

# Recitation 3

Computing Derivatives and Autograd

Talha Faiz, Fuyu Tang, and Zishen Wen

# Agenda

1. Autograd and Computational graphs
2. Back propagations: derivatives, gradients, and chain rules
3. Computing derivatives

# Autograd – HW1 Bonus

# Automatic differentiation

- Recall what we did in back propagation (will cover in details soon):
  - Express the computation as a series of computations of intermediate values
  - Repeatedly apply the chain rule of differentiation
- All computer functions can be rewritten in the form of nested differentiable operations
- We, thus, could use a framework, “Automatic Differentiation” (Autodiff), to calculate the derivatives of any arbitrarily complex function.

# Automatic differentiation

- In this bonus, we will build an alternative implementation of MyTorch (HW\*P1), based on a popular Automatic Differentiation framework – Autograd.
- By doing this bonus, you might find your time spent on part 1s is saved!
- Key components:
  - Autograd engine -> the core class for performing Autodiff
  - Functional scripts/ activation/ linear/ loss -> Similar to part 1s, but are expected to be decomposed into the most basic operation, in order to be recorded by autograd engine
  - Utils -> contains methods to store and update variables

# Automatic differentiation

- Key ideas:
  - All calculations are break down into several basic operations (e.g. add, div, matmul, etc.)
  - Use a list to track the sequence of operation
  - When performing back propagation, the list is evaluated in reverse order (i.e. calculate the gradient of inputs at each step and update them).

# Automatic differentiation

- Example:

- $y = Wx + b$

- We first break it down to two basic operations: matmul and add:

- $z = Wx$

- $y = z + b$

- For each of those operations, we add a spot in Memory buffer for each of the inputs, create an Operation object saving all information related to the operation and then append it to operation list of Autograd class
      - Iterate over the operation list in reverse order

# Automatic differentiation

Write down the formulas, derive the gradient by sum, product and chain rules

$$\frac{\partial(f(\theta)+g(\theta))}{\partial\theta} = \frac{\partial f(\theta)}{\partial\theta} + \frac{\partial g(\theta)}{\partial\theta} \quad \frac{\partial(f(\theta)g(\theta))}{\partial\theta} = g(\theta) \frac{\partial f(\theta)}{\partial\theta} + f(\theta) \frac{\partial g(\theta)}{\partial\theta} \quad \frac{\partial f(g(\theta))}{\partial\theta} = \frac{\partial f(g(\theta))}{\partial g(\theta)} \frac{\partial g(\theta)}{\partial\theta}$$

Naively do so can result in wasted computations

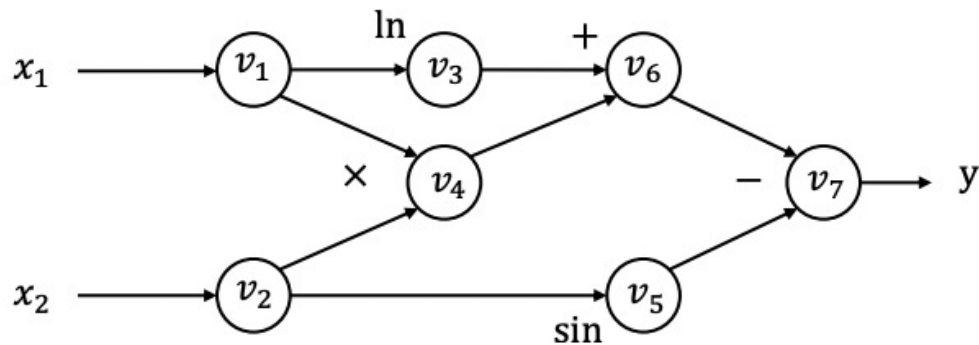
Example:  $f(\theta) = \prod_{i=0}^n \theta_i$        $\frac{\partial f(\theta)}{\partial \theta_k} = \prod_{j \neq k} \theta_j$

Cost  $n(n - 1)$  multiplies to compute all partial gradients



# Computational graph

$$y = f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin x_2$$



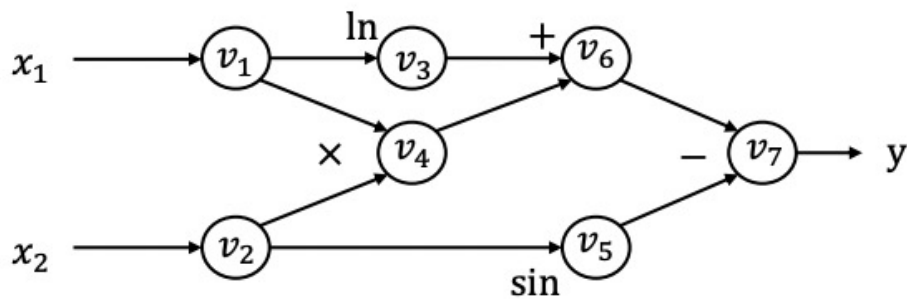
Forward evaluation trace

$$\begin{aligned}v_1 &= x_1 = 2 \\v_2 &= x_2 = 5 \\v_3 &= \ln v_1 = \ln 2 = 0.693 \\v_4 &= v_1 \times v_2 = 10 \\v_5 &= \sin v_2 = \sin 5 = -0.959 \\v_6 &= v_3 + v_4 = 10.693 \\v_7 &= v_6 - v_5 = 10.693 + 0.959 = 11.652 \\y &= v_7 = 11.652\end{aligned}$$

Each node represent an (intermediate) value in the computation. Edges present input output relations.

