

16720-B F20 Neural Networks for Recognition

Instructor: Kris Kitani

TAs: Anand B, Alireza G, Nadine C, Rawal K, Zhengyi L

Due October 19, 2018 11:59 PM

Total Points: 205 + 40 Extra Credit

Instructions

1. **Integrity and collaboration:** Students are encouraged to work in groups but each student must submit their own work. Include the names of your collaborators in your write up. Code should **NOT** be shared or copied. Please **DO NOT** use external code unless permitted. Plagiarism is prohibited and may lead to failure of this course.
2. **Start early!** This one is lengthy, training a network takes time. If you start late, you will be pressed for time while debugging. Also, speed read the entire assignment before you start so you get the gist of the assignment.
3. **Questions:** If you have any question, please look at piazza first.
4. **Submission:** The submission is on Gradescope, **you will be submitting both your writeup and code zip file**. The zip file, <andrew-id.zip> contains your code and any results files we ask you to save. **Note: You have to submit your writeup separately to Gradescope, and include results in the write-up.**
5. **Do not** submit anything from the `data/` folder in your submission.
6. For your code submission, **do not** use any libraries other than *numpy*, *scipy*, *scikit-image*, *matplotlib* and (in the appropriate section) *pytorch*. Including other libraries (for example, *cv2*, *ipdb*, etc.) **may lead to loss of credit** on the assignment.
7. To get the data, we have included some scripts in `scripts`.
8. Assignments that do not follow this submission rule will be **penalized upto 10% of the total score**.
9. Please make sure that the file paths that you use are relative and not absolute.

Contents

1	Theory (45 points)	2
2	Implement a Fully Connected Network (65 points)	4
2.1	Network Initialization	4
2.2	Forward Propagation	4
2.3	Backwards Propagation	4

2.4	Training Loop	5
2.5	Numerical Gradient Checker	5
3	Training Models (20 points)	6
4	Extract Text from Images (35 points)	8
5	Image Compression with Autoencoders (25 points Extra Credit)	10
5.1	Building the Autoencoder	10
5.2	Training the Autoencoder	10
5.3	Evaluating the Autoencoder	10
6	Comparing against PCA (15 points Extra Credit)	11
7	PyTorch (40 points)	12
7.1	Train a neural network in PyTorch	12
7.2	Fine Tuning	13
8	Appendix: Neural Network Overview	14
8.1	Basic Use	14
8.2	Backprop	15

1 Theory (45 points)

Q1.1 Theory [4 points] Prove that softmax is invariant to translation, that is

$$\text{softmax}(x) = \text{softmax}(x + c) \quad \forall c \in \mathbb{R}$$

Softmax is defined as below, for each index i in a vector x .

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

Often we use $c = -\max x_i$. Why is that a good idea? (Tip: consider the range of values that numerator will have with $c = 0$ and $c = -\max x_i$)

Q1.2 Theory [3 points] Softmax can be written as a three step processes, with $s_i = e^{x_i}$, $S = \sum s_i$ and $\text{softmax}(x_i) = \frac{1}{S}s_i$.

- As $x \in \mathbb{R}^d$, what are the properties of $\text{softmax}(x)$, namely what is the range of each element? What is the sum over all elements?
- One could say that "*softmax takes an arbitrary real valued vector x and turns it into a _____*".
- Can you see the role of each step in the multi-step process now? Explain them.

Q1.3 Theory [3 points] Show that multi-layer neural networks without a non-linear activation function are equivalent to linear regression.

Q1.4 Theory [4 points] Given the sigmoid activation function $\sigma(x) = \frac{1}{1+e^{-x}}$, derive the gradient of the sigmoid function and show that it can be written as a function of $\sigma(x)$ (without having access to x directly)

Q1.5 Theory [12 points] Given $y = Wx + b$ (or $y_j = \sum_{i=1}^d x_i W_{ij} + b_j$), and the gradient of some loss J with respect to y , show how to get $\frac{\partial J}{\partial W}$, $\frac{\partial J}{\partial x}$ and $\frac{\partial J}{\partial b}$. Be sure to do the derivatives with scalars and re-form the matrix form afterwards. Here are some notional suggestions.

$$\frac{\partial J}{\partial y} = \delta \in \mathbb{R}^{k \times 1} \quad W \in \mathbb{R}^{k \times d} \quad x \in \mathbb{R}^{d \times 1} \quad b \in \mathbb{R}^{k \times 1}$$

Q1.6 Theory [15 points] We will find the derivatives for Conv layers now. Since most Deep Learning frameworks such as Pytorch, Tensorflow use cross-correlation in their respective "convolution" functions ([here](#) and [here](#)), we will continue this abuse of notation. So the operation performed with the Conv Layer weights will be cross-correlation.

