

# Homework 3

CSCI 5525: Machine Learning

Due on Nov 7th 11:00am (before the class)

Please type in your info:

- **Name:** Zhongqi Zhao
- **Student ID:** 5229181
- **Email:** Zhao0680
- **Collaborators, and on which problems:**

**Homework Policy.** (1) You are encouraged to collaborate with your classmates on homework problems, but each person must write up the final solutions individually. You need to fill in above to specify which problems were a collaborative effort and with whom. (2) Regarding online resources, you should **not**:

- Google around for solutions to homework problems,
- Ask for help on online.
- Look up things/post on sites like Quora, StackExchange, etc.

**Submission.** Submit a PDF using this LaTeX template for written assignment part and submit .py files for all programming part. You should upload all the files on Canvas.

## Written Assignment

**Instruction.** For each problem, you are required to write down a full mathematical proof to establish the claim.

### Problem 1. VC Dimension.

(5 points) Prove this claim formally: If  $\mathcal{F}$  is finite, then  $\text{VCD}(\mathcal{F}) \leq \log(|\mathcal{F}|)$

**Your answer.** By VCD definition: The Vapnick-Chervonenkiss dimension,  $\text{VCD}(\mathcal{F})$ , of hypothesis space  $X$  defined over instance space  $X$  is the size of the largest finite subset of  $X$  shattered by  $\mathcal{F}$ . Assume  $\text{VCD}(\mathcal{F}) = n$  because of it is finite. Then  $\mathcal{F}$  need to shatter  $n$  data points with  $2^n$  label options. it can easily implies that  $|\mathcal{F}| \geq 2^n$ . Hence  $\text{VCD}(\mathcal{F}) \leq \log(|\mathcal{F}|)$ . The proof completes.

## Problem 2. Uniform convergence.

(10 points)

In this problem, we will prove a stronger generalization error bound for the binary classification that uses more information about the distribution. Let us say that  $P$  is a distribution over  $(X, Y)$  pairs where  $X \in \mathcal{X}$  and  $Y \in \{+1, -1\}$ . Let  $\mathcal{H} \subset \mathcal{X} \rightarrow \{+1, -1\}$  be a finite hypothesis class and let  $\ell$  denote the zero-one loss  $\ell(\hat{y}, y) = 1\{\hat{y} \neq y\}$ . Let  $R(h) = \mathbb{E} \ell(h(X), Y)$  denote the risk and let  $h^* = \min_{h \in \mathcal{H}} R(h)$ . Given  $n$  samples, let  $\hat{h}_n$  denote the empirical risk minimizer. Here, we want to prove a sample complexity bound of the form:

$$R(\hat{h}_n) - R(h^*) \leq c_1 \sqrt{\frac{R(h^*) \log(|\mathcal{H}|/\delta)}{n}} + c_2 \frac{\log(|\mathcal{H}|/\delta)}{n} \quad (1)$$

for constants  $c_1, c_2$ . If  $R(h^*)$  is small, this can be a much better bound than the usual excess risk bound. In particular, if  $R(h^*) = 0$ , this bound recovers the  $1/n$ -rate.

a) To prove the result, we will use Bernstein's inequality, which is a sharper concentration result.

**Theorem 1** (*Bernstein's inequality*). Let  $X_1, \dots, X_n$  be i.i.d real-valued random variables with mean zero, and such that  $|X_i| \leq M$  for all  $i$ . Then, for all  $t > 0$

$$\mathbb{P}\left[\sum_{i=1}^n X_i \geq t\right] \leq \exp\left(-\frac{t^2/2}{\sum_{i=1}^n \mathbb{E}[X_i^2] + Mt/3}\right) \quad (2)$$

Using this inequality, show that with probability at least  $1 - \delta$

$$|\bar{X}| \leq \sqrt{\left(\frac{2 \mathbb{E} X_1^2 \log(2/\delta)}{n}\right)} + \frac{2M \log(2/\delta)}{3n} \quad (3)$$

where  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  and  $X_i$ 's satisfy the conditions of Bernstein's inequality.

b) Using Eq. (3) and the union bound, show Eq. (1)

**Your answer.**

1. Proof:

The final result is about the absolute value of  $X$ . Replace the  $X_i$  with  $-X_i$  in equation (2). we could get

$$\mathbb{P}\left[\sum_{i=1}^n X_i \leq -t\right] \leq \exp\left(-\frac{t^2/2}{\sum_{i=1}^n \mathbb{E}[X_i^2] + Mt/3}\right) \quad (4)$$

take the union of both condition (2) and (4) we could get

$$\mathbb{P}\left[\left|\sum_{i=1}^n X_i\right| \geq t\right] \leq 2 \exp\left(-\frac{t^2/2}{\sum_{i=1}^n \mathbb{E}[X_i^2] + Mt/3}\right) \quad (5)$$

