

Homework 3

*Instructor: Vatsal Sharan**Due: March 27 by 11:59 pm PST*

A reminder on collaboration policy and academic integrity: Our goal is to maintain an optimal learning environment. You can discuss the homework problems at a high level with other groups, but you should not look at any other group's solutions. Trying to find solutions online or from any other sources for any homework or project is prohibited, will result in zero grade and will be reported. To prevent any future plagiarism, uploading any material from the course (your solutions, quizzes etc.) on the internet is prohibited, and any violations will also be reported. Please be considerate, and help us help everyone get the best out of this course.

Please remember the Student Conduct Code (Section 11.00 of the USC Student Guidebook). General principles of academic honesty include the concept of respect for the intellectual property of others, the expectation that individual work will be submitted unless otherwise allowed by an instructor, and the obligations both to protect one's own academic work from misuse by others as well as to avoid using another's work as one's own. All students are expected to understand and abide by these principles. Students will be referred to the Office of Student Judicial Affairs and Community Standards for further review, should there be any suspicion of academic dishonesty.

Total points: 99 points + Some bonus points to be released later

Notes on notation:

- Unless stated otherwise, scalars are denoted by small letter in normal font, vectors are denoted by small letters in bold font and matrices are denoted by capital letters in bold font.
- $\|\cdot\|$ means L2-norm unless specified otherwise, *i.e.*, $\|\cdot\| = \|\cdot\|_2$.

Instructions

We recommend that you use LaTeX to write up your homework solution. However, you can also scan handwritten notes. The homework will need to be submitted on Gradescope.

Theory-based Questions

Problem 1: Multi-class Perceptron (16pts)

Recall that a linear model for a multiclass classification problem with C classes is parameterized by C weight vectors $\mathbf{w}_1, \dots, \mathbf{w}_C \in \mathbb{R}^d$. In class, we derived the solution for multiclass logistic regression by minimizing the multiclass logistic loss. In this problem, you will derive the multiclass perceptron algorithm in a similar way. Specifically, the multiclass perceptron loss on a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times [C]$ is defined as

$$F(\mathbf{w}_1, \dots, \mathbf{w}_C) = \frac{1}{n} \sum_{i=1}^n F_i(\mathbf{w}_1, \dots, \mathbf{w}_C), \quad \text{where } F_i(\mathbf{w}_1, \dots, \mathbf{w}_C) = \max \left\{ 0, \max_{y \neq y_i} (\mathbf{w}_y^T \mathbf{x}_i - \mathbf{w}_{y_i}^T \mathbf{x}_i) \right\}.$$

1.1 (8pts) To optimize this loss function, we need to first derive its gradient. Specifically, for each $i \in [n]$ and $c \in [C]$, write down the partial derivative $\frac{\partial F_i}{\partial \mathbf{w}_c}$ (provide your reasoning). For simplicity, you can assume that for any i , $\mathbf{w}_1^T \mathbf{x}_i, \dots, \mathbf{w}_C^T \mathbf{x}_i$ are always C distinct values (so that there is no tie when taking the max over them, and consequently no non-differentiable points needed to be considered).

1.2 (4pts) Similarly to the binary case, multiclass perceptron is simply applying SGD with learning rate 1 to minimize the multiclass perceptron loss. Based on this information, fill in the missing details in the repeat-loop of the algorithm below (your solution cannot contain implicit quantities such as $\nabla F_i(\mathbf{w})$; instead, write down the exact formula based on your solution from the last question).

