

UNIVERSITY OF PENNSYLVANIA

ESE 546: PRINCIPLES OF DEEP LEARNING

FALL 2020

[09/09] HOMEWORK 1

DUE: 09/29 TUE, 1.30P ET

Changelog

- **09/24 11a:** Gabor Filterbank suggested parameters changed (see main text). You are also free to choose your own parameters.
 - **09/22 5p:** Due date extended by one day.
 - **09/15 9a:** Line 71-74. Computing the Gabor filters for 5,000 images is time consuming. You can do it only for 500 images (50 images/class).
 - **09/15 9a:** Line 299. No need to compute the validation loss/error for every weight update. Compute it every 1000 weight updates.
 - **09/14 7p:** Line 108: Image size 28×28 has been changed to 14×14 .
 - **09/17 3p:** Line 146, 218, 228: Fixed the dimension issue. $h^{(l+1)} = Wh^{(l)} + b$ has been changed to $h^{(l+1)} = h^{(l)}W^\top + b$.
-

Instructions

Read the following instructions carefully before beginning to work on the homework.

- You will submit solutions typeset in L^AT_EX on Gradescope (strongly encouraged). You can use hw_template.tex on Canvas in the “Homeworks” folder to do so. If your handwriting is *unambiguously legible*, you can submit PDF scans/tablet-created PDFs.
- Please start a new problem on a fresh page and mark all the pages corresponding to each problem. Failure to do so may result in your work not graded completely.
- Clearly indicate the name and Penn email ID of all your collaborators on your submitted solutions.
- For each problem in the homework, you should mention the total amount of time you spent on it.
- You can be informal while typesetting the solutions, e.g., if you want to draw a picture feel free to draw it on paper clearly, click a picture and include it in your solution. Do not spend undue time on typesetting solutions.
- You will see an entry of the form “HW 1 PDF” where you will upload the PDF of your solutions. You will also see entries like “HW 1 Problem 1 Code” and “HW 1 Problem 3 Code” where you will upload your solution for the respective problems. **For each programming problem, you should create a fresh Google Colab notebook.** This notebook should contain all the code to reproduce the results of the problem and you will upload the .ipynb

file obtained from Colab. Name your notebook to be “pennkey_hw1_problem1.ipynb”, e.g., I will name my code for Problem 1 as “pratikac_hw1_problem1.ipynb”. **Remember that you should also include the relevant plots in the PDF, without doing so you will not get full credit for the problem.**

- To be safe, make sure that the option “Omit code cell output when saving this notebook” is unchecked in the “Notebook Settings” tab in Colab. This way, the instructors may be able to give you partial credit for incorrect output/plots.
- **This is very important.** Note that the instructors will download your notebook and execute it on Colab themselves, so your notebook should be such that it can be executed independently without any errors to create all output/plots required in the problem.

Notice that even if Problem 1 contains multiple sub-parts, you should submit only one Colab notebook. In these cases, you will demarcate clearly within your notebook what is the output cell of each sub-problem. Jupyter notebooks can handle markdown very well so use it to write your solutions and narrative liberally. You should make sections/sub-sections in the Jupyter notebook for us to find your answers/plots easily. As an example, see Rahul Ramesh’s Colab notebook at https://github.com/rahul13ramesh/cis520-dataset/blob/master/Clean_CIS520_Project.ipynb for his CIS 520 project. Notice how the code cells are preceded by text cells with titles appropriately written in bold. You can have plots for your debugging in the notebook but please comment them when you submit it to Gradescope.

Credit The points for the problems add up to 140. You only need to solve for 100 points to get full credit, i.e., your final score will be $\min(\text{your total points}, 100)$.

1 **Problem 1 (60 points). Use Google Colab for this problem.** In this problem, we will fit the MNIST
2 dataset using a support vector machine (SVM) using the “scikit-learn” library. You can install it using

```
3  
4 [local] pip install scikit-learn scikit-image  
5 [colab] !pip install scikit-learn scikit-image
```

7 An SVM solves an optimization problem for maximizing the margin between two classes. Support
8 that we have a binary classification problem where (x_i, y_i) are the data and ground-truth labels
9 respectively and $y_i \in \{-1, 1\}$. We would like to find a hyper-plane that separates the data such that
10 all examples with labels $y_i = +1$ are on side and all examples with labels $y_i = -1$ are on the other
11 side. This involves solving the problem

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\theta\|^2 \\ & \text{subject to} \quad y_i(\theta^\top x_i + \theta_0) \geq 1 \quad \forall i = 1, \dots, n; \end{aligned} \tag{1}$$

12 here θ_0 is the offset parameter and θ is the hyper-plane. You can eliminate the offset parameter by
13 appending a 1 to the data, i.e., feeding in $x' = [x, 1]$ as the data with the same labels.

