

Neural Network Tuning and Advanced Features

1. Initializing and Training Models

- Simple question: What is the difference between initializing your model and training it in the same cell, or doing it in two different cells?
- Two ways to organize the code

i.

```
device = "cuda" if torch.cuda.is_available() else "cpu"
print(f"Using {device}")
print()

num_epochs = 500

spam_ham_model = SpamModel().to(device) # <=====

training_losses = np.zeros(num_epochs)
val_losses = np.zeros(num_epochs)
```

- Spam model and initialized the model
- Training the model in the same cell

ii.

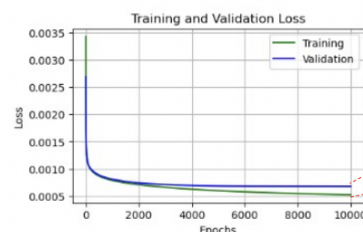
```
spam_ham_model = SpamModel().to(device)

device = "cuda" if torch.cuda.is_available() else "cpu"
print(f"Using {device}")
print()

patience = 20 # how many epoches to wait with no improvement in
               # loss score before termination
```

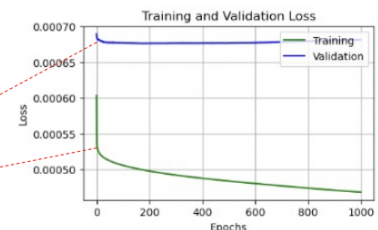
- Initialized the model in the above cell
 - Training the model in a different cell
- Answer: Nothing, as long as you always remember to run both cells for each training, especially if you change the hyperparameters!
 - If we do initialize and train, there is nothing different
 - The second picture → If you run it more than once, you are running the model for more epochs (continuing to run the trained model)
 - What if you don't?
 - Then you will be retraining an already-trained model! Fine if that is what you intend, but easy to forget, and you will get strange results:

Training for 10,000 epochs



Final Training Loss: 0.000521
Final Validation Loss: 0.000675

Retraining for 1000 more



Final Training Loss: 0.000468
Final Validation Loss: 0.000681

- g. Here is a nice way to avoid confusion and have both alternatives

```
: 1 # test if GPU is available
2
3 device = "cuda" if torch.cuda.is_available() else "cpu"
4 print(f"Using {device}")
5 print()
6
7 num_epochs = 10000
8
9 # Normally, will create model and train it in one run
10 # If want to retrain the model with more epoches, set next to True
11
12 # retrain = True
13 retrain = False
14
15 if not retrain:
16     spam_ham_model = SpamModel().to(device)
17     train_loss = np.zeros(num_epochs)
18     val_loss = np.zeros(num_epochs)
19
20     train_accuracy = np.zeros(num_epochs)
21     val_accuracy = np.zeros(num_epochs)
22
23 learning_rate = 0.1
24 learning_rate = 0.01
25 # learning_rate = 0.001
26 # learning_rate = 0.0001
27
28 # optimizer = torch.optim.SGD(spam_ham_model.parameters(),lr=learning_rate)
29 # optimizer = torch.optim.Adam(spam_ham_model.parameters(),lr=learning_rate)
30 optimizer = torch.optim.Adagrad(spam_ham_model.parameters(),lr=learning_rate)
31 # optimizer = torch.optim.RMSprop(spam_ham_model.parameters(),lr=learning_rate)
```

i.

2. Avoiding Redundant Computations

- Try to avoid redoing the same expensive operations over and over!
- We saw this with the Brown Corpus, which downloads the first time to your local disk, and thereafter checks to see if you already have it:

```
In [2]: 1 import numpy as np
2 import nltk
3 # First time you will need to download the corpus:
4 # Run the following and download the book collection
5
6 #nltk.download_shell()
7
8
```

```
In [3]: 1 from nltk.corpus import brown
2 nltk.download('brown')
3
```

```
[nltk_data] Downloading package brown to
[nltk_data] /Users/waynesnyder/nltk_data...
[nltk_data] Package brown is already up-to-date!
```

```
Out[3]: True
```

i.