Algorithm 1: Multiclass Perceptron

- 1 **Input:** A training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- 2 **Initialization:** $\mathbf{w}_1 = \dots = \mathbf{w}_C = \mathbf{0}$
- 3 **Repeat:**

└

1.3 (4pts) At this point, you should find that the parameters $\mathbf{w}_1, \dots, \mathbf{w}_C$ computed by Multiclass Perceptron are always linear combinations of the training points $\mathbf{x}_1, \dots, \mathbf{x}_n$, that is, $\mathbf{w}_c = \sum_{i=1}^n \alpha_{c,i} \mathbf{x}_i$ for some coefficients $\alpha_{c,i}$. Just as we saw for kernelized linear regression, this means that we can kernelize Multiclass Perceptron as well. Consider some kernel function $k(\cdot, \cdot)$ with a corresponding feature map $\phi(\cdot)$. Fill in the missing details in the repeat-loop of the algorithm below that maintains and updates the coefficient $\alpha_{c,i}$ for all c and i for the kernelized solution $\mathbf{w}_c = \sum_{i=1}^n \alpha_{c,i} \phi(\mathbf{x}_i)$. To do this, try to see how the weights $\alpha_{c,i}$ should be updated such that $\mathbf{w}_c = \sum_{i=1}^n \alpha_{c,i} \phi(\mathbf{x}_i)$ is always the same as what one would get by running Algorithm 1 with \mathbf{x}_i replaced by $\phi(\mathbf{x}_i)$ for all i .

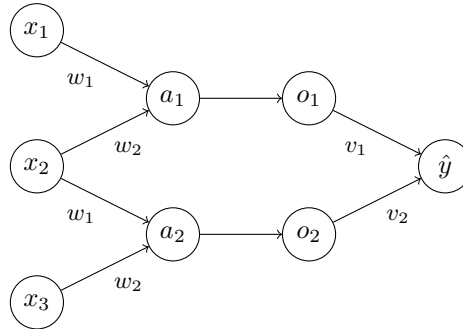
Algorithm 2: Multiclass Perceptron with kernel function $k(\cdot, \cdot)$

- 1 **Input:** A training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- 2 **Initialize:** $\alpha_{c,n} = 0$ for all $c \in [C]$ and $i \in [n]$
- 3 **Repeat:**

└

Problem 2: Backpropagation for CNN (18pts)

Consider the following mini convolutional neural net, where (x_1, x_2, x_3) is the input, followed by a convolution layer with a filter (w_1, w_2) , a ReLU layer, and a fully connected layer with weight (v_1, v_2) .



More concretely, the computation is specified by

$$\begin{aligned}
 a_1 &= x_1 w_1 + x_2 w_2 \\
 a_2 &= x_2 w_1 + x_3 w_2 \\
 o_1 &= \max\{0, a_1\} \\
 o_2 &= \max\{0, a_2\} \\
 \hat{y} &= o_1 v_1 + o_2 v_2
 \end{aligned}$$

For an example $(x, y) \in \mathbb{R}^3 \times \{-1, +1\}$, the logistic loss of the CNN is

$$\ell = \ln(1 + \exp(-y\hat{y})),$$

which is a function of the parameters of the network: w_1, w_2, v_1, v_2 .

2.1 (4pts) Write down $\frac{\partial \ell}{\partial v_1}$ and $\frac{\partial \ell}{\partial v_2}$ (show the intermediate steps that use chain rule). You can use the sigmoid function $\sigma(z) = \frac{1}{1+e^{-z}}$ to simplify your notation.

2.2 (6pts) Write down $\frac{\partial \ell}{\partial w_1}$ and $\frac{\partial \ell}{\partial w_2}$ (show the intermediate steps that use chain rule). The derivative of the ReLU function is $H(a) = \mathbb{I}[a > 0]$, which you can use directly in your answer.

2.3 (8pts) Using the derivations above, fill in the missing details of the repeat-loop of the Backpropagation algorithm below that is used to train this mini CNN.