14 (a) (5 point) It may not always be possible to classify a dataset cleanly into positive and negatively
15 labeled samples, i.e., there may not exist a θ that satisfies all constraints in (1). To handle such cases,
16 we relax the problem formulation. We create a “slack” variable that allows the constraint to be written
17 as

$$\text{subject to} \quad y_i(\theta^\top x_i + \theta_0) \geq 1 - \xi_i; \quad \xi_i \geq 0.$$

18 The variable ξ_i measures the degree to which we can violate the original constraint. We would like to
19 minimize the violation of the original constraints and the slack variable-based formulation of (1) will
20 use a different objective that does so. There can be many such objectives, write down one.

21 (b) (2 point) Define what are support samples in an SVM.

22 (c) (3 points) You can download the dataset using

```
23  
24 from sklearn.datasets import fetch_openml  
25 from sklearn.model_selection import train_test_split  
26  
27 ds = fetch_openml('mnist_784')  
28 x, y = ds.data, ds.target  
29  
30 x_train, x_val, y_train, y_val = train_test_split(x, y,  
31 test_size=0.2, random_state=42)
```

33 Check whether you have downloaded the data correctly; the images in x_{train} and x_{val} are in the
34 form of a vector of length 784, this is really the flattened matrix 28×28 . You can check it by plotting

```
35  
36 import matplotlib.pyplot as plt  
37 a = x_train[0].reshape((28, 28))  
38 plt.imshow(a)  
39  
40 # code for down-sampling  
41 import cv2  
42 b = cv2.resize(b, (14, 14))  
43 plt.imshow(b)
```

45 Construct training (80%) and validation (20%) datasets from the arrays x, y by sampling the images
46 and labels randomly. Why did we not construct a test dataset here?

47 (d) (15 points) Create the SVM classifier in scikit-learn using

```
48 classifier = svm.SVC(C=1.0, kernel='rbf', gamma='auto')
```

51 What do the parameters C and γ do? What are their default values? Fit the SVM classifier to the
52 data and predict the labels of the validation dataset using the trained classifier. Note that the input
53 data for an SVM is a vector of 784, not an image of size 28×28 . Provide the validation accuracy
54 and the 10-class confusion matrix. Note down the ratio of the number of support samples to the
55 total number of training samples for your trained classifier. **If training takes too long or runs out
56 of memory, you can down-sample the original 28×28 images to 14×14 (remember to reshape
57 it to 196 before training the SVM), and/or reduce the size of the training set.**

58 (e) (5 points) Read the manual of `svm.SVC` carefully. Identify all the options that you may not have
59 seen in your previous course on SVMs. Libraries that are used in production such as scikit-learn
60 will have numerous knobs to improve the performance; these knobs often implement state of the art
61 research and it is useful to know them. For instance, what does the parameter named “shrinking”
62 in `svm.SVC` do? Investigate and explain what optimization algorithm is used to fit the SVM in
63 scikit-learn.

64 (f) (5 points) The mathematical formulation of the SVM above is for a binary classifier. The MNIST
65 dataset consists of digits from 0-9 and has 10 classes in total. How does `svm.SVC` handle multiple
66 classes? Can you think of any alternative ways to use binary classifiers to perform multi-class
67 classification?

68 (g) (5 points) Use the `sklearn.model_selection.GridSearchCV` function to pick a better value than
69 the default one for the hyper-parameter C . Try at least 5 different hyper-parameters. Show all the
70 hyper-parameters tried by the method and their accuracies.

71 (h) **The following two parts are computationally intensive. Down-sample all images to 14×14
72 and create a training dataset using only 500 images from the full MNIST dataset. Make sure
73 that the training dataset is balanced, i.e., pick 50 images per digit. Similarly, pick an additional
74 500 images to form the validation set.**

75 The default kernel in `svm.SVC` is a radial basis function. The MNIST dataset consists of images and
76 since images have local regularities we can build a better classifier by exploiting them. It has been
77 found that the mammalian visual cortex consists of cells well-modeled by Gabor functions (named
78 after Dennis Gabor, a Hungarian physicist who invented holography). Let us represent each image as
79 a function $I(x, y)$, this function gives the intensity at pixel location (x, y) . A Gabor filter is given by
80 a function

$$g(x, y) = \exp(i 2\pi F(x \cos \omega + y \sin \omega)) \exp\left(-\pi \left(\frac{p^2}{\sigma_x^2} + \frac{q^2}{\sigma_y^2}\right)\right)$$

