

OSU CSE 5523 (2021 Spring)

Homework #5: Programming Set

Due Date: April 23th (23:59 ET), 2021
You can submit by April 25th (23:59 ET), 2021 with no penalty.

1 Starting Point

Scripts: From the GitHub repository, you will find:

- python script `dnn_misc.py`, which you will be amending by adding code for questions in Sect. 3.
- Various other python scripts: `dnn_mlp.py`, `dnn_mlp_nonlinear.py`, `dnn_cnn.py`, `hw5_dnn_check.py`, `dnn_im2col.py` and `data_loader.py`, which you are **not allowed** to modify.
- Various scripts: `q33.sh`, `q34.sh`, `q35.sh`, `q36.sh`, `q37.sh`, and `q38.sh`; you will use these to generate output files.

Data: You will use `mnist_subset.json`.

Submission: Please submit a single .zip file named `HW_5_name.number.zip` (e.g., `HW_5_chao.209.zip`) to Carmen. The following will constitute your submission:

- The python script `dnn_misc.py`, amended with the code you added for Sect. 3.
- Six .json files, which will be the output of the six scripts above. We reserve the right to run your code to regenerate these files, but you are expected to include them.

```
MLP_lr0.01_m0.0_w0.0_d0.0.json
MLP_lr0.01_m0.0_w0.0_d0.5.json
MLP_lr0.01_m0.0_w0.0_d0.95.json
LR_lr0.01_m0.0_w0.0_d0.0.json
CNN_lr0.01_m0.0_w0.0_d0.5.json
CNN_lr0.01_m0.9_w0.0_d0.5.json
```

Collaboration: You may discuss with your classmates at a high level. However, you need to complete your own code and submit it by yourself. Please consult the syllabus for what is and is not acceptable collaboration.

2 High-level descriptions

2.1 Dataset

We will use `mnist_subset` (images of handwritten digits from 0 to 9). The dataset is stored in a JSON-formatted file `mnist_subset.json`. You can access its training, validation, and test splits using the keys ‘train’, ‘valid’, and ‘test’, respectively. For example, suppose we load `mnist_subset.json` to the variable `x`. Then, `x['train']` refers to the training set of `mnist_subset`. This set is a list with two elements: `x['train'][0]` containing the features of size N (samples) $\times D$ (dimension of features), and `x['train'][1]` containing the corresponding labels of size N .

2.2 Tasks

You will implement neural networks (Sect. 3). Specifically, you will

- finish the implementation of all python functions in our template code.
- run your code by calling the specified scripts to generate output files.
- submit (1) all `*.py` files, and (2) all `*.json` files that **you have amended or created**.

In the next subsection, we will provide a **high-level** checklist of what you need to do. You are not responsible for loading/pre-processing data; we have done that for you. For specific instructions, please refer to text in Sect. 3, as well as corresponding python scripts.

2.2.1 Neural networks

Preparation: Read Sect. 3 as well as `dnn_mlp.py` and `dnn_cnn.py`.

Coding: In `dnn_misc.py`, finish implementing

- forward and backward functions in class `linear_layer`
- forward and backward functions in class `relu`
- backward function in class `dropout` (before that, please read forward function).

Refer to `dnn_misc.py` and Sect. 3 for more information.

Running your code Run the scripts `q33.sh`, `q34.sh`, `q35.sh`, `q36.sh`, `q37.sh`, and `q38.sh` after you finish your implementation. This will generate, respectively,

`MLP_lr0.01_m0.0_w0.0_d0.0.json`

`MLP_lr0.01_m0.0_w0.0_d0.5.json`

`MLP_lr0.01_m0.0_w0.0_d0.95.json`

LR_lr0.01_m0.0_w0.0_d0.0.json
CNN_lr0.01_m0.0_w0.0_d0.5.json
CNN_lr0.01_m0.9_w0.0_d0.5.json

What to submit Submit `dnn_misc.py` and the above six `.json` files.

