

UNIVERSITY OF PENNSYLVANIA  
ESE 546: PRINCIPLES OF DEEP LEARNING

---

1 **Changelog**

---

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

- 3 • You will submit solutions typeset in  $\text{\LaTeX}$  on Gradescope (strongly encouraged). You can  
4 use `hw_template.tex` on Canvas in the “Homeworks” folder to do so. If your handwriting is  
5 unambiguously legible, you can submit PDF scans/tablet-created PDFs.
- 6 • Clearly indicate the name and Penn email ID of all your collaborators on your submitted  
7 solutions.
- 8 • Start a new problem on a fresh page and mark all the pages corresponding to each problem.  
9 Failure to do so may result in your work not graded completely.
- 10 • For each problem in the homework, you should mention the total amount of time you spent  
11 on it. This helps us keep track of which problems most students are finding difficult.
- 12 • You can be informal while typesetting the solutions, e.g., if you want to draw a picture feel  
13 free to draw it on paper clearly, click a picture and include it in your solution. Do not spend  
14 undue time on typesetting solutions.
- 15 • You will see an entry of the form “HW 2 PDF” where you will upload the PDF of your  
16 solutions. You will also see entries like “HW 2 Problem 3 Code” where you will upload your  
17 solution for the respective problems.
- 18 • **For each programming problem/sub-problem, you should create a fresh .py file.** This  
19 file should contain **all** the code to reproduce the results of the problem/sub-problem, e.g., it  
20 should save the plot that is required (correctly with all the axes, title and legend) as a PDF  
21 in the same directory. You will upload the .py file as your solution for “HW 2 Problem  
22 3 Code”. Name your file as `pennkey_hw2_problem3.py`, e.g., I will name my code as  
23 `pratikac_hw2_problem3.py`. Note, we will not accept .ipynb files (i.e., Jupyter notebooks),  
24 you should only upload .py files. If you are using Google Colab to do your homework (and I  
25 suggest that you don’t...), you can export the notebook to a .py file.
- 26 • **This is very important.** Note that the instructors will download your code and execute it  
27 themselves, so your code should be such that it can be executed independently without any  
28 errors to create all output/plots required in the problem.

29 **Credit** The points for the problems add up to 110. You only need to solve for 100 points to get full  
30 credit, i.e., your final score will be  $\min(\text{your total points}, 100)$ .

---

- 31 **Problem 1 (25 points).** The torchvision library at <https://pytorch.org/vision/stable/models.html>  
 32 implements a number of popular architectures that you can use quickly in your code. In this problem,  
 33 you will take a deeper look at residual networks at  
 34 <https://github.com/pytorch/vision/blob/master/torchvision/models/resnet.py>. Understand how this  
 35 architecture is coded up.
- 36 (a) (2 points) Note down all the peculiarities that you notice in the code. E.g., which one is better:  
 37 using a batch-normalization layer before ReLU or after ReLU; what does this code do?
- 38 (b) (2 points) What do the calls “model.train()” and “model.eval()” do and where do you use them in  
 39 a typical training and validation code? Why did we not have them in HW 1 when we wrote our own  
 40 library for training deep networks?
- 41 (c) (3 points) Draw a rough picture of the Resnet-18 architecture and note down the number of  
 42 parameters in each layer; you will find it easier to write a function that computes the number of  
 43 parameters in each layer of the network.
- 44 (d) (3 points) Weight decay should not be applied to the biases of the different layers in the network,  
 45 argue why this is the case.
- 46 (e) (5 points) Write the code to iterate over all the network layers in Resnet-18 and separate out the  
 47 parameters in three groups: (i) batch-norm affine transform parameters, (ii) biases of convolutional  
 48 and fully-connected layers, and (iii) all the rest. There is no need to submit the code separately in this  
 49 case, since these are a few lines of code just copy them out into your PDF solutions.
- 50 (f) (10 points) Augmentations are increasingly becoming much more important for deep learn-  
 51 ing that what we initially believed. And therefore there are a number of sophisticated ways to  
 52 perform random augmentation. Go through the paper <https://arxiv.org/abs/1909.13719> (you don’t  
 53 have to read the details in order to do this question) and its implementation in torchvision at  
 54 <https://pytorch.org/vision/main/generated/torchvision.transforms.RandAugment.html>. The code of  
 55 this implementation is at [https://pytorch.org/vision/main/\\_modules/torchvision/transforms/autoaugment.html#RandAugment](https://pytorch.org/vision/main/_modules/torchvision/transforms/autoaugment.html#RandAugment).
- 56 Take any one image of your choice (you can use the same image of the astronaut as that of the  
 57 examples) and show the result of augmenting this image using the following 10 augmentations: (a)  
 58 ShearX, (b) ShearY, (c) TranslateX, (d) TranslateY, (e) Rotate, (d) Brightness, (e) Color, (f) Contrast,  
 59 (g) Sharpness, (h) Posterize, (i) Solarize and (j) Equalize. You will use existing functions from  
 60 torchvision to perform these augmentations (see the code of RandAugment linked here to understand  
 61 how to call each of them). Each of these augmentations has some parameters that you will have  
 62 to choose in order to run these corresponding augmentation functions. Any reasonable values are  
 63 okay, feel free to experiment. The goal of this problem is to understand how these augmentations  
 64 work. You don’t have to submit code in this case, pictures in the PDF correctly annotated with the  
 65 parameters of the augmentation that was used to create them are fine.
- 66 **Problem 2 (20 points).** Non-convex optimization problems are harder than convex optimization  
 67 problems. There are however a few special non-convex problems that are easy. We will look at one of  
 68 them here, namely unconstrained matrix factorization. Given a matrix  $X \in \mathbb{R}^{m \times n}$  we would like to  
 69 decompose it into two matrices of rank at most  $r$

