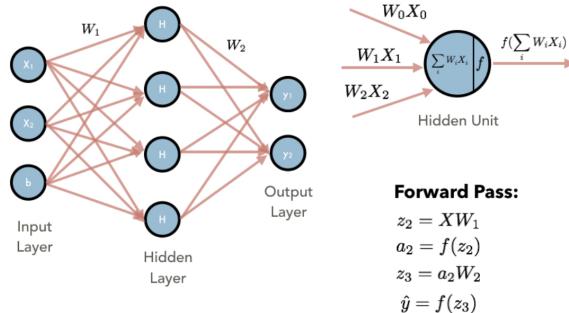
HOME • TUTORIALS • READINGS • CONTACT • ABOUT

VANILLA NEURAL NETWORK

[Edit](#)

I. OBJECTIVE:

Simple 1-layer multi layer perceptron (MLP). f is the activation function and introduces the non-linearity to our system.



II. LINEAR MODEL:

A simple linear model with a softmax layer on top. The main difference here is the lack of a non-linear activation function (ReLU, tanh, etc.). Thanks to Karpathy for the data and code structure, but we will break down the math behind the lines for better understanding. You can check out the code for loading the data on the Github repo but here we will focus on the main model operations.

 XW

```
1 # Class scores [NxC]
2 logits = np.dot(X, W)
```

$$\frac{e^{XW_y}}{\sum e^{XW}}$$

```
1 # Class probabilities
2 exp_logits = np.exp(logits)
3 probs = exp_logits / np.sum(exp_logits, axis=1, keepdims=True)
```

$$-\log(p) + \frac{\lambda}{2} \sum_k \sum_l W_{k,l}^2$$

```
1 # Loss
2 correct_class_logprobs = -np.log(probs[range(len(probs)), y])
3 loss = np.sum(correct_class_logprobs) / config.DATA_SIZE
4 loss += 0.5 * config.REG * np.sum(W*W)
```

$$\frac{\partial J}{\partial W_j} = XP, \frac{\partial J}{\partial W_y} = X(P-1)$$

```
1 # Backpropagation
2 dscores = probs
3 dscores[range(len(probs)), y] -= 1
4 dscores /= config.DATA_SIZE
```

$$\frac{\partial J}{\partial W_i} = \frac{\partial J}{\partial W_i} + \lambda W_i$$

```
1 dW = np.dot(X.T, dscores)
2 dW += config.REG*W
```

$$W \leftarrow W - \gamma \frac{\partial J}{\partial W}$$

- SEARCH -

Search ...

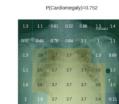
- FOLLOW ME ON TWITTER -

Tweets by @GokuMohandas

 **Goku Mohandas**

@GokuMohandas

Your model may be performing really well by incorrectly focusing on confounding features (extraneous influencers in the data that aren't accounted for). Check out [@johnzech's post](#) where x-ray stickers unintentionally influenced the classifications: medium.com/@jrzech/what-a...

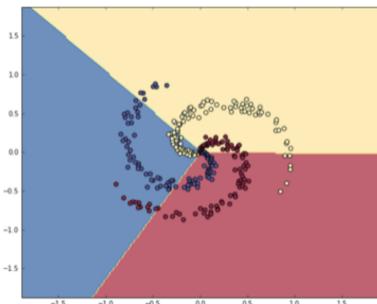
[Embed](#)[View on Twitter](#)**- RECENT POSTS -**[Update on Embeddings \(Spring 2017\)](#)[Exploring Sparsity in Recurrent Neural Networks](#)[Question Answering from Unstructured Text by Retrieval and Comprehension](#)[Overcoming Catastrophic Forgetting in Neural Networks](#)[Opening the Black Box of Deep Neural Networks via Information](#)

$$\nabla_i = \nabla_i - \alpha \frac{\partial J}{\partial W_i}$$

```
1 | W += -config.LEARNING_RATE * dW
```

Results:

We can see that the decision boundary of our classifier is linear and cannot adapt to the non-linear contortions of the data.