81 where $p = x \cos \theta + y \sin \theta$ and $q = -x \sin \theta + y \cos \theta$. First, note that this filter is a complex
82 function, this is different from a standard convolutional filter. Convolving the original image $I(x, y)$
83 with the filter $g(x, y)$ will result in two sets of co-efficients, one real and the other imaginary. The
84 parameters we will be concerned with are:

- 85 • F this is the spatial frequency of the filter,
- 86 • θ the rotation angle of the Gaussian,
- 87 • σ_x, σ_y : standard deviation of the kernel in the X and Y directions, and
- 88 • the parameter “bandwidth” in the code below is inversely related to the standard deviation
fixed the frequency.

90 You can read [this webpage](#) for a simple introduction to these filters (this is given in the OpenCV
91 format). You can also read this more mathematical [tutorial on Gabor filters](#) which is given in the
92 scikit-image format that we discussed above.

93 We will use the scikit-image library which implements a smaller machine learning-specific set of
94 image processing functions. Alternatively, you can also use the cv2.getGaborKernel function in
95 OpenCV.

```
96 from skimage.filters import gabor_kernel, gabor
97 import numpy as np
98
99 freq, theta, bandwidth = 0.1, np.pi/4, 1
100 gk = gabor_kernel(frequency=freq, theta=theta, bandwidth=bandwidth)
101 plt.figure(1); plt.clf(); plt.imshow(gk.real)
102 plt.figure(2); plt.clf(); plt.imshow(gk.imag)
103
104 # convolve the input image with the kernel and get co-efficients
105 # we will use only the real part and throw away the imaginary
106 # part of the co-efficients
107 image = x_train[0].reshape((14,14))
108 coeff_real, _ = gabor(image, frequency=freq, theta=theta,
109                         bandwidth=bandwidth)
110 plt.figure(1); plt.clf(); plt.imshow(coeff_real)
111
```

113 (j) (20 points) Run the above code a few times with different parameters for F, θ and bandwidth to
114 see how the filter changes in shape and size and the corresponding output after convolution. We will
115 create a filter bank that consists of multiple Gabor filters of fixed parameters. Instead of considering
116 the pixel intensities of the MNIST images as the features for training the SVM, the co-efficients of
117 the Gabor filter-bank will be used to train the SVM. You can pick

```
118 theta = np.arange(0, np.pi, np.pi/4)
119 frequency = np.arange(0.05, 0.5, 0.15)
120 bandwidth = np.arange(0.3, 1, 0.3)
```

123 This gives a total of 36 filters in the filter-bank. We therefore have converted a $14 \times 14 = 196$ pixel
124 image into a vector of length $196 \times 36 = 7056$. Plot the filter-bank to see that it gives you a good
125 spread of different filters. You want a diverse filter bank that can capture different rotations and scales.
126 Train the SVM on these features and report the training and validation accuracy.

127 Increase the number of filters next. You might have to use PCA to reduce the dimensionality of the
128 dataset to be able to fit the SVM in RAM; use scikit-learn to do so.

129 **Problem 2 (10 points).** Prove Jensen's inequality: for any random variable X with expectation μ
130 and a convex, finite function φ

$$\mathbb{E}_X[\varphi(X)] \geq \varphi(\mu).$$

131 You can assume that the random variable X takes values in a finite set. If you want to prove it in a
132 more general setting, you can assume that the function φ is differentiable.

133 **Problem 3 (70 points).** In this problem, you will write code to train a neural network completely
134 from scratch using only Numpy and basic Python (note no PyTorch/TF/other deep learning library).

135 **Work on this problem on your personal computer before moving to Colab, this will help during
136 debugging.**

137 (a) (5 points) Download the MNIST dataset using the following code.

```

138
139     import torchvision as thv
140     train = thv.datasets.MNIST('./', download=True, train=True)
141     val = thv.datasets.MNIST('./', download=True, train=False)
142     print(train.data.shape, len(train.targets))

```

144 The training dataset has 60,000 images while the validation dataset has 10,000 images spread roughly
 145 equally across 10 classes. Take 50% of the images *from each class* for training and validation, i.e.,
 146 about 30,000 training images and 5,000 validation images, almost evenly spread across all classes
 147 with a few minor differences. We will use this smaller dataset in this problem. **Plot the images of a**
 148 **few randomly chosen images from your dataset.**