2.3 Cautions

Please do not import packages that are not listed in the provided code. Follow the instructions in each section strictly to code up your solutions. **Do not change the output format. Do not modify the code unless we instruct you to do so.** A homework solution that does not match the provided setup, such as format, name, initializations, etc., **will not** be graded. It is your responsibility to **make sure that your code runs with the provided commands and scripts.**

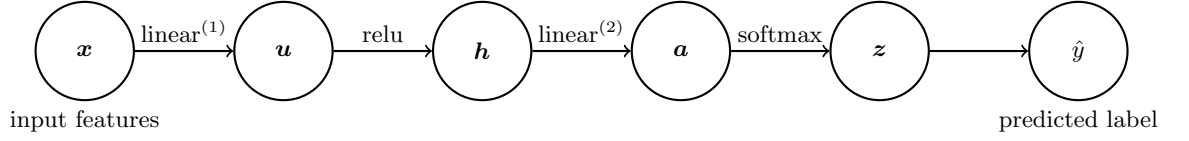


Figure 1: A diagram of a multi-layer perceptron (MLP). *The edges mean mathematical operations (modules), and the circles mean variables.* The term *relu* stands for rectified linear units.

3 Neural networks: multi-layer perceptrons (MLPs) and convolutional neural networks (CNNs)

3.1 Background

In recent years, neural networks have been one of the most powerful machine learning models. Many toolboxes/platforms (e.g., TensorFlow, PyTorch, Torch, Theano, MXNet, Caffe, CNTK) are publicly available for efficiently constructing and training neural networks. The core idea of these toolboxes is to treat a neural network as a combination of *data transformation (or mathematical operation) modules*.

For example, in Fig. 1 we provide a diagram of a multi-layer perceptron (MLP) for a K -class classification problem. The edges correspond to modules and the circles correspond to variables. Let $(\mathbf{x} \in \mathbb{R}^D, y \in \{1, 2, \dots, K\})$ be a labeled instance, such an MLP performs the following computations

$$\text{input features : } \mathbf{x} \in \mathbb{R}^D \quad (1)$$

$$\text{linear}^{(1)} : \mathbf{u} = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)} \quad , \mathbf{W}^{(1)} \in \mathbb{R}^{M \times D} \text{ and } \mathbf{b}^{(1)} \in \mathbb{R}^M \quad (2)$$

$$\text{relu : } \mathbf{h} = \max\{0, \mathbf{u}\} = \begin{bmatrix} \max\{0, u[1]\} \\ \vdots \\ \max\{0, u[M]\} \end{bmatrix} \quad (3)$$

$$\text{linear}^{(2)} : \mathbf{a} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)} \quad , \mathbf{W}^{(2)} \in \mathbb{R}^{K \times M} \text{ and } \mathbf{b}^{(2)} \in \mathbb{R}^K \quad (4)$$

$$\text{softmax : } \mathbf{z} = \begin{bmatrix} \frac{e^{a[1]}}{\sum_k e^{a[k]}} \\ \vdots \\ \frac{e^{a[K]}}{\sum_k e^{a[k]}} \end{bmatrix} \quad (5)$$

$$\text{predicted label : } \hat{y} = \arg \max_k z[k]. \quad (6)$$

For a K -class classification problem, one popular loss function for training (i.e., to learn $\mathbf{W}^{(1)}$, $\mathbf{W}^{(2)}$, $\mathbf{b}^{(1)}$, $\mathbf{b}^{(2)}$) is the cross-entropy loss,

$$\ell = - \sum_k \mathbf{1}[y == k] \log z[k], \quad (7)$$

$$\text{where } \mathbf{1}[\text{True}] = 1; \text{ otherwise, } 0. \quad (8)$$

For ease of notation, let us define the one-hot (i.e., 1-of- K) encoding

$$\mathbf{y} \in \mathbb{R}^K \text{ and } y[k] = \begin{cases} 1, & \text{if } y = k, \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

so that

$$\ell = - \sum_k y[k] \log z[k] = -\mathbf{y}^T \begin{bmatrix} \log z[1] \\ \vdots \\ \log z[K] \end{bmatrix} = -\mathbf{y}^T \log \mathbf{z}. \quad (10)$$

We can then perform error-back-propagation, a way to compute partial derivatives (or gradients) w.r.t the parameters of a neural network, and use gradient-based optimization to learn the parameters.