# Automatic differentiation (*Forwards*)

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$



**Define**  $\dot{v}_i = \frac{\partial v_i}{\partial x_1}$

We can then compute the  $\dot{v}_i$  iteratively in the forward topological order of the computational graph

Forward evaluation trace

$$\begin{aligned}v_1 &= x_1 = 2 \\v_2 &= x_2 = 5 \\v_3 &= \ln v_1 = \ln 2 = 0.693 \\v_4 &= v_1 \times v_2 = 10 \\v_5 &= \sin v_2 = \sin 5 = -0.959 \\v_6 &= v_3 + v_4 = 10.693 \\v_7 &= v_6 - v_5 = 10.693 + 0.959 = 11.652 \\y &= v_7 = 11.652\end{aligned}$$

Forward AD trace

$$\begin{aligned}\dot{v}_1 &= 1 \\\dot{v}_2 &= 0 \\\dot{v}_3 &= \dot{v}_1 / v_1 = 0.5 \\\dot{v}_4 &= \dot{v}_1 v_2 + \dot{v}_2 v_1 = 1 \times 5 + 0 \times 2 = 5 \\\dot{v}_5 &= \dot{v}_2 \cos v_2 = 0 \times \cos 5 = 0 \\\dot{v}_6 &= \dot{v}_3 + \dot{v}_4 = 0.5 + 5 = 5.5 \\\dot{v}_7 &= \dot{v}_6 - \dot{v}_5 = 5.5 - 0 = 5.5\end{aligned}$$

## Automatic differentiation (*Forwards*)

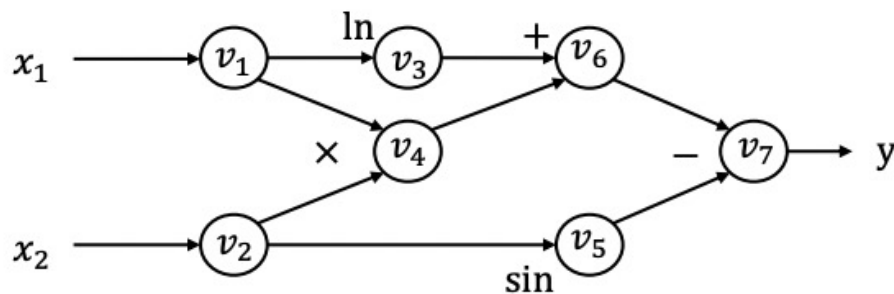
For  $f: \mathbb{R}^n \rightarrow \mathbb{R}^k$ , we need  $n$  forward AD passes to get the gradient with respect to each input.

We mostly care about the cases where  $k = 1$  and large  $n$ .

In order to resolve the problem efficiently, we need to use another kind of AD.

# Automatic differentiation (*Reversed*)

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$



Forward evaluation trace

$$v_1 = x_1 = 2$$

$$v_2 = x_2 = 5$$

$$v_3 = \ln v_1 = \ln 2 = 0.693$$

$$v_4 = v_1 \times v_2 = 10$$

$$v_5 = \sin v_2 = \sin 5 = -0.959$$

$$v_6 = v_3 + v_4 = 10.693$$

$$v_7 = v_6 - v_5 = 10.693 + 0.959 = 11.652$$

$$y = v_7 = 11.652$$

Define adjoint  $\bar{v}_i = \frac{\partial y}{\partial v_i}$

We can then compute the  $\bar{v}_i$  iteratively in the **reverse** topological order of the computational graph

Reverse AD evaluation trace

$$\bar{v}_7 = \frac{\partial y}{\partial v_7} = 1$$

$$\bar{v}_6 = \bar{v}_7 \frac{\partial v_7}{\partial v_6} = \bar{v}_7 \times 1 = 1$$

$$\bar{v}_5 = \bar{v}_7 \frac{\partial v_7}{\partial v_5} = \bar{v}_7 \times (-1) = -1$$

$$\bar{v}_4 = \bar{v}_6 \frac{\partial v_6}{\partial v_4} = \bar{v}_6 \times 1 = 1$$

$$\bar{v}_3 = \bar{v}_6 \frac{\partial v_6}{\partial v_3} = \bar{v}_6 \times 1 = 1$$

$$\bar{v}_2 = \bar{v}_5 \frac{\partial v_5}{\partial v_2} + \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_5 \times \cos v_2 + \bar{v}_4 \times v_1 = -0.284 + 2 = 1.716$$

$$\bar{v}_1 = \bar{v}_4 \frac{\partial v_4}{\partial v_1} + \bar{v}_3 \frac{\partial v_3}{\partial v_1} = \bar{v}_4 \times v_2 + \bar{v}_3 \frac{1}{v_1} = 5 + \frac{1}{2} = 5.5$$

# References

Recommend to take a look

- <https://dlsyscourse.org/slides/4-automatic-differentiation.pdf>
- [https://en.wikipedia.org/wiki/Automatic\\_differentiation](https://en.wikipedia.org/wiki/Automatic_differentiation)
- [https://deeplearning.cs.cmu.edu/S22/document/recitation/Recitation3/Recitation\\_3.pdf](https://deeplearning.cs.cmu.edu/S22/document/recitation/Recitation3/Recitation_3.pdf)

During Back Propagations, you will find  
we are doing the same thing....

# What is a loss function and loss?

“The function we want to minimize or maximize is called the **objective function** or **criterion**. When we are **minimizing** it, we may also call it the **cost function**, **loss function**, or **error function**.” [1]

Functions of loss:

**1. Monitor:** Loss evaluates the performance of the model. The lower the loss is, the better the model is.

**2. Part of the optimizer:**

Learning problem -> Optimization problem

Define loss function -> minimize the loss function

[1] Ian Goodfellow, Yoshua Bengio, Aaron Courville, *Deep Learning*, MIT Press, 2017

# Back propagation of loss

Loss is the starting point of the back propagation

Backpropagation aims to minimize the cost function by adjusting network's weights and biases. The level of adjustment is determined by the gradients of the cost function w.r.t. those parameters.



# Back propagation: Derivatives, Gradients, and the Chain Rule

Training a network:

1. Forward Propagation with current parameters
2. Calculate the loss
3. **Backward Propagation to calculate the gradients of the parameters**
4. Step to update the parameters with gradients

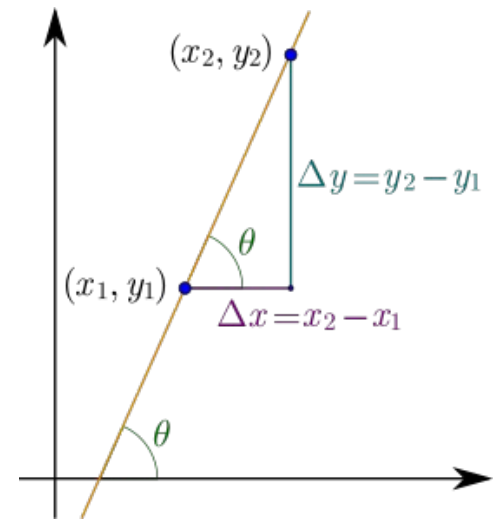
**The gradient is the transpose of the derivative**

# Derivatives

Mathematically, the derivative of a function  $f$  measures the sensitivity of change of the function value w.r.t. a change in its input value  $x$ .