Then, replace the  $t$  with  $nt$

$$\mathbb{P}\left[\frac{1}{n} \left|\sum_{i=1}^n X_i\right| \geq t\right] \leq 2 \exp\left(-\frac{n^2 t^2/2}{\sum_{i=1}^n \mathbb{E}[X_i^2] + Mnt/3}\right) \quad (6)$$

Then we let  $2 \exp \left( -\frac{n^2 t^2 / 2}{\sum_{i=1}^n \mathbb{E}[X_i^2] + Mnt/3} \right) = \delta$  and make  $t = f(\delta)$  delta is probability.

$$n^2 t^2 + \frac{2}{3} \log\left(\frac{\delta}{2}\right) Mnt + 2 \log \frac{\delta}{2} \sum_{i=1}^n \mathbb{E}[X_i^2] = 0 \quad (7)$$

with the condition that  $t \geq 0$ , we need to find the solution of the quadratic equation with one unknown.

$$t = \frac{\frac{M}{3} + \sqrt{\frac{M^2}{9} + \frac{2n \mathbb{E}[X_i^2]}{\log \frac{2}{\delta}}}}{\frac{n}{\log \frac{2}{\delta}}} \quad (8)$$

$\therefore \sqrt{x^2 + y^2} \leq x + y$  when  $x, y \geq 0$   $\therefore$  it is easy to split the second product.

$$t \leq \sqrt{\left( \frac{2 \mathbb{E} X_1^2 \log(2/\delta)}{n} \right)} + \frac{2M \log(2/\delta)}{3n} \quad (9)$$

2.  $\forall h$ , we can normalize  $X_i$  with mean value 0, which means  $X_i = \mathbf{1}\{\hat{y}_i \neq y_i\} - \mathbb{E}[\mathbf{1}\{\hat{y}_i \neq y_i\}]$  we could get

$$\mathbb{E}[X_i] \leq \mathbb{E}[(\mathbf{1}\{\hat{y}_i \neq y_i\})^2] \leq \mathbb{E}[\mathbf{1}\{\hat{y}_i \neq y_i\}] = R(h) \quad (10)$$

end  $M \leq 1$ . By using the equation(3) and applying the union bound, which the core is replacing the  $\delta$  with  $\delta/\mathcal{H}$  in (3) as teacher pointed out in the class. we could get  $\forall h \in \mathcal{H}$

$$|\hat{R}(h) - R(h)| \leq \sqrt{\left( \frac{2R(h) \log(2|\mathcal{H}|/\delta)}{n} \right)} + \frac{2 \log(2|\mathcal{H}|/\delta)}{3n} \quad (11)$$

Since  $\hat{h}$  will minimize  $\hat{R}$  we have  $\hat{R}(\hat{h}) - \hat{R}(h^*) \leq 0$   $\therefore$

$$\begin{aligned} \hat{R}(\hat{h}) - \hat{R}(h^*) &\leq (R(\hat{h}) - \hat{R}(\hat{h})) + (\hat{R}(h^*) - R(h^*)) \\ &\leq |R(\hat{h}) - \hat{R}(\hat{h})| + |\hat{R}(h^*) - R(h^*)| \\ &\leq \sqrt{\frac{2R(\hat{h}) \log(2|\mathcal{H}|/\delta)}{n}} + \sqrt{\frac{2R(h^*) \log(2|\mathcal{H}|/\delta)}{n}} \\ &\quad + \frac{4 \log(2|\mathcal{H}|/\delta)}{3n} \\ \therefore (\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}) &\leq \sqrt{\frac{2(R(\hat{h}) - R(h^*)) \log(2|\mathcal{H}|/\delta)}{n}} + 2\sqrt{\frac{2R(h^*) \log(2|\mathcal{H}|/\delta)}{n}} \\ &\quad + \frac{4 \log(2|\mathcal{H}|/\delta)}{3n} \end{aligned}$$

$\therefore (\sqrt{xy} \leq 1/2(x+y)$  when  $x, y \geq 0$ ) we do it on the 1st item in previous equation

$$\begin{aligned} &\leq \frac{1}{2}(R(\hat{h}) - R(h^*)) + 2\sqrt{\frac{2R(h^*) \log(2|\mathcal{H}|/\delta)}{n}} \\ &\quad + \frac{7 \log(2|\mathcal{H}|/\delta)}{3n} \end{aligned}$$