$$X = AB$$

70 where  $A \in \mathbb{R}^{m \times r}$  and  $B \in \mathbb{R}^{r \times n}$ . Think of arranging all your data as columns of  $X$ . Columns  
 71 of the matrix  $A$  are like elements of a dictionary, they correspond to different patterns in the data  
 72 and are called atoms. The matrix  $B$  chooses which patterns to collect together in order to create a  
 73 particular datum, i.e., column of  $X$ . Solving for factors  $A, B$  is usually done with constraints, e.g.,  $B$   
 74 is typically forced to be a sparse matrix which enables regenerating data  $X$  using as few atoms as  
 75 possible. We will solve a simpler problem:

$$A^*, B^* = \underset{A \in \mathbb{R}^{m \times r}, B \in \mathbb{R}^{r \times n}}{\operatorname{argmin}} \|X - AB\|_F^2.$$

76 where  $\|\cdot\|_F$  denotes the Frobenius norm.

77 (a) (4 points) Why is the above problem not convex?

78 (b) (8 points) The global optimum for this loss function can be obtained easily in spite of it being  
 79 non-convex; find it. You may find it useful to write down the SVD of  $X$ .

80 (c) (8 points) Is the solution to the above optimization problem unique? Given one solution  $A^*, B^*$   
 81 name one way using which you can obtain another solution.

82 **Problem 3 (45 points).** Use Google Colab to train the neural network, but after debugging  
 83 everything on your laptop first. If you have a new laptop then you should also be able to do  
 84 everything on it without Colab. Neural networks are high-dimensional classifiers. While we may be  
 85 able to train and regularize SVMs to have a large margin between two classes, it is often difficult  
 86 to do the same for neural networks. A consequence of this is that, roughly speaking, all samples in  
 87 typical training datasets lie close to the decision boundaries after training. This makes it quite easy  
 88 to make minor perturbations to the input image—perturbations that are imperceptible to the human  
 89 eye—and send the sample across the decision boundary so as to cause the network to mis-predict.  
 90 We will synthesize such adversarial perturbations in this problem. You can read more about it at  
 91 <https://arxiv.org/abs/1412.6572>; however be wary of the heuristic generalizations in this paper that  
 92 are incorrect.

93 We will find the best adversarial perturbation to a given image  $x$  and its target  $y$ , this amounts to  
 94 solving the optimization problem

$$\max_{\|x' - x\|_\infty \leq \epsilon} \ell(x', y) \quad (1)$$

95 where  $x'$  is the *variable of optimization* and it is the adversarially perturbed image corresponding to  $x$ ,  
 96 the quantity  $\ell(x', y)$  is the loss computed on the image  $x'$  for the label  $y$ . We have chosen to make the  
 97 parameters of the classifier  $w$  implicit for sake of clarity. This optimization problem searches for all  
 98 images within an  $\epsilon$ -ball of the original image  $x$ . We will use the  $\ell_\infty$ -norm

$$\|x\|_\infty = \max_k |x_k|$$

99 in this problem. Let us perform the Taylor expansion of the objective Eq. (1)

$$\ell(x', y) = \ell(x, y) + \epsilon d^\top \nabla \ell(x, y) + \mathcal{O}(\epsilon^2);$$

100 here  $d = (x' - x)/\epsilon$ . We can now write an approximate problem for finding adversarial perturbations  
 101 as

$$\max_{\|d\|_\infty \leq 1} d^\top \nabla \ell(x, y).$$

102 Notice that the constraint implies that any element of the vector  $d$  can be perturbed by at most 1; there  
 103 is no limit on the number of elements perturbed. The value of  $d$  that maximizes this objective is  
 104 therefore the “signed gradient”

$$d_k = \frac{\nabla \ell(x, y)_k}{|\nabla \ell(x, y)_k|};$$

105 where  $|\cdot|$  denotes the element-wise absolute value. If  $\nabla \ell(x, y)_k < 0$ , the corresponding  $d_k < 0$  and  
 106 vice versa. The maximal objective is  $\|\nabla \ell(x, y)\|_1$ . The perturbation  $d$  is what we want to compute.