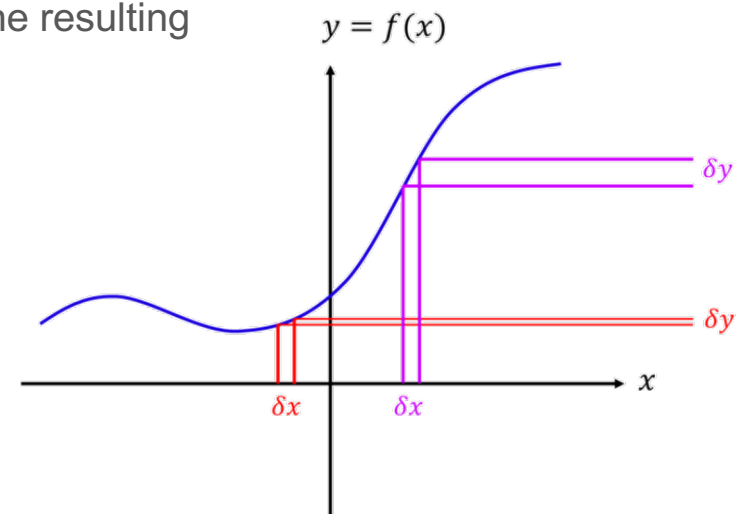
$$\frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}$$

Geometrically, the derivative of the  $f$  w.r.t.  $x$  at  $x_0$  is the slope of the tangent line to the graph of  $f$  at  $x_0$ .



# Derivatives of non-linear functions

Let  $y=f(x)$  be a relation between two variables,  $y$  and  $x$ . If  $f(x)$  is continuous and differentiable, any small perturbation of  $x$  will result in a small perturbation of  $y$ . We define the derivative of  $y$  with respect to  $x$  as the multiplier  $\alpha$  that relates a miniscule perturbation  $\delta x$  of  $x$  to the resulting perturbation  $\delta y$  of  $y$ .



# Derivatives

We note “the derivative of  $y$  with respect to  $x$ ” as

$$\Delta y = \nabla_x y \Delta x$$

The shape of the derivative for any variable will be transposed w.r.t that variable

Ex:

For a function with scalar input  $x$  and scalar output  $y$ ,  
its derivative is a scalar.

For a function with  $(D \times 1)$  vector input  $x$  and scalar output  $y$ ,  
its derivative is a  $(1 \times D)$  row vector.

For a function with  $(D \times 1)$  vector input  $x$  and  $(K \times 1)$  vector output  $y$ ,  
its derivative is a  $(K \times D)$  matrix.

# Derivatives

Scalar derivatives (scalar in, scalar out)

$$\Delta y = f'(x) \Delta x$$

Multivariable derivatives (vector in, scalar out)

$$\Delta y = \nabla_x y \Delta x = \left[ \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_D} \right] \begin{bmatrix} \Delta x_1 \\ \vdots \\ \Delta x_D \end{bmatrix}$$

Full derivative

Partial derivative



# Key Ideas about Derivatives

1. The derivative is the best linear approximation of  $f$  at a point
2. The derivative is a linear transformation (matrix multiplication)
3. The derivative describes the effect of each input on the output

# Computing Derivatives – Scalar Chain Rule

$$L = f(z)$$

$$z = g(x)$$

All terms are scalars

$\frac{\partial L}{\partial z}$  is given

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z} g'(x)$$

(given)



# Computing Derivatives – Scalar Addition

$$L = f(z)$$
$$z = x + y$$

All terms are scalars

$\frac{\partial L}{\partial z}$  is given

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z}$$
$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y} = \frac{\partial L}{\partial z}$$

# Computing Derivatives – Scalar Multiplication

$$L = f(z)$$

$$z = Wx$$

All terms are scalars

$\frac{\partial L}{\partial z}$  is given

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z} W$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial W} = x \frac{\partial L}{\partial z}$$

# Computing Derivatives – Scalar Generalized Chain Rule

$$L = f(z)$$

$$z = z_1 + z_2 + \cdots + z_n = g_1(x) + g_2(x) + \cdots + g_n(x)$$

All terms are scalars

$\frac{\partial L}{\partial z}$  is given

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x} = \frac{\partial L}{\partial z} \left( \frac{\partial g_1}{\partial x} + \frac{\partial g_2}{\partial x} + \cdots + \frac{\partial g_n}{\partial x} \right)$$

# Computing Derivatives – Multivariable Chain Rule

$$L = f(z)$$

$$z = g(x)$$

Here we assume that  $L$  is  $M \times 1$  vector

$x$  is  $D \times 1$  vector,  $z$  is  $K \times 1$  vector

$\nabla_z L$  is given  $(M \times K)$  matrix

$$\nabla_x L = \nabla_z L \nabla_x Z$$

$$M \times D \quad M \times K \quad K \times D$$

# Computing Derivatives – Multivariable Vector Addition

$$L = f(z)$$
$$z = x + y$$

$x, y, z$  are all  $D \times 1$  vectors  
 $\nabla_z L$  is given  $(M \times D)$  matrix

$$\nabla_x L = \nabla_z L \nabla_x Z = \nabla_z L$$

$$\nabla_y L = \nabla_z L \nabla_y Z = \nabla_z L$$

$$M \times D \quad M \times D \quad D \times D$$

# Computing Derivatives – Multivariable Vector Addition of derivatives

$$L = f_1(z) + f_2(y)$$

$$z = g(x)$$

$$y = h(x)$$

$x$  is  $D \times 1$  vector,  $z$  is  $K \times 1$  vector,  $y$  is  $M \times 1$  vector