149 (b) (10 points) We will next implement different parts of a typical neural network. First write a linear
 150 layer; this includes the forward function

$$h^{(l+1)} = h^{(l)} W^\top + b$$

151 and the corresponding backward function that takes the gradient $\overline{h^{(l+1)}}$ and outputs \overline{W} , \overline{b} and $\overline{h^{(l)}}$.
 152 Remember to write your function in such a way that it takes in a mini-batch of vectors $h^{(l)}$ as the
 153 input, i.e., if the feature vector $h^{(l)}$ is a -dimensional, for b images in the mini-batch, your forward
 154 function will take as input

$$h^{(l)} \in \mathbb{R}^{b \times a}$$

155 use

$$W \in \mathbb{R}^{c \times a}, \quad b \in \mathbb{R}^c$$

156 and output a mini-batch of feature vectors of size

$$h^{(l+1)} \in \mathbb{R}^{b \times c}.$$

157 Note that in this problem we have $a = 784$ because there are 28×28 pixels in MNIST images and
 158 $c = 10$ because there are 10 classes in MNIST. You should use numpy to write the forward function;
 159 do not use a for loop for computing the mini-batch-ed forward because it will be too slow for the next
 160 parts of the problem. You are advised to first write this function for $b = 1$ to understand the process
 161 and then you can extend it to $b > 1$. Some pseudo code is given below.

```

162
163 class linear_t:
164     def __init__(self):
165         # initialize to appropriate sizes, fill with Gaussian entires
166         # normalize to make the Frobenius norm of w, b equal to 1
167         self.w, self.b = ...
168
169     def forward(self, h^l):
170         h^{l+1} = ...
171         # cache h^l in forward because we will need it to compute
172         # dw in backward
173         self.hl = h^l
174         return h^{l+1}
175
176     def backward(self, dh^{l+1}):
177         dh^l, dw, db = ...
178         self.dw, self.db = dw, db
179         # notice that there is no need to cache dh^l
180         return dh^l
181
182     def zero_grad(self):
183         # useful to delete the stored backprop gradients of the
184         # previous mini-batch before you start a new mini-batch

```

```
185     self.dw, self.db = 0 * self.dw, 0 * self.db
```

187 (c) (5 points) Implement the rectified linear unit (ReLU) layer next. This will take the form of

$$h^{(l+1)} = \max(0, h^{(l)})$$

188 where the max is performed element-wise on the elements of $h^{(l)}$. Write the forward function and
189 the corresponding backward function.

190 (d) (10 points) Next we will write a combined softmax and cross-entropy loss layer. This is a layer
191 that first performs the operation

$$h_k^{(l+1)} = \frac{e^{h_k^{(l)}}}{\sum_{k'} e^{h_{k'}^{(l)}}}$$

192 where $h_k^{(l)}$ is the k^{th} element of the vector $h^{(l)}$. The input to this layer, i.e., $h^{(l)}$ are called the “logits”.
193 The output of this layer is a scalar, it is the negative log-probability of predicting the correct class, i.e.,

$$\ell(y) = -\log(h_y^{(l+1)}).$$

194 where y is the true label of the image. For a mini-batch with θ images, the average loss will be

$$\ell(\{y_i\}_{i=1,\dots,\theta}) = -\frac{1}{\theta} \sum_{i=1}^{\theta} \log(h_{y_i}^{(l+1)}).$$

195 You will again implement a forward function and a backward function for it yourself; remember to
196 implement both functions to take in a mini-batch of inputs. The pseudo-code for the log-softmax
197 layer is similar to that of the fully-connected layer. It does not have any parameters to initialize and
198 therefore does not need the zero_grad method.

```
199 class softmax_cross_entropy_t:
200     def __init__(self):
201         # no parameters, nothing to initialize
202
203     def forward(self, h^l, y):
204         h^{l+1} = ...
205         # compute average loss ell(y) over a mini-batch
206         ell = ...
207         error = ...
208         return ell, error
209
210     def backward(self):
211         # as we saw in the notes, the backprop input to the
212         # loss layer is 1, so this function does not take any
213         # arguments
214         dh^l = ...
215         return dh^l
```

218 We can also output the error of predictions in the forward function. It is computed as

$$\text{error} = \frac{1}{\theta} \sum_{i=1}^{\theta} \mathbf{1}_{\{y_i \neq \text{argmax}_k h_k^{(l+1)}\}}$$

219 and measures the number of mistakes the network makes.

