# Advanced Data Analysis
# Homework Week 7

Aswin Vijay

June 12, 2023

## 1] Math

We are asked to derive $\frac{\partial J}{\partial x_i}$ where $x_i$ is the $i'th$ units pre-activation value and $J$ is the loss. We are given $\frac{\partial J}{\partial u_i}$, where $u_i$ is the output after the batch normalization operation. We also have,

$$\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^{m} x_i \quad mini-batch mean \tag{1}$$

$$\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_\beta)^2 \quad mini-batch\ variance \tag{2}$$

$$\hat{x}_i = \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad normalized\ x_i \tag{3}$$

$$u_i = \gamma \hat{x}_i + \beta \quad Scale\ and\ shift \tag{4}$$

To derive $\frac{\partial J}{\partial x_i}$ we write out the partial derivatives using chain rule since $\hat{x}_i(x_i, \mu, \sigma^2)$, $\sigma^2(x_i, \mu)$ and $\mu(x_i)$, three variables are dependent on $x_i$ so:

$$\frac{\partial J}{\partial x_i} = \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial J}{\partial \mu} \cdot \frac{\partial \mu}{\partial x_i} + \frac{\partial J}{\partial \sigma^2} \cdot \frac{\partial \sigma^2}{\partial x_i}$$

The above derivatives are computed as follows,

$$\frac{\partial J}{\partial \hat{x}_i} = \frac{\partial J}{\partial u_i} \cdot \frac{\partial u_i}{\partial \hat{x}_i} = \gamma \cdot \frac{\partial J}{\partial u_i} \ From\ (4)$$

$$\frac{\partial \hat{x}_i}{\partial x_i} = \frac{1}{\sqrt{\sigma^2 + \epsilon}} \ From\ (3)$$

Computing $J$ derivatives w.r.t $\mu$,

$$\frac{\partial J}{\partial \mu} = \sum_{i=1}^{m} \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial \mu} + \frac{\partial J}{\partial \sigma^2} \cdot \frac{\partial \sigma^2}{\partial \mu} \; Summing \; over \; batch$$

$$\frac{\partial \hat{x}_i}{\partial \mu} = \frac{-1}{\sqrt{\sigma^2 + \epsilon}} \; From \; (3)$$

$$\frac{\partial \sigma^2}{\partial \mu} = \frac{-2}{m} \sum_{i=1}^{m} (x_i - \mu) = -2 \left( \frac{1}{m} \sum_{i=1}^{m} x_i - \mu \right) = 0 \; From \; (2,1)$$

Computing $J$ derivatives w.r.t $\sigma$,

$$\frac{\partial J}{\partial \sigma^2} = \sum_{i=1}^{m} \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial \sigma^2} \; Summing \; over \; batch$$

$$\frac{\partial \hat{x}_i}{\partial \sigma^2} = \frac{-1}{2} (\sigma^2 + \epsilon)^{-\frac{3}{2}} (x_i - \mu) \; From \; (3)$$

Putting everything together we have,

$$\frac{\partial J}{\partial \mu} = \sum_{i=1}^{m} \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma^2 + \epsilon}} + 0, \; then$$

$$\frac{\partial J}{\partial x_i} = \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma^2 + \epsilon}} + \sum_{i=1}^{m} \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma^2 + \epsilon}} \cdot \frac{\partial \mu}{\partial x_i} + \frac{-1}{2} (\sigma^2 + \epsilon)^{-\frac{3}{2}} \sum_{j=1}^{m} \frac{\partial J}{\partial \hat{x}_j} \cdot (x_j - \mu) \cdot \frac{\partial \sigma^2}{\partial x_i}$$

The remaining derivatives are calculated as below,

$$\frac{\partial \mu}{\partial x_i} = \frac{1}{m} \; From \; (1)$$

$$\frac{\partial \sigma^2}{\partial x_i} = \frac{2(x_i - \mu)}{m} \; From \; (2)$$