- You can do this with any data structure, such as tensors or numpy arrays.
- Here is a way to do that with HW 04, Problem 2:

```

def load_glove_model(file):
    ...

import os

if os.path.exists(data_dir+'texts_vector.pt'):
    texts_vector = torch.load(data_dir+'texts_vector.pt')
else:
    glove_model = load_glove_model(data_dir+'glove.6B/glove.6B.100d.txt')
    sp = spacy.load('en_core_web_sm')
    emails_raw = pd.read_csv(data_dir+'data_pa5/enron_spam_ham.csv').to_numpy()
    texts_vector = []

    for text,label in tqdm(emails_raw):
        text_vector = torch.tensor([0]*100,dtype=torch.float32) #size of the word vector
        document=sp(text.lower())
        count = 0

        for word in document:
            if str(word) in glove_model:
                str_word = str(word)
                text_vector = text_vector + glove_model[str_word]
                count += 1

        if count>0:
            text_vector /= count

        texts_vector.append((text_vector,torch.tensor(label,dtype=torch.int64)))

    torch.save(texts_vector,data_dir+'texts_vector.pt')

```

Thereafter, if you already have the file, just read it in using `torch.load(...)`

The first time, do the expensive computation, and save it to disk using `torch.save(...)`

- i.
- ii. Allocate the list and fill in the list instead of creating an empty list and appending it
- e. And of course you can save already-trained models, and not have to retrain them to use them later; again, we did this in HW 04 Problem 2:

```

1 # now save the best model found so far, defined by least validation loss
2
3 if epoch == 0:
4     least_val_loss = val_loss[epoch]
5
6 if val_loss[epoch] < least_val_loss:
7     least_val_loss = val_loss[epoch]
8     best_model = copy.deepcopy(spam_ham_model) # make deep copy to avoid I/O cost each time
9     best_epoch = epoch
10
11 .....
12
13 # testing using the best model found during training
14
15 best_model.eval()
16
17 testing_num_correct = 0
18
19 for X_test_batch,Y_test_batch in spam_ham_test_dl:
20     X_test_batch = X_test_batch.to(device)
21     Y_test_batch = Y_test_batch.to(device)
22
23     Y_hat_test = best_model(X_test_batch)
24
25     testing_num_correct += (torch.argmax(Y_hat_test,dim=1) == Y_test_batch).float().sum()
26
27 test_accuracy = testing_num_correct / N_test
28
29 .....
30
31
32 # Save the best model found during this training
33
34 torch.save(best_model, data_dir+'best_spam_ham_model.pt')
35
36
37 # You can load it any time to use it:
38
39 spam_ham_model = torch.load(data_dir+'best_spam_ham_model.pt')

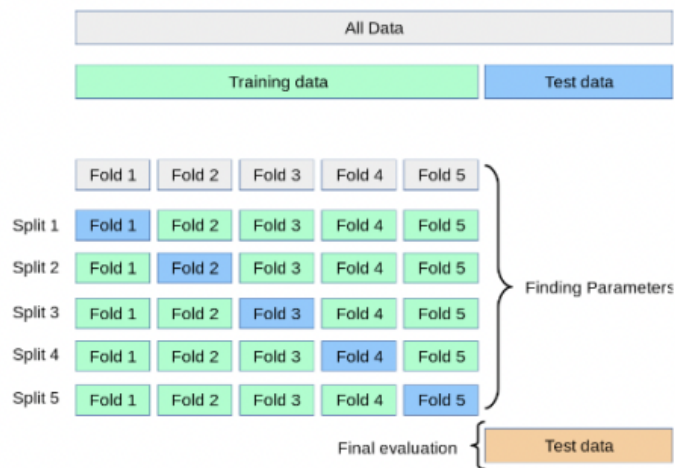
```

Use `copy.deepcopy(...)` to avoid doing expensive I/O in the main training loop!

Use the best model found during validation to do the testing!

Save and load the model as you wish!

- i.
 - ii. Use deepcopy to save the model instead of shallow copy since shallow copy simply shares the inner pointers
- ### 3. Cross Validation
- a. Cross-Validation is a dynamic alternative to choosing a fixed validation set (make sure that the data is shuffled, especially when the data is unbalanced):



- b.
- i. Dividing into 5 folds is generally most efficient (divide the validation set and test it using different validation)
 - ii. Split the data into a training set and a testing test, and hold out the testing set as usual;
 - iii. Now split the training set into K parts (“folds”) of approximately equal size;
 - iv. Training occurs in cycles of K epochs:
 - v. For each k in $\text{range}(K)$:
 - Train on all Folds except Fold k ;
 - Validate on Fold k
 - vi. At the conclusion of cycle of K epochs, take the mean of the loss and accuracy metrics;
 - vii. Report performance metrics for these means every K epochs.
- c. Cross-Validation and Static Validation have symmetric advantages and disadvantages:
- d. Static validation is
- i. Simpler and faster;
 - ii. Very dependent on quality of split, especially for small or unbalanced data sets:
 1. May overfit on that specific set;
 2. Performance metrics may be skewed.
- e. Cross-validation is
- i. More complex to implement, less efficient;
 - ii. Uses entire training set for validation, so exact split is less critical;
 1. Less possibility of overfitting;
 2. More accurate performance metrics
- f. Does not well as work for time-series data sets (e.g., stock prices, weather)
- g. Punchline: Static validation is fine for large datasets (always shuffle!!);