$\nabla_z L$  is given  $(N \times K)$  matrix

$\nabla_y L$  is given  $(N \times M)$  matrix

$$\nabla_x L = \nabla_z L \nabla_x Z + \nabla_y L \nabla_x Y$$

$$N \times D \quad N \times K \quad K \times D \quad N \times M \quad M \times D$$

# Computing Derivatives – Multivariable Matrix Multiplication

$$L = f(z)$$

$$z = Wx$$

$x$  is a  $D \times 1$  vector

$z$  is a  $K \times 1$  vector

$W$  is a  $K \times D$  matrix

$\nabla_z L$  is given  $(1 \times K)$  vector

$$\nabla_x L = \nabla_z L \nabla_x z = (\nabla_z L)W \quad 1 \times D$$

$$\nabla_W L = \nabla_z L \nabla_W z = x(\nabla_z L) \quad D \times K$$

## Computing Derivatives – Multivariable Generalized Chain Rule

$$L = f(z)$$

$$z = z_1 + z_2 + \cdots + z_n = g_1(x) + g_2(x) + \cdots + g_n(x)$$

Loss  $L$  is a  $M \times 1$  vector

$x$  is a  $D \times 1$  vector

$z$  is a  $K \times 1$  vector

$\nabla_z L$  is given  $(M \times K)$  matrix

$$\nabla_x L = \nabla_z L \nabla_x Z = \nabla_z L (\nabla_x Z_1 + \nabla_x Z_2 + \cdots + \nabla_x Z_n)$$

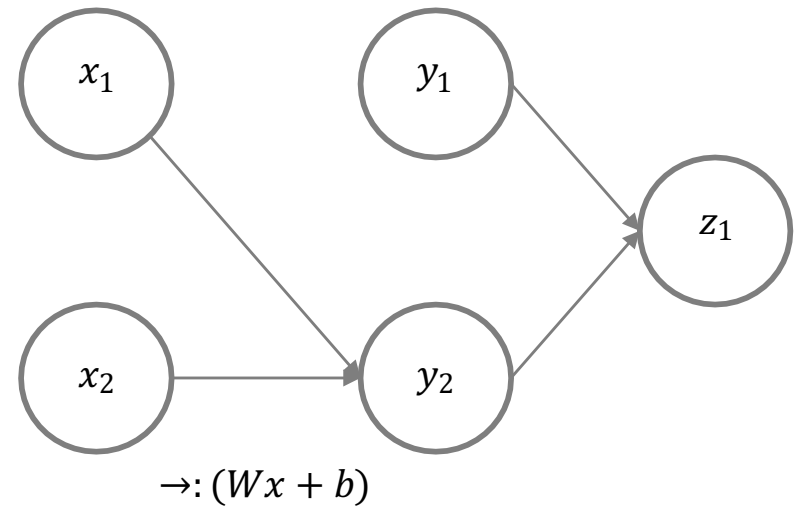


# Computing derivatives of complex functions

- We now are prepared to compute very complex derivatives
- Procedure:
  - Express the computation as a series of computations of intermediate values
  - Each computation must comprise either a unary or binary relation
    - Unary relation: RHS has one argument, e.g.  $y = g(x)$
    - Binary relation: RHS has two arguments  
e.g.  $z = x + y$  or  $z = xy$
  - Walk your way backward through the derivatives of the simple relations

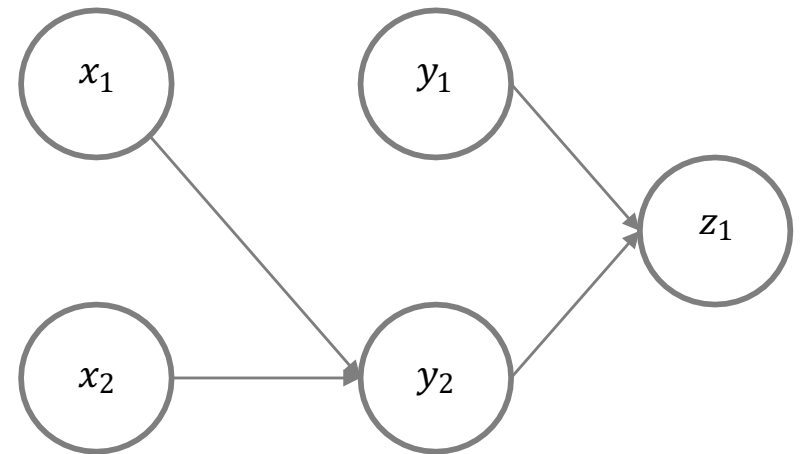
## Example:

- $y_2 = \tanh(W_{x_1}x_1 + b_{x_1} + W_{x_2}x_2 + b_{x_2})$
- $z_1 = \tanh(W_{y_1}y_1 + b_{y_1} + W_{y_2}y_2 + b_{y_2})$



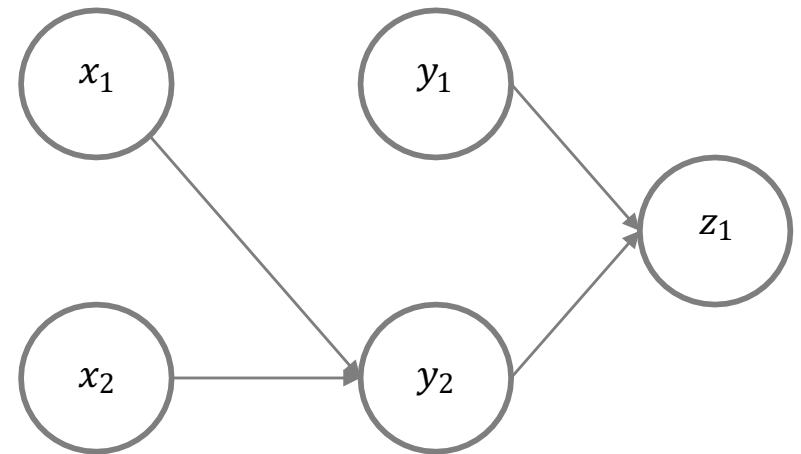
## Example:

