

## General Instructions

**Starting point** Your repository will now have a directory 'homework2/'. Please do not change the name of this repository or the names of any files we have added to it. `git pull` and you will find the following:

- Python scripts which you need to amend:
  - `logistic.py`
  - `dnn_misc.py`
  - `dnn_cnn_2.py`
- Python scripts which you are **not allowed** to modify:
  - `dnn_mlp.py`, `dnn_mlp_nonlinear.py`, `dnn_cnn.py`, `hw2_dnn_check.py`, `dnn_im2col.py` and `data_loader.py`
- Helper scripts that you will use to generate output files: `q33.sh`, `q34.sh`, `q35.sh`, `q36.sh`, `q37.sh`, `q38.sh`, `q310.sh`, `logistic_binary.sh` and `logistic_multiclass.sh`

**Environment** This assignment has to be done in Python 3.5.2. Make sure you have the correct version installed. As in previous assignments, there are multiple ways you can install Python 3, for example:

- Install **Virtualbox** and import this **virtual machine**. Everything is setup already in this VM. Your submission will eventually be graded within this VM too.
- Alternatively you can also use **virtualenv** or **miniconda** to create a Python 3.5.2 environment for this programming assignment.
- [For testing purposes only] Other online coding platforms such as **Google Colab** which pre-installs all the packages used in this assignment (but with different versions).

**Python Packages** You are allowed to use the following Python packages:

- all built-in packages in Python 3.5.2, such as `sys`.
- **numpy** (1.13.1)
- `scipy` (0.19.1)
- **matplotlib** (2.0.2)

You will use Numpy mostly; in contrast, Scipy is usually not needed unless you have special needs. As for Matplotlib, you can use it to visualize the results, but it is not required for implementing this assignment.

You will also need the following package for testing your code, but do not import it yourself or use any of its functions in your implementation:

- `sklearn` (0.19.0)

**Download the data** Please download `mnist_subset.json` from Piazza/Resource/Homework. **DO NOT push it into your repository** when you submit your results; otherwise, you will get 20% deduction of your score of this assignment.

**Submission Instructions** The following will constitute your submission:

- The three Python scripts that you need to amend according to Sect. 2 and Sect. 3. Make sure that you committed your changes.
- Seven `.json` files and two `.out` files. These are the outputs from the eight helper scripts.

```
logistic_binary.out  
logistic_multiclass.out  
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  
CNN2_lr0.001_m0.9_w0.0_d0.5.json
```

## Problem 1 High-level descriptions

**1.1 Dataset** (Same as in Homework 1.) We will use `mnist_subset` (images of handwritten digits from 0 to 9). The dataset is stored in a JSON-formated 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  $N$ .

Besides, for logistic regression in Sect. 2, you will be using synthetic datasets with two, three and five classes.

**1.2 Tasks** You will be asked to implement binary and multiclass classification (Sect. 2) and neural networks (Sect. 3). Specifically, you will

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

In the next two subsections, 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. 2 and Sect. 3, as well as corresponding python scripts.

### 1.2.1 Logistic regression

**Coding** In `logistic.py`, finish implementing the following functions: `binary_train`, `binary_predict`, `multinomial_train`, `multinomial_predict`, `ovr_train` and `ovr_predict`. Refer to `logistic.py` and Sect. 2 for more information.

**Running your code** Run the scripts `logistic_binary.sh` and `logistic_multiclass.sh` after you finish your implementation. This will output:

- `logistic_binary.out`
- `logistic_multiclass.out`

**What to submit** Submit `logistic.py`, `logistic_binary.out`, `logistic_multiclass.out`.

### 1.2.2 Neural networks

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

**Coding** First, 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.

Second, in `dnn_cnn_2.py`, finish implementing the `main` function. There are five TODO items. Refer to `dnn_cnn_2.py` and Sect. 3 for more information.