The input, x is of shape $M \times N$ with C channels. This will be *convolved* (actually cross-correlation) with D number of $K \times K$ filters, each with a bias term. The stride is 1 and there will be no padding. We know the gradient of some loss J with respect to the output y , which will have D channels. Show how to get $\frac{\partial J}{\partial W}$, $\frac{\partial J}{\partial x}$ and $\frac{\partial J}{\partial b}$.

The dimensions and notation are as follows:

$$\frac{\partial J}{\partial y} = \delta \in \mathbb{R}^{D \times M_o \times N_o} \quad M_o = M - K + 1 \quad N_o = N - K + 1$$

$$x \in \mathbb{R}^{C \times M \times N} \quad W \in \mathbb{R}^{D \times C \times K \times K} \quad b \in \mathbb{R}^D$$

$x_{c,i,j}$: The element at the i^{th} row, the j^{th} column and the c^{th} channel of the input

$y_{c,i,j}$: The element at the i^{th} row, the j^{th} column and the c^{th} channel of the output

$W_{d,c,i,j}$: The element at the i^{th} row, the j^{th} column, the c^{th} channel of the kernel of the d^{th} filter

For this question, you may compute the derivatives with scalars only. You don't need to re-form the matrix

Q1.7 Theory [4 points] When the neural network applies the elementwise activation function (such as sigmoid), the gradient of the activation function scales the back-propagation update. This is directly from the chain rule, $\frac{d}{dx} f(g(x)) = f'(g(x))g'(x)$.

1. Consider the sigmoid activation function for deep neural networks. Why might it lead to a "vanishing gradient" problem if it is used for many layers (consider plotting Q1.3)?
2. Often it is replaced with $\tanh(x) = \frac{1-e^{-2x}}{1+e^{-2x}}$. What are the output ranges of both \tanh and sigmoid? Why might we prefer \tanh ?
3. Why does $\tanh(x)$ have less of a vanishing gradient problem? (plotting the derivatives helps! for reference: $\tanh'(x) = 1 - \tanh(x)^2$)
4. \tanh is a scaled and shifted version of the sigmoid. Show how $\tanh(x)$ can be written in terms of $\sigma(x)$. (*Hint: consider how to make it have the same range*)

2 Implement a Fully Connected Network (65 points)

All of these functions should be implemented in `python/nm.py`

2.1 Network Initialization

Q2.1.1 Theory [3 points] Why is it not a good idea to initialize a network with all zeros? If you imagine that every layer has weights and biases, what can a zero-initialized network output after training?

Q2.1.2 Code [3 points] Implement a function to initialize neural network weights with Xavier initialization [1], where $Var[w] = \frac{2}{n_{in} + n_{out}}$ where n is the dimensionality of the vectors and you use a **uniform distribution** to sample random numbers (see eq 16 in [1]).

Q2.1.3 Theory [2 points] Why do we initialize with random numbers? Why do we scale the initialization depending on layer size (see near Fig 6 in [1])?

2.2 Forward Propagation

The appendix (sec 8) has the math for forward propagation, we will implement it here.

Q2.2.1 Code [12 points] Implement sigmoid, along with forward propagation for a single layer with an activation function, namely $y = \sigma(XW + b)$, returning the output and intermediate results for an $N \times D$ dimension input X , with examples along the rows, data dimensions along the columns.

Q2.2.2 Code [5 points] Implement the softmax function. Be sure the use the numerical stability trick you derived in Q1.1 softmax.

