

Alex Walczak | HW6 | Neural Networks

Part I Backpropagation

Please see attached written portion.

Part II Training

I used 50,000 training points and 10,000 validation points in my neural net testing.

Preprocessing:

All the images were first scaled by 255 and then the mean was subtracted. The mean of the training data was also subtracted the validation set.

Labels adjusted from $\{0, 1\}$ to $\{0.02, 0.98\}$ since these target reachable by logistic and tanh functions.

Weight initialization:

As suggested here <http://cs231n.github.io/neural-networks-2/> and other places, it is recommended to adjust the variances of the weights. This is done so that the output of each neuron has the same variance (and the network can converge faster). The variance of the output of the neuron depends on the variance of the input, so our goal can be to make the variance of the output equal the variance of the input. For my weight initialization, I sampled weights from the unit Gaussian, and scaled this matrix by $1/\sqrt{\text{number of input features}}$ to achieve that goal.

Learning Rate:

Different layers learn at different rates (cs.cmu.edu/afs/cs/academic/class/15883-f15/slides/backprop.pdf).

Input-to-hidden (in code: i2h) weights learn more slowly than hidden-to-output (h2o) weights.

I made the learning rate variable, so that decayed over time, yet in later iterations it wouldn't jump out of the local minimum of the cost function.

Stopping Conditions:

None, basically. I'm a little impatient, so whenever the calculated loss was the smallest ever seen by my neural net, I pickled it, naming the files using the unique ID generated for each instance of the neural net (a pickled NN is included with the submission). Luckily, my convergence was fast (less than 15 mins). When I was tuning parameters, I would run the neural net for 20×10^4 iterations at a time. I suppose I would never let my net run for longer than 60 mins, so that would be my only true stopping condition.

Momentum:

I implemented the common learning rate technique of adding a momentum term (a fraction of the previous weight update) every fixed number of iterations (10^4 , usually). This let gradient updates from previous iterations to persist in later iterations.

Part III

Kaggle:
0.97640

Runtime:

The neural net took 10 seconds to run 1,000 iterations. Generating the two plots directly below took 200×10 seconds \sim 33 minutes.

However, in about a quarter of that time (8.3 minutes), the loss was already flat-lining, and this corresponded to 2.7% error on the validation set.

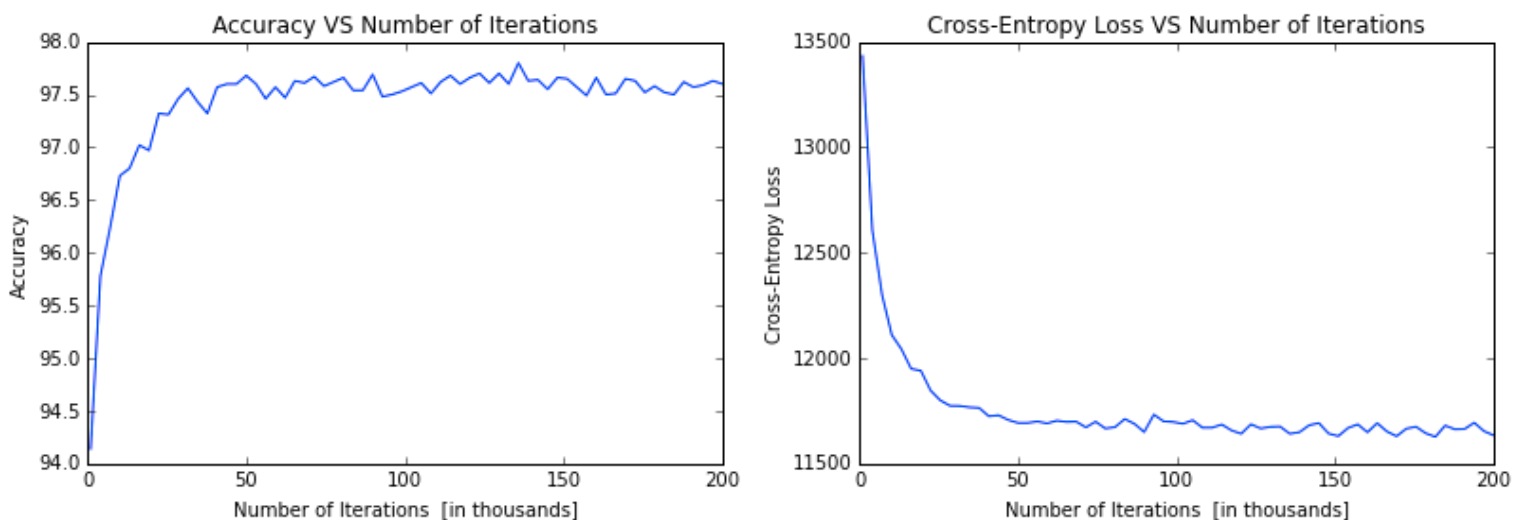
Accuracy:

On the best neural net, the validation set error was 2.4%, i.e., 97.6% accuracy (same as Kaggle). The training set error was 0.02% (99.98% accuracy). Knowing this, overfitting must have occurred, even though the loss on the validation set was minimized for this net.

