

Advanced Computer Vision HW4

Ricky Yuan
rickyy@andrew.cmu.edu

November 7, 2024

1 Theory

1.1

Adding a constant c to the vector \mathbf{x} gives us:

$$\text{softmax}(x_i + c) = \frac{e^{x_i + c}}{\sum_j e^{x_j + c}} = \frac{e^c \cdot e^{x_i}}{\sum_j e^c \cdot e^{x_j}} = \frac{e^c \cdot e^{x_i}}{e^c \cdot \sum_j e^{x_j}} = \frac{e^{x_i}}{\sum_j e^{x_j}} = \text{softmax}(x_i)$$

for all x_i in \mathbf{x} holds, hence proved.

Since x_i could be large, making e^x exponentially large, resulting in a huge number divided by a huge number, it can lead to numerical instability. By setting $c = -\max x_i$, the largest term becomes 0, and the range of the exponential function becomes manageable.

1.2

The output of the softmax function lies within the range $(0, 1)$, and the sum over all the elements is 1. Therefore, we can view the output of the softmax function as a **probability distribution**. The exponential operation makes the input value positive and amplifies larger values than smaller values. The summation normalizes the values of the next step, and ensures that the final softmax output sums to 1. The last operation then calculates the corresponding probabilities.

1.3

Consider a multi-layer neural network with input vector \mathbf{x} and parameters W_i, b_i , where $i \in (1, 2, \dots, N)$, and N is the number of layer. The output y is

$$y = W_N W_{N-1} \dots W_1 \mathbf{x} + b_N + W_N b_{N-1} + W_N W_{N-1} b_{N-1} + \dots$$

We can combine and express the neural network as $W\mathbf{x} + b$, where

$$\begin{aligned} W &= W_N W_{N-1} \dots W_1 \\ b &= b_N + W_N b_{N-1} + W_N W_{N-1} b_{N-1} + \dots \end{aligned}$$

Therefore it is equivalent to linear regression.

1.4

$$\begin{aligned} \frac{\partial \sigma(x)}{\partial x} &= \frac{\partial}{\partial x} \frac{1}{1 + e^{-x}} = -(1 + e^{-x})^{-2} \cdot (-e^{-x}) = \frac{e^{-x}}{(1 + e^{-x})^2} \\ &= \frac{1}{1 + e^{-x}} \cdot \left(1 - \frac{1}{1 + e^{-x}}\right) = \sigma(x) \cdot (1 - \sigma(x)) \end{aligned}$$

1.5

$$y \in \mathbb{R}^{k \times 1} \quad W \in \mathbb{R}^{k \times d} \quad x \in \mathbb{R}^{d \times 1} \quad b \in \mathbb{R}^{k \times 1}$$

The loss J is a function of y , so we can write:

$$\frac{\partial J}{\partial W} = \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial W}$$

Since $y_i = \sum_{j=1}^d W_{ij} x_j + b_i$, then the partial derivative of y_i w.r.t W_{ij} is:

$$\frac{\partial y_i}{\partial W_{ij}} = x_j$$

We then have the matrix form $\frac{\partial y}{\partial W} = x^T$. Thus the answer is:

$$\frac{\partial J}{\partial W} = \delta x^T$$

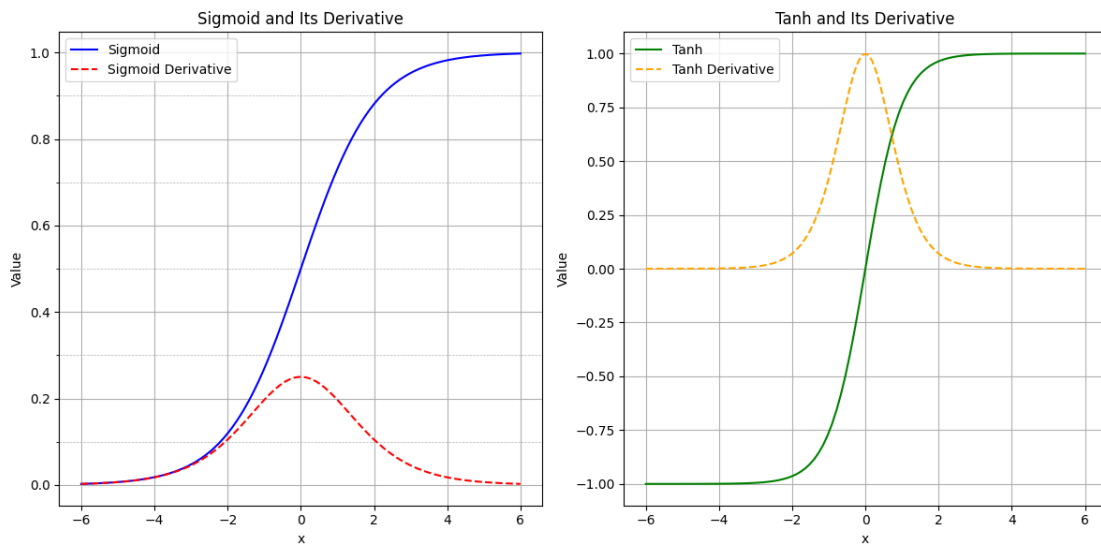
where $\delta \in \mathbb{R}^{k \times 1}$ and $x^T \in \mathbb{R}^{1 \times d}$. The result is a $k \times d$ matrix.

We use a similar technique to derive other partial derivatives.

$$\frac{\partial J}{\partial x} = \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial x} = W^T \delta \quad \frac{\partial J}{\partial b} = \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial b} = \delta$$

where $W^T \in \mathbb{R}^{d \times k}$ and $\delta \in \mathbb{R}^{k \times 1}$. The result of $\frac{\partial J}{\partial x}$ is a $d \times 1$ vector.