**Running your code** Run the scripts `q33.sh`, `q34.sh`, `q35.sh`, `q36.sh`, `q37.sh`, `q38.sh`, `q310.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
CNN2_lr0.001_m0.9_w0.0_d0.5.json
```

**What to submit** Submit `dnn_misc.py`, `dnn_cnn_2.py`, and the above seven `.json` files.

### 1.3 Cautions

- Do not import packages that are not listed above (See Python Packages section).
- 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 mismatches 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 Python 3.5.2 in the VM.

**1.4 Advice** We are extensively using softmax and sigmoid function in this homework. To avoid numerical issues such as overflow and underflow caused by `numpy.exp()` and `numpy.log()`, please use the following implementations:

- Let  $x$  be a input vector to the softmax function. Use  $\tilde{x} = x - \max(x)$  instead of using  $x$  directly for the softmax function  $f$ . That is, if you want to compute  $f(x)_i$ , compute  $f(\tilde{x})_i = \frac{\exp(\tilde{x}_i)}{\sum_{j=1}^D \exp(\tilde{x}_j)}$  instead, which is clearly mathematically equivalent but numerically more stable.
- If you are using `numpy.log()`, make sure the input to the log function is positive. Also, there may be chances that one of the outputs of softmax, e.g.  $f(\tilde{x})_i$ , is extremely small but you need the value  $\ln(f(\tilde{x})_i)$ . In this case you should convert the computation equivalently into  $\tilde{x}_i - \ln(\sum_{j=1}^D \exp(\tilde{x}_j))$ .

We have implemented and run the code ourselves without problems, so if you follow the instructions and settings provided in the python files, you should not encounter overflow or underflow.

## Problem 2 Logistic Regression

(20 Points)

For this assignment you are asked to implement Logistic Regression for binary and multiclass classification.

### Q2.1 (6 Points)

In lecture 3 we discussed logistic regression for binary classification. In this problem, you are given a training set  $\mathcal{D} = \{(x_n, y_n)_{n=1}^N\}$ , where  $y_i \in \{0, 1\} \forall i = 1 \dots N$ . **Important:** note that here the binary labels are not  $-1$  or  $+1$  as used in the lecture, so be very careful about applying formulas from the lecture notes.

Your task is to learn the linear model specified by  $w^T x + b$  that minimizes the logistic loss. Note that we do not explicitly append the feature 1 to the data, so you need to explicitly learn the **bias**/intercept term  $b$  too. Specifically you need to implement function `binary_train` in `logistic.py` which uses gradient descent (*not* stochastic gradient descent) to find the optimal parameters (recall logistic regression does not admit a closed-form solution).

In addition you need to implement function `binary_predict` in `logistic.py`. We discuss two ways of making predictions in logistic regression in lecture 4: deterministic prediction or randomized prediction. Here you need to use the *deterministic prediction*.

After finishing implementation, please run `logistic_binary.sh` which generates `logistic_binary.out`.

What to submit:

- `logistic.py`
- `logistic_binary.out`

**Q2.2 (7 Points)** In the lectures you learned several methods to perform multiclass classification. One of them was one-versus-rest or one-versus-all approach.

For one-versus-rest classification in a problem with  $K$  classes, we need to train  $K$  classifiers using a black-box. Classifier  $k$  is trained on a binary problem, where the two labels corresponds to belonging or not belonging to class  $k$ . After that, the multiclass prediction is made based on the combination of all predictions from  $K$  binary classifiers.

In this problem you will implement one-versus-rest using binary logistic regression (that you have implemented in Q2.1) as the black-box. **Important:** the way to predict discussed in the lecture is to randomized over the classifiers that say “yes”; however, here since binary logistic regression naturally predicts a probability for each class (recall the sigmoid model), we will simply predict the class with the highest probability (using numpy `argmax`).

To sum up, you need to complete functions `OVR_train` and `OVR_predict` to perform one-versus-rest classification. After you finished implementation, please run `logistic_multiclass.sh` script, which will produce `logistic_multiclass.out`.

What to submit: `logistic.py` and `logistic_multiclass.out`.