Loss Functions:

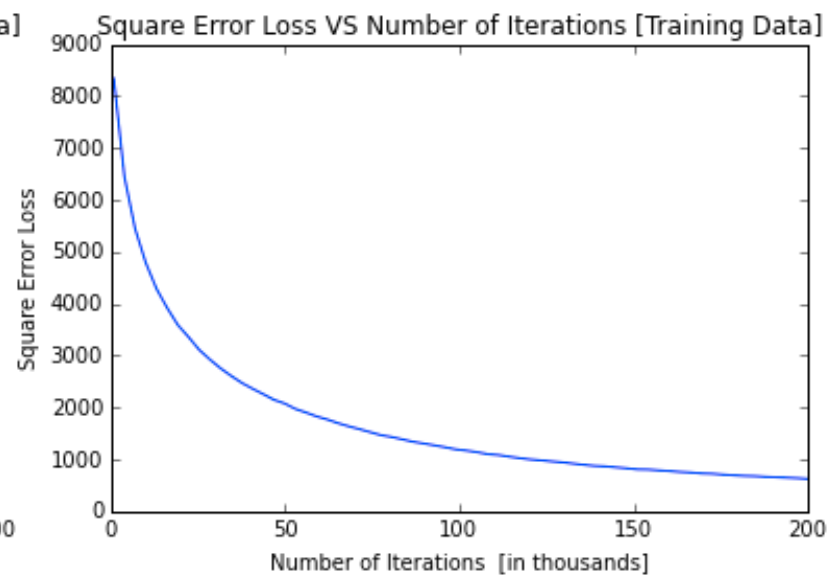
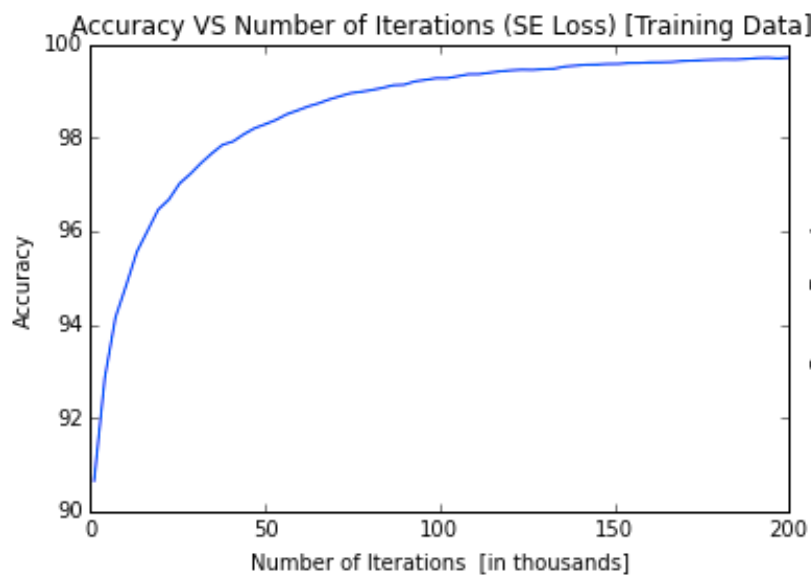
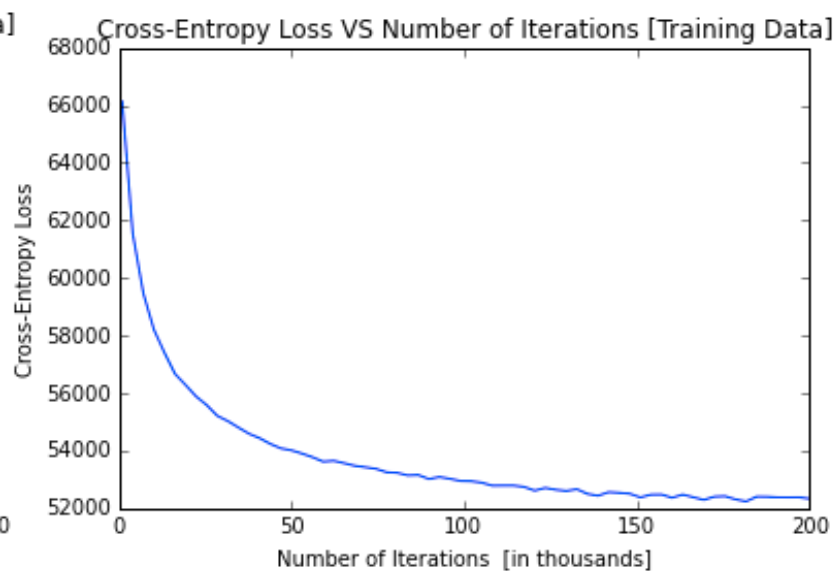
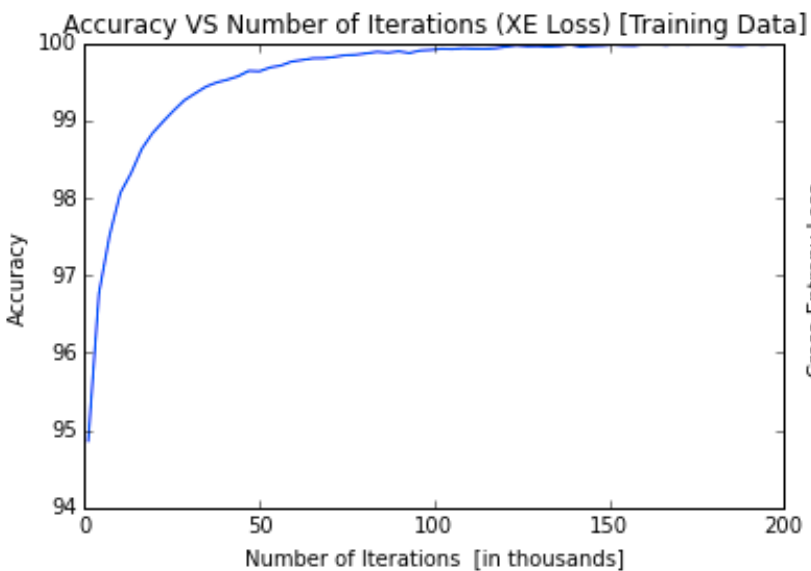
Cross-entropy (XE) loss resulted in faster convergence than square error (SE) loss on the same neural net. The rate was very roughly double. On the validation set, the XE loss-trained net started with 5.9% error at 1000 iterations in, and then settled to 2.5% after 50,000 iterations.