Hence,  $\frac{1}{2}(R(\hat{h}) - R(h^*)) \leq 2\sqrt{\frac{2R(h^*) \log(2|\mathcal{H}|/\delta)}{n}} + \frac{7 \log(2|\mathcal{H}|/\delta)}{3n} = c_1 \sqrt{\frac{R(h^*) \log(|\mathcal{H}|/\delta)}{n}} + c_2 \frac{\log(|\mathcal{H}|/\delta)}{n}$   
proved

### Problem 3. Model selection.

(10 points) In problem 2, we proved the  $R(h^*)$  bound. The goal in this problem is to do this simultaneously while doing structural risk minimization. Specifically, given a family of hypothesis classes  $\mathcal{H}_1 \subset \mathcal{H}_2 \dots \subset \mathcal{H}_L$  of sizes  $N_1 \leq N_2 \leq \dots \leq N_L < \infty$ , a loss function bounded on  $[0,1]$  and a sample of size  $n$ , design an algorithm that guarantees

$$R(\hat{h}) \leq \min_{i \in [L]} \min_{h^* \in \mathcal{H}_i} \left\{ R(h^*) + c_1 \sqrt{\frac{R(h^*) \log(LN_i/\delta)}{n}} + c_2 \frac{\log(LN_i/\delta)}{n} \right\}$$

for  $n \geq 2$ . Your algorithm may use ERM (so need not be efficient) and your constants may vary.

You may find it useful to use the following *empirical Bernstein inequality*.

**Theorem 2** Let  $X_1, \dots, X_n$  be iid random variables from a distribution  $P$  supported on  $[0,1]$  and define the sample variance  $V_n = \frac{1}{n(n-1)} \sum_{1 \leq i < j \leq n} (X_i - X_j)^2$ . Then for any  $\delta \in (0,1)$  with probability at least  $1 - \delta$

$$\mathbb{E} X - \frac{1}{n} \sum_{i=1}^n X_i \leq \sqrt{\frac{2V_n \log(2/\delta)}{n}} + \frac{7 \log(2/\delta)}{3(n-1)} \quad (12)$$

**Your answer.**  $\forall h \in \mathcal{H}_k$ , we can do the same trick we did in Problem 2, let  $X_i = \mathbf{1}\{\hat{y}_i \neq y_i\} - \mathbb{E}[\mathbf{1}\{\hat{y}_i \neq y_i\}]$  The Var of the sample is

$$\begin{aligned} V_{N_k} &= \frac{1}{N_k(N_k-1)} \sum_{i,j} (\mathbf{1}\{\hat{y}_i \neq y_i\} - \mathbf{1}\{\hat{y}_j \neq y_j\})^2 \\ &\leq \frac{1}{N_k(N_k-1)} \sum_{i,j} (\mathbf{1}\{\hat{y}_i \neq y_i\})^2 + (\mathbf{1}\{\hat{y}_j \neq y_j\})^2 \\ &\leq \frac{1}{N_k(N_k-1)} \sum_{i,j} \mathbf{1}\{\hat{y}_i \neq y_i\} + \mathbf{1}\{\hat{y}_j \neq y_j\} \\ &= 2 \sum_i \mathbf{1}\{\hat{y}_i \neq y_i\} \\ &= 2\hat{R}(h) \end{aligned}$$

$\therefore$  theorem 2, we a inequality.

$$\begin{aligned} R(h) &\leq \hat{R}(h) + \sqrt{\frac{4\hat{R}(h) \log(2LN_k/\delta)}{n}} + \frac{7 \log(2LN_k/\delta)}{3(n-1)} \\ \hat{R}(h) &\leq R(h) + \sqrt{\frac{2R(h) \log(LN_k/\delta)}{n}} + \frac{3 \log(LN_k/\delta)}{3n} \end{aligned}$$

Then Assume  $\exists i$  such that it could minimize  $\hat{R}(\hat{h}_i) + \sqrt{\frac{4\hat{R}(\hat{h}_i) \log(2LN_k/\delta)}{n}} + \frac{7 \log(2LN_k/\delta)}{3(n-1)}$   
Then we has a good estimator  $\hat{h} = \hat{h}_i$ .  $\forall$  class  $m$

$$\begin{aligned} R(\hat{h}) &\leq \hat{R}(\hat{h}) + \sqrt{\frac{4\hat{R}(\hat{h}) \log(2L|\mathcal{H}_{\hat{h}}|/\delta)}{n}} + \frac{7 \log(2L|\mathcal{H}_{\hat{h}}|/\delta)}{3(n-1)} \\ &\leq \hat{R}(\hat{h}_m) + \sqrt{\frac{4\hat{R}(\hat{h}_m) \log(2L|\mathcal{H}_{\hat{h}_m}|/\delta)}{n}} + \frac{7 \log(2L|\mathcal{H}_{\hat{h}_m}|/\delta)}{3(n-1)} \end{aligned}$$