1.6



1. The value of the derivative of sigmoid ranges from 0 to 0.25. When we do backpropagation in a very deep network, multiple values that range from 0 to 0.25 propagate back to the input, causing the gradient to vanish.
2. The output ranges for sigmoid and tanh are $(0, 1)$, $(-1, 1)$, respectively. The mean of tanh is 0 so that the model would not learn bias compared to softmax.
3. The derivative of tanh ranges from $(0, 1)$, which has a larger domain compared to the derivative of softmax, making the model less likely to cause gradient vanishing.
4. The definition of tanh can be written as:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{1 - e^{-2x}}{1 + e^{-2x}} = \frac{2 - 1 - e^{-2x}}{1 + e^{-2x}} = \frac{2}{1 + e^{-2x}} - 1 = 2\sigma(2x) - 1$$

2 Implement a Fully Connected Network

2.1 Network Initialization

2.1.1

Initializing a network with all zeros would make the output of each neuron to be 0, causing the gradient to be all zeros. The network would then learn nothing since the gradient is zero.

2.1.2

```
1 def initialize_weights(in_size, out_size, params, name=""):
2     init_range = np.sqrt(6 / (in_size + out_size))
3     W = np.random.uniform(-init_range, init_range, (in_size, out_size))
4     b = np.zeros(out_size)
5
6     params["W" + name] = W
7     params["b" + name] = b
```

2.1.3

Initialization with random numbers helps different neurons better capture good and different features. If we initialize all the weights to be the same, effectively all the neurons would make no difference at all. Scaling the initialization properly makes the variance of the output reasonable, and it prevents the gradient vanishing or gradient exploding.

2.2 Forward Propagation

2.2.1

```
1 def sigmoid(x):
2     res = 1 / (1 + np.exp(-x))
3     return res
4
5 def forward(X, params, name="", activation=sigmoid):
6     # get the layer parameters
7     W = params["W" + name]
8     b = params["b" + name]
9
10    pre_act = X @ W + b
11    post_act = activation(pre_act)
12
13    # store the pre-activation and post-activation values
14    # these will be important in backprop
15    params["cache_" + name] = (X, pre_act, post_act)
16
17    return post_act
```

2.2.2

```
1 def softmax(x):
2     _x = x - np.max(x, axis=1, keepdims=True)
3     res = np.exp(_x) / np.sum(np.exp(_x), axis=1, keepdims=True)
4     return np.array(res)
```

2.2.3

```
1 def compute_loss_and_acc(y, probs):
2     loss = -np.sum(y * np.log(probs)) # prob = f(x), y is one-hot
3     acc = np.mean(np.argmax(y, axis=1) == np.argmax(probs, axis=1))
4     return loss, acc
```

2.3 Backwards Propagation

```
1 def backwards(delta, params, name="", activation_deriv=sigmoid_deriv):
2     """
3     Do a backwards pass
4
5     Keyword arguments:
6     delta -- errors to backprop
7     params -- a dictionary containing parameters
8     name -- name of the layer
9     activation_deriv -- the derivative of the activation_func
10    """
11    grad_X, grad_W, grad_b = None, None, None
12    # everything you may need for this layer
13    W = params["W" + name]
14    b = params["b" + name]
15    X, pre_act, post_act = params["cache_" + name]
16
17    # do the derivative through activation first
18    # (don't forget activation_deriv is a function of post_act)
19    # then compute the derivative W, b, and X
20    delta = delta * activation_deriv(post_act)
21    grad_W = X.T @ delta
22    grad_b = np.sum(delta, axis=0)
23    grad_X = delta @ W.T
24
25    # store the gradients
26    params["grad_W" + name] = grad_W
27    params["grad_b" + name] = grad_b
28    return grad_X
```

2.4 Training Loop

```
1 def get_random_batches(x, y, batch_size):
2     batches = []
3     num_samples = x.shape[0]
4     indices = np.random.permutation(num_samples)
5
6     for i in range(0, num_samples, batch_size):
7         batch_indices = indices[i : i + batch_size]
8         batches.append((x[batch_indices], y[batch_indices]))
9
10    return batches
```

```
1 batches = get_random_batches(x, y, 5)
2 batch_num = len(batches)
3
4 # WRITE A TRAINING LOOP HERE
5 max_iters = 500
6 learning_rate = 1e-3
7 # with default settings, you should get loss < 35 and accuracy > 75%
8 for itr in range(max_iters):
9     total_loss = 0
10    avg_acc = 0
11    for xb, yb in batches:
12        # forward
13        h1 = forward(xb, params, "layer1")
14        probs = forward(h1, params, "output", softmax)
15
16        # loss
17        loss, acc = compute_loss_and_acc(yb, probs)
18        total_loss += loss
19        avg_acc += acc
20
21        # backward
22        delta1 = probs - yb
23        delta2 = backwards(delta1, params, "output", linear_deriv)
24        backwards(delta2, params, "layer1", sigmoid_deriv) # No need to
25        store the final grad_X
26
27        # apply gradient
28        # gradients should be summed over batch samples
29        for k in params.keys():
30            if "grad" in k:
31                name = k.split("_")[1] # layer name
32                params[name] -= learning_rate * params[k]
33
34    avg_acc /= batch_num
35    if itr % 100 == 0:
36        print(
37            "itr: {:02d} \t loss: {:.2f} \t acc : {:.2f}".format(
38                itr, total_loss, avg_acc
39            )
40        )
```

2.5 Numerical Gradient Checker

```
1 # compute gradients using finite difference
2 eps = 1e-6
3 for k, v in params.items():
4     if "_" in k:
5         continue
6     # for each value inside the parameter
7     # add epsilon
8     # run the network
9     # get the loss
10    # subtract 2*epsilon
11    # run the network
12    # get the loss
13    # restore the original parameter value
14    # compute derivative with central diffs
15
16    # Create an iterator to go through each element of the parameter
17    it = np.nditer(v, flags=['multi_index'], op_flags=['readwrite'])
18    while not it.finished:
19        idx = it.multi_index # Current index, 2D for W, 1D for b
20
21        # Compute f(x + eps)
22        v[idx] += eps
23        h1 = forward(x, params, "layer1")
24        probs = forward(h1, params, "output", softmax)
25        loss1, _ = compute_loss_and_acc(y, probs)
26
27        # Compute f(x - eps)
28        v[idx] -= 2 * eps
29        h1 = forward(x, params, "layer1")
30        probs = forward(h1, params, "output", softmax)
31        loss2, _ = compute_loss_and_acc(y, probs)
32
33        v[idx] += eps # Restore original value
34
35        # Compute the numerical gradient
36        grad = (loss1 - loss2) / (2 * eps)
37        params["grad_" + k][idx] = grad
38
39        it.iternext() # Move to the next element
```