Algorithm 3: Backpropagation for the above mini CNN

1 **Input:** A training set $(x_1, y_1), \dots, (x_n, y_n)$, learning rate η

2 **Initialize:** set w_1, w_2, v_1, v_2 randomly

3 **Repeat:**

4 randomly pick an example (x_i, y_i)

5 Forward propagation:

6 Backward propagation:

Programming-based Questions

As in previous homeworks, you need to have your coding environment setup for this part. We use python3 (version ≥ 3.7) in our programming-based questions. There are multiple ways you can install python3, for example:

- You can use **conda** to configure a python3 environment for all programming assignments.
- Alternatively, you can also use **virtualenv** to configure a python3 environment for all programming assignments

After you have a python3 environment, you will need to install the following python packages:

- numpy
- tqdm (a new package used in this HW, please install it to monitor training progress)
- matplotlib (for plotting figures)

Note: You are **not allowed** to use other packages such as *tensorflow*, *pytorch*, *keras*, *scikit-learn*, *scipy*, etc. for 3.1-3.2. If you have other package requests, please ask first before using them. You are **allowed** to use any packages for 3.3-3.5 and 4.

Download the files for the programming part from <https://vatsalsharan.github.io/spring24/hw3.zip>.

Problem 3: Fashion MNIST (65pts)

We will use a subset of the Fashion MNIST dataset (<https://github.com/zalandoresearch/fashion-mnist>) dataset in this part. The entire dataset contains 70,000 grayscale images in 10 categories. In this HW, we will use 10,000 examples for training, 2,000 examples for validation and 10,000 examples for testing. The datapoints in the Fashion MNIST dataset are images of individual articles of clothing at low resolution (28 by 28 pixels), as seen in Fig. 1:



Figure 1: Fashion-MNIST samples¹

Label	Class
0	T-shirt/top
1	Trouser
2	Pullover
3	Dress
4	Coat
5	Sandal
6	Shirt
7	Sneaker
8	Bag
9	Ankle boot

Figure 2: Fashion MNIST class name

Fashion MNIST is intended as a drop-in replacement for the classic MNIST dataset. MNIST is often regarded as the “Hello, World” of machine learning algorithms for computer vision. The MNIST dataset contains images of handwritten digits (0, 1, 2, etc.), and is similar in format to the Fashion MNIST dataset. This problem uses Fashion MNIST instead of MNIST to add some variety (and fashion), and also because it’s slightly more challenging than regular MNIST. Both datasets are relatively small and are often used to verify that an algorithm works as expected. They’re good starting points to test and debug code, and to explore ideas.

¹Figure from <https://medium.com/@ipylypenko/exploring-neural-networks-with-fashion-mnist-b0a8214b7b7b>

The images from Fashion MNIST are stored as 28x28 arrays, with pixel values ranging from 0 to 255. The labels are an array of integers, ranging from 0 to 9. The labels correspond to the item of clothing the image represents, as in Fig. 2. Each image is mapped to a single label.

Part I - Multi-Layer Perceptron (MLP) (38pts)

General instructions

- In this part you will implement neural nets. We provide the bootstrap code and you are expected to complete the functions.
- Do not import libraries other than those already imported in the original code.
- Please follow the data type annotations. You have to make the function's return values match the required type.
- Only make modifications in `{neural_networks.py}`.
- For parts where you will be running the code, we include the expected running time on Google Colab notebook in brackets (such as [5min]). If the running time of your code is much larger, it is likely that it has a bug. Your runtime maybe much faster though (indeed, the code ran faster on our laptops than on Colab).

3.1 Neural nets (24pts) In this task, you are asked to implement neural networks! You will later use this neural network to classify Fashion MNIST images into articles of clothing (0-9). The architecture of the neural network you will implement is based on the *multi-layer perceptron* (also known as MLP, just another term for fully connected feedforward networks we discussed in class). It is designed for a K -class classification problem. Let (\mathbf{x}, y) where $\mathbf{x} \in \mathbb{R}^d, y \in \{1, 2, \dots, K\}$ be a labeled instance. A MLP predicts the label by carrying out the following computations:

