

Tutorial September 26, 2019

This document contains rough code outline that describes several key steps in training and evaluating a neural network model, using PyTorch. It describes the dataset class, the trainloader class, and the training and evaluation loop.

Please note that this code HAS NOT been compiled and run, and so may contain errors.

```
# Assume that, at this point, we have a DataFrame object, called allData,
# that contains both the input data/features (all in numerical form
# at this point, having converted any categorical features
# to 1-hot and normalized the continuous features) and the
# labels for each sample (= row in the original csv file)
# To be clear, to get to this point you will need to have done the
# processing described in Assignment 3 from Section 3.1 through 3.6
```

```
# Let's keep it simple, and assume there are only 4 numerical
# input features and one (binary) output *label* which is
# either 0 or 1; (in assignment 3 you should have over 100 input
# features)
```

```
# This data needs to be split in two ways:
# 1. We need to split the set into training and validation sets
# 2. We need to separate the input features and output labels
# We should also convert these into numpy array types, as that is the
# input to the Dataloader class requires
```

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
```

```
import torch
from torch.utils.data import DataLoader
```

```
# Below is also rough outline of the model definition
# Previous descriptions in class should help you with this part
```

```
class MultiLayerPerceptron(nn.Module):
```

```
    def __init__(self, input_size, Other_NN_HyperParameters):
```

```
        super(MultiLayerPerceptron, self).__init__()
```

```
        self.input_size = input_size
```

```
        self.output_size = 1
```

```
        # YOU WILL NEED OTHER FUNCTIONS NEEDED TO DEFINE THE MODEL
        # HERE, as discussed previously
```

```
    def forward(self, features):
```

```

55         x = SOMETHING (features)
56         x = SOMETHINGELSE (x)
57
58     # CONTINUE TO DEFINE MODEL as discussed previously in Lecture 8
59
60     return x
61
62 # Separate features from labels
63
64 # 'labelcol' is the heading in the csv file of the binary encoded label
65 allLabels = allData["labelcol"]
66
67 #convert DataFrame object to numpy array because dataloader requires
68 #a numpy array type as its input
69
70 allLabels = allLabels.values
71
72 # remove label column from DataFrame object, leaving just the features
73
74 features = allData.drop(columns=["labelcol"])
75
76 features = features.values # also convert DataFrame to numpy 2D array
77
78
79 # now, separate into training and validation set, randomly,
80 # With 20% going to validation set
81 # Should set random seed so that program is the same each execution
82
83 seed = 0
84
85 feat_train, feat_valid, label_train, label_valid =
86
87     train_test_split(features, labels, test_size=0.2,random_state=seed)
88
89
90 # Need to set up the DataLoader - which feeds the
91 # training and validation loops with the data
92 #
93 # DataLoader requires the Dataset object:
94 # A Dataset contains both the features and the labels
95 # it requires the methods: __init__, __len__ (number of features)
96 # and __getitem__ - get one sample
97
98 import torch.utils.data as data
99
100 class myDataset(data.Dataset):
101
102     def __init__(self, features, label):
103
104         # these must be numpy array as mentioned above
105         # this gives methods in this class access to features and labels
106
107         self.features = features
108         self.label = label

```

```

109
110     def __len__(self):
111         return len(self.features)
112
113     # getitem returns an individual sample's features and label
114     # if there is more than one feature, it should be a numpy array
115
116     def __getitem__(self, index):
117         features = self.features[index]
118
119         label = self.label[index]
120
121         return features, label
122
123     # The code below sets up the separate training and validation datasets
124     # putting them into a DataLoader object
125
126     # First, instantiate an object in the dataset class defined above, and
127     # fill it with the training data
128
129     train_dataset = myDataset(feats_train, label_train)
130
131     # create a callable object that will provide 'batches' of the samples
132     # later on, when asked for the training data
133     # this converts (invisibly) the data to be PyTorch-compatible tensors
134     # This is also where the batch_size is set
135     # Shuffle = True re-orders data every Epoch
136
137     train_loader = DataLoader(train_dataset, batch_size=batch_size,
138                               shuffle=True)
139
140     # similarly, do the same for the validation data set.
141
142     valid_dataset = myDataset(feats_valid, label_valid)
143
144     valid_loader = DataLoader(valid_dataset, batch_size=batch_size,
145                               shuffle=False)
146
147     # skip validation function below for the moment,
148     # will return to it after training loop
149
150     # a function to run the validation data set through the model
151     # This will be done every so often during the training loop below
152
153     def evaluate(model, valid_loader):
154         total_corr = 0
155
156         for i, vbatch in enumerate(valid_loader):
157             feats, label = vbatch      #feats will have shape (batch_size,4)
158
159             # run the neural network model on the data!
160             # note that there are batch_size samples being run through
161             # the model, not just one sample
162

```