The SE loss-trained net began at 9.5% error and after 50,000 iterations reached 3.0% error. Interestingly, the SE loss saw much greater improvements in cost — almost 50% in 50,000 iterations. However, the XE loss was only reduced by 13%. This leads me think that the topology of the SE loss function was steeper, whereas the XE loss function had smaller gradients and could only decrease slightly with each iteration.



The plots above were generated using the validation set.

The plots below were generated with the training set.



Alex Walczak | CS189 | HW6 | PART I.

#1. MSE Loss.

$$J = \frac{1}{2} \sum_k^n (y_k - f(x_k))^2$$

Let:

① $z_k(x) = \tanh(W_k x)$, ② $W = \begin{bmatrix} -w_1^T \\ \vdots \\ -w_i^T \end{bmatrix}$,
(V , same arrangement)

③ $z_i = \tanh(Wx)$, ← applied to each

④ $f(x_k) = \sigma(Vz_i)$. ← entry.

So (i) $\frac{d}{dV} J = \frac{d}{dV} \left[\frac{1}{2} \|Y - f(x)\|_2^2 \right]$

$$= \frac{d}{dV} \left[\frac{1}{2} \|Y - \sigma(Vz)\|_2^2 \right]$$

$$= \frac{1}{2} \left[\frac{d}{dV} \left[\|\sigma(Vz)\|_2^2 \right] + 2 \frac{d}{dV} \left[\sigma(Vz)^T Y \right] \right] = \frac{d}{dV} \left[\left(\begin{bmatrix} \vdots \\ \sigma(V_i^T z) \\ \vdots \end{bmatrix} \right)^T \right] + 2 \frac{d}{dV} \left[\left(\begin{bmatrix} \vdots \\ \sigma(V_i^T z) \\ \vdots \end{bmatrix} \right)^T Y \right] \left(\frac{1}{2} \right)$$

$$= \frac{2}{2} \left[-z^T \begin{bmatrix} \vdots \\ \sigma(V_i^T z) \sigma'(V_i^T z) \\ \vdots \end{bmatrix} + \frac{2}{2} \left[-y_i \sigma'(V_i^T z) z^T \right] \right] \Rightarrow \left[\frac{d}{dV} J \right]_{\text{ith row}} = \left[-(y - \sigma(Vz)) * \sigma'(Vz) \right] z^T$$

multiply across in vector

(ii) $\frac{d}{dW} J = \frac{1}{2} \frac{d}{dW} \|Y - \sigma(Vz)\|_2^2 = \frac{1}{2} \frac{d}{dW} \left[\left\| \begin{bmatrix} \vdots \\ \sigma(V_i^T \tanh(Wx)) \\ \vdots \end{bmatrix} \right\|_2^2 \right] = \frac{1}{2} \frac{d}{dW} \left[1 + \sigma(V_i^T \tanh(Wx))^2 + \dots \right]$

Now, $\frac{d}{dW_i} [\text{prev}] = \frac{2}{2} \sigma(V_i^T \tanh(Wx)) \cdot \sigma'(V_i^T \tanh(Wx)) \cdot \frac{d}{dW_i} [V_i^T \tanh(Wx)]$

$$\Rightarrow \frac{d}{dW_i} = \sigma(V_i^T \tanh(Wx)) \cdot \sigma'(V_i^T \tanh(Wx)) \cdot \frac{d}{dW_i} [V_i^T \tanh(Wx)]$$

$$\Rightarrow \frac{d}{dW_j} J = \left[-\sum_i (y_i - \sigma(V_i^T z)) \cdot \sigma'(V_i^T z) V_j \right] \tanh(W_j^T x) x^T$$

SGD: $V_i^T \leftarrow V_i^T + \lambda \frac{d}{dV_i} J$
 $W_j^T \leftarrow W_j^T + \lambda \frac{d}{dW_j} J$

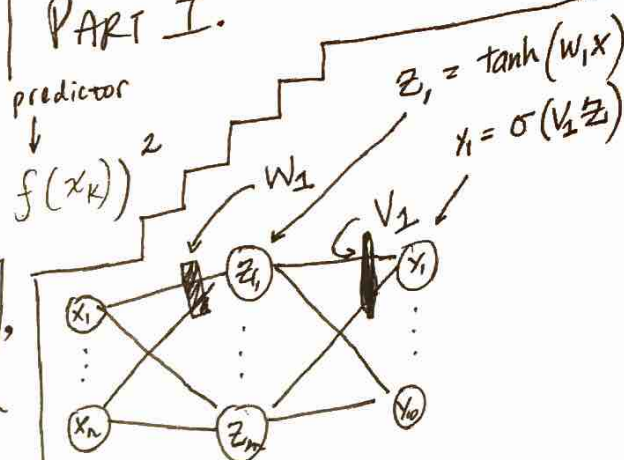


Figure 1.
The conventions used in deriving SGD.