3.2 Modules

Now we will provide more information on modules for this assignment. Each module has its own parameters (but note that a module may have no parameters). Moreover, each module can perform a *forward pass* and a *backward pass*. The forward pass performs the computation of the module, given the input to the module. The backward pass computes the partial derivatives of the loss function w.r.t. the input and parameters, given the partial derivatives of the loss function w.r.t. the output of the module. Consider a module `<module_name>`. Let `<module_name>.forward` and `<module_name>.backward` be its forward and backward passes, respectively.

For example, the linear module may be defined as follows.

$$\begin{aligned} \text{forward pass:} \quad \mathbf{u} &= \text{linear}^{(1)}.forward(\mathbf{x}) = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}, \\ \text{where } \mathbf{W}^{(1)} \text{ and } \mathbf{b}^{(1)} &\text{ are its parameters.} \end{aligned} \quad (11)$$

$$\text{backward pass:} \quad \left[\frac{\partial \ell}{\partial \mathbf{x}}, \frac{\partial \ell}{\partial \mathbf{W}^{(1)}}, \frac{\partial \ell}{\partial \mathbf{b}^{(1)}} \right] = \text{linear}^{(1)}.backward(\mathbf{x}, \frac{\partial \ell}{\partial \mathbf{u}}). \quad (12)$$

Let us assume that we have implemented all the desired modules. Then, getting $\hat{\mathbf{y}}$ for \mathbf{x} is equivalent to running the forward pass of each module in order, given \mathbf{x} . All the intermediated variables (i.e., \mathbf{u} , \mathbf{h} , etc.) will all be computed along the forward pass. Similarly, getting the partial derivatives of the loss function w.r.t. the parameters is equivalent to running the backward pass of each module in a reverse order, given $\frac{\partial \ell}{\partial \mathbf{z}}$.

In this question, we provide a Python environment based on the idea of modules. Every module is defined as a class, so you can create multiple modules of the same functionality by creating multiple object instances of the same class. Your work is to finish the implementation of several modules, where these modules are elements of a multi-layer perceptron (MLP) or a convolutional neural network (CNN). We will apply these models to the 10-class classification problem in mnist. We will train the models using stochastic gradient descent with mini-batch, and explore how different hyperparameters of optimizers and regularization techniques affect training and validation accuracies over training epochs. For deeper understanding, check out, the seminal work of Yann LeCun et al. “Gradient-based learning applied to document recognition,” 1998.

We give a specific example below. Suppose that, at iteration t , you sample a mini-batch of N examples $\{(\mathbf{x}_i \in \mathbb{R}^D, \mathbf{y}_i \in \mathbb{R}^K)\}_{i=1}^N$ from the training set ($K = 10$). Then, the loss of such a mini-batch given by Fig. 1 is

$$\ell_{mb} = \frac{1}{N} \sum_{i=1}^N \ell(\text{softmax.forward}(\text{linear}^{(2)}.forward(\text{relu.forward}(\text{linear}^{(1)}.forward(\mathbf{x}_i))), \mathbf{y}_i) \quad (13)$$

$$= \frac{1}{N} \sum_{i=1}^N \ell(\text{softmax.forward}(\text{linear}^{(2)}.forward(\text{relu.forward}(\mathbf{u}_i))), \mathbf{y}_i) \quad (14)$$

$$= \dots \quad (15)$$

$$= \frac{1}{N} \sum_{i=1}^N \ell(\text{softmax.forward}(\mathbf{a}_i), \mathbf{y}_i) \quad (16)$$

$$= -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_i[k] \log z_i[k]. \quad (17)$$