$\because \frac{1}{3(n-1)} \leq \frac{1}{n}$  when  $n > 1$  we can change the third product to  $\frac{1}{n}$  form. Then,we implies  $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$  with  $x = \hat{R}(\hat{h}_m)$   $y = Rest$  to 2nd product. And we implies we apply the 2nd inequality to the first product to 2nd product. In the end

$$R(\hat{h}) \leq \min_{i \in [L]} \min_{h^* \in \mathcal{H}_i} \left\{ R(h^*) + c_1 \sqrt{\frac{R(h^*) \log(L|\mathcal{H}_{\hat{h}_m}|/\delta)}{n}} + c_2 \frac{\log(L|\mathcal{H}_{\hat{h}_m}|/\delta)}{n} \right\} \quad (13)$$

#### Problem 4. Neural network.

(10 points) Consider a neural neural network with input  $x \in R^{d_i,1}$  and the number of classes being  $d_o$ . The network can be given by  $\hat{y} = \text{softmax}(W_2(\text{ReLU}(W_1x + b_1)) + b_2)$  where  $W_1 \in R^{d_1,d_i}$ ,  $b_1 \in R^{d_1,1}$ ,  $W_2 \in R^{d_o,d_1}$ ,  $b_2 \in R^{d_o,1}$  and Loss  $L = \text{cross\_entropy}(y, \hat{y})$

Cross entropy loss between  $y, \hat{y}$  is given by  $-\sum_i^m y_i \log(\hat{y}_i)$  and softmax clamps  $x$  to  $[0,1]$  range using  $\frac{e^{x_i}}{\sum_i e^{x_i}}$ .

Derive expressions for partial derivatives (be precise and clear)  $\frac{\partial L}{\partial W_2}, \frac{\partial L}{\partial b_2}, \frac{\partial L}{\partial W_1}, \frac{\partial L}{\partial b_1}$

**Your answer.** We have to make several marks to nodes in the NN.

$$\begin{aligned} f_1 &= W_1x + b_1 \\ g_1 &= \text{ReLU}(f_1) \\ f_2 &= W_2g_1 + b_2 \\ g_2 &= \text{Softmax}(f_2) \\ \hat{y} &= g_2 \end{aligned}$$

1.  $\frac{\partial L}{\partial W_2}$

$$\begin{aligned} \frac{\partial L}{\partial W_2} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial g_2} \frac{\partial g_2}{\partial f_2} \frac{\partial f_2}{\partial W_2} \\ \because \left[ \frac{\partial L}{\partial \hat{y}} \right]_i &= -\frac{y_i}{\hat{y}_i} \\ \left[ \frac{\partial g_2}{\partial f_2} \right]_{i,j} &= \begin{cases} \frac{e^{x_i} \sum_k e^{x_k} - e^{2x_i}}{(\sum_k e^{x_k})^2}, & i = j \\ -\frac{e^{x_i+x_j}}{(\sum_k e^{x_k})^2}, & i \neq j \end{cases} \\ \frac{\partial f_2}{\partial W_2} &= g_1^T \end{aligned}$$

Hence, The final result is

$$\frac{\partial L}{\partial W_2} = \left(\frac{\partial g_2}{\partial f_2}\right)^\top \frac{\partial L}{\partial \hat{y}} g_1^\top$$

2.  $\frac{\partial L}{\partial b_2}$

Hence, The final result is:

$$\begin{aligned} \frac{\partial L}{\partial b_2} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial g_2} \frac{\partial g_2}{\partial f_2} \frac{\partial f_2}{\partial b_2} \\ &= \left(\frac{\partial g_2}{\partial f_2}\right)^\top \frac{\partial L}{\partial \hat{y}} \mathbf{1} \end{aligned}$$

3.  $\frac{\partial L}{\partial W_1}$

$$\begin{aligned} \frac{\partial L}{\partial W_2} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial g_2} \frac{\partial g_2}{\partial f_2} \frac{\partial f_2}{\partial g_1} \frac{\partial g_1}{\partial f_1} \frac{\partial f_1}{\partial W_1} \\ \left[\frac{\partial g_1}{\partial f_1}\right]_{i,j} &= \begin{cases} 1, & i = j, \text{ and } (f_1)_i \geq 0 \\ 0, & \text{otherwise} \end{cases} \\ \frac{\partial f_1}{\partial W_1} &= x^\top \end{aligned}$$