#1 XE - Loss SGD derivation.

Now, $J = - \sum_{k=1}^n \left(y_k \ln(\sigma(\overbrace{v_k^T \tanh(x)}^{z_k(x)})) + (1-y_k) \ln(1-\sigma(z_k(x))) \right)$

$$\frac{dJ}{dv} \left(y_k \ln(\sigma(z_k(x))) + (1-y_k) \ln(1-\sigma(z_k(x))) \right)$$

$$= \frac{y_k}{\sigma(z_k(x))} \cdot \sigma'(z_k(x)) \cdot \frac{d}{dv} [v_k^T \tanh(wx)] + \frac{(1-y_k)}{1-\sigma(z_k(x))} \cdot (-\sigma'(z_k(x)) \cdot \frac{d}{dv} [v_k^T \tanh(wx)])$$

$$= y_k (1-\sigma(z_k(x))) \cdot \frac{d}{dv} (v_k^T \tanh(wx)) + (1-y_k) \cdot (-\sigma(z_k(x))) \cdot \frac{d}{dv} [v_k^T \tanh(wx)],$$

$$\Rightarrow (y_k - \sigma(z_k(x))) \cdot \tanh(wx) = \frac{d}{dv_k} J$$

$\sigma(z_k) \triangleq \sigma(z_k(x))$
 \downarrow
 $z_k \triangleq z_k(x)$

$$\frac{dJ}{dw_i} = \frac{d}{dw_i} \left[- \sum (y_k \ln \sigma(z_k) + (1-y_k) \ln (1-\sigma(z_k))) \right],$$

$$\frac{d}{dw_i} (y_k \ln \sigma(z_k) + (1-y_k) \ln (1-\sigma(z_k))) = \frac{y_k}{\sigma(z_k)} \sigma'(z_k) \frac{d}{dw_i} z_k + \frac{(1-y_k)}{1-\sigma(z_k)} (-\sigma'(z_k)) \frac{d}{dw_i} z_k$$

$$= (1-\sigma(z_k)) y_k \frac{d}{dw_i} z_k + \sigma(z_k) (1-y_k) \cdot \frac{d}{dw_i} z_k. \text{ Since } \frac{d}{dw_i} z_k = \frac{d}{dw_i} (v_k^T \tanh(wx))$$

$$= v_k^T \tanh'(wx) \cdot x, \text{ then putting it all together, we get:}$$

$$\frac{dJ}{dw_i} = \left[V^T (y - \sigma(V \tanh(wx))) \tanh'(wx) \right] x^T_{i \text{th row (entry)}}$$

So SGD equations are:

$$v_i^T \leftarrow \lambda \left(\frac{dJ}{dv_i} \right) + v_i$$

$$w_j^T \leftarrow \lambda \left(\frac{dJ}{dw_j} \right) + w_j$$


```

1 # coding: utf-8
2
3 # In[2484]:
4
5 # Alex Walczak | CS 189 | Homework 6: Neural Networks
6
7 # Import functions and libraries
8 from __future__ import division, print_function
9 import cPickle as pickle
10 import numpy as np, matplotlib.pyplot as plt
11 from matplotlib.pyplot import *
12 import scipy.io
13 from pylab import rcParams
14 from sklearn.cross_validation import KFold
15 rcParams['figure.figsize'] = 7, 7
16 import csv
17 get_ipython().magic(u'matplotlib inline')
18 from scipy import ndimage
19
20
21 # In[2996]:
22
23 # # # Serialize object
24 # f = open('pick_name.pickle', 'wb')
25 # pickle.dump(NN, f, protocol=pickle.HIGHEST_PROTOCOL)
26 # f.close()
27
28
29 # In[2997]:
30
31 # # Load object
32 # f = open('pick_name.pickle', 'rb')
33 # NN = pickle.load(f)
34 # f.close()
35
36
37 # In[2999]:
38
39 # rand_inds = np.random.permutation(60000)
40
41
42 # In[3000]:
43
44 # Load & Preprocess (training and validation sets)
45 # To enable bias, I append a 1 to front of each feature. x0 = 1. Then later, w0 = TBD.
46 data = scipy.io.loadmat("dataset/train.mat")
47 imgs = np.reshape( data["train_images"], (-1, 60000)).T[rand_inds]
48 labels0 = data["train_labels"][rand_inds]
49
50 ##### goooooooooo #####
51 imgs = imgs/255
52 imgs_mean = np.mean(imgs, axis = 0)
53 imgs = imgs - imgs_mean # zero-center the data (important)
54 #####
55
56 # 0.1/.9 target reachable by logistic and tanh functions.
57 # cs.cmu.edu/afs/cs/academic/class/15883-f15/slides/backprop.pdf
58 # Compare to: normal (unreachable) 0/1 targets.
59
60 all_labels = []
61 for label in labels0:
62     l = np.zeros(10) + 0.02
63     l[label] = 0.98
64     all_labels += [l]