- $y_2 = \tanh(W_{x_1}x_1 + b_{x_1} + W_{x_2}x_2 + b_{x_2})$
- ~~$z_1 = \tanh(W_{y_1}y_1 + b_{y_1} + W_{y_2}y_2 + b_{y_2})$~~
- $i_1 = W_{x_1}x_1$
- $i_2 = W_{x_2}x_2$
- $i_3 = i_1 + b_{x_1} + i_2 + b_{x_2}$
- $y_2 = \tanh(i_3)$



## Example:

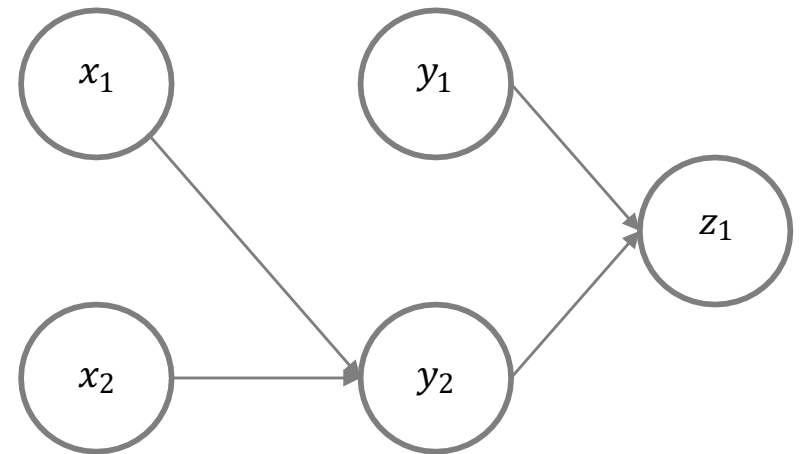
- ~~$y_2 = \tanh(W_{x_1}x_1 + b_{x_1} + W_{x_2}x_2 + b_{x_2})$~~
- $z_1 = \tanh(W_{y_1}y_1 + b_{y_1} + W_{y_2}y_2 + b_{y_2})$
- $i_4 = W_{y_1}y_1$
- $i_5 = W_{y_2}y_2$
- $i_6 = i_4 + b_{y_1} + i_5 + b_{y_2}$
- $z_1 = \tanh(i_6)$



## Example:

- $y_2 = \tanh(W_{x_1}x_1 + b_{x_1} + W_{x_2}x_2 + b_{x_2})$
- $z_1 = \tanh(W_{y_1}y_1 + b_{y_1} + W_{y_2}y_2 + b_{y_2})$

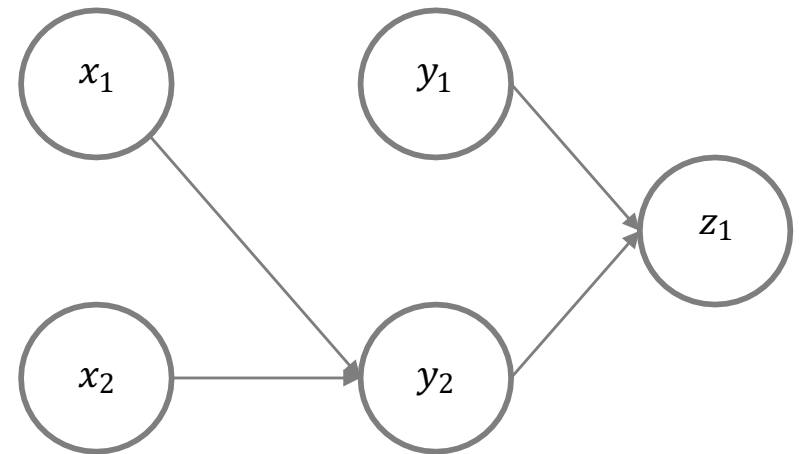
- $i_1 = W_{x_1}x_1$
- $i_2 = W_{x_2}x_2$
- $i_3 = i_1 + b_{x_1} + i_2 + b_{x_2}$
- $y_2 = \tanh(i_3)$
- $i_4 = W_{y_1}y_1$
- $i_5 = W_{y_2}y_2$
- $i_6 = i_4 + b_{y_1} + i_5 + b_{y_2}$
- $z_1 = \tanh(i_6)$



## Example:

- $y_2 = \tanh(W_{x_1}x_1 + b_{x_1} + W_{x_2}x_2 + b_{x_2})$
- $z_1 = \tanh(W_{y_1}y_1 + b_{y_1} + W_{y_2}y_2 + b_{y_2})$
- Given  $\frac{dL}{dz_1}$  ( $\nabla_{z_1} L$ )

- |   |   |
|---|---|
| • $i_1 = W_{x_1}x_1$                    | • $i_4 = W_{y_1}y_1$                    |
| • $i_2 = W_{x_2}x_2$                    | • $i_5 = W_{y_2}y_2$                    |
| • $i_3 = i_1 + b_{x_1} + i_2 + b_{x_2}$ | • $i_6 = i_4 + b_{y_1} + i_5 + b_{y_2}$ |
| • $y_2 = \tanh(i_3)$                    | • $z_1 = \tanh(i_6)$                    |



## Example:

- Given  $\frac{dL}{dz_1} (\nabla_{z_1} L)$
- $\nabla_{i_6} L = \nabla_{z_1} L \nabla_{i_6} z_1 = \nabla_{z_1} L (1 - \tanh^2(i_6))$
- $z_1 = \tanh(i_6)$

## Example:

- Given  $\frac{dL}{dz_1} (\nabla_{z_1} L)$
  - $\nabla_{i_6} L = \nabla_{z_1} L \nabla_{i_6} z_1 = \nabla_{z_1} L (1 - \tanh^2(i_6))$
  - $\nabla_{i_4} L = \nabla_{i_6} L \nabla_{i_4} i_6 = \nabla_{i_6} L$
  - $\nabla_{b_{y_1}} L = \nabla_{i_6} L \nabla_{b_{y_1}} i_6 = \nabla_{i_6} L$
  - $\nabla_{i_5} L = \nabla_{i_6} L \nabla_{i_5} i_6 = \nabla_{i_6} L$
  - $\nabla_{b_{y_2}} L = \nabla_{i_6} L \nabla_{b_{y_2}} i_6 = \nabla_{i_6} L$
- $z_1 = \tanh(i_6)$
  - $i_6 = i_4 + b_{y_1} + i_5 + b_{y_2}$