3 Training Models

3.1

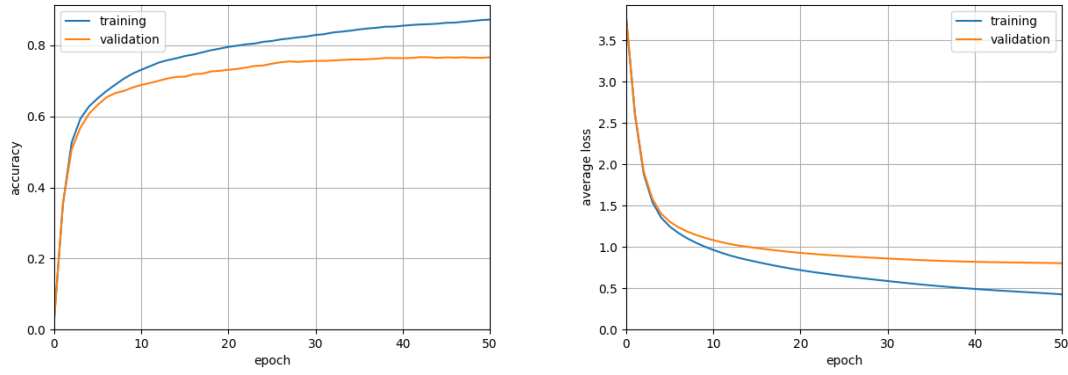


Figure 1: Accuracy & Loss with Batch = 64, Learning Rate = 0.005, Random Seed = 1487

3.2

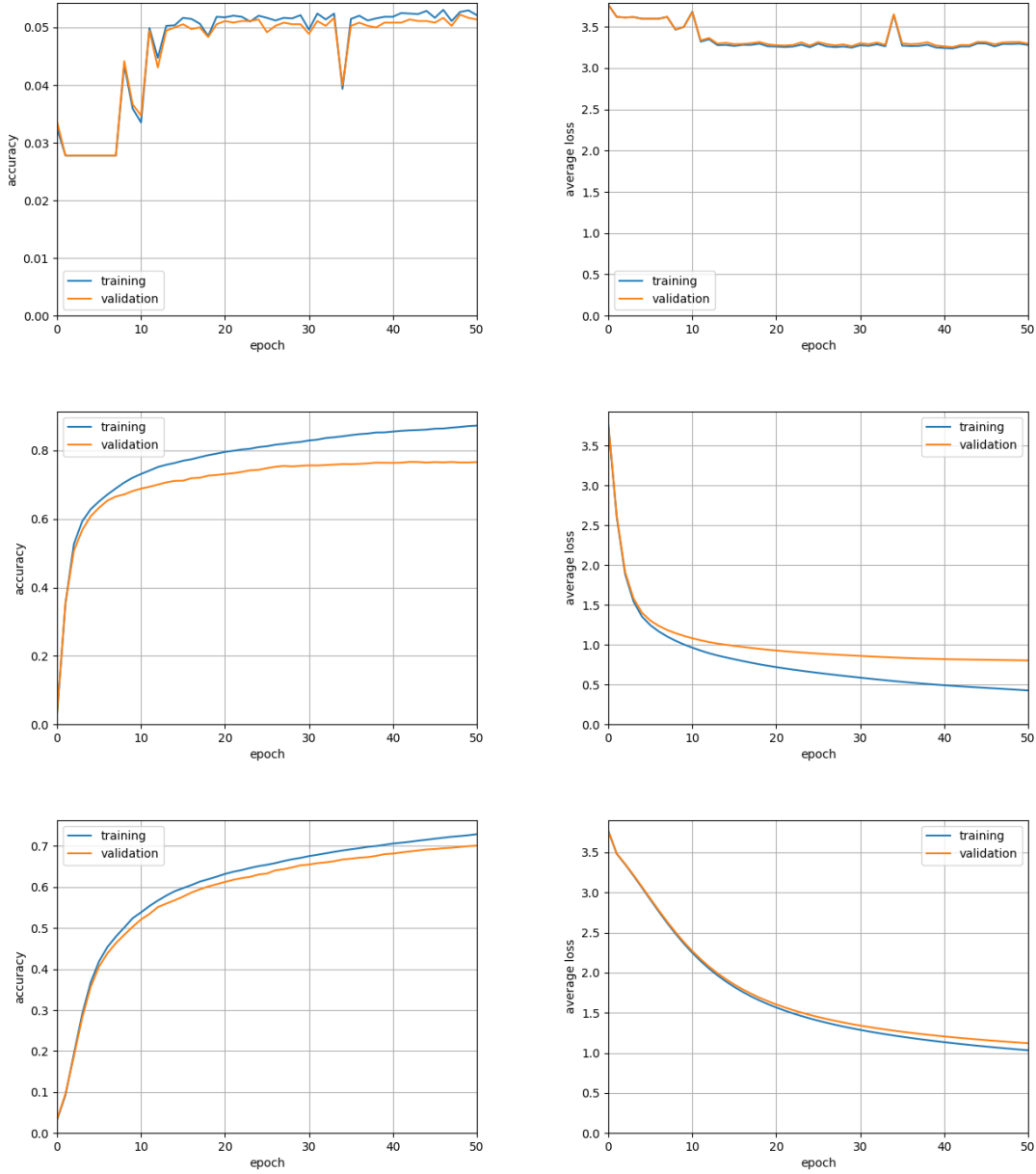


Figure 2: Accuracy & Loss with Learning Rates = $[0.05, 0.005, 0.0005]$, respectively

From the figures above, we see that the result for learning rate = 0.05 is bad because the update steps are too large, causing the loss oscillating and not able to converge to a minimum. On the other hand, the result for learning rate = 0.0005 have the validating accuracy around 0.71, which is still worse than the original experiment settings because the learning rate is too small. It needs more training iterations until fully convergence. The best network's accuracy reaches 0.766 for validation data and 0.773 for test data.