```

```

65 all_labels = np.array(all_labels)
66
67 images_with_concatd1 = np.c_[ np.ones(len(imgs)) , imgs ]
68
69 images = images_with_concatd1[10000:]
70 labels = all_labels[10000:]
71
72 # Create validation set, size 10,000.
73 vimages = images_with_concatd1[:10000]
74 vlabels = all_labels[:10000]
75
76 # Uncomment and run for sanity check on preprocessing.
77 # I = 35;
78 # plt.figure(); showme(images[I]); print('First image (train):', labels[I].argmax())
79 # plt.figure(); showme(vimages[I]); print('Second image (validation):', vlabels[I].argmax())
80
81
82 # In[3001]:
83
84 # Load & Preprocess (test set)
85 test_images = scipy.io.loadmat("dataset/test.mat")['test_images']
86 test_imgs = []
87 # Load and norm test images.
88 for i in range(len(test_images)):
89     ti = test_images[i]
90     ti = np.ndarray.flatten(np.fliplr(((np.rot90(np.reshape(ti, (28,-1))))).T)))
91     # ^ Reshape data
92     test_imgs += [ti]
93
94 test_imgs = np.array(test_imgs)
95
96 # test_imgs = np.reshape( test_images, (-1, 10000)).T
97 test_imgs_scaled = test_imgs/255
98 test_imgs = test_imgs_scaled - imgs_mean
99 test_imgs = np.c_[ np.ones(len(test_imgs)) , test_imgs ]
100
101
102 # In[3004]:
103
104 # Computation
105
106 def sigmoid(gamma):
107     # activation function for output units
108     return 1/(1+np.exp(-gamma))
109
110 def deriv_sigmoid(gamma):
111     # d/d(gamma) sigmoid
112     return sigmoid(gamma)*(1-sigmoid(gamma))
113
114 def tanh(gamma):
115     # activation function for hidden units
116     return np.tanh(gamma)
117
118 def deriv_tanh(gamma):
119     # d/d(gamma) tanh
120     return -tanh(gamma)**2 + 1
121
122 def square_error(x,y,W,V):
123     return -1*-1*0.5*np.sum((y - sigmoid(V.dot(tanh(W.dot(x)))))**2)
124
125 def cross_entropy(x,y,W,V):
126     return -1*-1*-np.sum(y*np.log(sigmoid(V.dot(tanh(W.dot(x))))) + (1-y)*np.log(1-
127         sigmoid(V.dot(tanh(W.dot(x)))))
128
129 def evaluate_grad(gW,gV,W,V):

```

```

129 wnum, wden = np.linalg.norm(gW-W), np.linalg.norm(gW+W)
130 vnum, vden = np.linalg.norm(gV-V), np.linalg.norm(gV+V)
131 if vnum == vden == 0:
132     vden = 1
133 if wnum == wden == 0:
134     wden = 1
135 return wnum/wden, vnum/vden
136
137 # Utilities
138
139 def showme(img):
140     # skip first entry of array because it's always a one.
141     plt.figure()
142     plt.imshow(np.reshape(img[1:], (28,-1)))
143
144 def save_labels(labels, fname):
145     # Ex. save_labels(spam_labels, 'kaggle_spam_mean.csv')
146     f1 = open(fname, 'w+')
147     print('Id,Category', file = f1)
148     for i in range(len(labels)):
149         print(str(i+1)+" "+str(int(labels[i])), file = f1)
150
151
152 # In[3015]:
153
154 class Neural_Network():
155     # (785-200-10)
156     # with a single hidden layer.
157     # classifies MNIST images.
158
159     def __init__(self, learn_weights, cost_fn=None, Lambda=0.1, momentum=0.5,
160 MAX_ITERS=10e6):
161         self.ID = str(np.random.randint(10e4,10e8))
162         self.cost_fn = cost_fn
163         self.Lambda = Lambda
164         self.momentumW = momentum
165         self.momentumV = momentum
166         self.stddW = 1/np.sqrt(785)
167         self.stddV = 1/np.sqrt(200)
168         self.W = self.stddW*np.random.randn(200, 785)
169         self.V = self.stddV*np.random.randn(10, 200)
170
171         self.prev_gradV = 0
172         self.prev_gradW = 0
173
174 # Different layers learn at different rates.
175 # This was also mentioned in the same CMU ppt as in the early slides.
176 # Input-to-hidden (i2h) weights learn more slowly than hidden-to-output (h2o) weights.
177
178 # learn_weights = (i2h_lr, i2h_inc, i2h_dec, h2o_lr, h2o_inc, h2o_dec)
179 # 7.75% error for (0.1, 1.02, 0.6, 0.02, 1.002, 0.4)
180
181 self.i2h_lr = learn_weights[0]
182 self.i2h_inc = learn_weights[1]
183 self.i2h_dec = learn_weights[2]
184 self.h2o_lr = learn_weights[3]
185 self.h2o_inc = learn_weights[4]
186 self.h2o_dec = learn_weights[5]
187
188 self.iters = 1
189 self.total_iters = 1
190 self.MAX_ITERS = MAX_ITERS
191
192 self.W_dot_x = 'Nothing yet!'
193 self.tanh_W_dot_x = 'Nothing yet!'