That is, in the forward pass, we can perform the computation of a certain module to all the N input examples, and then pass the N output examples to the next module. This is the same case for the backward pass. For example, according to Fig. 1, given the partial derivatives of the loss w.r.t. $\{\mathbf{a}_i\}_{i=1}^N$

$$\frac{\partial \ell_{mb}}{\partial \{\mathbf{a}_i\}_{i=1}^N} = \begin{bmatrix} \left(\frac{\partial \ell_{mb}}{\partial \mathbf{a}_1}\right)^T \\ \left(\frac{\partial \ell_{mb}}{\partial \mathbf{a}_2}\right)^T \\ \vdots \\ \left(\frac{\partial \ell_{mb}}{\partial \mathbf{a}_{N-1}}\right)^T \\ \left(\frac{\partial \ell_{mb}}{\partial \mathbf{a}_N}\right)^T \end{bmatrix} \in \mathbb{R}^{N \times K}, \quad (18)$$

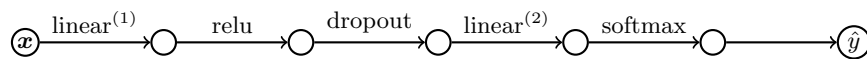


Figure 2: The diagram of the MLP implemented in `dnn_mlp.py`. The circles mean variables and edges mean modules.

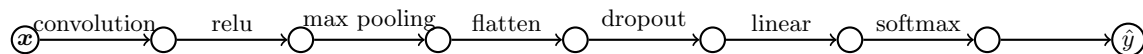


Figure 3: The diagram of the CNN implemented in `dnn_cnn.py`. The circles correspond to variables and edges correspond to modules. Note that the input to CNN may not be a vector (e.g., in `dnn_cnn.py` it is an image, which can be represented as a 3-dimensional tensor). The flatten layer is to reshape its input into vector.

`linear(2).backward` will compute $\frac{\partial \ell_{mb}}{\partial \{\mathbf{h}_i\}_{i=1}^N}$ and pass it back to `relu.backward`.

We note that, $\mathbb{R}^{N \times K}$ means that $\frac{\partial \ell_{mb}}{\partial \{\mathbf{a}_i\}_{i=1}^N}$ is represented in the row-wise fashion for Python code. For your problem set, you will still use the column-wise fashion where each column is for a data instance.

3.3 Preparation

Q3.1 Please read through `dnn_mlp.py` and `dnn_cnn.py`. Both files will use modules defined in `dnn_misc.py` (which you will modify). Your work is to understand how modules are created, how they are linked to perform the forward and backward passes, and how parameters are updated based on gradients (and momentum). The architectures of the MLP and CNN defined in `dnn_mlp.py` and `dnn_cnn.py` are shown in Fig. 2 and Fig. 3, respectively.

What to submit: Nothing.

3.4 Coding: Modules

[38 points: 7+7+8+8+8]

Q3.2 You will modify `dnn_misc.py`. This script defines all modules that you will need to construct the MLP and CNN in `dnn_mlp.py` and `dnn_cnn.py`, respectively. You have three tasks. First, finish the implementation of `forward` and `backward` functions in `class linear_layer`. Please follow Equation 2 for the forward pass and derive the partial derivatives accordingly. Second, finish the implementation of `forward` and `backward` functions in `class relu`. Please follow Equation 3 for the forward pass and derive the partial derivatives accordingly. Third, finish the the implementation of `backward` function in `class dropout`. We define the forward and the backward passes of dropout as follows.

$$\text{forward: } \mathbf{s} = \text{dropout.forward}(\mathbf{q} \in \mathbb{R}^J) = \frac{1}{1-r} \times \begin{bmatrix} \mathbf{1}[p[1] \geq r] \times q[1] \\ \vdots \\ \mathbf{1}[p[J] \geq r] \times q[J] \end{bmatrix}, \quad (19)$$

where $p[j]$ is sampled uniformly from $[0, 1)$, $\forall j \in \{1, \dots, J\}$,
and $r \in [0, 1)$ is a pre-defined scalar named dropout rate. (20)

$$\text{backward: } \frac{\partial \ell}{\partial \mathbf{q}} = \text{dropout.backward}(\mathbf{q}, \frac{\partial \ell}{\partial \mathbf{s}}) = \frac{1}{1-r} \times \begin{bmatrix} \mathbf{1}[p[1] \geq r] \times \frac{\partial \ell}{\partial s[1]} \\ \vdots \\ \mathbf{1}[p[J] \geq r] \times \frac{\partial \ell}{\partial s[J]} \end{bmatrix}. \quad (21)$$

Note that $p[j], j \in \{1, \dots, J\}$ and r are not learned so we do not need to compute the derivatives w.r.t. to them. Moreover, $p[j], j \in \{1, \dots, J\}$ are re-sampled every forward pass (for every data instance), and are kept for the following backward pass. The dropout rate r is set to 0 during testing.