$$\begin{aligned} \text{input features : } \mathbf{x} &\in \mathbb{R}^d \\ \text{linear}^{(1)} : \mathbf{u} &= \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}, \mathbf{W}^{(1)} \in \mathbb{R}^{M \times d} \text{ and } \mathbf{b}^{(1)} \in \mathbb{R}^M \\ \text{relu : } \mathbf{h} &= \max\{0, \mathbf{u}\} = \begin{bmatrix} \max\{0, u_1\} \\ \vdots \\ \max\{0, u_M\} \end{bmatrix} \\ \text{linear}^{(2)} : \mathbf{a} &= \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)} \\ \text{softmax : } \mathbf{z} &= \begin{bmatrix} \frac{e^{a_1}}{\sum_k e^{a_k}} \\ \vdots \\ \frac{e^{a_K}}{\sum_k e^{a_k}} \end{bmatrix} \\ \text{predicted label : } \hat{y} &= \underset{k}{\operatorname{argmax}} z_k \end{aligned}$$

You might find it helpful to first go through what the following parts are asking, and then read through the `main` function to see how we will be using the functions you will write. Note that to better utilize the validation set, early-stopping is provided in the starter code. Early-stopping has a hyperparameter which we call *Patience*. If the patience parameter is set to k epochs, then training will terminate if there is no improvement in the validation accuracy for k epochs in a row. An epoch is just one pass through the training set.

3.1.1 Linear layer (6pts) First, you need to implement the linear layer of the MLP by implementing three python functions in `class linear_layer`. This layer has two parameters \mathbf{W} and \mathbf{b} .

- In the function `def __init__(self, input_D, output_D)`, you need to randomly initialize the entries of \mathbf{W} and \mathbf{b} with mean 0 and standard deviation 0.1 using `np.random.normal`. You also need to initialize the gradients to zeroes in the same function.
- In `def forward(self, X)`, implement the forward pass of this layer. Note that the input \mathbf{X} contains several examples, each of which needs to be passed through this layer. Try to use matrix operation instead of for loops to speed up your code. (As you might have realized with previous homeworks, matrix operations are much faster than for loops in Python. Please read this nice guide on *vectorized* operations to learn more: <https://www.askpython.com/python-modules/numpy/vectorization-numpy>).

- In `def backward(self, X, grad)`, implement the backward pass of this layer. Here, `grad` is the gradient with respect to the output of this layer, and you need to find and store the gradients with respect to `W` and `b`, and also find and return the gradients with respect to the input `X` of this layer.

3.1.2 ReLU (4pts) Next, you need to implement the ReLU activation by implementing 2 python functions in `class relu`. There are no parameters to be learned in this module.

- In `def forward(self, X)`, implement the forward pass of ReLU activation.
- In `def backward(self, X, grad)`, implement the backward pass of ReLU activation. Here, `X` is the gradient with respect to the output of this layer, and you need to return the gradient with respect to the input `X`.

3.1.3 Mini-batch stochastic gradient descent (2pts) Next, implement mini-batch version of stochastic gradient descent to learn the parameters of the neural network. Recall that for a general optimization problem with a parameter w , SGD iteratively computes $w_t = w_{t-1} - \eta \nabla F(w)$, where η is the step size, and $\nabla F(w)$ is the stochastic gradient (at w_{t-1} w.r.t loss function $F(\cdot)$). You need to complete `def miniBatchGradientDescent(model, _learning_rate)`, where we update the parameters of each layer. Note that this function is executed after the backward pass of each layer has been called and the gradients have been stored properly.

3.1.4 Backpropagation (4pts) Your network is almost ready to be trained. The only function you need to implement in this part is `backward_pass`, which calls the backward pass of each layer in the correct order and with the correct inputs.