## Example:

- Given  $\frac{dL}{dz_1} (\nabla_{z_1} L)$
  - $\nabla_{i_6} L = \nabla_{z_1} L \nabla_{i_6} z_1 = \nabla_{z_1} L (1 - \tanh^2(i_6))$
  - $\nabla_{i_4} L = \nabla_{i_6} L \nabla_{i_4} i_6 = \nabla_{i_6} L$
  - $\nabla_{b_{y_1}} L = \nabla_{i_6} L \nabla_{b_{y_1}} i_6 = \nabla_{i_6} L$
  - $\nabla_{i_5} L = \nabla_{i_6} L \nabla_{i_5} i_6 = \nabla_{i_6} L$
  - $\nabla_{b_{y_2}} L = \nabla_{i_6} L \nabla_{b_{y_2}} i_6 = \nabla_{i_6} L$
  - $\nabla_{W_{y_2}} L = \nabla_{i_5} L \nabla_{W_{y_2}} i_5 = y_2 \nabla_{i_5} L$
  - $\nabla_{y_2} L = \nabla_{i_5} L \nabla_{y_2} i_5 = \nabla_{i_5} L W_{y_2}$
- $z_1 = \tanh(i_6)$
  - $i_6 = i_4 + b_{y_1} + i_5 + b_{y_2}$
  - $i_5 = W_{y_2} y_2$

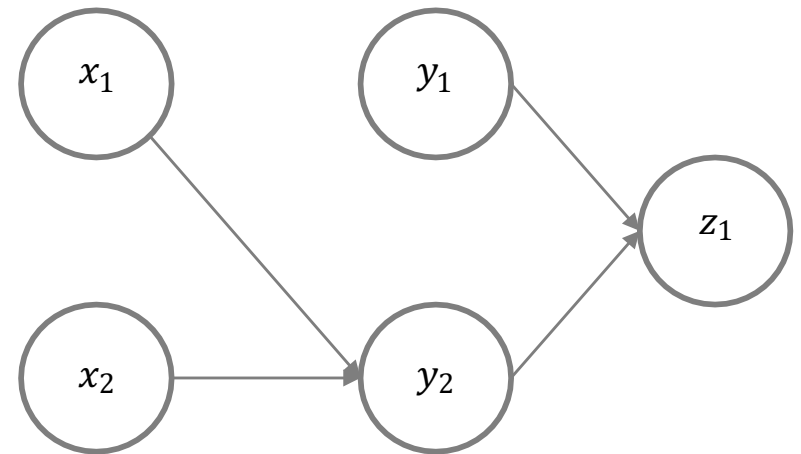
## Example:

- Given  $\frac{dL}{dz_1} (\nabla_{z_1} L)$
  - $\nabla_{i_6} L = \nabla_{z_1} L \nabla_{i_6} z_1 = \nabla_{z_1} L (1 - \tanh^2(i_6))$
  - $\nabla_{i_4} L = \nabla_{i_6} L \nabla_{i_4} i_6 = \nabla_{i_6} L$
  - $\nabla_{b_{y_1}} L = \nabla_{i_6} L \nabla_{b_{y_1}} i_6 = \nabla_{i_6} L$
  - $\nabla_{i_5} L = \nabla_{i_6} L \nabla_{i_5} i_6 = \nabla_{i_6} L$
  - $\nabla_{b_{y_2}} L = \nabla_{i_6} L \nabla_{b_{y_2}} i_6 = \nabla_{i_6} L$
  - $\nabla_{W_{y_2}} L = \nabla_{i_5} L \nabla_{W_{y_2}} i_5 = y_2 \nabla_{i_5} L$
  - $\nabla_{y_2} L = \nabla_{i_5} L \nabla_{y_2} i_5 = \nabla_{i_5} L W_{y_2}$
  - $\nabla_{W_{y_1}} L = \nabla_{i_4} L \nabla_{W_{y_1}} i_4 = y_1 \nabla_{i_4} L$
  - $\nabla_{y_1} L = \nabla_{i_4} L \nabla_{y_1} i_4 = \nabla_{i_4} L W_{y_1}$
- $z_1 = \tanh(i_6)$
  - $i_6 = i_4 + b_{y_1} + i_5 + b_{y_2}$
  - $i_5 = W_{y_2} y_2$
  - $i_4 = W_{y_1} y_1$

## Example:

- $y_2 = \tanh(W_{x_1}x_1 + b_{x_1} + W_{x_2}x_2 + b_{x_2})$
- $z_1 = \tanh(W_{y_1}y_1 + b_{y_1} + W_{y_2}y_2 + b_{y_2})$
- Given  $\frac{dL}{dz_1}$  ( $\nabla_{z_1} L$ )

- |   |   |
|---|---|
| • $i_1 = W_{x_1}x_1$                    | • $i_4 = W_{y_1}y_1$                    |
| • $i_2 = W_{x_2}x_2$                    | • $i_5 = W_{y_2}y_2$                    |
| • $i_3 = i_1 + b_{x_1} + i_2 + b_{x_2}$ | • $i_6 = i_4 + b_{y_1} + i_5 + b_{y_2}$ |
| • $y_2 = \tanh(i_3)$                    | • $z_1 = \tanh(i_6)$                    |



## Example:

- Given  $\frac{dL}{dy_2} (\nabla_{y_2} L)$
- $\nabla_{i_3} L = \nabla_{y_2} L \nabla_{i_3} y_2 = \nabla_{y_2} L (1 - \tanh^2(i_3))$
- $y_2 = \tanh(i_3)$