```

```

193 self.forward = 'Nothing yet!'
194
195 self.last_error = 10e90
196 self.oldW = None
197 self.oldV = None
198 self.alternateVW = 1
199
200 self.errors = np.array([10,20,30])
201
202 def train(self, images, labels, learn_weights, epoch, error_inc=1.0001, error_func=None,
203 iters=None, reset_iters=False, vimages=None, vlabels=None):
204
205     self.epoch = epoch #4*10e3
206     self.error_inc = error_inc # if the error increased by more than this percentage of
207     last epoch, then we did worse. Change back weights.
208
209     self.i2h_lr = learn_weights[0]
210     self.i2h_inc = learn_weights[1]
211     self.i2h_dec = learn_weights[2]
212     self.h2o_lr = learn_weights[3]
213     self.h2o_inc = learn_weights[4]
214     self.h2o_dec = learn_weights[5]
215
216 self.iters = 1
217
218 if reset_iters:
219     self.total_iters = 1
220
221 if error_func == None:
222     print('Please choose an error function.')
223     return
224 if error_func == 'square_error':
225     gradVfunc = self.gradV_square_error
226     gradWfunc = self.gradW_square_error
227     _error = square_error
228 if error_func == 'cross_entropy':
229     gradVfunc = self.gradV_cross_entropy_error
230     gradWfunc = self.gradW_cross_entropy_error
231     _error = cross_entropy
232
233 if iters == None:
234     iters = self.MAX_ITERS
235
236 while self.iters <= iters:
237     idx = np.random.randint(1,len(images))
238     x = images[idx]
239
240     self.W_dot_x = self.W.dot(x)
241     self.tanh_W_dot_x = tanh(self.W_dot_x)
242     self.forward = sigmoid(self.V.dot(self.tanh_W_dot_x))
243
244     _gradV = self.h2o_lr*gradVfunc(x, labels[idx])
245     _gradW = self.i2h_lr*gradWfunc(x, labels[idx])
246
247     self.V += (_gradV + self.momentumV*self.prev_gradV) # does sign/order matter?
248     self.W += (_gradW + self.momentumW*self.prev_gradW) # ''
249
250     self.prev_gradV = _gradV
251     self.prev_gradW = _gradW
252
253     self.iters += 1
254     self.total_iters += 1
255
256 # Adjust learning rate.

```

```

256     if (self.total_iters % self.epoch) == 0:
257
258         self.i2h_inc = self.i2h_inc*(1 - (self.total_iters/self.MAX_ITERS))
259         self.i2h_dec = self.i2h_dec*(1 - (self.total_iters/self.MAX_ITERS))
260         self.h2o_inc = self.h2o_inc*(1 - (self.total_iters/self.MAX_ITERS))
261         self.h2o_dec = self.h2o_dec*(1 - (self.total_iters/self.MAX_ITERS))
262
263         data_error, XXX = self.check_error(vimages, vlabels)
264
265         if error_func == 'cross_entropy':
266             error = self.cost_xe(vimages, vlabels)
267         if error_func == 'square_error':
268             error = self.cost_se(vimages, vlabels)
269
270         self.errors = np.append(self.errors,error)
271
272         self.alternateVW += 1
273
274         if error < self.last_error:
275
276             if (error == self.errors[3:].min()):
277                 # Pickle if best.
278                 # # # Serialize object
279                 f = open('BEST_ENCOUNTERED_'+self.ID+'_.pickle', 'wb')
280                 pickle.dump(self, f, protocol=pickle.HIGHEST_PROTOCOL)
281                 f.close()
282                 print('pickled!')
283
284             print('better!')
285
286             if (error < 1.0):
287                 print('good enough!')
288                 break
289
290             # Doing well? Speed up!
291             if self.alternateVW % 2 == 0:
292                 self.i2h_lr = self.i2h_lr*self.i2h_inc
293                 self.oldW = np.copy(self.W)
294                 self.momentumW = self.momentumW*1.1
295
296             else:
297                 self.h2o_lr = self.h2o_lr*self.h2o_inc
298                 self.oldV = np.copy(self.V)
299                 self.momentumV = self.momentumV*1.05
300
301             if error >= self.last_error*self.error_inc:
302
303                 print('worse')
304
305                 # Getting worse? Slow down!
306                 if self.alternateVW % 2 == 1: # (sometimes) trying 1 not 0 because if I
just change W, that was a bad change bc error increased.
307                     # experiment with this!
308                     self.i2h_lr = self.i2h_dec*self.i2h_lr # i2h_dec should be at least
309 bigger than h2o_dec because don't want i2h_lr to slow down too much (takes longer to get up)
310                     self.W = np.copy(self.oldW)
311                     self.momentumW = self.momentumW*0.6
312
313             else:
314                 self.h2o_lr = self.h2o_dec*self.h2o_lr
315                 self.V = np.copy(self.oldV)
316                 self.momentumV = self.momentumV*0.3
317
318         self.last_error = error