3.3

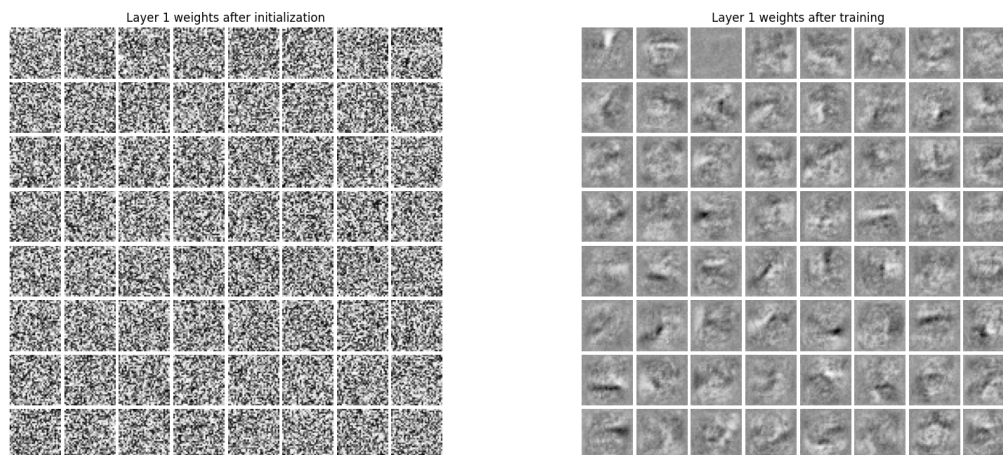


Figure 3: Visualization of the first layer before and after training

The first layer before training looks noisy and random, whereas the one after training looks there are patterns like strokes or lines.

3.4

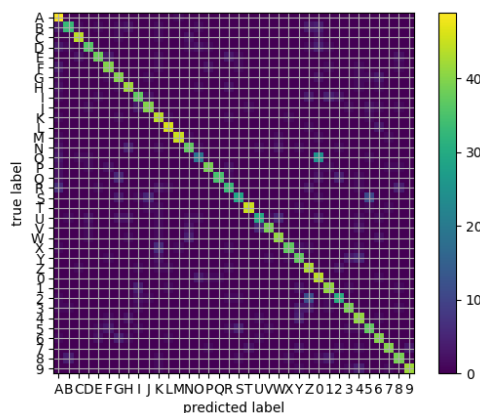


Figure 4: Confusion Matrix of the test data

From the figure, we can see that the character "O" and number "0" are the most easily to be confused by the model. Other pairs like "Z" and "2", "S" and "5" are easy to be confused as well.

4 Extract Text from Images

4.1

The sample method relies on several assumptions. One key assumption is that characters are not connected, allowing us to draw a bounding box around each character and feed them into our neural network one at a time. However, as shown in the example below, when characters are connected, it becomes difficult to separate them during the bounding box stage. This can lead to poor classification results.



Another important assumption is that the testing data is sampled from the same domain as the training data. The model is trained solely on English characters, and as a result, it cannot recognize characters from other languages. For example, the image below shows two Traditional Chinese characters meaning “hello.” Since these characters are not part of the model’s output classes, the model fails to classify them correctly.



4.2

```
1 def findLetters(image):
2     # insert processing in here
3     # one idea estimate noise -> denoise -> greyscale -> threshold ->
4     # morphology -> label -> skip small boxes
5     # this can be 10 to 15 lines of code using skimage functions
6
7     # Denoise the image using a denoising filter
8     # Ref: https://scikit-image.org/docs/stable/auto\_examples/filters/plot\_denoise.html
9     # https://medium.com/@betulmesci/image-processing-tutorial-using-scikit-image-noise-20e181c541a1
10    denoised = skimage.restoration.denoise_bilateral(image, channel_axis=-1)
11
12    # denoised = skimage.filters.gaussian(image, sigma=1)
13
14    # 2. Convert to grayscale
```

```

13 gray = skimage.color.rgb2gray(denoised)
14
15 # 3. Apply a threshold to get a binary (black-and-white) image
16 threshold = skimage.filters.threshold_otsu(gray)
17 bw = gray < threshold # Invert the binary image if needed
18
19 # 4. Apply morphological operations to clean up small noise
20 bw = skimage.morphology.remove_small_objects(bw, min_size=100)
21 bw = skimage.morphology.closing(bw, skimage.morphology.square(6))
22
23 # 5. Label connected components, 8-connectivity
24 labeled = skimage.measure.label(bw, connectivity=2)
25
26 # 6. Find regions and filter small components (bounding boxes)
27 regions = skimage.measure.regionprops(labeled)
28
29 # 7. Filter regions based on area
30 areas = [region.area for region in regions]
31 mean, std = np.mean(areas), np.std(areas)
32 # Drop all regions with area less than mean - 2.5 * std
33 bboxes = [region.bbox for region in regions if region.area > mean -
34            2.5 * std]
35
36 # Make the stroke of the letters thicker
37 bw = skimage.morphology.dilation(bw, footprint=skimage.morphology.
38                                   square(5))
39 return bboxes, bw.astype(np.float64)

