

Homework 1: Backpropagation

The goal of homework 1 is to help you understand how to update network parameters by the using backpropagation algorithm.

For part 1, you need to answer the questions with mathematics equations. You should put all your answers in a PDF file and we will not accept any scanned hand-written answers. It is recommended to use \LaTeX .

For part 2, you need to program with Python. It requires you to implement your own forward and backward pass without using autograd. You need to submit your `mlp.py` file for this part.

The due date of homework 1 is 23:55 of 02/25. Submit the following files in a zip file `your_net_id.zip` through NYU classes:

- `theory.pdf`
- `mlp.py`

The following behaviors will result in penalty of your final score:

1. 5% penalty for submitting your files without using the correct format. (including naming the zip file, PDF file or python file wrong, or adding extra files in the zip folder, like the testing scripts from part 2).
2. 20% penalty for late submission within the first 24 hours. We will not accept any late submission after the first 24 hours.
3. 20% penalty for code submission that cannot be executed using the steps we mentioned in part 2. So please test your code before submit it.

1 Theory (50pt)

To answer questions in this part, you need some basic knowledge of linear algebra and matrix calculus. Also, you need to follow the instructions:

1. Every vector is treated as column vector.
2. You need to use the numerator-layout notation for matrix calculus. Please refer to Wikipedia about the notation. *
3. You are only allowed to use vector and matrix. You cannot use tensor in any of your answer.
4. Missing transpose are considered as wrong answer.

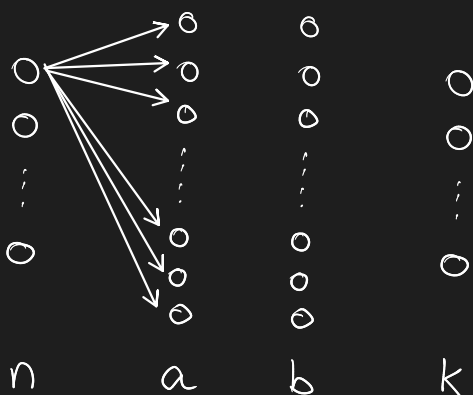
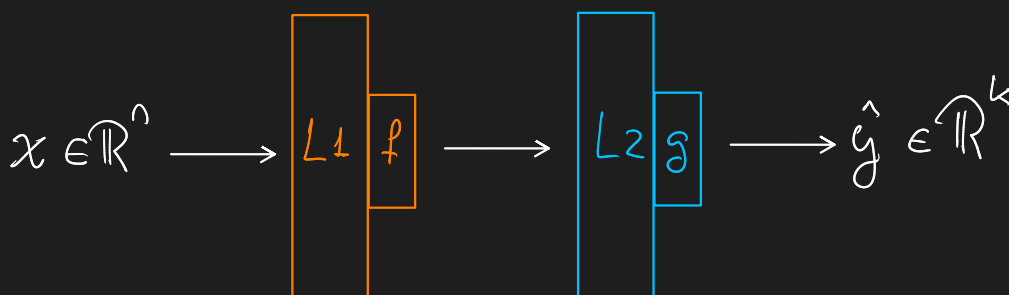
* https://en.wikipedia.org/wiki/Matrix_calculus#Numerator-layout_notation

1.1 Two-Layer Neural Nets

You are given the following neural net architecture:

$$\text{Linear}_1 \rightarrow f \rightarrow \text{Linear}_2 \rightarrow g$$

where $\text{Linear}_i(x) = \mathbf{W}^{(i)}\mathbf{x} + \mathbf{b}^{(i)}$ is the i -th affine transformation, and f, g are element-wise nonlinear activation functions. When an input $\mathbf{x} \in \mathbb{R}^n$ is fed to the network, $\hat{\mathbf{y}} \in \mathbb{R}^k$ is obtained as the output.



1.2 Regression Task

We would like to perform regression task. We choose $f(\cdot) = (\cdot)^+ = \text{ReLU}(\cdot)$ and g to be the identity function. To train this network, we choose MSE loss function $\ell_{\text{MSE}}(\hat{\mathbf{y}}, \mathbf{y}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2$, where \mathbf{y} is the target output.

(a) Name and mathematically describe the 5 programming steps you would take to train this model with PyTorch using SGD on a single batch of data.

1. Feed forward pass
2. Compute Loss (or free energy)
3. Delete/Zero past gradients
4. Backward pass: Compute and accumulate gradients
5. Update parameters

```
# Full-batch training loop
```

```
for t in range(1_000):
```

```
    # Feed forward to get the logits
```

```
    l = model(X)
```

$$\hat{y} := g \left(L_2 \left(f \left(L_1(x) \right) \right) \right)$$

```
    # Compute the free energy F
```

```
    F = C(l, y)
```

```
    L = F.mean()
```

$$\mathcal{L} := \|\hat{y} - y\|$$

```
    # Zero the gradients
```

```
    optimiser.zero_grad()
```

```
    # Backward pass to compute and accumulate the gradient
```

```
    # of the free energy w.r.t our learnable params
```

```
    L.backward()
```

```
    # Update params
```

```
    optimiser.step()
```

```
    # Display epoch, L, and accuracy
```

```
    overwrite(f'[EPOCH]: {t}, [LOSS]: {L.item():.6f}, [ACCURACY]: {acc(l, y):.3f}')
```

- (b) For a single data point (x, y) , write down all inputs and outputs for forward pass of each layer. You can only use variable $\mathbf{x}, \mathbf{y}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}$ in your answer. (note that $\text{Linear}_i(\mathbf{x}) = \mathbf{W}^{(i)}\mathbf{x} + \mathbf{b}^{(i)}$).

Layer	Input	Output
Linear ₁		
f		
Linear ₂		
g		
Loss		

- (c) Write down the gradient calculated from the backward pass. You can only use the following variables: $\mathbf{x}, \mathbf{y}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}, \frac{\partial \ell}{\partial \hat{\mathbf{y}}}, \frac{\partial z_2}{\partial z_1}, \frac{\partial \hat{\mathbf{y}}}{\partial z_3}$ in your answer, where $z_1, z_2, z_3, \hat{\mathbf{y}}$ are the outputs of $\text{Linear}_1, f, \text{Linear}_2, g$.

Parameter	Gradient
$\mathbf{W}^{(1)}$	
$\mathbf{b}^{(1)}$	
$\mathbf{W}^{(2)}$	
$\mathbf{b}^{(2)}$	

- (d) Show us the elements of $\frac{\partial z_2}{\partial z_1}$, $\frac{\partial \hat{\mathbf{y}}}{\partial z_3}$ and $\frac{\partial \ell}{\partial \hat{\mathbf{y}}}$ (be careful about the dimensionality)?

1.3 Classification Task

We would like to perform multi-class classification task, so we set both $f, g = \sigma$, the logistic sigmoid function $\sigma(z) \doteq (1 + \exp(-z))^{-1}$.

- (a) If you want to train this network, what do you need to change in the equations of (b), (c) and (d), assuming we are using the same MSE loss function.
- (b) Now you think you can do a better job by using a *Binary Cross Entropy* (BCE) loss function $\ell_{\text{BCE}}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{K} \sum_{i=1}^K -[y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$. What do you need to change in the equations of (b), (c) and (d)?
- (c) Things are getting better. You realize that not all intermediate hidden activations need to be binary (or soft version of binary). You decide to use $f(\cdot) = (\cdot)^+$ but keep g as σ . Explain why this choice of f can be beneficial for training a (deeper) network.