```

```

319         # If std dev of most recent errors is low, assume we are stuck in a local
minimum.
320
321         if np.std(self.errors[-3:]) < 0.10:
322             print('Errors not changing! Time to shake things up.')
323             self.momentumW = np.random.randint(5,20)/150
324             self.momentumV = np.random.randint(5,100)/1000
325             self.i2h_lr += 0.0001*np.random.randint(1,61)
326             self.h2o_lr += 0.0001*np.random.randint(1,61)
327             self.errors[-1]+=10e9
328
329             if self.i2h_lr < 10e-8:
330                 print('i2h too small!')
331                 self.i2h_lr += 10e-6
332
333             if self.h2o_lr < 10e-14:
334                 print('h2o too small!')
335                 self.h2o_lr += 10e-10
336
337             print(self.i2h_lr, self.h2o_lr, self.momentumW, self.momentumV)
338
339         def cost_xe(self, vimages, vlabels):
340             tot=0
341             for i in range(len(vimages)):
342                 tot+=
np.sum(vlabels[i]*np.log(sigmoid(self.V.dot(tanh(self.W.dot(vimages[i]))))) + (1-
vlabels[i])*np.log(1-sigmoid(self.V.dot(tanh(self.W.dot(vimages[i]))))))
343             print('XE Cost', tot)
344             return tot
345
346         def cost_se(self, vimages, vlabels):
347             tot=0
348             for i in range(len(vimages)):
349                 tot+=np.sum( (vlabels[i] - sigmoid(self.V.dot(tanh(self.W.dot(vimages[i])))))**2
350
351             print('SE Cost', tot)
352             return tot
353
354         def predict(self, image):
355             return sigmoid(self.V.dot(tanh(self.W.dot(image))))
356
357         def predict_all(self, images):
358             return(np.array(map(self.predict, images)))
359
360         def check_error(self, images, labels):
361             # labels are assumed to be each of size 10 as in preprocessing.
362             pdn = self.predict_all(images)
363             pdn_labels = np.array([p.argmax() for p in pdn])
364             labels = np.array([l.argmax() for l in labels])
365             error = 100*np.sum(pdn_labels != labels)/len(labels)
366             print('Error', error,'%')
367             return error, pdn_labels
368
369         def check_grad(self, error_func, x, y, eps=10e-5):
370             # Call check grad to set self.forward before comparing to grad_XXX_error
371             self.forward = sigmoid(self.V.dot(tanh(self.W.dot(x))))
372
373             ogW, ogV = np.copy(self.W), np.copy(self.V)
374             gradW, gradV = np.zeros(self.W.shape),np.copy( np.zeros(self.V.shape))
375
376             for j in range(len(self.W)):
377                 for i in range(len(self.W[j])):
378                     og = self.W[j][i]
379                     self.W[j][i] = og + eps
380                     res1 = error_func(x,y,self.W,self.V)
381                     self.W[j][i] = og - eps

```



```

380         res2 = error_func(x,y,self.W,self.V)
381         gradW[j][i] = (res1 - res2)/(2*eps)
382         self.W[j][i] = og
383
384     for j in range(len(self.V)):
385         for i in range(len(self.V[j])):
386             og = self.V[j][i]
387             self.V[j][i] = og + eps
388             res1 = error_func(x,y,self.W,self.V)
389             self.V[j][i] = og - eps
390             res2 = error_func(x,y,self.W,self.V)
391             gradV[j][i] = (res1 - res2)/(2*eps)
392             self.V[j][i] = og
393
394     self.W, self.V = ogW, ogV
395     return gradW, gradV
396
397 def gradW_square_error(self, x,y):
398     a = ((y - self.forward)*deriv_sigmoid(self.V.dot(self.tanh_W_dot_x)))
399     return np.outer(self.V.T.dot(a)*deriv_tanh(self.W_dot_x), x.T)
400
401 def gradV_square_error(self, x,y):
402     return np.outer((y - self.forward)*deriv_sigmoid(self.V.dot(self.tanh_W_dot_x)),
self.tanh_W_dot_x.T)
403
404 def gradV_cross_entropy_error(self, x,y):
405     return np.outer((y - self.forward), self.tanh_W_dot_x.T)
406
407 def gradW_cross_entropy_error(self,x,y):
408     return np.outer(((self.V.T.dot(y - self.forward))*deriv_tanh(self.W_dot_x)), x.T)
409
410 def grad_cross_entropy(self, x,y):
411     return self.gradW_cross_entropy_error(x,y), self.gradV_cross_entropy_error(x,y)
412
413 def grad_square_error(self, x,y):
414     return self.gradW_square_error(x,y), self.gradV_square_error(x,y)
415
416
417 # In[3038]:
418
419 # Training and testing
420
421
422 # In[ ]:
423
424 # Values I had good results using.
425 _i2h_lr = 0.1
426 _i2h_inc = 1.015
427 _i2h_dec = 0.8
428 _h2o_lr = 0.001
429 _h2o_inc = 1.0001
430 _h2o_dec = 0.45
431 _epoch = 3*10e3
432 _error_inc = 1+10e-6
433 _momentum = 0.001
434
435 learning_rate_weights = (_i2h_lr, _i2h_inc, _i2h_dec, _h2o_lr, _h2o_inc, _h2o_dec)
436
437
438 # In[3020]:
439
440 # Training with XE Loss
441
442 NN = Neural_Network(learn_weights=learning_rate_weights, momentum=_momentum)
443