- h. Cross-Validation should be used for small or unbalanced data sets.
- i. In Pytorch, you can simply create K different DataLoaders, and DIY as just described; sklearn also has a popular library KFold to make it simple but inefficient:

```
from sklearn.model_selection import KFold
from torch.utils.data import DataLoader, Subset

dataset = MyDataset()
k_fold = KFold(n_splits=5)
for train_indices, val_indices in k_fold.split(dataset):

    # Using Subset to create datasets for training and validation
    train_subset = Subset(dataset, train_indices)
    val_subset = Subset(dataset, val_indices)

    train_loader = DataLoader(train_subset, batch_size=32, shuffle=True)
    val_loader = DataLoader(val_subset, batch_size=32)

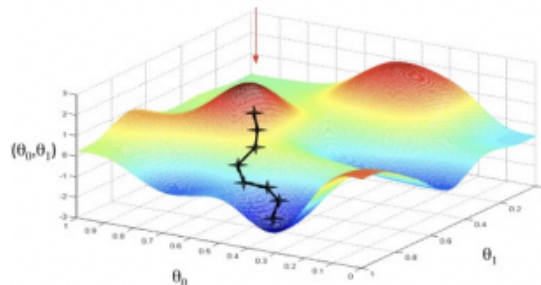
    # Now you can use train_loader and val_loader in your training and
    # validation loops
```

4. Optimizers: SGD, Adam, Adagrad, RMSProp

- a. SGD (Stochastic Gradient Descent):
 - i. Classic optimizer that updates the weights by taking a step in the direction of the negative gradient of the loss function
 - ii. Calculate the gradient in each direction
 - iii. Function depends on the learning rate as well
 - 1. Higher the learning rate, you reach the optimal point faster but are not guaranteed to reach the point due to the big step
 - 2. Lower learning rate, you reach the optimal point slower but is more guaranteed to reach the optimal point

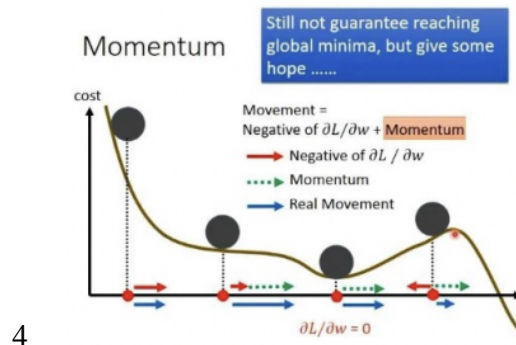
$$w_{\text{new}} = w_{\text{old}} - \eta \times \nabla_w \text{Loss}$$

iv.



v.

- b. Can set various parameters such as
- i. Learning Rate Schedules (changing the learning rate in the middle by starting with large step and gradually decreasing the learning rate) with
 1. Step decay: reduce the learning rate by some factor each epoch
 2. Exponential decay: Decrease the learning rate exponentially over the epochs;
 3. 1/t decay: reduce the lr as the inverse of the square root of the number of epochs;
 - ii. Momentum:
 1. Add a fraction (between 0 and 1) of the previous weight update to the current update
 2. Recommended is point 0.9
 3. Helps accelerate in the relevant direction and dampen oscillations



- a. May help to avoid local optimal point to reach the global optimal point
- iii. Weight decay:
 1. Equivalent to L2 regularization
- c. Adagrad (Adaptive Gradient Algorithm):
 - i. Stored a running sum of the squares of past gradients and divides the learning rate by the square root of the running sum:

$$w_{\text{new}} = w_{\text{old}} - \frac{\eta}{\sqrt{s_t + \epsilon}} \times \nabla_w \text{Loss}$$

$$s_t = s_{t-1} + \nabla_w \text{Loss} \odot \nabla_w \text{Loss}$$
 - ii.
 - iii. It is similar to the momentum idea but uses gradients
 - iv. Pro: Adapts to size of gradients.
 - v. Con: Can adapt too much and stop learning!
- d. RMSProp (Root Mean Square Propagation):
 - i. More effective version of Adagrad, using a moving average of squared past gradients:

$$s_t = \beta s_{t-1} + (1 - \beta) \nabla_w \text{Loss} \odot \nabla_w \text{Loss}$$
 - ii.
 - iii. RMSProp tends to work better for very deep neural networks

