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

```
[local] pip install scikit-learn scikit-image
[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

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

12 here $\theta_0$ is the offset parameter and $\theta$ 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 $\theta$ 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

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

18 The variable $\xi_i$ 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

```
from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split

ds = fetch_openml('mnist_784')
x, y = ds.data, ds.target

x_train, x_val, y_train, y_val = train_test_split(x, y,
                    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

```
import matplotlib.pyplot as plt
a = x_train[0].reshape((28,28))
plt.imshow(a)

# code for down-sampling
import cv2
b = cv2.resize(b, (14,14))
plt.imshow(b)
```

45 Construct training (80%) and validation (20%) datasets from the arrays $x, y$ by sampling the images
46 and labels randomly. You should make sure that each class has an equal number of samples. Why did
47 we not construct a test dataset here?

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

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

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

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

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

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

72 (h) **The following two parts are computationally intensive. Down-sample all images to 14×14**
73 **and create a training dataset using only 500 images from the full MNIST dataset. Make sure**
74 **that the training dataset is balanced, i.e., pick 50 images per digit. Similarly, pick an additional**
75 **500 images (50 images/digit) to form the validation set.**

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

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

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

86     • $F$ this is the spatial frequency of the filter,
87     • $\theta$ the rotation angle of the Gaussian,

88    • $\sigma_x, \sigma_y$: standard deviation of the kernel in the X and Y directions, and

89    • the parameter "bandwidth" in the code below is inversely related to the standard deviation

90      fixed the frequency.

91    You can read this webpage for a simple introduction to these filters (this is given in the OpenCV

92    format). You can also read this more mathematical tutorial on Gabor filters which is given in the

93    scikit-image format that we discussed above.

94    We will use the scikit-image library which implements a smaller machine learning-specific set of

95    image processing functions. Alternatively, you can also use the cv2.getGaborKernel function in

96    OpenCV.

```python
97
98   from skimage.filters import gabor_kernel, gabor
99   import numpy as np
100
101  freq, theta, bandwidth = 0.1, np.pi/4, 1
102  gk = gabor_kernel(frequency=freq, theta=theta, bandwidth=bandwidth)
103  plt.figure(1); plt.clf(); plt.imshow(gk.real)
104  plt.figure(2); plt.clf(); plt.imshow(gk.imag)
105
106  # convolve the input image with the kernel and get co-efficients
107  # we will use only the real part and throw away the imaginary
108  # part of the co-efficients
109  image = x_train[0].reshape((14,14))
110  coeff_real, _ = gabor(image, frequency=freq, theta=theta,
111                        bandwidth=bandwidth)
112  plt.figure(1); plt.clf(); plt.imshow(coeff_real)
113
```

114    (j) (20 points) Run the above code a few times with different parameters for $F, \theta$ and bandwidth to

115    see how the filter changes in shape and size and the corresponding output after convolution. We will

116    create a filter bank that consists of multiple Gabor filters of fixed parameters. Instead of considering

117    the pixel intensities of the MNIST images as the features for training the SVM, the co-efficients of

118    the Gabor filter-bank will be used to train the SVM. You can pick

```python
119
120  theta = np.arange(0,np.pi,np.pi/4)
121  frequency = np.arange(0.05,0.5,0.15)
122  bandwidth = np.arange(0.3,1,0.3)
123
```

124    This gives a total of 36 filters in the filter-bank. We therefore have converted a $14 \times 14 = 196$ pixel

125    image into a vector of length $196 \times 36 = 7056$. Plot the filter-bank to see that it gives you a good

126    spread of different filters. You want a diverse filter bank that can capture different rotations and scales.

127    Train the SVM on these features and report the training and validation accuracy.

128    Increase the number of filters next. You might have to use PCA to reduce the dimensionality of the

129    dataset to be able to fit the SVM in RAM; use scikit-learn to do so.

deep learning library except for downloading the data). **Work on this problem on your personal computer before moving to Colab, this will help during debugging.**

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

```python
import torchvision as thv
train = thv.datasets.MNIST('./', download=True, train=True)
val = thv.datasets.MNIST('./', download=True, train=False)
print(train.data.shape, len(train.targets))
```

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

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

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

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

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

use

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

and output a mini-batch of feature vectors of size

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

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

```python
class linear_t:
    def __init__(self):
        # initialize to appropriate sizes, fill with Gaussian entires
        # normalize to make the Frobenius norm of w, b equal to 1
        self.w, self.b = ...

    def forward(self, h^l):
        h^{l+1} = ...
        # cache h^l in forward because we will need it to compute
        # dw in backward
        self.hl = h^l
        return h^{l+1}

    def backward(self, dh^{l+1}):
        dh^l, dw, db = ...
        self.dw, self.db = dw, db
        # notice that there is no need to cache dh^l
        return dh^l
```

```
182
183    def zero_grad(self):
184        # useful to delete the stored backprop gradients of the
185        # previous mini-batch before you start a new mini-batch
186        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

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

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

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

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

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

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

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

$$\ell(\{y_i\}_{i=1,\ldots,b}) = -\frac{1}{b}\sum_{i=1}^{b}\log\left(h_{y_i}^{(l+1)}\right).$$

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

```
200
201  class softmax_cross_entropy_t:
202      def __init__(self):
203          # no parameters, nothing to initialize
204
205      def forward(self, h^l, y):
206          h^{l+1} = ...
207          # compute average loss ell(y) over a mini-batch
208          ell = ...
209          error = ...
210          return ell, error
211
212      def backward(self):
213          # as we saw in the notes, the backprop input to the
214          # loss layer is 1, so this function does not take any
215          # arguments
216          dh^l = ...
217          return dh^l
218
```

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

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

7

220  and measures the number of mistakes the network makes.

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

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

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

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

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

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

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

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

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

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

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

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

```
# load dataset
...

# initialize all the layers
l1, l2, l3 = linear_t(), relu_t(), softmax_cross_entropy_t()
net = [l1, l2, l3]

# train for at least 1000 iterations
for t in range(1000):
    # 1. sample a mini-batch of size = 32
    # each image in the mini-batch is chosen uniformly randomly from the
    # training dataset
    x, y = ...

    # 2. zero gradient buffer
    for l in net:
```

```
262            l.zero_grad()
263
264        # 3. forward pass
265        h1 = l1.forward(x)
266        h2 = l2.forward(h1)
267        ell, error = l3.forward(h2, y)
268
269        # 4. backward pass
270        dh2 = l3.backward()
271        dh1 = l2.backward(dh2)
272        dx = l1.backward(dh1)
273
274        # 5. gather backprop gradients
275        dw, db = l1.dw, l1.db
276
277        # 6. print some quantities for logging
278        # and debugging
279        print(t, ell, error)
280        print(t, np.linalg.norm(dw/l1.w), np.linalg.norm(db/l1.b))
281
282        # 7. one step of SGD
283        l1.w = l1.w - lr*dw
284        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.

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

```
293  def validate(w, b):
294      # 1. iterate over mini-batches from the validation dataset
295      # note that this should not be done randomly, we want to check
296      # every image only once
297
298      loss, tot_error = 0, 0
299      for i in range(0, 5000, 32):
300          x, y = val.data[i:i+32], val.targets[i:i+32]
301
302          # 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**.

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

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