```

4.3

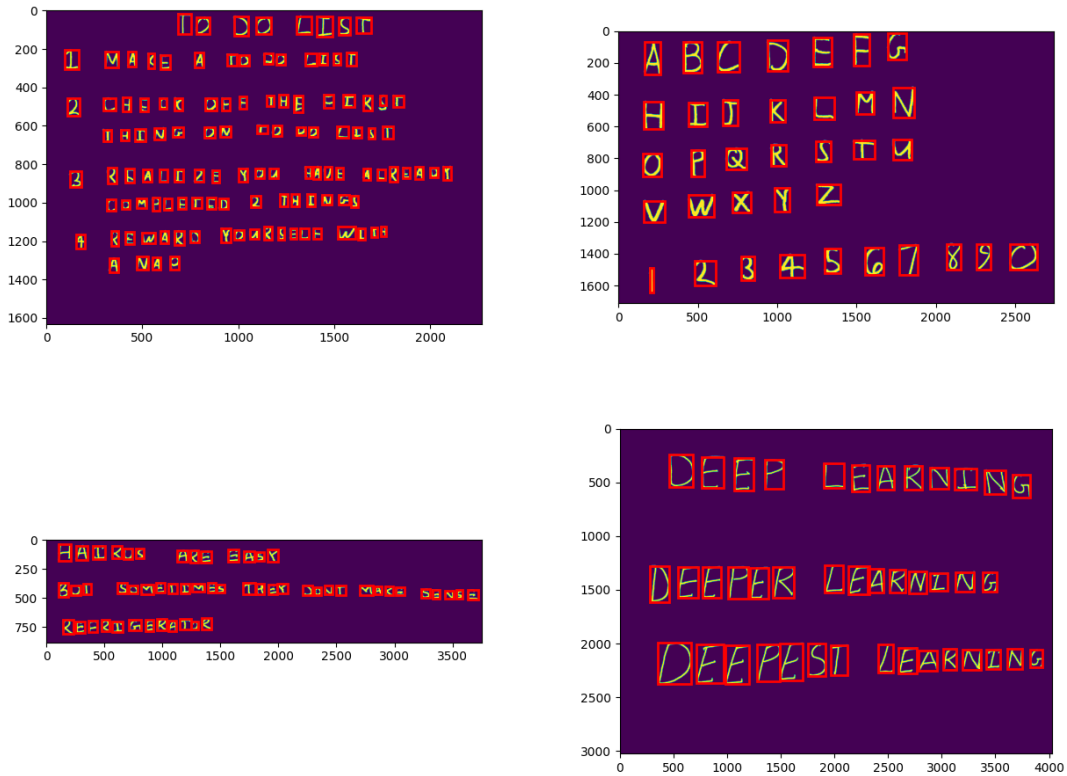


Figure 5: Result of `findLetters()`

4.4

list

TODOLIST
IMAXEATDDULIST
2LHF5KDFF7HEFIRFT
7HINGQNTQWQLIST
3RIALIZEYQUHAVEALR8ADT
CDMPLFTIDZTHINGS
AREWARDYOURSELFWITH
ANAP

haiku

HAIKUSAREHAS
BUTSQMETIMESTHEYDDNTMAK2SGNGE
RBFRIGERATQR

letters

ABCDEFGH
IIJKLMN
QPRSTU
VWXYZ
1234567890

deep

DEEFLEARNIHG
DE8PERLEARAING
D55PE5TLEARNIHG

The total accuracy is 193/246 (78.46%).

5 (Extra Credit) Image Compression with Autoencoders

5.1 Building the Autoencoder

5.1.1

```
1 # Initialize the layers for the autoencoder
2 input_size = train_x.shape[1] # 1024 for the NIST36 dataset
3 initialize_weights(input_size, hidden_size, params, "layer1")
4 initialize_weights(hidden_size, hidden_size, params, "layer2")
5 initialize_weights(hidden_size, hidden_size, params, "layer3")
6 initialize_weights(hidden_size, input_size, params, "output")
7 # Initialize the momentum terms
8 for k in ["layer1", "layer2", "layer3", "output"]:
9     params["m_W" + k] = np.zeros_like(params["W" + k])
10    params["m_b" + k] = np.zeros_like(params["b" + k])
11
12 # should look like your previous training loops
13 losses = []
14 for itr in range(max_iters):
15     total_loss = 0
16     for xb, _ in batches:
17         # forward pass
18         h1 = forward(xb, params, "layer1", relu)
19         h2 = forward(h1, params, "layer2", relu)
20         h3 = forward(h2, params, "layer3", relu)
21         reconstructed_x = forward(h3, params, "output", sigmoid)
22
23         # loss
24         loss = np.sum((xb - reconstructed_x) ** 2) # total squared error
25         total_loss += loss
26
27         # backward
28         delta = 2 * (reconstructed_x - xb)
29         delta = backwards(delta, params, "output", sigmoid_deriv)
30         delta = backwards(delta, params, "layer3", relu_deriv)
31         delta = backwards(delta, params, "layer2", relu_deriv)
32         backwards(delta, params, "layer1", relu_deriv)
33
34         # apply gradient, remember to update momentum as well
35         for k in params.keys():
36             if "grad" in k:
37                 layer_name = k.split('_')[1] # Layer name, e.g. Wlayer1
38                 m_key = "m_" + layer_name # Momentum key, e.g. m_Wlayer1
39                 params[m_key] = 0.9 * params[m_key] - learning_rate *
40                     params[k]
41                 params[layer_name] += params[m_key]
42
43     losses.append(total_loss/train_x.shape[0])
44     if itr % 2 == 0:
45         print("itr: {:02d} \t loss: {:.2f}".format(itr, total_loss))
46     if itr % lr_rate == lr_rate-1:
47         learning_rate *= 0.9
```

5.1.2

See the code above.