3.1.5 Gradient checker (4pts) As you probably realized as you wrote the code for the previous parts, backpropagation is a subtle algorithm and it is easy to make mistakes. In this part, you will implement gradient checking to give you confidence that you are calculating the gradients of your model correctly. The idea behind gradient checking is as follows. As mentioned in the lecture, one naive way to approximate gradient is by

$$\frac{dF(w)}{dw} = \lim_{\epsilon \rightarrow 0} \frac{F(w + \epsilon) - F(w - \epsilon)}{2\epsilon}$$

We can use this relation to check that we are computing gradients correctly using backpropagation. For example, to check the gradient of the $(0,0)$ entry of $\mathbf{W}^{(1)}$, we apply a small $\epsilon = 1e-3$ perturbation to the $(0,0)$ entry of $\mathbf{W}^{(1)}$. Then, we compare the approximate gradient of the objective function with respect to the $(0,0)$ entry of $\mathbf{W}^{(1)}$ computed by the above equation, to the value obtained by backpropagation. These two values should be close to each other if backpropagation is implemented correctly. Please implement the `gradient_checker` function to do this. Note that to enable this function at runtime, you should turn on the `check_gradient` flag when you run the code later.

Apart from the gradient checker, we also provide a `magnitude_checker` function to help you debug the code. The gradient on any mini-batch provides an unbiased estimate of the true gradient, therefore we expect the gradient on a batch size of 50 to be similar in magnitude to the gradient on a batch size of 5,000. The `magnitude_checker` function (which is provided to you) checks that the ℓ_1 -norm of $\mathbf{W}^{(1)}$ is of the same magnitude on these two batch sizes.

3.1.6 Screenshots and plots (4pts) [40sec] After finishing all five modules above, run the below command to get a result file `MLP_lr0.01_b5.json`. Take a screenshot of your output log, like in `screenshot_run_b5.png` (included in the zip file). (2pts) Plot the training accuracy vs. epochs using `plot_train_process.py`. (2pts) For the magnitude check, you should expect the two numbers are of the same magnitude (however, there can be small differences).. For the four gradient checks, you should expect the number from the backpropagation to be the same as that from the approximation. For training accuracy, the curve should grow as the number of epochs increases. Some fluctuation in the middle is acceptable.

```
python neural_networks.py --minibatch_size 5 --check_gradient --check_magnitude
```

3.2 Batch size (14pts) [6min] We now explore the effect of different batch size on training. Run the command from the previous part with 5 different batch sizes: 1, 5, 50, 500, 5000.

3.2.1 Accuracy and time (4pts) Plot testing accuracy and training time *w.r.t.* batch size using `plot_batch.py`.

3.2.2 Epochs and gradient updates (4pts) What is the number of training epochs required to get the best model for each batch size? How many gradient updates are required get the best model for each batch size? (We define a gradient update to be a single stochastic gradient descent step, where the gradient may have been computed over some mini-batch).

3.2.3 Analysis (6pts) Answer the following question based on the previous plots and results:

- (i). Does smaller batch size guarantee faster training? Why do you think this is the case?
- (ii). Does larger batch size imply higher test accuracy? Why do you think this is the case?
- (iii). Does larger batch size imply less gradient updates to converge? Why do you think this is the case?

Part II - Explore on Colab (12pts)

The two-layer MLP from the previous part is our base model, and in the rest of this HW we will explore some variations on it. We will do this on a Colab notebook `fashion_mnist.ipynb`. To run the notebook, upload the `fashion_mnist.ipynb` to your USC Google Drive. Then, add Google Colab to your Google App by New → More → Connect more apps → Type in Google Colab and install it, as in Fig. 3 and Fig. 4. After that, you can run the notebook with your browser. The notebook is based on *TensorFlow*, a popular library for neural networks. TensorFlow enables us to build and train neural networks with a few simple lines of code. We recommend that you go through the TensorFlow code we provide and understand what's happening, but for the purpose of this homework you will mainly have to run the code and understand the results.