## Example:

- Given  $\frac{dL}{dy_2} (\nabla_{y_2} L)$
  - $\nabla_{i_3} L = \nabla_{y_2} L \nabla_{i_3} y_2 = \nabla_{y_2} L (1 - \tanh^2(i_3))$
  - $\nabla_{i_2} L = \nabla_{i_3} L \nabla_{i_2} i_3 = \nabla_{i_3} L$
  - $\nabla_{b_{x_1}} L = \nabla_{i_3} L \nabla_{b_{x_1}} i_3 = \nabla_{i_3} L$
  - $\nabla_{i_1} L = \nabla_{i_3} L \nabla_{i_1} i_3 = \nabla_{i_3} L$
  - $\nabla_{b_{x_2}} L = \nabla_{i_3} L \nabla_{b_{x_2}} i_3 = \nabla_{i_3} L$
- $y_2 = \tanh(i_3)$
  - $i_3 = i_1 + b_{x_1} + i_2 + b_{x_2}$

## Example:

- Given  $\frac{dL}{dy_2} (\nabla_{y_2} L)$
- $\nabla_{i_3} L = \nabla_{y_2} L \nabla_{i_3} y_2 = \nabla_{y_2} L (1 - \tanh^2(i_3))$
- $\nabla_{i_2} L = \nabla_{i_3} L \nabla_{i_2} i_3 = \nabla_{i_3} L$
- $\nabla_{b_{x_1}} L = \nabla_{i_3} L \nabla_{b_{x_1}} i_3 = \nabla_{i_3} L$
- $\nabla_{i_1} L = \nabla_{i_3} L \nabla_{i_1} i_3 = \nabla_{i_3} L$
- $\nabla_{b_{x_2}} L = \nabla_{i_3} L \nabla_{b_{x_2}} i_3 = \nabla_{i_3} L$
- $\nabla_{W_{x_2}} L = \nabla_{i_2} L \nabla_{W_{x_2}} i_2 = x_2 \nabla_{i_2} L$
- $\nabla_{x_2} L = \nabla_{i_2} L \nabla_{x_2} i_2 = \nabla_{i_2} L W_{x_2}$
- $y_2 = \tanh(i_3)$
- $i_3 = i_1 + b_{x_1} + i_2 + b_{x_2}$
- $i_2 = W_{x_2} x_2$

## Example:

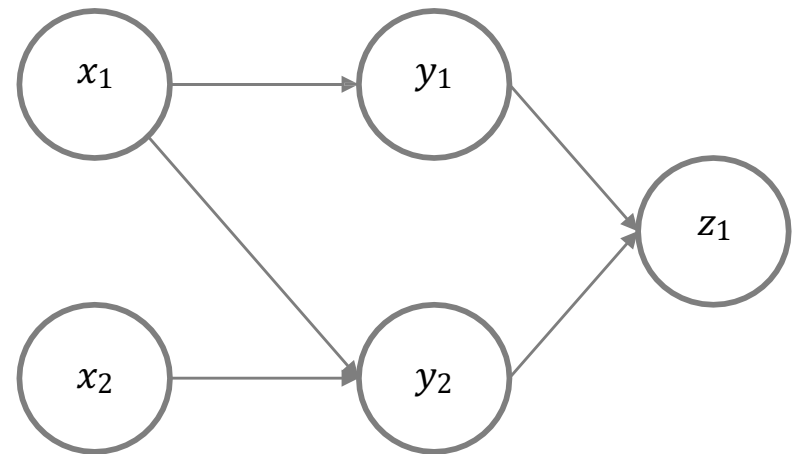
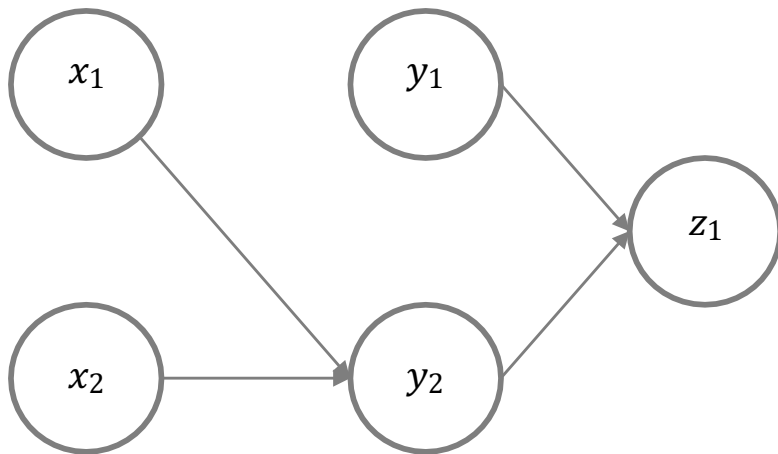
- Given  $\frac{dL}{dy_2} (\nabla_{y_2} L)$
- $\nabla_{i_3} L = \nabla_{y_2} L \nabla_{i_3} y_2 = \nabla_{y_2} L (1 - \tanh^2(i_3))$
- $\nabla_{i_2} L = \nabla_{i_3} L \nabla_{i_2} i_3 = \nabla_{i_3} L$
- $\nabla_{b_{x_1}} L = \nabla_{i_3} L \nabla_{b_{x_1}} i_3 = \nabla_{i_3} L$
- $\nabla_{i_1} L = \nabla_{i_3} L \nabla_{i_1} i_3 = \nabla_{i_3} L$
- $\nabla_{b_{x_2}} L = \nabla_{i_3} L \nabla_{b_{x_2}} i_3 = \nabla_{i_3} L$
- $\nabla_{W_{x_2}} L = \nabla_{i_2} L \nabla_{W_{x_2}} i_2 = x_2 \nabla_{i_2} L$
- $\nabla_{x_2} L = \nabla_{i_2} L \nabla_{x_2} i_2 = \nabla_{i_2} L W_{x_2}$
- $\nabla_{W_{x_1}} L = \nabla_{i_1} L \nabla_{W_{x_1}} i_1 = x_1 \nabla_{i_1} L$
- $\nabla_{x_1} L = \nabla_{i_1} L \nabla_{x_1} i_1 = \nabla_{i_1} L W_{x_1}$
- $y_2 = \tanh(i_3)$
- $i_3 = i_1 + b_{x_1} + i_2 + b_{x_2}$
- $i_2 = W_{x_2} x_2$
- $i_1 = W_{x_1} x_1$