5.2 Training the Autoencoder

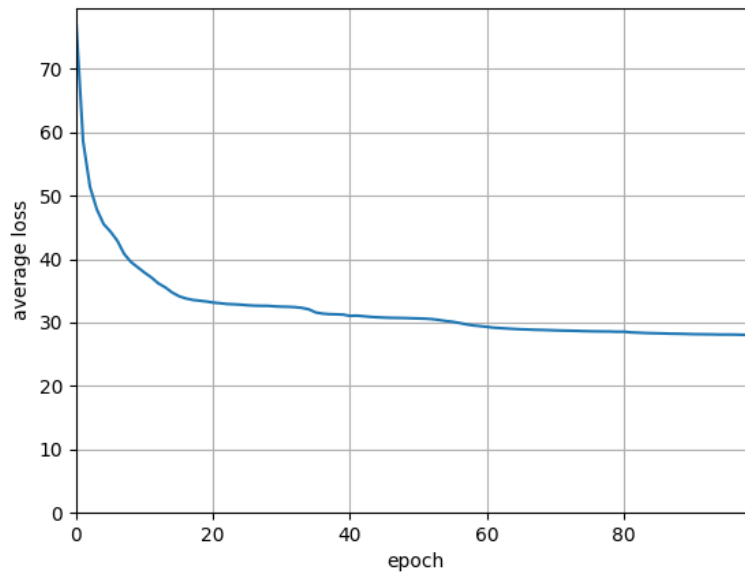


Figure 6: Training loss curve for the autoencoder

From the figure above, we observe that the loss decreases rapidly during early training. After that, the loss converges to value around 30.

5.3 Evaluating the Autoencoder

5.3.1



Figure 7: Reconstruction result for character R, I, C, K, Y

From the figure above, we observe that the reconstructed characters retain general shape but are more blurry and noisy. Some characters, such as "R" or "K" suffer from noticeable distortion and lack of details. This suggests that the autoencoder captures basic features but struggles with finer details due to the limited latent space.

5.3.2

With default hyper-parameters, the average PSNR in my experiment is 16.00.

6 PyTorch

6.1 Train a neural network in PyTorch

6.1.1

Hyper-parameters:

```
1 batch_size = 64
2 learning_rate = 0.005
3 num_epochs = 100
4 hidden_size = 64
5 momentum = 0.9
```

Network Structure:

```
1 class FullyConnectedNetwork(nn.Module):
2     def __init__(self):
3         super(FullyConnectedNetwork, self).__init__()
4         self.hidden_size = hidden_size
5         self.model = nn.Sequential(
6             nn.Linear(1024, self.hidden_size),
7             nn.Sigmoid(),
8             nn.Linear(self.hidden_size, 36),
9         )
10
11     def forward(self, x):
12         return self.model(x)
```

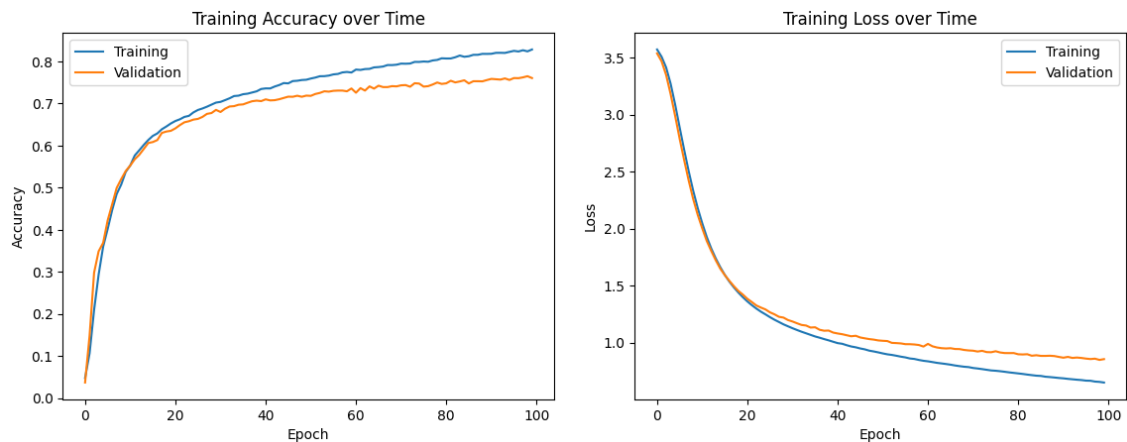


Figure 8: Accuracy and Loss of Training and Validation for fully-connected network

The network's accuracy reaches 0.7642 for validation data and 0.7656 for test data.

6.1.2

The hyper-parameters are the same as the previous setting, except here I used `num_epoch=30` since it converges much faster.

Network Structure:

```
1 class ConvNetwork(nn.Module):
2     def __init__(self):
3         super(ConvNetwork, self).__init__()
4         self.model = nn.Sequential(
5             nn.Conv2d(1, 4, kernel_size=3, stride=1, padding=1),
6             nn.ReLU(),
7             nn.MaxPool2d(kernel_size=2, stride=2, padding=0),
8             nn.Conv2d(4, 16, kernel_size=3, stride=1, padding=1),
9             nn.ReLU(),
10            nn.MaxPool2d(kernel_size=2, stride=2, padding=0),
11            nn.Flatten(),
12            nn.Linear(16 * 8 * 8, 256),
13            nn.ReLU(),
14            nn.Linear(256, 36)
15        )
16
17    def forward(self, x):
18        x = x.reshape(-1, 1, 32, 32)
19        return self.model(x)
```

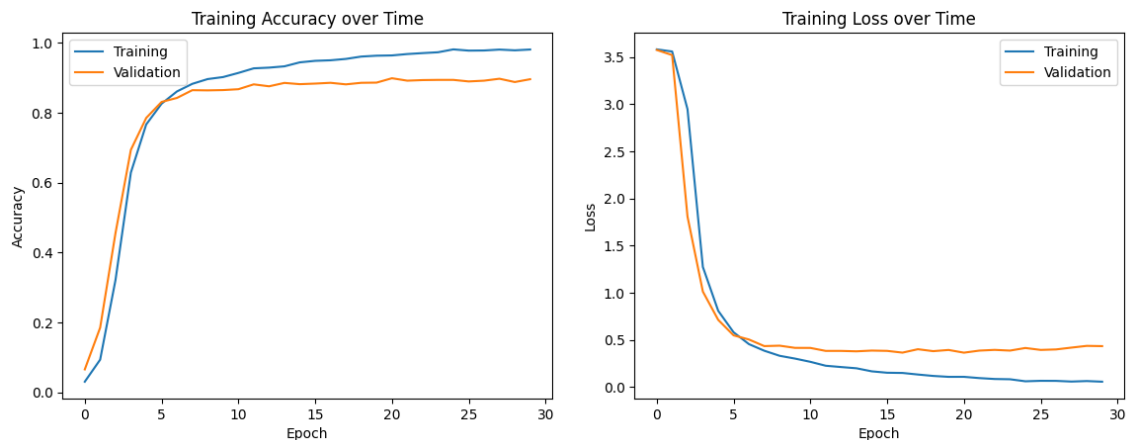


Figure 9: Accuracy and Loss of Training and Validation for convolution neural network

The network's accuracy reaches 0.8864 for validation data and 0.8917 for test data. Compared to the fully connected network, CNN converges much faster and has better performance.