Finally we get,

$$\frac{\partial J}{\partial x_i} = \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma^2 + \epsilon}} + \frac{1}{m} \sum_{i=1}^{m} \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma^2 + \epsilon}} - \frac{(\sigma^2 + \epsilon)^{-\frac{3}{2}}}{m} \sum_{j=1}^{m} \frac{\partial J}{\partial \hat{x}_j} \cdot (x_j - \mu) \cdot (x_i - \mu)$$

$$= \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma^2 + \epsilon}} + \frac{1}{m} \sum_{i=1}^{m} \frac{\partial J}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma^2 + \epsilon}} - \frac{1}{m\sqrt{\sigma^2 + \epsilon}} \cdot \hat{x}_i \cdot \sum_{j=1}^{m} \frac{\partial J}{\partial \hat{x}_j} \hat{x}_j$$

$$= \frac{1}{m\sqrt{\sigma^2 + \epsilon}} \left( m \frac{\partial J}{\partial \hat{x}_i} - \sum_{i=1}^{m} \frac{\partial J}{\partial \hat{x}_i} - \hat{x}_i \sum_{j=1}^{m} \frac{\partial J}{\partial \hat{x}_j} \hat{x}_j \right)$$

## 2] Architecture

In Convolutional Neural Networks the bias parameter is associated with each filter and its value is added to the filter output. Such a bias is called a tied

bias where the bias remains constant for each location of a feature map, here the learnable parameters is lesser. Untied biases can also be used in which case each location on the input map can have its own bias, the number of learnable parameters drastically increases but the network would now be able to learn location specific information allowing for more fine tuning. Their use depends on the training data, if the data is shifted uniformly over all features a tied bias would work well, but if it has location specific shifts an untied bias can be used for correction.