III. NEURAL NETWORK:

Now we introduce a neural net with a softmax on the last layer for class probabilities. We use a ReLU unit to introduce non-linearity. Our network will have two layers, where the shape of the input will be manipulated as follows:

$$\begin{array}{ccccccc} \mathbf{x} & & \mathbf{x} & \mathbf{w_1} & & \mathbf{x} & \mathbf{w_2} & \hat{\mathbf{y}} \\ [\text{N} \times \text{D}] & \longrightarrow & [\text{N} \times \text{D}] & [\text{D} \times \text{H}] & \longrightarrow & [\text{N} \times \text{H}] & [\text{N} \times \text{H}] & [\text{H} \times \text{C}] \longrightarrow [\text{N} \times \text{C}] \end{array}$$

Once again, let's break down the code.

$$z_2 = XW_1$$

```
1 | z_2 = np.dot(X, W_1)
a2 = max(0, XW1)

1 | a_2 = np.maximum(0, z_2) # ReLU
```

$$\text{logits} = a_2 W_2$$

```
1 | logits = np.dot(a_2, W_2)
```

$$P = \frac{e^{\text{logits}_y}}{\sum(e^{\text{logits}})}$$

```
1 | # Class probabilities
2 | exp_logits = np.exp(logits)
3 | probs = exp_logits / np.sum(exp_logits, axis=1, keepdims=True)
```

$$\text{loss} = -\ln(P) + \frac{\lambda}{2} \sum \sum W_1 + \frac{\lambda}{2} \sum \sum W_2$$

```
1 | # Loss
2 | correct_class_logprobs = -np.log(probs[range(len(probs)), y])
3 | loss = np.sum(correct_class_logprobs) / config.DATA_SIZE
4 | loss += 0.5 * config.REG * np.sum(W_1 * W_1)
5 | loss += 0.5 * config.REG * np.sum(W_2 * W_2)
```

$$\frac{\partial J}{\partial W_{2j}} = a_2 P, \frac{\partial J}{\partial W_{2y}} = a_2 (P - 1)$$

```
1 | # Backpropagation
2 | dscores = probs
3 | dscores[range(len(probs)), y] -= 1
4 | dscores /= config.DATA_SIZE
5 | dW2 = np.dot(a_2.T, dscores)
```

$$\frac{\partial J}{\partial W_1} = \frac{\partial J}{\partial P} \frac{\partial P}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial W_1} = W_2(dscores) \partial(\text{ReLU}) X$$

```
1 | dhiddens = np.dot(dscores, W_2.T)
2 | dhiddens[a_2 > 0] *= 1
3 | dW1 = np.dot(X.T, dhiddens)
```

$$\frac{\partial J}{\partial W_i} = \frac{\partial J}{\partial P} + \lambda W_i$$

∂W_i

```
1 | dW2 += config.REG * W_2
2 | dW1 += config.REG * W_1
```

$$W_i = W_i - \alpha \frac{\partial J}{\partial W_i}$$

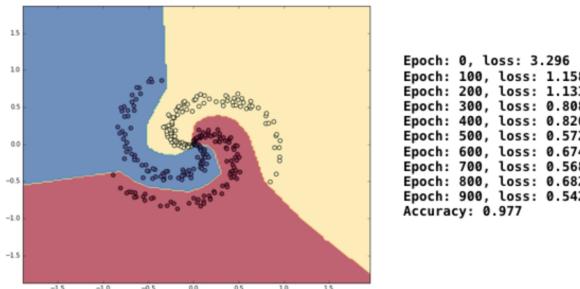
```
1 | W_1 += -config.LEARNING_RATE * dW1
2 | W_2 += -config.LEARNING_RATE * dW2
```

Our accuracy is very simple and involves doing the forward pass and then comparing the predicted class with the target class.

```
1 | def accuracy(X, y, W_1, W_2=None):
2 |     logits = np.dot(X, W_1)
3 |     if W_2 is None:
4 |         predicted_class = np.argmax(logits, axis=1)
5 |         print "Accuracy: %.3f" % (np.mean(predicted_class))
6 |     else:
7 |         z_2 = np.dot(X, W_1)
8 |         a_2 = np.maximum(0, z_2)
9 |         logits = np.dot(a_2, W_2)
10 |        predicted_class = np.argmax(logits, axis=1)
11 |        print "Accuracy: %.3f" % (np.mean(predicted_class))
```

Results:

The resulting decision boundary is able to classify the non-linear data really well.