Detailed descriptions/instructions about each pass (i.e., what to compute and what to return) are included in `dnn_misc.py`. Please do read carefully.

Note that in this script we do `import numpy as np`. Thus, to call a function `XX` from `numpy`, please use `np.XX`.

What to do and submit: Finish the implementation of 5 functions specified above in `dnn_misc.py`. Submit your completed `dnn_misc.py`. **We do provide a checking code `hw5_dnn_check.py` to check your implementation.** Just simply do `python3 hw5_dnn_check.py`.

3.5 Testing `dnn_misc.py` with multi-layer perceptron (MLP)

[8 points: 2+2+2+2]

Q3.3 *What to do and submit:* run script `q33.sh`. It will output `MLP_lr0.01_m0.0_w0.0_d0.0.json`.

What it does: `q33.sh` will run `python3 dnn_mlp.py` with learning rate 0.01, no momentum, no weight decay, and dropout rate 0.0. The output file stores the training and validation accuracies over 30 training epochs.

Q3.4 *What to do and submit:* run script `q34.sh`. It will output `MLP_lr0.01_m0.0_w0.0_d0.5.json`.

What it does: `q34.sh` will run `python3 dnn_mlp.py --dropout_rate 0.5` with learning rate 0.01, no momentum, no weight decay, and dropout rate 0.5. The output file stores the training and validation accuracies over 30 training epochs.

Q3.5 *What to do and submit:* run script `q35.sh`. It will output `MLP_lr0.01_m0.0_w0.0_d0.95.json`.
What it does: `q35.sh` will run `python3 dnn_mlp.py --dropout_rate 0.95` with learning rate 0.01, no momentum, no weight decay, and dropout rate 0.95. The output file stores the training and validation accuracies over 30 training epochs.

You will observe that the model in Q3.4 will give better validation accuracy (at epoch 30) compared to Q3.3. Specifically, dropout is widely-used to prevent over-fitting. However, if we use a too large dropout rate (like the one in Q3.5), the validation accuracy (together with the training accuracy) will be relatively lower, essentially under-fitting the training data.

Q3.6 *What to do and submit:* run script `q36.sh`. It will output `LR_lr0.01_m0.0_w0.0_d0.0.json`.
What it does: `q36.sh` will run `python3 dnn_mlp_nonlinear.py` with learning rate 0.01, no momentum, no weight decay, and dropout rate 0.0. The output file stores the training and validation accuracies over 30 training epochs.

The network has the same structure as the one in Q3.3, except that we remove the relu (nonlinear) layer. You will see that the validation accuracies drop significantly (the gap is around 0.03). Essentially, without the nonlinear layer, the model is learning multinomial logistic regression similar to Q4.2.

3.6 Testing `dnn_misc.py` with convolutional neural networks (CNN)

[4 points: 2+2]

Q3.7 *What to do and submit:* run script `q37.sh`. It will output `CNN_lr0.01_m0.0_w0.0_d0.5.json`.
What it does: `q37.sh` will run `python3 dnn_cnn.py` with learning rate 0.01, no momentum, no weight decay, and dropout rate 0.5. The output file stores the training and validation accuracies over 30 training epochs.

Q3.8 *What to do and submit:* run script `q38.sh`. It will output `CNN_lr0.01_m0.9_w0.0_d0.5.json`.
What it does: `q38.sh` will run `python3 dnn_cnn.py --alpha 0.9` with learning rate 0.01, momentum 0.9, no weight decay, and dropout rate 0.5. The output file stores the training and validation accuracies over 30 training epochs.

You will see that Q3.8 will lead to faster convergence than Q3.7 (i.e., the training/validation accuracies will be higher than 0.94 after 1 epoch). That is, using momentum will lead to more stable updates of the parameters.