

# Homework 1 Bonus

Adam, AdamW and Dropout

11-785: Introduction to Deep Learning (Fall 2023)

## Start Here

- **Collaboration policy:**

- You are expected to comply with the [University Policy on Academic Integrity and Plagiarism](#).
- You are allowed to help your friends debug
- You are allowed to look at your friends code
- You are allowed to copy math equations from any source that are not in code form
- You are not allowed to type code for your friend
- You are not allowed to look at your friends code while typing your solution
- You are not allowed to copy and paste solutions off the internet
- You are not allowed to import pre-built or pre-trained models
- You can share ideas but not code, you must submit your own code. All submitted code will be compared against all code submitted this semester and in previous semesters using [MOSS](#).

We encourage you to meet regularly with your study group to discuss and work on the homework. You will not only learn more, you will also be more efficient that way. However, as noted above, the actual code used to obtain the final submission must be entirely your own.

- **Directions:**

- You are required to do this assignment in the Python (version 3) programming language. Do not use any auto-differentiation toolboxes (PyTorch, TensorFlow, Keras, etc) - you are only permitted and recommended to vectorize your computation using the Numpy library.
- We recommend that you look through all of the problems before attempting the first problem. However we do recommend you complete the problems in order, as the difficulty increases, and questions often rely on the completion of previous questions.

- **Overview:**

- MyTorch
- Adam and AdamW optimizers
- Dropout

## Homework objective

After this homework, you would ideally have learned:

- How to write code to implement different optimizers from scratch
  - How to implement Adam
  - How to implement AdamW
- How to write code to implement dropout from scratch
  - How to implement Forward Propagation with Dropout
  - How to implement Back Propagation with Dropout

# Contents

1	MyTorch	4
2	Setup and Submission	4
3	ADAM [5 points]	5
4	AdamW [5 points]	6
5	Dropout [10 points]	7

# 1 MyTorch

The culmination of all of the Homework Part 1's will be your own custom deep learning library, which we are calling MyTorch. It will act similar to other deep learning libraries like PyTorch or Tensorflow. The files in your homework are structured in such a way that you can easily import and reuse modules of code for your subsequent homeworks.

## 2 Setup and Submission

- **Extract** the handout *hw1p1\_bonus.tar* by running the following command in the same directory<sup>1</sup>

```
tar -xvf hw1p1_bonus_handout.tar
```

This will create a directory called HW1P1 with the following file structure: <sup>2</sup>

```
HW1P1_bonus
├── mytorch
│   ├── models
│   │   └── mlp.py (Copy your file from HW1P1)
│   ├── nn
│   │   ├── activation.py (Copy your file from HW1P1)
│   │   ├── batchnorm.py (Copy your file from HW1P1)
│   │   ├── linear.py (Copy your file from HW1P1)
│   │   ├── loss.py (Copy your file from HW1P1)
│   │   └── dropout.py
│   └── optim
│       ├── adam.py
│       ├── adamW.py
│       └── sgd.py (Copy your file from HW1P1)
└── hw1p1_bonus_autograder.py
```

- **Install** Activate the same conda environment created for HW1P1:

```
conda activate idlF23
```

- **Autograde** your code by executing the following in terminal from top level directory:

```
python3 hw1p1_bonus_autograder.py
```

- **Hand-in** your code by running the following command from the top level directory, then **SUBMIT** the created *handin.tar* file to autolab:

```
tar cvf handin.tar mytorch
```

- **DO**

- We strongly recommend that you review the [ADAM](#) and [Dropout](#) papers.

- **DO NOT**

- Import other external libraries other than **numpy** in your submission, as extra packages that do not exist in autolab will cause submission failures.
  - Add, move, or remove any files or change any filenames.

---

<sup>1</sup>The handout might have an extension like *handout.tar.112*. In such a case, you will have to first rename downloaded file as *handout.tar* by removing the *.112* extension and then untar the file.

<sup>2</sup>For using code from Homework 1, ensure that you received all 100 autograded points.

### 3 ADAM [5 points]

Adam is a per-parameter adaptive optimizer that considers both the first and second moment of the current gradient. Implement the `adam` class in `mytorch/optim/adam.py`.

At any time step  $t$ , Adam keeps moving averages of the **biased first moment**  $m_t$  for each parameter (in the case of a linear layer we have parameters  $W$  and  $b$ ) and the **biased second moment**  $v_t$  for each parameter.  $t$  is initialized as 0 and  $m_t, v_t$  are initialized as 0 tensors. They are updated via:

$$m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$$

$$v_t = \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$$

where  $g_t$  are the current gradients for the parameters.  $\beta_1$  and  $\beta_2$  are hyper-parameters which control the exponential decay rates of these moving averages, their default values are  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ .

We call  $m_t$  and  $v_t$  as **biased** moments, because they are initialized as 0 tensors. As such, they are biased towards 0, especially in earlier steps. As such, Adam corrects this with:

$$\hat{m}_t = m_t / (1 - \beta_1^t)$$

$$\hat{v}_t = v_t / (1 - \beta_2^t)$$

Lastly, the parameters are updated via

$$\theta_t = \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$$

where  $\alpha$  is the learning rate. Recall that we intuitively divide  $\hat{m}_t$  by  $\sqrt{\hat{v}_t}$  in the last step to normalize out the magnitude of the running average gradient and to incorporate the second moment estimate. Note:  $\theta_t$  here represents relevant parameters, such as weight and bias in the case of a linear layer.

For more detailed explanations, we recommend that you reference the original [paper](#).

You need to keep a running estimate for some parameters related to  $W$  and  $b$  of the linear layer.

## 4 AdamW [5 points]

AdamW is an optimizer which uses weight decay regularization with Adam. Implement the `adamW` class in `mytorch/optim/adamW.py`.

If you implemented Adam, then the only additional parameter in AdamW is the weight decay, a regularization method which reduces the chance of overfitting. In this, we reduce the network parameter a portion of the model parameter at each time step.

$$\begin{aligned}W_t &= W_t - W_{t-1} * \alpha * \lambda \\b_t &= b_t - b_{t-1} * \alpha * \lambda\end{aligned}$$

Where  $W_t$  and  $b_t$  are your parameters after the standard Adam update in iteration  $t$ ,  $\lambda$  is weight decay and  $\alpha$  is the learning rate. Alternatively, you can *first* perform the following two weight decay updates:

$$\begin{aligned}W_t &= W_{t-1} - W_{t-1} * \alpha * \lambda \\b_t &= b_{t-1} - b_{t-1} * \alpha * \lambda\end{aligned}$$

and then, add the standard Adam update to  $W_t$  and  $b_t$  obtained above using  $W_{t-1}$  and  $b_{t-1}$ . (**NOTE:** Observe the subscripts indicating the iteration number in the two presented ways).

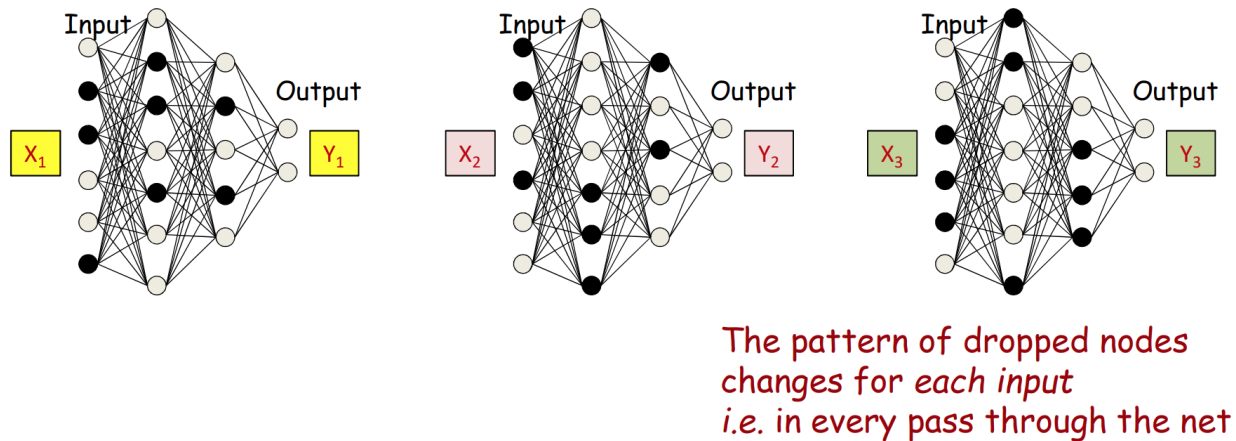
**Tips:** for using AdamW for homework part 2s: Optimal weight decay depends on the total number of batch passes/weight updates. The paper's empirical analysis of SGD and Adam suggests that the larger the runtime/number of batch passes to be performed, the smaller the optimal weight decay. If the value of weight decay is too high, a lot of entries will have values close to 0, and we will get a sparse model that fails to capture a lot of information.

For more detailed explanations, we recommend that you reference the original [paper](#)

## 5 Dropout [10 points]

Dropout is a regularization method that approximates ensemble learning of networks by randomly "turning off" neurons in a network during training. Implement the `Dropout` class in `mytorch/nn/modules/dropout.py`.

For every input, the neurons that are "turned off" are randomly chosen via the dropout rate  $p$ . The probability of zero-ing out a neuron output is then  $p$ .



We can implement this by generating and applying a binary mask to the output tensor of a layer. As dropout zeros out a portion of the tensor, we need to re-scale the remaining numbers so the total expected "intensity" of the output is same as in testing, where we don't apply dropout.

For more detailed explanations, we recommend that you reference the [paper](#).

### Implementation Notes:

- You should use `np.random.binomial`
- You should scale during training and not during testing