Q2.2.3 Code [5 points] Write a function to compute the accuracy of a set of labels, along with the scalar loss across the data. The loss function generally used for classification is the cross-entropy loss.

$$L_f(\mathbf{D}) = - \sum_{(\mathbf{x}, \mathbf{y}) \in \mathbf{D}} \mathbf{y} \cdot \log(\mathbf{f}(\mathbf{x}))$$

Here \mathbf{D} is the full training dataset of data samples \mathbf{x} ($N \times 1$ vectors, N = dimensionality of data) and labels \mathbf{y} ($C \times 1$ one-hot vectors, C = number of classes).

2.3 Backwards Propagation

Q2.3.1 Code [10 points] Compute back-propagation for a single layer, given the original weights, the appropriate intermediate results, and given gradient with respect to the loss. You should return the gradient with respect to X so you can feed it into the next layer. As a size check, your gradients should be the same dimensions as the original objects.

2.4 Training Loop

You will tend to see gradient descent in three forms: “normal”, “stochastic” and “batch”. “Normal” gradient descent aggregates the updates for the entire dataset before changing the weights. Stochastic gradient descent applies updates after every single data example. Batch gradient descent is a compromise, where random subsets of the full dataset are evaluated before applying the gradient update.

Q2.4.1 Code [10 points] Write a training loop that generates random batches, iterates over them for many iterations, does forward and backward propagation, and applies a gradient update step.

2.5 Numerical Gradient Checker

Q2.5.1 [15 points] Implement a numerical gradient checker. Instead of using the analytical gradients computed from the chain rule, add ϵ offset to each element in the weights, and compute the numerical gradient of the loss with central differences. Central differences is just $\frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$. Remember, this needs to be done for each scalar dimension in all of your weights independently.

3 Training Models (20 points)

First, be sure to run the script, from inside the scripts folder, `get_data.sh`. This will use `wget` and `unzip` to download

http://www.cs.cmu.edu/~lkeselma/16720a_data/data.zip

http://www.cs.cmu.edu/~lkeselma/16720a_data/images.zip

and extract them to `data` and `image` folders

Since our input images are 32×32 images, unrolled into one 1024 dimensional vector, that gets multiplied by $\mathbf{W}^{(1)}$, each row of $\mathbf{W}^{(1)}$ can be seen as a weight image. Reshaping each row into a 32×32 image can give us an idea of what types of images each unit in the hidden layer has a high response to.

We have provided you three data `.mat` files to use for this section. The training data in `nist36_train.mat` contains samples for each of the 26 upper-case letters of the alphabet and the 10 digits. This is the set you should use for training your network. The cross-validation set in `nist36_valid.mat` contains samples from each class, and should be used in the training loop to see how the network is performing on data that it is not training on. This will help to spot over fitting. Finally, the test data in `nist36_test.mat` contains testing data, and should be used for the final evaluation on your best model to see how well it will generalize to new unseen data.

Q3.1.1 Code [8 points] Train a network from scratch. Use a single hidden layer with 64 hidden units, and train for at least 30 epochs. **Modify** the script to generate two plots: one showing the **accuracy on both the training and validation** set over the epochs, and the other showing the **cross-entropy loss averaged** over the data. The x-axis should represent the epoch number, while the y-axis represents the accuracy or loss. With these settings, you should see an accuracy on the validation set of at least 75%.

Q3.1.2 Writeup [3 points] Use your modified training script to train three networks, one with your best learning rate, one with 10 times that learning rate and one with one tenth that learning rate. Include all 6 plots in your writeup. Comment on how the learning rates affect the training, and report the final accuracy of the best network on the test set.

Q3.1.3 Writeup [2 points] Visualize the first layer weights that your network learned. Compare these to the network weights immediately after initialization. Include both visualizations in your writeup. Comment on the learned weights. Do you notice any patterns?

Q3.1.4 Writeup [3 points] To visualize the deeper second layer we use the following method:

Given y_1 (size 64) as the output of the first layer and y_2 (size 36) as the output of the second layer. We will take a **one-hot encoding** for any number i in range $[0, 36)$ and multiply it with the transposed weight matrix of layer 2, W_2^T . Similarly multiply the output with W_1^T to get back an array of the size of the input image (1024). Reshape and visualize the result. Add images for 3 different values of i in your write-up. Also mention how they are different from the visualization in the previous part.