- e. Adam (Adaptive Moment Estimation):
 - i. Keeps the best features of Adagrad and RMSprop
 - f. Improves on Adagrad and RMSprop by combining both approaches with regard to past gradients:
 - i. Keep a moving weighted average of both the past gradients and the squared past gradients (called first and second moments), and adjust the learning rate accordingly.
 - ii. Corrects for initial bias in the moving averages, so tends to have more stable starts than other algorithms.
 - g. Punchlines:
 - i. SGD has many parameters which can be tuned for excellent performance, and may lead to better performance.
 - ii. Adam is the default optimizer for many tasks because it tends to “work well out of the box” without a lot of tuning.
5. Regularization: L1, L2, Dropout
- a. Regularization attempts to prevent overfitting by preventing models from becoming too complex. There is a large variety of ways to accomplish this:
 - b. Adding noise:
 - i. Produce random fluctuations in the data through augmentation;
 - 1. Kind of see the best outline (general shape)
 - ii. The network generalizes instead of focusing on the details.
 - c. L1 Regularization (Lasso Regression):
 - i. Adds a penalty proportional to the absolute value of the coefficients:

$$\lambda \sum |w_i|$$
 - ii. This prevents the parameters from becoming too large, limiting their range.
 - d. L2 Regularization (Ridge Regression):
 - i. Adds a penalty proportional to the square of the magnitude of the coefficients:

$$\lambda \sum w_i^2$$
 - e. L2 is generally preferred, since L1 can force some parameters to 0.
 - f. L2 Regularization is accomplished in Pytorch using the `weight_decay` parameter in the optimizer:

```

39 # weight_decay = 0.02 is the strength of the L2 regularization
40
41 optimizer = torch.optim.Adam(spam_ham_model.parameters(), lr=0.01, weight_decay = 0.02)
42

```

- g. The effect is to add the following penalty term to the loss L calculated during training:

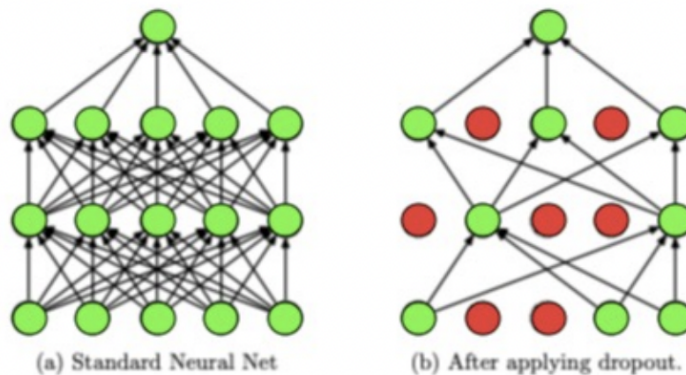
$$L' = L + \frac{\text{weight_decay}}{2} \sum w^2$$

Parameters of model

- h. (L1 regularization is not implemented in Pytorch and you would have to DIY.)

i. Dropout

- i. During training, parameters are set to 0 with some probability p
- ii. This prevents parameters from co-evolving and effectively memorizing the data
- iii. The knowledge implicit in the data is generalized throughout the network and not localized in specific parameters