6.1.3

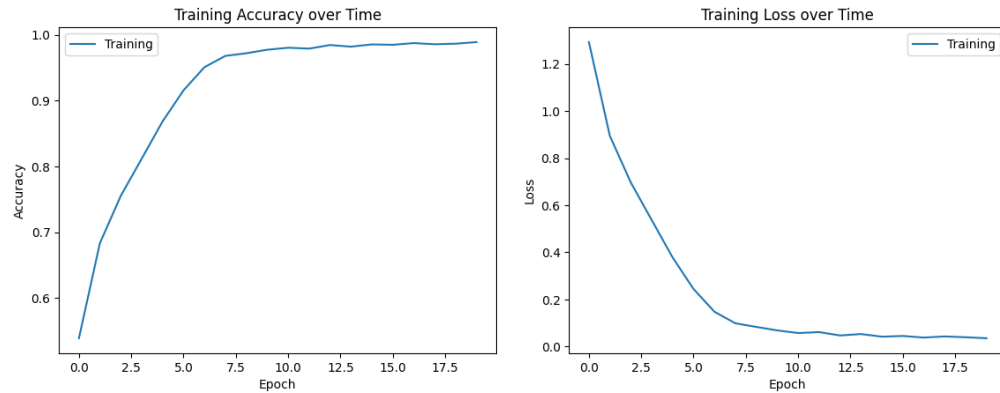


Figure 10: Accuracy and Loss of Training for CIFAR10 dataset

The network reaches around 0.72 accuracy over the test data.

6.2 Fine Tuning

I chose to finetune on the flowers 17 dataset, and trained only 10 epochs. Below are my results.

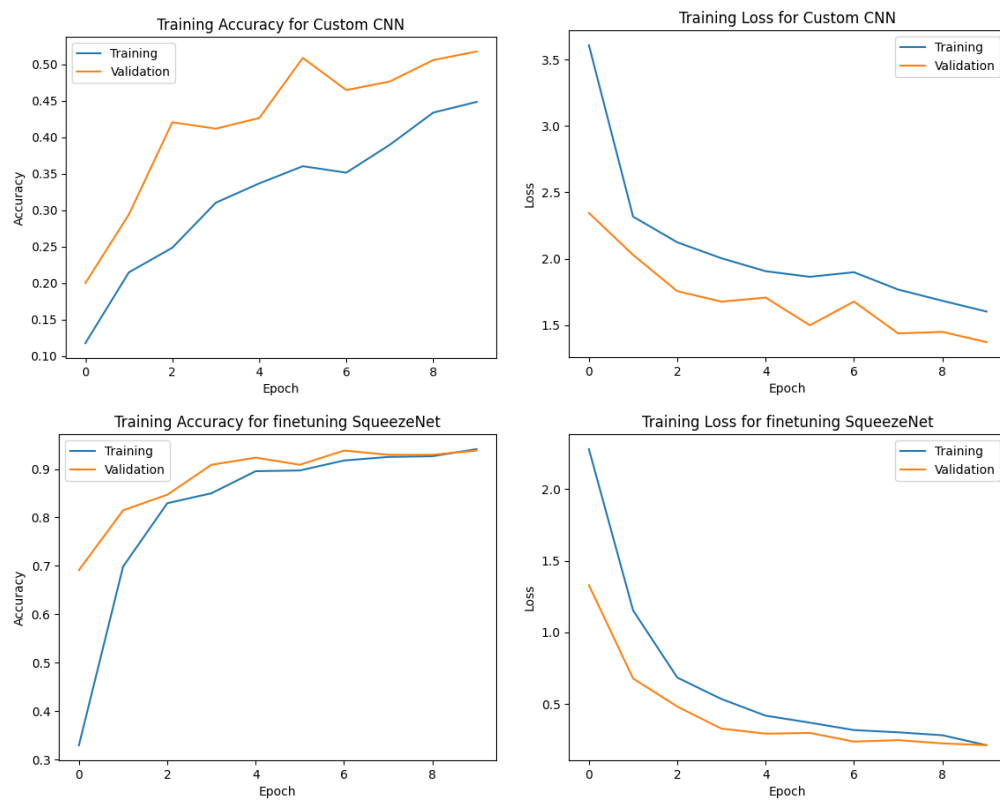


Figure 11: Accuracy and Loss of finetuning and training the whole model

As we can see, finetuning converges much faster and has a better performance. The performance of finetuning SqueezeNet for 10 epoch reaches 0.94 accuracy on validation set whereas training the whole model only reaches 0.52 accuracy. Below are my code and model structures:

```

1 model_squeezenet = models.squeezenet1_1(weights=models.
    SqueezeNet1_1_Weights.DEFAULT)
2
3 num_classes = len(class_names)
4 # Replace the classifier part of the model
5 model_squeezenet.classifier[1] = nn.Conv2d(512, num_classes, kernel_size
    =(1, 1))
6 model_squeezenet.num_classes = num_classes
7
8 # Freeze all layers in the feature extractor
9 for param in model_squeezenet.features.parameters():
10     param.requires_grad = False