For these values of i , additionally include the image in the validation set for which $y_2(i)$ attains it's maximum value.

Q3.1.5 Writeup [4 points] Visualize the confusion matrix for your best model for the test set. Comment on the top few pairs of classes that are most commonly confused.

4 Extract Text from Images (35 points)

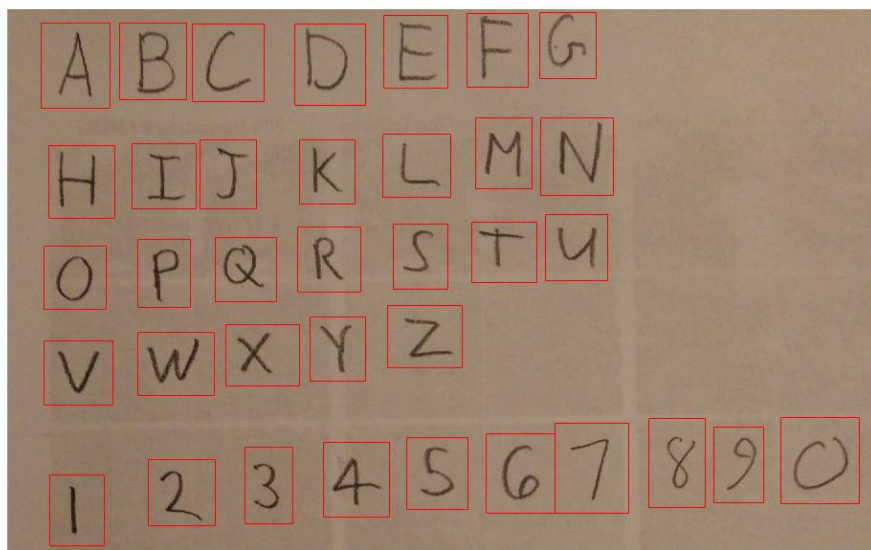


Figure 1: Sample image with handwritten characters annotated with boxes around each character.

Now that you have a network that can recognize handwritten letters with reasonable accuracy, you can now use it to parse text in an image. Given an image with some text on it, our goal is to have a function that returns the actual text in the image. However, since your neural network expects a binary image with a single character, you will need to process the input image to extract each character. There are various approaches that can be done so feel free to use any strategy you like.

Here we outline one possible method, another is that given in a [tutorial](#)

1. Process the image ([blur](#), [threshold](#), [opening morphology](#), etc. (perhaps in that order)) to classify all pixels as being part of a character or background.
2. Find connected groups of character pixels (see [skimage.measure.label](#)). Place a bounding box around each connected component.
3. Group the letters based on which line of the text they are a part of, and sort each group so that the letters are in the order they appear on the page.
4. Take each bounding box one at a time and resize it to 32×32 , classify it with your network, and report the characters in order (inserting spaces when it makes sense).

Since the network you trained likely does not have perfect accuracy, you can expect there to be some errors in your final text parsing. Whichever method you choose to implement for the character detection, you should be able to place a box on most of these characters in the image. We have provided you with `01_list.jpg`, `02_letters.jpg`, `03_haiku.jpg` and `04_deep.jpg` to test your implementation on.

Q4.1 Theory [3 points] The method outlined above is pretty simplistic, and makes several assumptions. What are two big assumptions that the sample method makes. In your writeup, include two example images where you expect the character detection to fail (either miss valid letters, or respond to non-letters).

Q4.2 Code [13 points] Find letters in the image. Given an RGB image, this function should return bounding boxes for all of the located handwritten characters in the image, as well as a binary black-and-white version of the image `im`. Each row of the matrix should contain `[y1,x1,y2,x2]` the positions of the top-left and bottom-right corners of the box. The black and white image should be floating point, 0 to 1, with the characters in black and background in white.

Q4.3 Writeup [6 points] Run `findLetters(..)` on all of the provided sample images in `images/`. Plot all of the located boxes on top of the image to show the accuracy of your `findLetters(..)` function. Include all the result images in your writeup.