```

```

444 for i in range(10):
445     print(i)
446     NN.train(images, labels, epoch=_epoch, error_inc=_error_inc,
learn_weights=learning_rate_weights, error_func='cross_entropy', iters=10e4,
reset_iters=False, vimages=vimages, vlabels=vlabels)
447     print()
448
449 # This gets printed: (error, self.i2h_lr, self.h2o_lr, self.momentumW, self.momentumV)
450
451
452 # In[3036]:
453
454 # Training with SE Loss
455
456 # Run this line if not creating new instance of NN
457 # learning_rate_weights = (NN.i2h_lr, _i2h_inc, _i2h_dec, NN.h2o_lr, _h2o_inc, _h2o_dec)
458
459 NN = Neural_Network(learn_weights=learning_rate_weights, momentum=_momentum)
460
461 for i in range(10):
462     print(i)
463     NN.train(images, labels, epoch=_epoch, error_inc=_error_inc,
learn_weights=learning_rate_weights, error_func='square_error', iters=10e4,
reset_iters=False, vimages=vimages, vlabels=vlabels)
464     print()
465
466
467 # In[3037]:
468
469 ti = vimages[900]
470 pdn = NN.predict(ti)
471 print(pdn.argmax())
472 showme(ti)
473
474
475 # In[ ]:
476
477 # Generate test data labels.
478
479
480 # In[2974]:
481
482 # __, predicted_labels = best_NN.check_error(test_imgs, vlabels)
483
484 # # Sanity check.
485 # _t = np.array([test_imgs[1], test_imgs[22], test_imgs[333], test_imgs[444]])
486 # _l = np.array([predicted_labels[1], predicted_labels[22], predicted_labels[333],
predicted_labels[444]])
487 # for i in range(len(_t)):
488 #     showme(_t[i])
489 #     print('Label:', _l[i])
490
491 # Save labels.
492 # save_labels(predicted_labels, '784_200_10_NN_simple_preproc_labels.csv')
493
494 # Kaggle Score: 0.97640
495 # NN ID: 354600521
496
497
498 # In[ ]:
499
500 # Below, data saved for reference
501 #
502 # iterations: 10e3 to 20*10e4 with 67 pts.
503 # loss = array([ 13432.93569313, 12606.51591719, 12301.06718922, 12110.30162897,

```

```

504 #      12042.57279805, 11949.50139832, 11940.63769764, 11846.70304254,
505 #      11799.9777952 , 11774.20877548, 11773.6444034 , 11767.0068726 ,
506 #      11764.23110766, 11724.4665161 , 11729.1478896 , 11706.98543717,
507 #      11693.03176399, 11693.04117115, 11699.62146882, 11690.84224179,
508 #      11703.65743183, 11697.59636369, 11699.06362548, 11671.76060895,
509 #      11698.92654993, 11666.41520784, 11673.01219267, 11711.40579961,
510 #      11688.67719221, 11649.71231889, 11732.35222877, 11700.34949357,
511 #      11697.32432718, 11688.9760188 , 11705.92981412, 11670.40111919,
512 #      11670.51814198, 11685.40041635, 11657.34283291, 11641.43742866,
513 #      11686.59015803, 11666.59176936, 11674.29694423, 11675.79256358,
514 #      11641.35108495, 11648.90007011, 11682.20514718, 11692.19700681,
515 #      11642.26152735, 11630.38092284, 11669.01027341, 11686.12114843,
516 #      11648.40395379, 11692.13425494, 11652.37114204, 11629.60390472,
517 #      11665.57532849, 11674.92947925, 11643.88537853, 11626.71857277,
518 #      11680.50829383, 11662.69287411, 11663.85694344, 11693.98101753,
519 #      11653.85801366, 11633.78618821))
520 # accuracy = np.array([5.86 , 4.22 , 3.76 , 3.27 , 3.2 , 2.98 , 3.03 , 2.68 , 2.69 , 2.54 ,
      2.44 , 2.57 , 2.68 , 2.43 , 2.4 , 2.4 , 2.32 , 2.4 , 2.54 , 2.43 , 2.53 , 2.37 , 2.39 , 2.33
      , 2.42 , 2.38 , 2.34 , 2.46 , 2.46 , 2.31 , 2.52 , 2.5 , 2.47 , 2.43 , 2.39 , 2.49 , 2.38 ,
      2.32 , 2.4 , 2.34 , 2.3 , 2.39 , 2.3 , 2.4 , 2.2 , 2.37 , 2.36 , 2.45 , 2.34 , 2.35 , 2.43 ,
      2.51 , 2.34 , 2.5 , 2.49 , 2.35 , 2.37 , 2.48 , 2.42 , 2.48 , 2.5 , 2.38 , 2.43 , 2.41 ,
      2.37 , 2.4])
521
522
523 # In[3039]:
524
525 # plt.plot(np.linspace(10e3,20*10e4,66)/10e3,NN.errors[3:])
526 # plt.title('Cross-Entropy Loss VS Number of Iterations')
527 # plt.xlabel('Number of Iterations [in thousands]')
528 # plt.ylabel('Cross-Entropy Loss')
529
530
531 # In[3040]:
532
533 # plt.plot(np.linspace(10e3,20*10e4,66)/10e3,accuracy)
534 # plt.title('Accuracy VS Number of Iterations')
535 # plt.xlabel('Number of Iterations [in thousands]')
536 # plt.ylabel('Accuracy')

```