107 (a) (25 points) We will first train a convolutional neural network for this problem with all the bells  
 108 and whistles. You can use the code for the model and training provided on Canvas. This is a small  
 109 model with about 1.6M parameters. Train this model on the CIFAR-10 dataset for 100 epochs, you  
 110 should try to get a validation error below 12%. You should use a GPU on Colab for training; else  
 111 your code will be very slow. Roughly speaking, running for 100 epochs will take 1-2 hours, so be  
 112 patient. You can use data augmentation such as mirror flips and brightness and contrast changes to  
 113 improve your validation accuracy. Plot the training and validation losses and errors as a function of  
 114 the number of epochs. Some hints for choosing hyper-parameters:

- 115 • Learning rate of 0.1 for the first 40 epochs, then 0.01 for the next 40 epochs and then 0.001  
 116 for the final 20 epochs.
- 117 • Weight decay of  $10^{-3}$ .
- 118 • No need to perform data augmentation, although you can do so if you wish.
- 119 • Use SGD with Nesterov’s momentum of 0.9 to train the network.

120 Make sure you save the parameters of the network because we will need them for the next part.

121 (b) (10 points) We will next compute the backprop gradient of the loss with respect to the input. We  
 122 know that code of the form

```
123 ...
124
125 yh = net.forward(x)
126 loss = loss.forward(yh, y)
127
128 loss.backward()
```

130 computes the gradient of the loss with respect to the weights, i.e., it computes  $\bar{w}$  in our notation. You  
 131 can get the gradient  $\bar{x}$  very easily by adding the following line after `loss.backward()`.

```
132
133 dx = x.grad.data.clone()
```

135 Plot this gradient  $dx$  for a few input images which the network classifies correctly and also for a few  
 136 images which the network misclassifies. Comment on the similarities or the differences.

137 Note that each pixel of the RGB image  $x$  lies in  $[0, 255]$ , we will pick  $\epsilon = 8$ . Pick a particular  
 138 mini-batch  $\{x_1, \dots, x_\ell\}$  with  $\ell = 100$ . For every image in this mini-batch perform the “5-step  
 139 signed gradient attack”, i.e., perturb that image 5 times using the signed gradient, at each step you  
 140 feed in the perturbed image from the previous step and perturb it a bit more. Your pseudo-code will  
 141 look as follows.

```
142
143 xs, ys = mini-batch of inputs and targets
144 for x, y in zip(xs, ys):
```

```

145     for k in range(5):
146         # forward propagate x through the network
147         # backprop the loss
148         dx = ...
149         x += eps*sign(dx)
150         # record loss on the perturbed image
151     ell = loss(x, y)
152

```

153 Plot the loss on the perturbed images as a function of the number of steps in the attack averaged  
154 across your mini-batch.

155 (c) (10 points) Compute the accuracy of the network on 1-step perturbed images, i.e., for every image  
156 in the validation set, perturb the image using a 1-step attack and check the prediction of the network.  
157 How does this accuracy on adversarially perturbed images compare with the accuracy on the clean  
158 validation set?

159 (d) (0 points) The human brain also has a lot of neurons and is likely a high-dimensional classifier.  
160 Are humans susceptible to adversarial perturbations? Can you give examples of images that fool the  
161 human visual system? Are these “small”, i.e., is  $\|\epsilon/x\|_\infty$  small for these examples?

162 **Problem 4 (20 points).** Training of recurrent neural networks (RNNs) is often difficult because of  
163 the vanishing or exploding gradient problem. We will study this in a simple setting without any  
164 nonlinearities.

165 (a) (5 points) If the input to an RNN at the  $t^{\text{th}}$  timestep is  $x^t \in \mathbb{R}^d$  and the hidden vector is  $z^t \in \mathbb{R}^p$ ,  
166 the hidden vector  $z^{t+1}$  is given by

$$z^{t+1} = \sigma(w_x x^t + w_z z^t)$$

167 where  $w_x \in \mathbb{R}^{p \times d}$  and  $w_z \in \mathbb{R}^{p \times p}$  are weights and  $\sigma$  is a nonlinearity. If  $\sigma(z) = z$ , i.e., there is no  
168 nonlinearity, and  $w_x = 0$  then the update boils down to

$$z^{t+1} = w_z z^t.$$

169 Write down the back-propagation gradient for  $w_z$  if the loss function is only a function of the hidden  
170 vector at time  $T$ , i.e., the loss function is  $\ell(z^T)$ . Compute the conditions on the weight matrix  $w_z$   
171 under which the gradient explodes and vanishes.

172 (b) (2 points) Argue how the nonlinearity  $\sigma(\cdot)$  affects the exploding/vanishing gradients. Which  
173 nonlinearities are well-suited for training RNNs?

174 (c) (3 points) Updates to the weights are computed using SGD as

$$w_z^{\text{new}} = w_z^{\text{old}} - \eta \frac{d\ell}{dw_z}.$$

175 We would like to protect the weights  $w_z^{\text{new}}$  from blowing up even if the gradient  $\overline{w_z}$  explodes. Can you  
176 think of a way to do this? Similarly, can you modify these updates to handle the vanishing gradient  
177 problem?

178 (d) (5 points) Explain how an LSTM solves the problem of vanishing gradients.

179 (e) (5 points) Explain how a self-attention layer works for solving the problem of vanishing gradients.