Q4.4 Code/Writeup [13 points] Now you will load the image, find the character locations, classify each one with the network you trained in **Q3.2.1**, and return the text contained in the image. Be sure you try to make your detected images look like the images from the training set. Visualize them and act accordingly.

Run your `run_q4` on all of the provided sample images in `images/`. Include the extracted text in your writeup.

Note: We would recommend finishing section 7 (pytorch) before attempting the extra credit sections (5, 6)

5 Image Compression with Autoencoders (25 points Extra Credit)

An autoencoder is a neural network that is trained to attempt to copy its input to its output, but it usually allows copying only approximately. This is typically achieved by restricting the number of hidden nodes inside the autoencoder; in other words, the autoencoder would be forced to learn to *represent* data with this limited number of hidden nodes. This is a useful way of learning compressed representations. In this section, we will continue using the NIST36 dataset you have from the previous questions.

5.1 Building the Autoencoder

Q5.1.1 Code [10 points] Due to the difficulty in training auto-encoders, we have to move to the $\text{relu}(x) = \max(x, 0)$ activation function. It is provided for you in `util.py`. Implement a 2 hidden layer autoencoder where the layers are

- 1024 to 32 dimensions, followed by a ReLU
- 32 to 32 dimensions, followed by a ReLU
- 32 to 32 dimensions, followed by a ReLU
- 32 to 1024 dimensions, followed by a sigmoid (this normalizes the image output for us)

The loss function that you're using is total squared error for the output image compared to the input image (they should be the same!).

Q5.1.2 Code [5 points] To help even more with convergence speed, we will implement [momentum](#). Now, instead of updating $W = W - \alpha \frac{\partial J}{\partial W}$, we will use the update rules $M_W = 0.9M_W - \alpha \frac{\partial J}{\partial W}$ and $W = W + M_W$. To implement this, populate the parameters dictionary with zero-initialized momentum accumulators, one for each parameter. Then simply perform both update equations for every batch.

5.2 Training the Autoencoder

Q5.2 Writeup/Code [3 points] Using the provided default settings, train the network for 100 epochs. What do you observe in the plotted training loss curve as it progresses?

5.3 Evaluating the Autoencoder

Q5.3.1 Writeup/Code [4 points] Now let's evaluate how well the autoencoder has been trained. Select 5 classes from the total 36 classes in your dataset and for each selected class include in your report 2 validation images and their reconstruction. What differences do you observe that exist in the reconstructed validation images, compared to the original ones?

Q5.3.2 Writeup [3 points] Let's evaluate the reconstruction quality using Peak Signal-to-noise Ratio (PSNR). PSNR is defined as

$$\text{PSNR} = 20 \times \log_{10}(\text{MAX}_I) - 10 \times \log_{10}(\text{MSE}) \quad (1)$$

where MAX_I is the maximum possible pixel value of the image, and MSE (mean squared error) is computed across all pixels. You may use [skimage.measure.compare_psnr](#) for convenience. Report the average PSNR you get from the autoencoder across all validation images.

6 Comparing against PCA (15 points Extra Credit)

As a baseline for comparison, we will use one of the most popular methods for data dimensionality reduction - Principle Component Analysis (PCA). PCA allows one to find the best low-rank approximation of the data by keeping only a specified number of principle components. To perform PCA, we will use a factorization method called Singular Value Decomposition (SVD).

Run SVD on the training data. One of the matrices will be an orthonormal matrix that indicates the components of your data, sorted by their importances. Extract the first 32 principle components and form a projection matrix; you will need to figure out how to do these from the U, S, V matrices.

Q6.1 Writeup [4 points] What is the size of your projection matrix? What is its rank?

This projection matrix was “trained” from our training data. Now let’s “test” it on our test data. Use the projection matrix on test data to obtain the reconstructed test images. Note that these reconstructions can also be represented with only 32 numbers.

Q6.2 Writeup [4 points] Use the classes you selected in Q5.3.1, and for each of these 5 classes, include in your report 2 test images and their reconstruction. You may use test labels to help you find the corresponding classes. What differences do you observe that exist in the reconstructed test images, compared to the original ones? How do they compare to the ones reconstructed from your autoencoder?

Q6.3 Writeup [4 points] Report the average PSNR you get from PCA. Is it better than your autoencoder? Why?