```

And the custom CNN model structure:

```

1 class CustomCNN(nn.Module):
2     def __init__(self, num_classes):
3         super(CustomCNN, self).__init__()
4         self.features = nn.Sequential(
5             nn.Conv2d(3, 64, kernel_size=3, stride=1, padding=1),
6             nn.ReLU(),
7             nn.MaxPool2d(kernel_size=2, stride=2),
8             nn.Conv2d(64, 128, kernel_size=3, stride=1, padding=1),
9             nn.ReLU(),
10            nn.MaxPool2d(kernel_size=2, stride=2),
11            nn.Conv2d(128, 256, kernel_size=3, stride=1, padding=1),
12            nn.ReLU(),
13            nn.MaxPool2d(kernel_size=2, stride=2)
14        )
15        self.classifier = nn.Sequential(
16            nn.Linear(256 * 28 * 28, 512),
17            nn.ReLU(),
18            nn.Dropout(0.5),
19            nn.Linear(512, num_classes)
20        )
21
22    def forward(self, x):
23        x = self.features(x)
24        x = x.view(x.size(0), -1) # Flatten the features
25        x = self.classifier(x)
26        return x

```

Training logics:

```
1 def train_model(model, criterion, optimizer, num_epochs=10):
2     train_losses, valid_losses = [], []
3     train_accs, valid_accs = [], []
4     for epoch in range(num_epochs):
5         model.train()
6         train_acc, train_loss = 0, 0
7         for inputs, labels in dataloaders['train']:
8             optimizer.zero_grad()
9             outputs = model(inputs)
10            loss = criterion(outputs, labels)
11
12            loss.backward()
13            optimizer.step()
14
15            train_loss += loss.item()
16            _, predicted = torch.max(outputs, 1)
17            train_acc += (predicted == labels).sum().item()
18
19        train_acc /= dataset_sizes['train']
20        train_loss /= len(dataloaders['train'])
21        train_losses.append(train_loss)
22        train_accs.append(train_acc)
23
24        # Validation loop
25        model.eval()
26        valid_acc, valid_loss = 0, 0
27        with torch.no_grad():
28            for inputs, labels in dataloaders['val']:
29                outputs = model(inputs)
30                _, predicted = torch.max(outputs, 1)
31                valid_acc += (predicted == labels).sum().item()
32                valid_loss += criterion(outputs, labels).item()
33
34        valid_acc /= dataset_sizes['val']
35        valid_loss /= len(dataloaders['val'])
36        valid_losses.append(valid_loss)
37        valid_accs.append(valid_acc)
38
39        print(f"Epoch {epoch+1}, Train Loss: {train_loss:.4f}, Train Acc:
40              {train_acc:.4f}, Valid Loss: {valid_loss:.4f}, Valid Acc: {
41              valid_acc:.4f}")
42
43    return model, train_losses, valid_losses, train_accs, valid_accs
```

6.3 Neural Networks in the Real World

I chose class index 291 (Lion) for my experiment and used a video from National Geographic on YouTube (link). On the validation set, the pretrained model achieves perfect accuracy (1.00), but for the video, the accuracy drops to 0.9595. This makes sense since the video comes from a slightly different domain. To improve performance, we could apply techniques like domain adaptation or data augmentation. Below is my code for downloading and processing the video sequence. Note that I trimmed the video to a few seconds before feeding to my model.

```
1 import torch
2 import torchvision.models as models
3 from torchvision.models.resnet import ResNet50_Weights
4 import torchvision.transforms as transforms
5 from PIL import Image
6 import numpy as np
7 from tqdm import tqdm
8 from matplotlib import pyplot as plt
9
10 model = models.resnet50(weights=ResNet50_Weights.IMAGENET1K_V2)
11 model.eval()
12 imagenet_transform = transforms.Compose([
13     transforms.Resize(256),
14     transforms.CenterCrop(224),
15     transforms.ToTensor(),
16     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224,
17         0.225]),
18 ])
19
20 from torchvision.datasets import ImageNet
21
22 val_dir = '../data/'
23 imagenet_data = ImageNet(val_dir, split='val', transform=
24     imagenet_transform)
25 val_loader = torch.utils.data.DataLoader(imagenet_data, batch_size=32,
26     shuffle=False)
27 class_idx = 291
28 specific_class_data = [(img, label) for img, label in imagenet_data if
29     label == class_idx]
30 # Evaluate model on specific class
31 correct = 0
32 for inputs, labels in tqdm(specific_class_data[:]):
33     inputs = inputs.unsqueeze(0) # Add batch dimension
34     outputs = model(inputs)
35     _, preds = torch.max(outputs, 1)
36     correct += (preds == labels).sum().item()
37
38 print(f"Validation accuracy for class {class_idx}: {correct / 50:.4f}")
39
40 from pytubefix import YouTube
41
42 yt = YouTube('https://www.youtube.com/watch?v=OMkEVX23BdM')
43 stream = yt.streams.get_highest_resolution()
44 stream.download(output_path='../data/', filename='lion.mp4')
```

```

6 print('Download complete!')
7
8 import cv2
9 import os
10
11 video_path = '../data/lion_trimed.mp4'
12 output_dir = '../data/lion_frames'
13 os.makedirs(output_dir, exist_ok=True)
14
15 # Extract frames from the video
16 cap = cv2.VideoCapture(video_path)
17 frame_count = 0
18
19 while cap.isOpened():
20     ret, frame = cap.read()
21     if not ret:
22         break
23     frame_path = os.path.join(output_dir, f"frame_{frame_count:04d}.jpg")
24     cv2.imwrite(frame_path, frame)
25     frame_count += 1
26
27 cap.release()
28 print(f"Extracted {frame_count} frames.")
29 frame_paths = sorted([os.path.join(output_dir, fname) for fname in os.
30     listdir(output_dir)])
31 frame_preds = []
32
33 for frame_path in frame_paths:
34     frame = Image.open(frame_path).convert('RGB')
35     input_tensor = imagenet_transform(frame).unsqueeze(0) # Add batch
36     dimension
37     outputs = model(input_tensor)
38     _, preds = torch.max(outputs, 1)
39     frame_preds.append(preds.item())
40
41 correct_frames = sum([1 for pred in frame_preds if pred == class_idx])
42 print(f"Accuracy on video frames: {correct_frames / len(frame_preds):.4f}")
43 )

```

7 Collaboration

For this homework, I discussed with Daniel Yang (danielya).