**Q2.3 (7 Points)** Yet another multiclass classification method you learned was multinomial logistic regression. Complete the functions `multinomial_train` and `multinomial_predict` to perform multinomial logistic regression, following the same notes as in Q2.1, that is, 1) explicitly learn the biased term; 2) perform gradient descent instead of stochastic gradient descent; 3) make deterministic predictions.

After you finished implementation, please run `logistic_multiclass.sh` script, which will produce `logistic_multiclass.out`.

What to submit: `logistic.py` and `logistic_multiclass.out`.

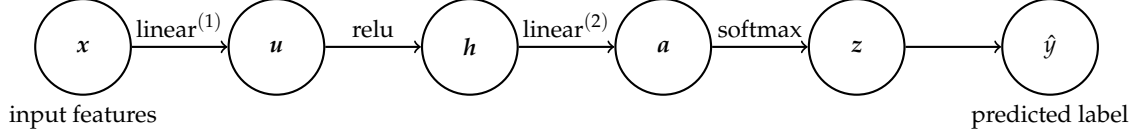


Figure 1: A diagram of a multi-layer perceptron (MLP). The edges mean mathematical operations (modules), and the circles mean variables. The term *relu* stands for rectified linear units.

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

(30 Points)

#### 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, just another term for fully connected feedforward networks we discussed in the lecture) for a  $K$ -class classification problem. The edges correspond to modules and the circles correspond to variables. Let  $(x \in \mathbb{R}^D, y \in \{1, 2, \dots, K\})$  be a labeled instance, such an MLP performs the following computations

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

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

$$\text{relu} : h = \max\{0, u\} = \begin{bmatrix} \max\{0, u_1\} \\ \vdots \\ \max\{0, u_M\} \end{bmatrix} \quad (3)$$

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

$$\text{softmax} : z = \begin{bmatrix} \frac{e^{a_1}}{\sum_k e^{a_k}} \\ \vdots \\ \frac{e^{a_K}}{\sum_k e^{a_k}} \end{bmatrix} \quad (5)$$

$$\text{predicted label : } \hat{y} = \arg \max_k z_k. \quad (6)$$

For a  $K$ -class classification problem, one popular loss function for training (i.e., to learn  $W^{(1)}, W^{(2)}, b^{(1)}, b^{(2)}$ ) is the cross-entropy loss. Specifically we denote the cross-entropy loss with respect to the training example  $(x, y)$  by  $l$ :

$$l = -\log(z_y) = \log \left( 1 + \sum_{k \neq y} e^{a_k - a_y} \right)$$

Note that one should look at  $l$  as a function of the parameters of the network, that is,  $W^{(1)}, b^{(1)}, W^{(2)}$  and  $b^{(2)}$ . For ease of notation, let us define the one-hot (i.e., 1-of- $K$ ) encoding of a class  $y$  as

$$\mathbf{y} \in \mathbb{R}^K \text{ and } y_k = \begin{cases} 1, & \text{if } y = k, \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

so that

$$l = -\sum_k y_k \log z_k = -\mathbf{y}^T \begin{bmatrix} \log z_1 \\ \vdots \\ \log z_K \end{bmatrix} = -\mathbf{y}^T \log \mathbf{z}. \quad (8)$$

We can then perform error-backpropagation, 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.

## 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.

$$\begin{aligned} \text{forward pass: } \mathbf{u} &= \text{linear}^{(1)}.forward(\mathbf{x}) = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}, \\ \text{where } \mathbf{W}^{(1)} \text{ and } \mathbf{b}^{(1)} &\text{ are its parameters.} \end{aligned} \quad (9)$$

$$\text{backward pass: } \left[ \frac{\partial l}{\partial \mathbf{x}}, \frac{\partial l}{\partial \mathbf{W}^{(1)}}, \frac{\partial l}{\partial \mathbf{b}^{(1)}} \right] = \text{linear}^{(1)}.backward(\mathbf{x}, \frac{\partial l}{\partial \mathbf{u}}). \quad (10)$$

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  $\frac{\partial l}{\partial \mathbf{z}}$ .

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 same 10-class classification problem introduced in Sect. 2. 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, e.g., the seminal work of Yann LeCun et al. “Gradient-based learning applied to document recognition,” written in 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

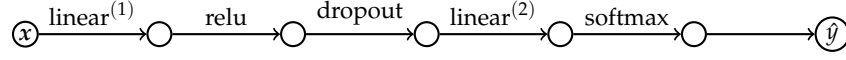


Figure 2: The diagram of the MLP implemented in `dnn_mlp.py`. The circles mean variables and edges mean modules.

$$l_{mb} = \frac{1}{N} \sum_{i=1}^N l(\text{softmax.forward}(\text{linear}^{(2)}.forward(\text{relu.forward}(\text{linear}^{(1)}.forward(x_i))), y_i) \quad (11)$$

$$= \frac{1}{N} \sum_{i=1}^N l(\text{softmax.forward}(\text{linear}^{(2)}.forward(\text{relu.forward}(u_i))), y_i) \quad (12)$$

$$= \dots \quad (13)$$

$$= \frac{1}{N} \sum_{i=1}^N l(\text{softmax.forward}(a_i), y_i) \quad (14)$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log z_{ik}. \quad (15)$$

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.  $\{a_i\}_{i=1}^N$

$$\frac{\partial l_{mb}}{\partial \{a_i\}_{i=1}^N} = \begin{bmatrix} (\frac{\partial l_{mb}}{\partial a_1})^T \\ (\frac{\partial l_{mb}}{\partial a_2})^T \\ \vdots \\ (\frac{\partial l_{mb}}{\partial a_{N-1}})^T \\ (\frac{\partial l_{mb}}{\partial a_N})^T \end{bmatrix}, \quad (16)$$

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

## 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.

## Coding: Modules



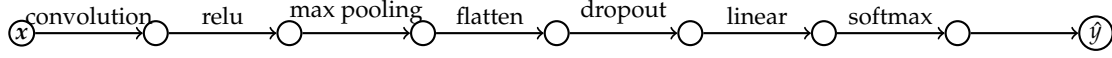