Q6.4 Writeup [3 points] Count the number of learned parameters for both your autoencoder and the PCA model. How many parameters do the respective approaches learn? Why is there such a difference in terms of performance?

7 PyTorch (40 points)

While you were able to derive manual back-propagation rules for sigmoid and fully-connected layers, wouldn't it be nice if someone did that for lots of useful primitives and made it fast and easy to use for general computation? Meet [automatic differentiation](#). Since we have high-dimensional inputs (images) and low-dimensional outputs (a scalar loss), it turns out **forward mode AD** is very efficient. Popular autodiff packages include [pytorch](#) (Facebook), [tensorflow](#) (Google), [autograd](#) (Boston-area academics). Autograd provides its own replacement for numpy operators and is a drop-in replacement for numpy, except you can ask for gradients now. The other two are able to act as shim layers for [cuDNN](#), an implementation of auto-diff made by Nvidia for use on their GPUs. Since GPUs are able to perform large amounts of math much faster than CPUs, this makes the former two packages very popular for researchers who train large networks. Tensorflow asks you to build a computational graph using its API, and then is able to pass data through that graph. PyTorch builds a dynamic graph and allows you to mix autograd functions with normal python code much more smoothly, so it is currently more popular among CMU students.

We will use [PyTorch](#) as a framework. Many computer vision projects use neural networks as a basic building block, so familiarity with one of these frameworks is a good skill to develop. Here, we basically replicate and slightly expand our handwritten character recognition networks, but do it in PyTorch instead of doing it ourselves. Feel free to use any tutorial you like, but we like [the official one](#) or [this tutorial](#) (in a jupyter notebook) or [these slides](#) (starting from number 35).

For this section, you're free to implement these however you like. All of the tasks required here are fairly small and don't require a GPU if you use small networks. Including 7.2.

7.1 Train a neural network in PyTorch

Q7.1.1 Code/Writeup [10 points] Re-write and re-train your fully-connected network on NIST36 in PyTorch. Plot training accuracy and loss over time.

Q7.1.2 Code/Writeup [3 points] Train a convolutional neural network with PyTorch on MNIST. Plot training accuracy and loss over time.

Q7.1.3 Code/Writeup [2 points] Train a convolutional neural network with PyTorch on the included NIST36 dataset.

Q7.1.4 Code/Writeup [15 points] Train a convolutional neural network with PyTorch on the EMNIST Balanced dataset (available in `torchvision.datasets`, use `balanced` split) and evaluate it on the findLetters bounded boxes from the images folder.

7.2 Fine Tuning

When training from scratch, a lot of epochs and data are often needed to learn anything meaningful. One way to avoid this is to instead initialize the weights more intelligently.

These days, it is most common to initialize a network with weights from another deep network that was trained for a different purpose. This is because, whether we are doing image classification, segmentation, recognition etc..., most real images share common properties. Simply copying the weights from the other network to yours gives your network a head start, so your network does not need to learn these common weights from scratch all over again. This is also referred to as fine tuning.

Q7.2.1 Code/Writeup [10 points] Fine-tune a single layer classifier using pytorch on the [flowers 17](#) (or [flowers 102!](#)) dataset using [squeezenet1_1](#), as well as an architecture you've designed yourself (*3 conv layers, followed 2 fc layers, it's standard [slide 6](#)*) and trained from scratch. How do they compare?

We include a script in `scripts/` to fetch the flowers dataset and extract it in a way that [PyTorch ImageFolder](#) can consume it, see [an example](#), from `data/oxford-flowers17`. You should look at how SqueezeNet is [defined](#), and just replace the classifier layer. There exists a pretty good example for [fine-tuning](#) in PyTorch.

References

- [1] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. 2010. <http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>.
- [2] P. J. Grother. Nist special database 19 – handprinted forms and characters database. <https://www.nist.gov/srd/nist-special-database-19>, 1995.

8 Appendix: Neural Network Overview

Deep learning has quickly become one of the most applied machine learning techniques in computer vision. Convolutional neural networks have been applied to many different computer vision problems such as image classification, recognition, and segmentation with great success. In this assignment, you will first implement a fully connected feed forward neural network for hand written character classification. Then in the second part, you will implement a system to locate characters in an image, which you can then classify with your deep network. The end result will be a system that, given an image of hand written text, will output the text contained in the image.