- iv.
 - v. Note: Due to the random nature of dropout, different neurons will be dropped out for each data sample.
 - vi. This changes over time constantly
- j. Dropout in Pytorch is easily accomplished with a dropout layer build into the network geometry:

```

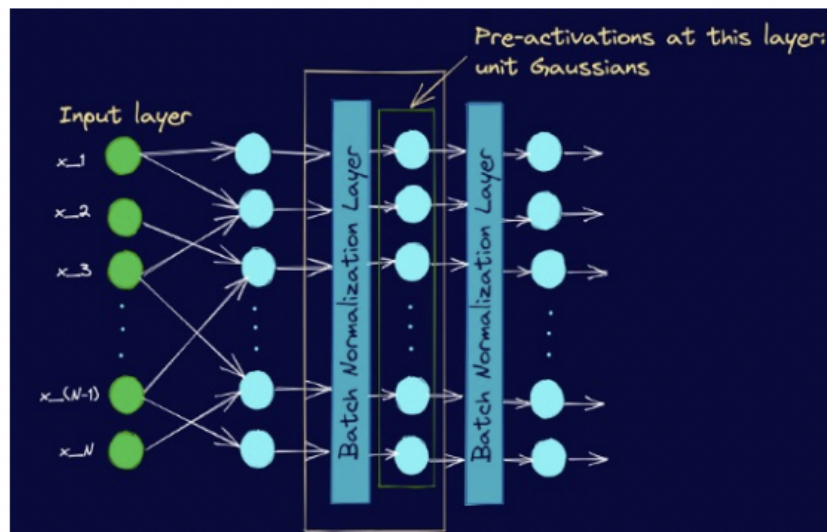
1 class SpamModel(torch.nn.Module):
2     def __init__(self):
3         super(SpamModel, self).__init__()
4         self.linear1 = torch.nn.Linear(100,15)
5         self.activation1 = torch.nn.ReLU()
6         self.linear2 = torch.nn.Linear(15,2)
7         self.dropout = nn.Dropout(0.4) # dropout neurons with probability 0.4
8
9     def forward(self, x):
10        x = self.linear1(x)
11        x = self.dropout(x)
12        x = self.activation1(x)
13        x = self.linear2(x)
14        return x

```


6. Layer and Batch Normalization

- a. Layer Normalization: For each output value from a layer:
 - i. Compute its mean μ and standard deviation σ
 - ii. Normalize to mean 0 and standard deviation 1; and then
 - iii. Scale and shift it by two parameters learned during training.
- b. This is done after every individual data sample.

$$y_i = \gamma \left(\frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}} \right) + \beta$$



- c.
- d. Batch Normalization is the same process, but applied to all layer outputs for a whole mini-batch.
- e. Normalization is also considered to be a form of regularization, because it limits the range of parameters.
- f. Why normalize layer and batch outputs?
 - i. Helps gradient flow by avoiding disappearing or exploding gradients;
 - ii. Acts as a regularizer to avoid overfitting by introducing “useful noise” into the parameters;
 - iii. Smooths the gradient landscape:
 - 1. Allows for higher learning rates and faster convergence;
 - 2. Makes weight initialization strategy less critical.
- g. Normalization is particularly effective with deep networks

h. Batch normalization in PyTorch

```
class BatchNormNet(nn.Module):
    def __init__(self, input_dim, hidden_dim, output_dim):
        super(BatchNormNet, self).__init__()
        self.fc1 = nn.Linear(input_dim, hidden_dim)
        self.bn1 = nn.BatchNorm1d(hidden_dim)
        self.fc3 = nn.Linear(hidden_dim, output_dim)

    def forward(self, x):
        x = self.fc1(x)
        x = self.bn1(x)          # note that normalization is
        x = F.relu(x)           # done before relu or sigmoid
        x = self.fc3(x)
        return x
```

7. Early Stopping

- a. Your goal is to find the best possible model for your task, typically measured by
 - i. Minimum loss score
 - ii. Maximum accuracy score
 - iii. Optimal value of some other metric (specificity, F1, etc.)
- b. The consensus view is that “it depends” but as a first approximation, loss is better than accuracy.
- c. There is no reason to overfit by continuing past this point, and there are simple ways to implement early stopping:
 - i. Stop when you reach some threshold of loss;
 - ii. Define a parameter patience, and stop training if your model does not improve (i.e., the loss does not decrease) after patience epochs. Typical values are 20 – 50 epochs.
- d. Early Stopping in PyTorch

```
8
9 patience = 50 # how many epoches to wait with no improvement in
10 # loss score before termination
11
12 early_use_stopping = True
13 # use_early_stopping = False
14
```

```
87 # now save the best model found so far, defined by least validation loss
88
89 if epoch == 0:
90     least_val_loss = val_loss[epoch]
91     num_epoches_no_improvement = 0
92
93 if val_loss[epoch] < least_val_loss:
94     least_val_loss = val_loss[epoch]
95     best_model = copy.deepcopy(span_ham_model) # make deep copy to avoid I/O cost each time
96     best_epoch = epoch
97     num_epoches_no_improvement = 0
98 else:
99     num_epoches_no_improvement += 1
100
101 last_epoch = epoch # save if early stopping
102
103 # early stopping
104
105 if use_early_stopping and num_epoches_no_improvement == patience:
106     print(f"Early stopping at epoch {last_epoch}, no improvement in validation loss after {patience} epoches.")
107     break
108
```