IV. TENSORFLOW IMPLEMENTATION:

We will start by setting up our tensorflow model but we will have an extra function called `summarize()` which will store the progress as we training through the epochs. We will decide which values to store with `tf.scalar_summary()` so we can see the changes later.

```
1 | def create_model(sess, FLAGS):
2 |
3 |     model = mlp(FLAGS.DIMENSIONS,
4 |                  FLAGS.NUM_HIDDEN_UNITS,
5 |                  FLAGS.NUM_CLASSES,
6 |                  FLAGS.REG,
7 |                  FLAGS.LEARNING_RATE)
8 |     sess.run(tf.initialize_all_variables())
9 |     return model
10 |
11 class mlp(object):
12 |
13     def __init__(self,
14                  input_dimensions,
15                  num_hidden_units,
16                  num_classes,
17                  regularization,
18                  learning_rate):
19 |
20     # Placeholders
21     self.X = tf.placeholder("float", [None, None])
22     self.y = tf.placeholder("float", [None, None])
23 |
24     # Weights
25     W1 = tf.Variable(tf.random_normal(
26                     [input_dimensions, num_hidden_units], stddev=0
27     W2 = tf.Variable(tf.random_normal(
28                     [num_hidden_units, num_classes], stddev=0.01),
29 |
30     with tf.name_scope('forward_pass') as scope:
31         z_2 = tf.matmul(self.X, W1)
32         a_2 = tf.nn.relu(z_2)
33         self.logits = tf.matmul(a_2, W2)
34 |
35     # Add summary ops to collect data
36     W_1 = tf.histogram_summary("W1", W1)
37     W_2 = tf.histogram_summary("W2", W2)
38 |
39     with tf.name_scope('cost') as scope:
40         self.cost = tf.reduce_mean(
41             tf.nn.softmax_cross_entropy_with_logits(
42                 + 0.5 * regularization * tf.reduce_sum(
43                     + 0.5 * regularization * tf.reduce_sum(
44                         tf.scalar_summary("cost", self.cost)
45 |
46     with tf.name_scope('train') as scope:
47         self.optimizer = tf.train.AdamOptimizer(
```

```

49     learning_rate=learning_rate).minimize(self
50
51
52
53     def step(self, sess, batch_X, batch_y):
54
55         input_feed = {self.X: batch_X, self.y: batch_y}
56         output_feed = {self.logits, self.cost, self.optimizer}
57
58         outputs = sess.run(output_feed, input_feed)
59         return outputs[0], outputs[1], outputs[2]
60
61
62     def summarize(self, sess, batch_X, batch_y):
63         # Merge all summaries into a single operator
64         merged_summary_op = tf.merge_all_summaries()
65
66         return sess.run(merged_summary_op,
67                         feed_dict={self.X:batch_X, self.y:batch_y})

```

Then we will train for several epochs and save the summary each time.

```

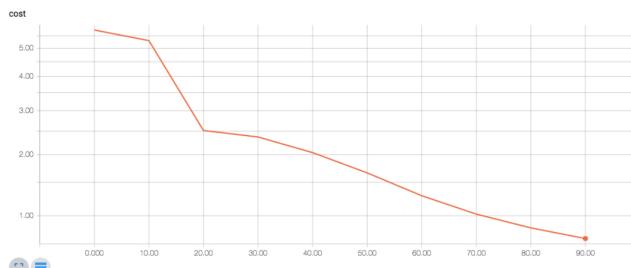
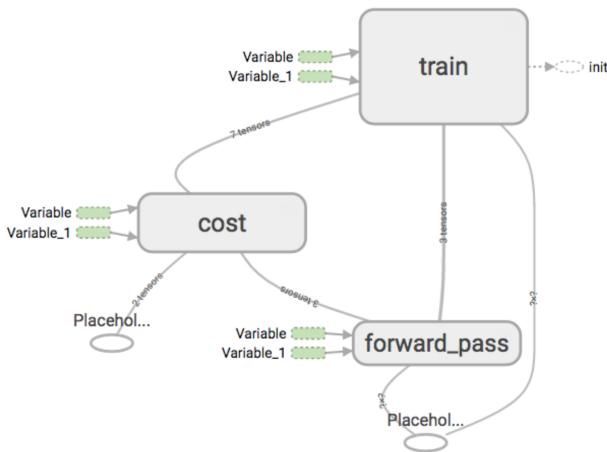
1  def train(FLAGS):
2
3      # Load the data
4      FLAGS, X, y = load_data(FLAGS)
5
6      with tf.Session() as sess:
7
8          model = create_model(sess, FLAGS)
9          summary_writer = tf.train.SummaryWriter(
10              FLAGS.TENSORBOARD_DIR, graph=sess.graph)
11
12          # y to categorical
13          Y = tf.one_hot(y, FLAGS.NUM_CLASSES).eval()
14
15          for epoch_num in range(FLAGS.NUM_EPOCHS):
16              logits, training_loss, _ = model.step(sess, X,
17                  # Display
18                  if epoch_num%FLAGS.DISPLAY_STEP == 0:
19                      print "EPOCH %i: \n Training loss: %.3f, A
20                      % (epoch_num,
21                          training_loss,
22                          np.mean(np.argmax(logits, 1) == y))
23
24              # Write logs for each epoch_num
25              summary_str = model.summarize(sess, X, Y)
26              summary_writer.add_summary(summary_str, ep
27
28      if __name__ == '__main__':
29          FLAGS = parameters()
30          train(FLAGS)