If there is batch normalization operation after the convolution then biasing would be point less as the bias would also get normalized in the case of tied biases. The bias term thus becomes redundant. In case of untied biases the effect will remain as each neuron gets biased differently.

```python
1   # Implementation of a Three-layer neural network for MNIST classification
2   # with manual gradient calculation and manual optimization.
3
4   import torch
5   import torch.nn.functional as F
6   from torch.utils.data import DataLoader
7   from torchvision import datasets, transforms
8   import matplotlib.pyplot as plt
9
10  torch.manual_seed(0)
11  lr = 0.005
12  hidden_dim = 15
13  batch_size = 64
14  epochs = 20
15  plot = True
16
17  transform = transforms.Compose([transforms.ToTensor(), transforms.Normalize((0.1307,),
    (0.3081,))]) # mnist mean and std
18  train_dataset = datasets.MNIST(root='./data', train=True, download=True,
    transform=transform)
19  test_dataset = datasets.MNIST(root='./data', train=False, download=True,
    transform=transform)
20  train_loader = DataLoader(dataset=train_dataset, batch_size=batch_size, shuffle=True)
21  test_loader = DataLoader(dataset=test_dataset, batch_size=1000, shuffle=False)
22
23  class ThreeLayerNet:
24      def __init__(self, input_size, hidden_size1, hidden_size2, output_size):
25          self.W1 = torch.randn(input_size, hidden_size1) * 0.1
26          self.b1 = torch.randn(hidden_size1) * 0.1
27          self.W2 = torch.randn(hidden_size1, hidden_size2) * 0.1
28          self.b2 = torch.randn(hidden_size2) * 0.1
29          self.W3 = torch.randn(hidden_size2, output_size) * 0.1
30          self.b3 = torch.randn(output_size) * 0.1
31
32      def forward(self, x):
33          # Input
34          self.x = x
35          # Hidden Layer 1
36          self.z1 = x @ self.W1 + self.b1
37          self.a1 = F.relu(self.z1)
38          # Hidden Layer 2
39          self.z2 = self.a1 @ self.W2 + self.b2
40          self.a2 = F.relu(self.z2)
41          # Output Layer
42          self.z3 = self.a2 @ self.W3 + self.b3
43          return self.z3
44
45
46  model = ThreeLayerNet(784, hidden_dim, hidden_dim, 10)
47
48  train_loss_values = []
49  test_loss_values = []
50  train_accuracy_values = []
51  test_accuracy_values = []
52
53  def train(epoch):
54      for batch_idx, (data, target) in enumerate(train_loader):
```

```python
        data = data.view(-1, 784)  # flatten the input

        # === FORWARD PASS ===
        output = model.forward(data)
        log_softmax = F.log_softmax(output, dim=1)
        loss = - torch.mean(log_softmax[range(len(target)), target]) # Equivalent to
NLLLoss

        # === BACKWARD PASS ===
        # gradient of the loss w.r.t. output of model
        grad_z3 = F.softmax(output, dim=1)
        grad_z3[range(len(target)), target] -= 1
        grad_z3 /= len(target) # recall that loss is average over batch

        # gradient of the loss w.r.t. the output after the second hidden layer ReLU
        grad_a2 = grad_z3 @ model.W3.T

        # gradient of the loss w.r.t. the output before the second hidden layer ReLU
        grad_z2 = grad_a2.clone()
        grad_z2[model.z2 < 0] = 0

        # gradient of the loss w.r.t. the output after the first hidden layer ReLU
        grad_a1 = grad_z2 @ model.W2.T

        # gradient of the loss w.r.t. the output before the first hidden layer ReLU
        grad_z1 = grad_a1.clone()
        grad_z1[model.z1 < 0] = 0

        # gradient of the loss w.r.t. the model parameters
        model.W3.grad = model.a2.T @ grad_z3
        model.b3.grad = grad_z3.sum(axis=0)
        model.W2.grad = model.a1.T @ grad_z2
        model.b2.grad = grad_z2.sum(axis=0)
        model.W1.grad = model.x.T @ grad_z1
        model.b1.grad = grad_z1.sum(axis=0)

        # === PARAM UPDATES === # Vanilla gradient descent
        model.W3 -= lr * model.W3.grad
        model.b3 -= lr * model.b3.grad
        model.W2 -= lr * model.W2.grad
        model.b2 -= lr * model.b2.grad
        model.W1 -= lr * model.W1.grad
        model.b1 -= lr * model.b1.grad

        # === PRINT EVERY 200 ITERATIONS ===
        if batch_idx % 200 == 0:
            print('Train Epoch: {} [{}/{} ({:.0f}%)]\tLoss: {:.6f}'.format(
                epoch, batch_idx * len(data), len(train_loader.dataset),
                100. * batch_idx / len(train_loader), loss.item()))

@torch.no_grad()
def check(train_or_test, loader):
    loss, correct = 0, 0
    for data, target in loader:
        data = data.view(-1, 784) # flatten the input
        output = model.forward(data)
        log_softmax = F.log_softmax(output, dim=1)
        loss += F.nll_loss(log_softmax, target, reduction='sum').item() # sum up batch
loss
        pred = output.argmax(dim=1, keepdim=True) # get the index of the max log-
probability
```

```python
113        correct += pred.eq(target.view_as(pred)).sum().item() # pred is batch_size x 1,
    target is batch_size

115    loss /= len(loader.dataset) # note that reduction is 'sum' instead of 'mean' in
    F.nll_loss
116    accuracy = 100. * correct / len(loader.dataset)
117    print('{} set: Average loss: {:.4f}, Accuracy: {}/{} ({:.0f}%)'.format(
118        train_or_test, loss, correct, len(loader.dataset), accuracy))
119    return loss, accuracy

121 for epoch in range(1, 1+epochs):
122    train(epoch)
123    train_loss, train_acc = check(train_or_test='train', loader=train_loader)
124    train_loss_values.append(train_loss)
125    train_accuracy_values.append(train_acc)
126    test_loss, test_acc = check(train_or_test='test', loader=test_loader)
127    test_loss_values.append(test_loss)
128    test_accuracy_values.append(test_acc)
129    print('', end='\n')
130    # Notice how we are doing a full forward pass again at the end of each epoch to
    compute the train in loss and accuracy,
131    # but to save time, we could keep track of the loss (or number of correct predictions)
    for each mini-batch
132    # and take an average at the end of the epoch instead.

134 if plot:
135    plt.figure(figsize=(12, 5))
136    plt.subplot(1, 2, 1)
137    plt.plot(range(1, 1+epochs), train_loss_values, label='Training Loss')
138    plt.plot(range(1, 1+epochs), test_loss_values, label='Test Loss')
139    plt.xlabel('Epochs')
140    plt.ylabel('Loss')
141    plt.legend()

143    plt.subplot(1, 2, 2)
144    plt.plot(range(1, 1+epochs), train_accuracy_values, label='Training Accuracy')
145    plt.plot(range(1, 1+epochs), test_accuracy_values, label='Test Accuracy')
146    plt.xlabel('Epochs')
147    plt.ylabel('Accuracy (%)')
148    plt.legend()

150    plt.show()
```