8.1 Basic Use

Here we will give a brief overview of the math for a single hidden layer feed forward network. For a more detailed look at the math and derivation, please see the class slides.

A fully-connected network \mathbf{f} , for classification, applies a series of linear and non-linear functions to an input data vector \mathbf{x} of size $N \times 1$ to produce an output vector $\mathbf{f}(\mathbf{x})$ of size $C \times 1$, where each element i of the output vector represents the probability of \mathbf{x} belonging to the class i . Since the data samples are of dimensionality N , this means the input layer has N input units. To compute the value of the output units, we must first compute the values of all the hidden layers. The first hidden layer *pre-activation* $\mathbf{a}^{(1)}(\mathbf{x})$ is given by

$$\mathbf{a}^{(1)}(\mathbf{x}) = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$$

Then the *post-activation* values of the first hidden layer $\mathbf{h}^{(1)}(\mathbf{x})$ are computed by applying a non-linear activation function \mathbf{g} to the *pre-activation* values

$$\mathbf{h}^{(1)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(1)}(\mathbf{x})) = \mathbf{g}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

Subsequent hidden layer ($1 < t \leq T$) pre- and post activations are given by:

$$\begin{aligned}\mathbf{a}^{(t)}(\mathbf{x}) &= \mathbf{W}^{(t)}\mathbf{h}^{(t-1)} + \mathbf{b}^{(t)} \\ \mathbf{h}^{(t)}(\mathbf{x}) &= \mathbf{g}(\mathbf{a}^{(t)}(\mathbf{x}))\end{aligned}$$

The output layer *pre-activations* $\mathbf{a}^{(T)}(\mathbf{x})$ are computed in a similar way

$$\mathbf{a}^{(T)}(\mathbf{x}) = \mathbf{W}^{(T)}\mathbf{h}^{(T-1)}(\mathbf{x}) + \mathbf{b}^{(T)}$$

and finally the *post-activation* values of the output layer are computed with

$$\mathbf{f}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(T)}(\mathbf{x})) = \mathbf{o}(\mathbf{W}^{(T)}\mathbf{h}^{(T-1)}(\mathbf{x}) + \mathbf{b}^{(T)})$$

where \mathbf{o} is the output activation function. Please note the difference between \mathbf{g} and \mathbf{o} ! For this assignment, we will be using the sigmoid activation function for the hidden layer, so:

$$\mathbf{g}(y) = \frac{1}{1 + \exp(-y)}$$



Figure 2: Samples from NIST Special 19 dataset [2]

where when \mathbf{g} is applied to a vector, it is applied element wise across the vector.

Since we are using this deep network for classification, a common output activation function to use is the softmax function. This will allow us to turn the real value, possibly negative values of $\mathbf{a}^{(T)}(\mathbf{x})$ into a set of probabilities (vector of positive numbers that sum to 1). Letting \mathbf{x}_i denote the i^{th} element of the vector \mathbf{x} , the softmax function is defined as:

$$\mathbf{o}_i(\mathbf{y}) = \frac{\exp(\mathbf{y}_i)}{\sum_j \exp(\mathbf{y}_j)}$$

Gradient descent is an iterative optimisation algorithm, used to find the local optima. To find the local minima, we start at a point on the function and move in the direction of negative gradient (steepest descent) till some stopping criteria is met.

8.2 Backprop

The update equation for a general weight $W_{ij}^{(t)}$ and bias $b_i^{(t)}$ is

$$W_{ij}^{(t)} = W_{ij}^{(t)} - \alpha * \frac{\partial L_{\mathbf{f}}}{\partial W_{ij}^{(t)}}(\mathbf{x}) \quad b_i^{(t)} = b_i^{(t)} - \alpha * \frac{\partial L_{\mathbf{f}}}{\partial b_i^{(t)}}(\mathbf{x})$$

α is the learning rate. Please refer to the back-propagation slides for more details on how to derive the gradients. Note that here we are using softmax loss (which is different from the least square loss in the slides).