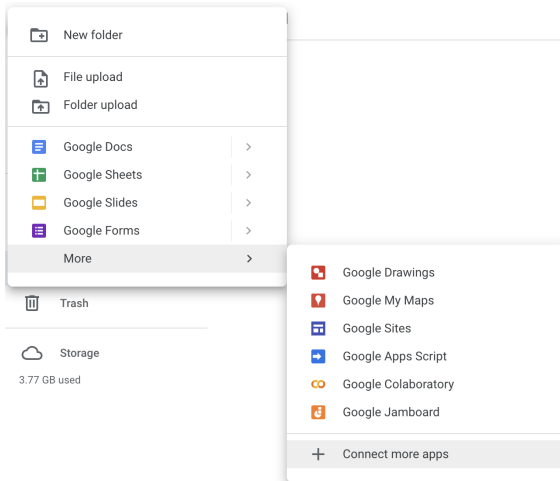


Figure 3: Colab install 1

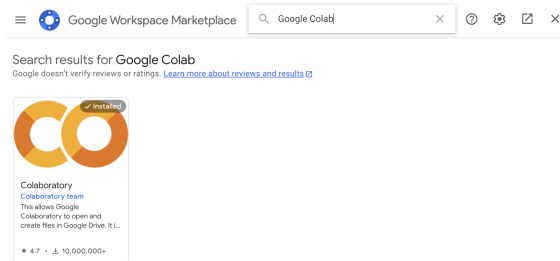


Figure 4: Colab install 2

3.3 The effect of early-stopping (6pts) [3.5min] We provided built-in early-stopping in the starter code, but did not explore the effect that it has on training. In this question, we will explore the effect of early-stopping. To do this, we evaluate the training, validation and test accuracy for each training epoch. We then plot the training, validation and test accuracy throughout the training process to see how early-stopping works. We train 30 epochs on a batch size of 5 without early-stopping.

- Plot the training accuracy vs. the number of epochs. On the same graph, plot the validation and test accuracy vs. the number of epochs. Also, mark the point on the test accuracy curve where we previously early-stopped. (We provide the plotting code, you only need to run it). (2pts)
- What is the trend of training and test accuracy after the early-stopped point? (2pts)
- Based on the plot, what do you think could go wrong if the patience parameter for early-stopping is too small? (Recall that if the patience parameter is set to k epochs, then training will terminate if there is no improvement in the validation accuracy for k epochs in a row.) (2pts)

3.4 SGD with momentum (6pts) [14min] We mentioned in class that adding a “momentum” term, which encourages the model to continue along the previous gradient direction helps the network to converge. Concretely, with an initial velocity $\mathbf{v} = \mathbf{0}$, we update the gradient by

$$\begin{aligned}\mathbf{v} &\leftarrow \alpha \mathbf{v} + \nabla F(\mathbf{w}) \\ \mathbf{w} &\leftarrow \mathbf{w} - \eta \mathbf{v}\end{aligned}$$

where η is the learning rate and α is the factor describing how much weight we put on the previous gradients. $\alpha = 0$ is equivalent to gradient update without momentum.

We provide the `sgd_with_momentum` function with argument α . In this question, we will explore how the momentum factor α effects training loss and test accuracy.

- (i). Run the code to call the function `sgd_with_momentum` with $\alpha = 0.1, 0.3, 0.5, 0.7, 0.9$. We store the training loss and test accuracy returned by the function. We visualize the training loss by plotting 5 curves for the 5 different values of α on the same graph. Each curve has the epoch on the x -axis and the training loss on the y -axis. (4pts)
- (ii). Based on this, what is a suitable value of α ? Therefore, how should training ideally rely on previous gradients for better convergence? (2pts)

Part III- Exploring Out-of-Distribution (OOD) generalization on Colab (15pts) [15min]

A major research direction at present is to ensure that our ML models not only do well on test data drawn from the same distribution as training data, but that they also do well when they get data from distributions different from the original distribution. This is known as Out-of-Distribution (OOD) generalization. In the Fashion MNIST dataset, the training and test set are drawn from the same distribution. Let's explore how our models do on test data coming from a slightly different distribution.