```

163         prediction = model(feats)
164
165         # if a prediction is OVER 0.5 - that is considered to be 1
166         # otherwise answer is 0.
167         # Squeeze: shape is (batchsize,1) results - don't want
168         # that 1 dimension
169         # use long because that is the basic integer type
170
171         corr = (prediction > 0.5).squeeze().long() == label
172
173         # sum up the number of correct predictions
174
175         total_corr += int(corr.sum())
176
177     return float(total_corr)/len(valid_loader.dataset)
178
179 # NOW, heading towards the training loop!!!
180
181 # Choose the loss function to be Mean Squared Errir
182 # it is a callable object; will describe intuition later.
183
184 loss_function = torch.nn.MSELoss()
185
186 # Instantiate the callable object that is the neural NETWORK model
187 # defined in the model section above
188 #
189
190
191 model = MultiLayerPerceptron(4, OTHERPARAMETERS....)
192
193 # Choose the optimization method - Stochastic Gradient Descent
194 # model.parameters() contains all of the weights and biases defined in
195 # the model definition
196
197 # lr is the value of the learning rate that you're setting
198
199 optimizer = torch.optim.SGD(model.parameters(), lr=args.learningrate)
200
201
202 #
203 # THE TRAINING OPTIMIZING LOOP
204 #
205
206     t = 0    # used to count batch number putting through the model
207
208
209     for epoch in range(MaxEpochs): # Epoch: one pass through all training
210         accum_loss = 0
211         tot_corr = 0
212
213         for i, batch in enumerate(train_loader):    #from DataLoader
214             # this gets one "batch" of data;p
215
216             feats, label = batch    #feats will have shape (batch_size,4)

```

```

217
218     # need to send batch through model and do a gradient opt step;
219     # first set all gradients to zero
220
221     optimizer.zero_grad()
222
223     # Run the neural network model on the batch, and get answers
224
225     predictions = model(feats)    # has shape (batch_size,1)
226
227     # compute the loss function (MSE as above) using the
228     # correct answer for the entire batch
229     # label was an int, needs to become a float
230
231     batch_loss = loss_fnc(input=predictions.squeeze(),
232 target=label.float())
233
234     accum_loss += batch_loss
235
236     # computes the gradient of loss with respect to the parameters
237     # pytorch keeps all kinds of information in the Tensor object
238     # to make this possible; uses back-propagation
239
240     batch_loss.backward()
241
242     # Change the parameters in the model with one 'step' guided by
243     # the learning rate. Recall parameters are the weights & bias
244
245     optimizer.step()
246
247     # calculate number of correct predictions
248
249     corr = (predictions > 0.5).squeeze().long() == label
250
251     tot_corr += int(corr.sum())
252
253
254     # evaluate model on the validation set every eval_every steps
255
256     if (t+1) % args.eval_every == 0:
257         valid_acc = evaluate(model, valid_loader)
258
259         print("Epoch: {}, Step {} | Loss: {} | Valid acc:
260 {}".format(epoch+1, t+1, accum_loss / eval_every, valid_acc))
261
262
263         accum_loss = 0
264
265         t = t + 1
266     # output final, post-training training accuracy
267
268     print("Train acc: {}".format(float(tot_corr)/len(train_loader.dataset)))
269
270

```