220 (e) (10 points) Before moving on to training, let us check whether we have implemented the forward
 221 and backward correctly for all the three layers. Consider the function for the linear layer. **Use a**
 222 **batch-size $b = 1$ for this part.** The forward function for the linear layer implements

$$h^{(l+1)} = h^{(l)} W^\top + b$$

223 which is easy enough. However, we would like to check our implementation of the backward function.

```
224 def backward(self, dh^{l+1}):
225     dh^l, self.dw, self.db = ...
226     return dh^l
```

229 Think carefully about your implementation of the backward function. Notice that if you call the
 230 backward function with the argument $\bar{h}^{l+1} = [0, 0, \dots, 0, 1, 0, 0, \dots]$, i.e., there is a 1 at the k^{th}
 231 element, the function is going to calculate the quantities

$$\text{self}.dw = \frac{\partial h_k^{(l+1)}}{\partial W}, \quad \text{self}.db = \frac{\partial h_k^{(l+1)}}{\partial b}, \quad dh^{(l)} = \frac{\partial h_k^{(l+1)}}{\partial h^{(l)}}.$$

232 We now compute the estimate of the derivative using finite-differences, e.g.,

$$\frac{\partial h_k^{(l+1)}}{\partial W_{ij}} \approx \frac{\left(h^{(l)}(W + \epsilon)^\top\right)_k - \left(h^{(l)}(W - \epsilon)^\top\right)_k}{2\epsilon_{ij}}$$

233 where ϵ is a matrix with a Gaussian random variable as the $(ij)^{\text{th}}$ entry and zero everywhere else. In
 234 simple words, you can perturb the $(ij)^{\text{th}}$ element of weight W by ϵ_{ij} , compute the right hand-side of
 235 the finite-difference estimate above and compare it with the $(ij)^{\text{th}}$ element of your variable `self.dw`.

236 This idea checks the gradient with respect to only one element of W , namely W_{ij} . Do this for about
 237 10 randomly chosen elements of W and a few (5 should be enough) different entries k of $h_k^{(l+1)}$ and
 238 check if the answer matches `self.dw` that you have implemented in the backward function. Repeat
 239 this process for the other two gradients.

240 **Do not move on to the next part until you are convinced your implementation of forward/back-**
 241 **ward is correct for all the three layers. It is essential that the gradient is implemented correctly,**
 242 **your training will not work if the gradient is wrong.**

243 (f) (10 points) You will now train your neural network. The pseudo-code looks as follows:

```
244 # load dataset
245 ...
246
247
248 # initialize all the layers
249 l1, l2, l3 = linear_t(), relu_t(), softmax_cross_entropy_t()
250 net = [l1, l2, l3]
251
252 # train for at least 1000 iterations
253 for t in range(1000):
254     # 1. sample a mini-batch of size = 32
255     # each image in the mini-batch is chosen uniformly randomly from the
256     # training dataset
257     x, y = ...
258
259     # 2. zero gradient buffer
260     for l in net:
261         l.zero_grad()
```

```

263     # 3. forward pass
264     h1 = l1.forward(x)
265     h2 = l2.forward(h1)
266     ell, error = l3.forward(h2, y)
267
268     # 4. backward pass
269     dh2 = l3.backward()
270     dh1 = l2.backward(dh2)
271     dx = l1.backward(dh1)
272
273     # 5. gather backprop gradients
274     dw, db = l1.dw, l1.db
275
276     # 6. print some quantities for logging
277     # and debugging
278     print(t, ell, error)
279     print(t, np.linalg.norm(dw/l1.w), np.linalg.norm(db/l1.b))
280
281     # 7. one step of SGD
282     l1.w = l1.w - lr*dw
283     l1.b = l1.b - lr*db

```

285 You can pick the learning rate to be $lr = 0.1$. **Plot the training loss and training error as a function of the number of weight updates.** Make sure that the training loss decreases with the number of updates. You should try to get better than/around 15% error on the training dataset after 10,000-50,000 updates.

289 (g) (5 points) We have implemented the training loop. Write the corresponding code for computing the validation loss and error.

```

291 def validate(w, b):
292     # 1. iterate over mini-batches from the validation dataset
293     # note that this should not be done randomly, we want to check
294     # every image only once
295
296     loss, tot_error = 0, 0
297     for i in range(0, 5000, 32):
298         x, y = val.data[i:i+32], val.targets[i:i+32]
299
300         # 2. compute forward pass and error

```

303 **Plot the validation loss and validation error as a function of the number of weight updates, every 1000 weight updates.**

305 If everything works as expected, congratulations! You have implemented your own little library for training neural networks, completely from scratch!

307 (h) (15 points) Repeat the entire process in parts (b)-(g) using the pre-built functions inside PyTorch.
308 You will take help of the code provided in the recitation sessions for this purpose. Train the network
309 for at least 10,000 weight updates this time. Plot the training loss, training error, validation loss and
310 the validation error as a function of the number of weight updates.