3.5 Translation and rotation (15pts) We modify the original test datapoints, by moving up the images of the test set by 4 pixels. By doing this, we create a new translated test set. Run the code to see the original and modified images. You will notice that to a human eye, the new test set is not any more difficult than the original test set.

- (i). But can our MLP model still do well on the new test set? What's the test accuracy on the two-layer MLP? (1pt)
- (ii). What if we try a different model architecture? For example, in class we saw that convolutional neural networks are good at dealing with image translation. We replace the first linear layer with a convolutional layer of 64 kernels with size 7×7 , followed by max-pooling. This can be done with just one or two lines of code in TensorFlow. Calculate the number of parameters of the CNN model and the original 2-layer MLP model (show your calculation), and verify that the CNN has fewer parameters than the MLP model. (3pts)
- (iii). Train the 2-layer CNN on the original training data and test it on both the original and the translated test sets. What is the in-domain test accuracy and the translated test accuracy of the CNN model? Can you provide some intuition behind these numbers? (2pts) [6min]
- (iv). Going one step further, we can make the CNN deeper by adding one more convolutional layer of 64×128 (64 input channels and 128 output channels) kernels with size 2×2 between the two existing layers, followed by max-pooling. Verify that the deeper 3-layer CNN model has fewer parameters than the 2-layer CNN model (show your calculation). (3pts)
- (v). What is the in-domain and the translated test accuracy of the deeper 3-layer CNN model? Provide some intuition on why it is better than the 2-layer CNN on the translated set. (3pts) [9min] You will notice that the deeper model takes some time to train. If we trained on GPUs instead of CPUs, training would be much faster. This is because GPUs are much more efficient at matrix computations, which is exactly what neural network training demands.

Next, we create 3 more OOD test sets by rotation. We rotate the images in the original test set by 90, 180 and 270 degrees.

- (vi). Provide the test accuracy of the 2-layer MLP model, the 2-layer CNN model and the 3-layer CNN model on the three rotation test sets. Are the 2-layer CNN and the 3-layer CNN still doing well? (3pts)

Deliverables for Problem 3: Code for 3.1 as a separate Python file `neural_networks.py`. Screenshot for 3.1. Plots for part 3.1, 3.2, 3.3 and 3.4. Number of epochs and gradient updates for part 3.2. Analysis and explanation for parts 3.2, 3.3, 3.4. For question 3.5, submit the test accuracy on in-domain and OOD test sets of the three models, the number of parameters of the three models, and the explanations we ask for.

Problem 4: Bonus Question on Transformers (10pts)

The Transformer architecture is important for ML because it revolutionized NLP by introducing the self-attention mechanism, enabling effective modeling of long-range dependencies in sequential data. This led to state-of-the-art performance in various NLP tasks such as machine translation, text generation, and sentiment analysis. In this problem, you will run a tiny decoder-only Transformer architecture and gain insight into its basic structure. We will do this using a Colab notebook, `transformer.ipynb`. To run the notebook, upload the `transformer.ipynb` to your USC Google Drive and follow the instructions in Problem 3 Part II. The notebook is based on PyTorch, a popular library for neural networks. Similar to TensorFlow, PyTorch enables us to build and train neural networks with just a few simple lines of code.

The main purpose of this assignment is to familiarize you with the PyTorch framework and the Transformer architecture. To do this, you will train a small decoder-only Transformer model with the tiny Shakespeare dataset and observe the effects of some hyperparameters.