## When to use “=” vs “+=”

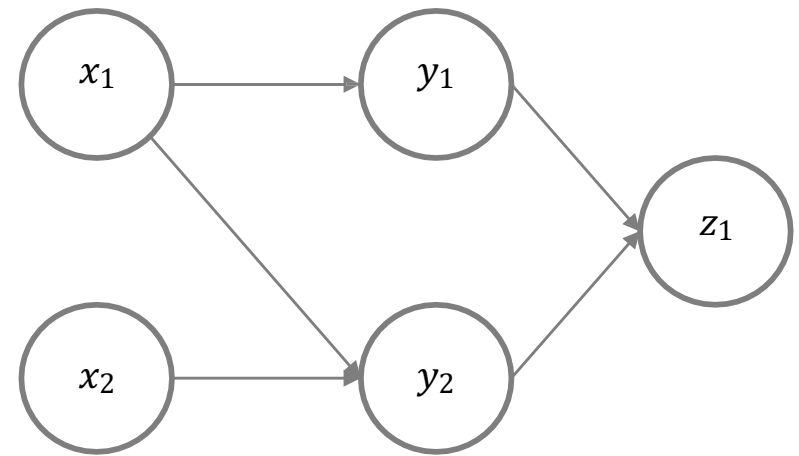
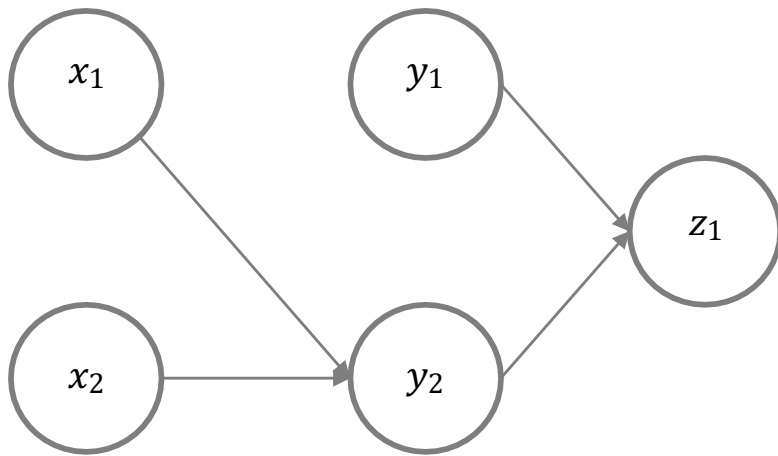
- In the forward computation a variable may be used multiple times to compute other intermediate variables
- During backward computations, the first time the derivative is computed for the variable, then we will use “=”
- In subsequent computations we use “+=”
- It may be difficult to keep track of when we first compute the derivative for a variable
  - When to use “=” vs when to use “+=”
- Cheap trick:
  - Initialize all derivatives to 0 during computation
  - *Always* use “+=”
  - You will get the correct answer (why?)



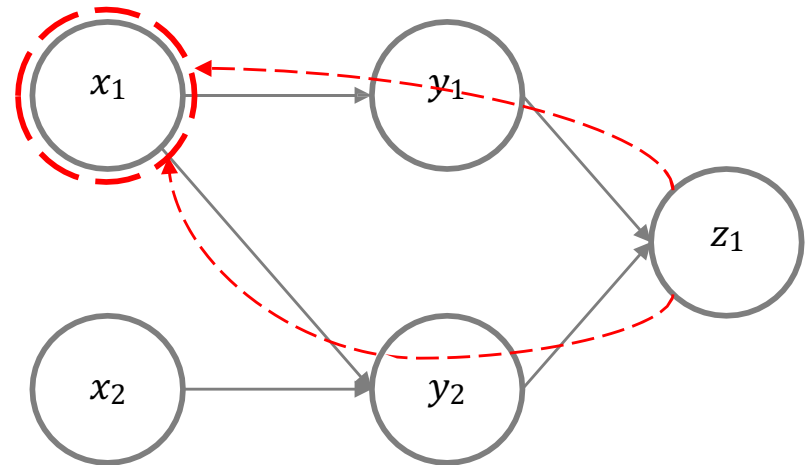
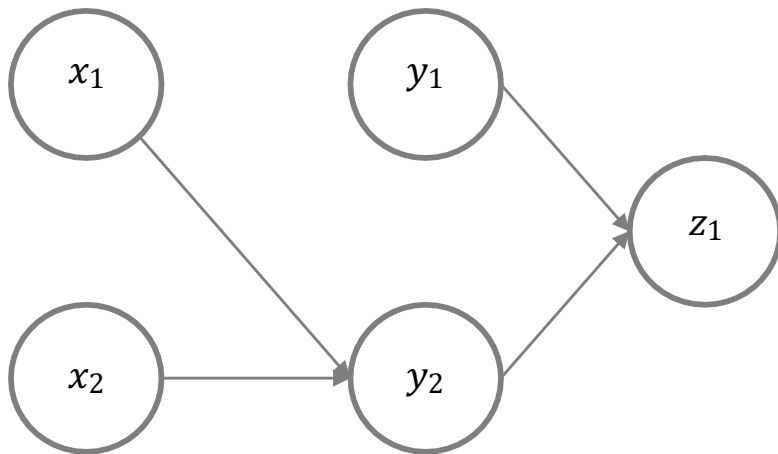
- In the figures below which example do you think uses “+=”?



- In the example (left figure) we showed before, we kept using "=", think about why it worked
- In the new example (right figure), which variable requires "+=" ?



- In the example (left figure) we showed before, we kept using "=", think about why it worked
- In the new example (right figure), which variable requires "+=" ?



# References

- <https://deeplearning.cs.cmu.edu/S21/document/recitation/Recitation2.pdf>
- <https://deeplearning.cs.cmu.edu/F20/document/recitation/recitation2.1.pdf>
- <https://deeplearning.cs.cmu.edu/F20/document/recitation/recitation2.2.pdf>
- <https://deeplearning.cs.cmu.edu/S20/document/recitation/recitation-2.pdf>
- <https://pytorch.org/docs/stable/nn.html#loss-functions>
- <https://towardsdatascience.com/understanding-backpropagation-algorithm-7bb3aa2f95fd>