Hence, The final result is

$$\frac{\partial L}{\partial W_2} = \left(\frac{\partial g_2}{\partial f_2} \frac{\partial f_2}{\partial g_1} \frac{\partial g_1}{\partial f_1}\right)^\top \frac{\partial L}{\partial \hat{y}} x^\top$$

4.  $\frac{\partial L}{\partial b_1}$

Hence, The final result is

$$\frac{\partial L}{\partial W_2} = \left(\frac{\partial g_2}{\partial f_2} \frac{\partial f_2}{\partial g_1} \frac{\partial g_1}{\partial f_1}\right)^\top \frac{\partial L}{\partial \hat{y}} \mathbf{1}$$

## Programming Assignment

**Instruction.** For each problem, you are required to report descriptions and results in the PDF and submit code as python file (.py) (as per the question).

- **Python** version: Python 3.
- Please follow PEP 8 style of writing your Python code for better readability in case you are wondering how to name functions & variables, put comments and indent your code
- **Packages allowed:** pytorch, numpy, pandas, matplotlib
- **Submission:** Submit report, description and explanation of results in the main PDF and code in .py files.
- Please PROPERLY COMMENT your code in order to have utmost readability

### Programming Common

This programming assignment focuses on implementing neural network for handwritten digits. This problem is about multi class classification and you will use MNIST dataset. You already have an idea of the dataset by now.

You will use pytorch to implement neural network. The implementation has to be for the CPU version only - No GPU's or MPI parallel programming is required for this assignment. Follow the installation instructions at <https://pytorch.org> in case you want to use your local machine, we recommend using conda environment.

Here is the pytorch tutorial. If you are new to pytorch, we recommend you to go through the tutorial here. [https://pytorch.org/tutorials/beginner/deep\\_learning\\_60min\\_blitz.html](https://pytorch.org/tutorials/beginner/deep_learning_60min_blitz.html)

### Problem 5. Neural Network Implementation.

(25 Points)

1. **(No Points)** You can directly download MNIST in pytorch using `torchvision.datasets.MNIST`. It will allow you to

```
trainset = torchvision.datasets.MNIST(root='./data', train=True, download=True, transform=transform)
testset = torchvision.datasets.MNIST(root='./data', train=False, download=True, transform=transform)
```
2. **(8 Points)** First, implement a multi-layer fully connected network:
  - Input: 1-channel input, size 28x28
  - Keep batch size as 32.
  - Fully connected layer 1: Input with bias; output - 128
  - ReLu Layer
  - Fully connected layer 2: input - 128; output - 10
  - Softmax layer
  - Use cross entropy as loss function

- Use SGD as optimizer.

Train using mini-batches of the given batch size. Plot loss and training accuracy for every epoch. At the end of the training, save your model with the name "mnist-fc". Load the saved model and report testing accuracy on the reloaded model. We should be able to reload your trained model and test.

**epoch:** An epoch tells you the number of times the learning algorithm sees the whole training data set. When it has seen all the training samples, you say 1 epoch is over and you start iterating in the next epoch.

3. **(10 Points)** Implement a convolutional neural network with the following specifications.

- Input: 1-channel input, size 28x28
- Keep batch size as 32.
- Convolution layer: Convolution kernel size is (3, 3) with stride as 1. Input channels - 1; Output channels - 20
- Max-pool: 2x2 max pool
- ReLu Layer
- Flatten input for feed to fully connected layers
- Fully connected layer 1: flattened input with bias; output - 128
- ReLu Layer
- Fully connected layer 2: input - 128; output - 10
- Softmax layer as above
- Use cross entropy as loss function
- Use SGD as optimizer.

Train using mini-batches of the given batch size. Plot loss and training accuracy for every epoch. At the end of the training, save your model with the name "mnist-cnn". Load the saved model and report testing accuracy on the reloaded model. We should be able to reload your trained model and test.

4. **(4 Points)** For the convolutional network implemented above, vary the batch sizes as [32, 64, 96, 128] and plot the convergence run time vs batch sizes.

5. **(3 Points)** "torch optim" provides many optimizers. Try the following optimizers in your last implementation - "SGD", "ADAM", "ADAGRAD". (SGD is already covered in the class, for ADAM and ADAGRAD refer to <https://arxiv.org/abs/1412.6980> and <http://www.jmlr.org/papers/volume12/duchi11a/duchi11a.pdf> respectively). Plot loss vs epochs for each optimizer. Briefly describe.

**Submission:** Submit all plots requested and explanation in latex PDF. Submit your program in a file named (hw3\_mnistfc.py and hw3\_mnistcnn.py).