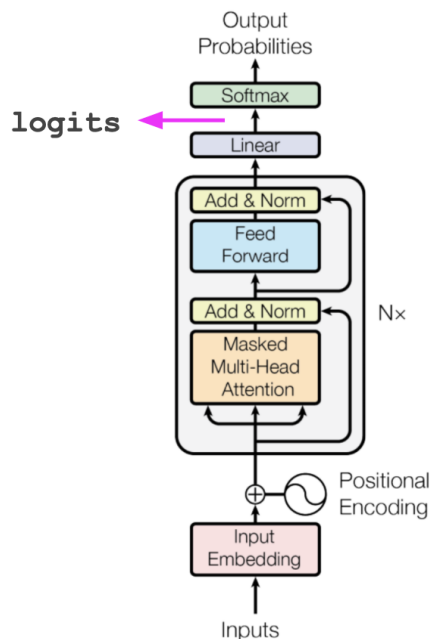


Figure 5: Decoder-only transformer architecture

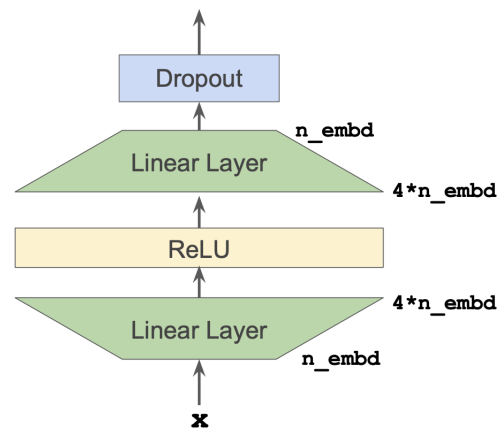


Figure 6: Feed Forward Module of Transformer

Transformer Structure First, you are expected to understand the code of the transformer which composed of multiple components. This section is not graded.

- (a) **Attention Layer:** We divide the attention layer into two main components: Head, and MultiHeadAttention. Inside Head module, the attention mechanism is implemented. We provide you the implemented version, however you are highly encouraged to understand the code. MultiHeadAttention module is responsible for combining the results of Heads. It first gets all results from Heads, then concatenates them. Finally, this is followed by a linear layer and dropout regularization.
- (b) **Feed Forward (FF) Module:** FF module in our network contains 2 linear layers with the ReLU activation function applied in between. After the second linear layer, dropout regularization is applied. FF module first increases the embedding dimension 4 times via the first linear layer. Then decreases it to its original embedding dimension via the second linear layer. You can see the diagram in Figure 6.
- (c) **Transformer Block:** A block contains 4 components (See also Figure 5, grey shaded block): Self-Attention module, Feed Forward module and 2 Layer-Normalization modules.

- (d) Model: `BigramLanguageModel` module brings everything together and creates the whole transformer architecture (Figure 5). We also calculate the loss inside this function. We also provide the implementation of the `generate` function.

4.1 Training the Model (10pts) Now, we will train our model. In this section, you are not expected to implement anything.

- (i). Before training the model, make it generate some sequence. Then, you will train the model with tiny Shakespeare dataset. We provided you the initial hyper-parameters. With Colab, the training session should take around 10 minutes. After training the model, make it generate some sequence again. Compare it with the first generation. Also, compare the generation with training data. Do they look similar? Explain the reason why the generation is not like a Wikipedia article, for example. (3pts)
- (ii). Provide the loss curve for train and validation during training. (1pt)
- (iii). For different values of block-size (sequence-length/context-length) [2, 16, 64], train the model from scratch. Provide the train-validation loss curves and model's generations for different sequence-lengths. Discuss how the sequence-length changes the generation of the model. (3pts)
- (iv). For different number of heads [2, 4, 8], train the model from scratch. Provide the train-validation loss curves and model's generations. Discuss how the number of heads changes the train-validation loss curves. (3pts)

Deliverables for Problem 4: In general you should include all explanations we asked for in the question. For 4.1 (i), discussion on generations and screenshots of the generations. Plots for 4.1 (ii), (iii) (iv). Discussion and analysis for 4.1 (iii) and (iv). Screenshots of the generations for 4.1 (iii).