```

Finally, we can view our training progress using:

```
$ tensorboard --logdir=logs
```

and then heading over to <http://localhost:6006> on your browser to view the results. Here are a few:

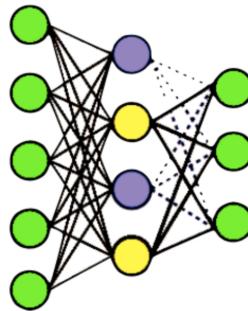


EXTRAS (DROPOUT AND DROPCONNECT):

There are many add on techniques to this vanilla neural network that works to increase optimization, robustness and overall performance. We will be covering many of them in future posts but I will briefly talk about a very common regularization technique: dropout.

What is it? Dropout is a regularization technique that allows us to nullify the outputs of certain neurons to zero. This will effectively be the same as the neuron not existing in the network. We will do this for p% of the total neurons in each layer and for each batch, a new p% of the neurons in each layer are “dropped”.

Why do we do this? It works out to be a great regularization technique because for each input batch, we are sampling from a different neural net since a whole new set of neurons are dropped. By repeating this, we are preventing the units from co-adapting too much to the data. The original paper describes each iteration as a “thinned” network because p% of the neurons are dropped. **Note:** Dropout is only for training time. At test time, we will not be dropping any neurons.

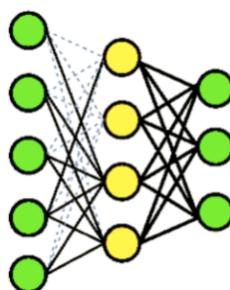


In the image above, the layer has p=0.5 which means half of its units are dropped. In another iteration, a different set of 1/2 of the neurons will be dropped. Let's take a look at masking code to really understand what's happening.

$$\begin{aligned} r_j^{(l)} &\sim \text{Bernoulli}(p), \\ \tilde{\mathbf{y}}^{(l)} &= \mathbf{r}^{(l)} * \mathbf{y}^{(l)}, \\ z_i^{(l+1)} &= \mathbf{w}_i^{(l+1)} \tilde{\mathbf{y}}^l + b_i^{(l+1)}, \\ y_i^{(l+1)} &= f(z_i^{(l+1)}). \end{aligned}$$

We use a Bernoulli distribution to generate 0/1 with probability p for 0. We apply this mask to the outputs from our layer. The parts that are multiplied by zero are our “dropped” neurons since they will yield an output of 0 when multiplied by the next set of weights.

Another regularization method, which is an extension of dropout, is dropconnect. It also involves a similar mechanism but is applied to the weights instead.



Notice that here, a set of weights are dropped instead of the neurons.

$$\begin{aligned} \text{mask}_{ij}^{(l)} &\sim \text{Bernoulli}(p) \\ \tilde{W}^{(l)} &= W^{(l)} * \text{mask}^{(l)} \\ z_i^{(l+1)} &= \sum_{j=1}^n \tilde{W}_{ij}^{(l)} y_j^{(l)} + b_i^{(l)} \\ y^{(l+1)} &= f(z^{(l+1)}) \end{aligned}$$

We apply a similar bernoulli mask to the weights and we use those weights for the layers. Any inputs that are dot produced with the zeroed weights will result in 0. You can see the similarity with dropout and so, empirically, both techniques offer similar results. Dropconnect was proposed because you always have more weights than neurons, so there are more ways to create “thinned” models thus results in more robust training. However, in more papers you will see mostly dropout being utilized and very rarely drop connect since results are similar.

You can read more about dropout [here](#) and dropconnect [here](#).

V. RAW CODE:

[GitHub Repo](#) (Updating all repos, will be back up soon!)



Posted in: Uncategorized

Tagged: Tutorials

← Logistic Regression

Convolutional Neural Networks (CNN) →

2 THOUGHTS ON “VANILLA NEURAL NETWORK”

OLGA August 17, 2017 at 9:30 pm

[EDIT](#) [REPLY →](#)



One of the most succinct and clearest NN guides I've seen so far, great material!

★ Liked by you

GOKUMOHANDAS August 17, 2017 at 10:15 pm

[EDIT](#) [REPLY →](#)



Thanks 😊 made it for myself initially just to be able to look at it and recall everything quickly

★ Like

LEAVE A REPLY

Enter your comment here...