Figure 3: The diagram of the CNN implemented in `dnn_cnn.py`. The circles correspond to variables and edges correspond to modules. Note that the input to CNN may not be a vector (e.g., in `dnn_cnn.py` it is an image, which can be represented as a 3-dimensional tensor). The flatten layer is to reshape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

**Q3.2 (14 Points)** 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 Eqn. (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 Eqn. (3) for the forward pass and derive the partial derivatives accordingly. Third, finish the implementation of `backward` function in class `dropout`. We define the forward and the backward passes as follows.

$$\text{forward pass: } s = \text{dropout.forward}(q \in \mathbb{R}^J) = \frac{1}{1-r} \times \begin{bmatrix} \mathbf{1}[p_1 \geq r] \times q_1 \\ \vdots \\ \mathbf{1}[p_J \geq r] \times q_J \end{bmatrix}, \quad (17)$$

$$\begin{aligned} &\text{where } p_j \text{ is sampled uniformly from } [0, 1), \forall j \in \{1, \dots, J\}, \\ &\text{and } r \in [0, 1) \text{ is a pre-defined scalar named dropout rate.} \end{aligned} \quad (18)$$

$$\text{backward pass: } \frac{\partial l}{\partial q} = \text{dropout.backward}(q, \frac{\partial l}{\partial s}) = \frac{1}{1-r} \times \begin{bmatrix} \mathbf{1}[p_1 \geq r] \times \frac{\partial l}{\partial s_1} \\ \vdots \\ \mathbf{1}[p_J \geq r] \times \frac{\partial l}{\partial s_J} \end{bmatrix}. \quad (19)$$

Note that  $p_j, j \in \{1, \dots, J\}$  and  $r$  are not be 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, 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 `hw2_dnn_check.py` to check your implementation.**

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

**Q3.3 (2 Points)** *What to do and submit:* run script `q33.sh`. It will output `MLP_lr0.01_m0.0_w0.0_d0.0.json`. Add, commit, and push this file before the due date.

*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 (2 Points)** *What to do and submit:* run script `q34.sh`. It will output `MLP_lr0.01_m0.0_w0.0_d0.5.json`. Add, commit, and push this file before the due date.

*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 (2 Points)** *What to do and submit:* run script `q35.sh`. It will output `MLP_lr0.01_m0.0_w0.0_d0.95.json`. Add, commit, and push this file before the due date.

*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 (2 Points)** *What to do and submit:* run script `q36.sh`. It will output `LR_lr0.01_m0.0_w0.0_d0.0.json`. Add, commit, and push this file before the due date.

*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 Q2.3.

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

**Q3.7 (2 Points)** *What to do and submit:* run script `q37.sh`. It will output `CNN_lr0.01_m0.0_w0.0_d0.5.json`. Add, commit, and push this file before the due date.

*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 (2 Points)** *What to do and submit:* run script `q38.sh`. It will output `CNN_lr0.01_m0.9_w0.0_d0.5.json`. Add, commit, and push this file before the due date.

*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.

### Coding: Building a deeper architecture

**Q3.9 (2 Points)** The CNN architecture in `dnn_cnn.py` has only one convolutional layer. In this question, you are going to construct a two-convolutional-layer CNN (see Fig. 4 using the modules you implemented in Q3.2. Please modify the `main` function in `dnn_cnn_2.py`. The code in `dnn_cnn_2.py` is similar to that in `dnn_cnn.py`, except that there are a few parts marked as `TODO`. You need to fill in your code so as to construct the CNN in Fig. 4.

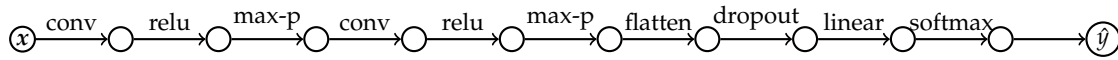


Figure 4: The diagram of the CNN you are going to implement in `dnn_cnn_2.py`. The term *conv* stands for convolution; *max-p* stands for max pooling. The circles correspond to variables and edges correspond to modules. Note that the input to CNN may not be a vector (e.g., in `dnn_cnn_2.py` it is an image, which can be represented as a 3-dimensional tensor). The flatten layer is to reshape its input into vector.

*What to do and submit:* Finish the implementation of the `main` function in `dnn_cnn_2.py` (search for `TODO` in `main`). Submit your completed `dnn_cnn_2.py`.

### Testing `dnn_cnn_2.py`

**Q3.10 (2 Points)** *What to do and submit:* run script `q310.sh`. It will output `CNN2_lr0.001_m0.9_w0.0_d0.5.json`. Add, commit, and push this file before the due date.

*What it does:* `q310.sh` will run `python3 dnn_cnn_2.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 you can achieve slightly higher validation accuracies than those in Q3.8.