

Recitation 4

Computing Derivatives

Agenda

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

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

Commonly used loss functions

Mean absolute error loss: `nn.L1Loss`

Mean squared error loss: `nn.MSELoss`

Cross Entropy loss (Classification): `nn.CrossEntropyLoss` (hw1, hw2)

Connectionist Temporal Classification loss: `nn.CTCLoss` (hw3)

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

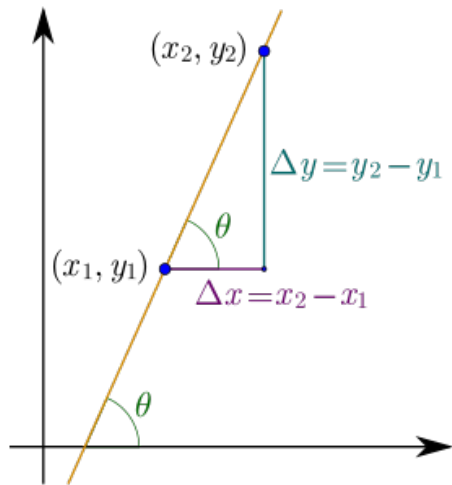
The gradient is the transpose of the derivative

Derivatives

Mathematically, the derivative of a function f measures the sensitivity of change of the function value y 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

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

Derivatives

Multivariable derivatives (vector in, vector out)

$$\text{Input } x = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix}, \text{ Output } y = \begin{bmatrix} y_1 \\ \vdots \\ y_K \end{bmatrix}$$

This is also referred to as the Jacobian of $f(x)$ and designated as $J_x f(x)$.

$$\begin{bmatrix} \Delta y_1 \\ \vdots \\ \Delta y_K \end{bmatrix} = \nabla_x y \Delta x = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_K}{\partial x_1} & \dots & \frac{\partial y_K}{\partial x_D} \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \vdots \\ \Delta x_D \end{bmatrix}$$

$K \times 1$ $K \times D$ $D \times 1$

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

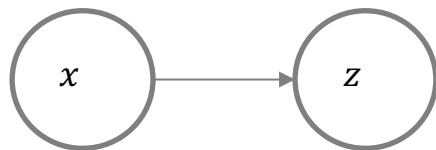
$$z = g(x)$$

All terms are scalars

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

Target: calculate $\frac{\partial L}{\partial x}$

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



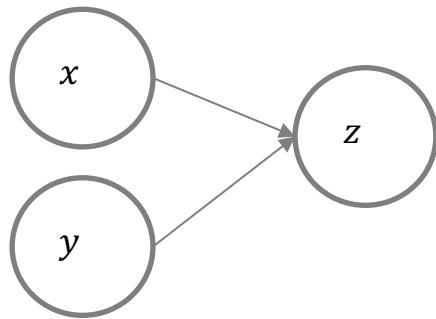
Computing Derivatives – Scalar Addition

$$z = x + y$$

All terms are scalars

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

Target: calculate $\frac{\partial L}{\partial x}$, $\frac{\partial L}{\partial y}$



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

Computing Derivatives – Scalar Multiplication

$$z = Wx$$

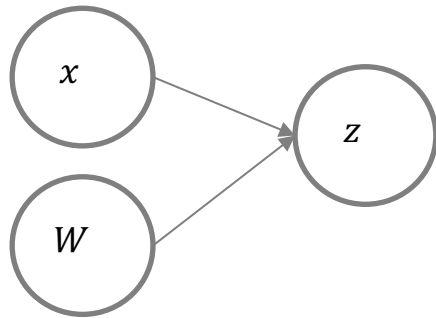
All terms are scalars

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

Target: calculate $\frac{\partial L}{\partial x}$, $\frac{\partial L}{\partial W}$

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

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



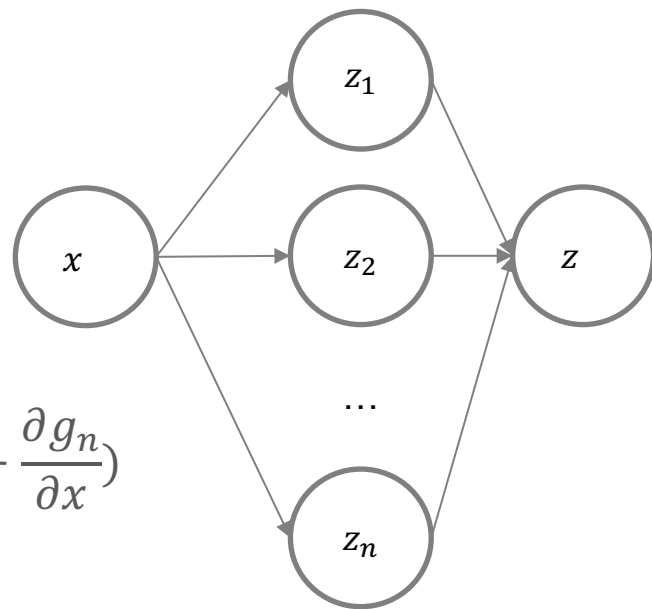
Computing Derivatives – Scalar Generalized Chain Rule

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

All terms are scalars

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

Target: calculate $\frac{\partial L}{\partial x}$



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

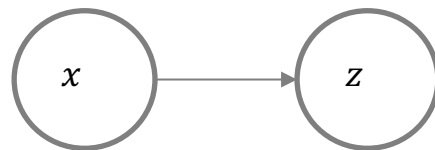
Computing Derivatives – Multivariable Chain Rule

$$z = g(x)$$

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

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

Target: calculate $\nabla_x L$



$$\nabla_x L = \nabla_z L \nabla_x z$$

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

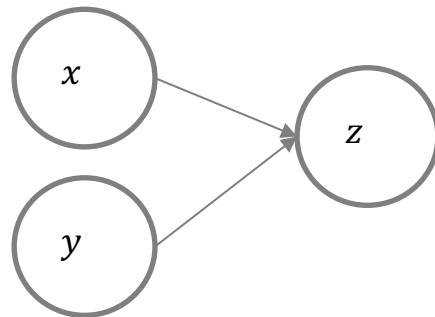
Computing Derivatives – Multivariable Vector Addition

$$z = x + y$$

x, y, z are all $D \times 1$ vectors

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

Target: calculate $\nabla_x L, \nabla_y L$



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

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

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

Computing Derivatives – Multivariable Vector Addition of derivatives

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

$$z = g(x)$$

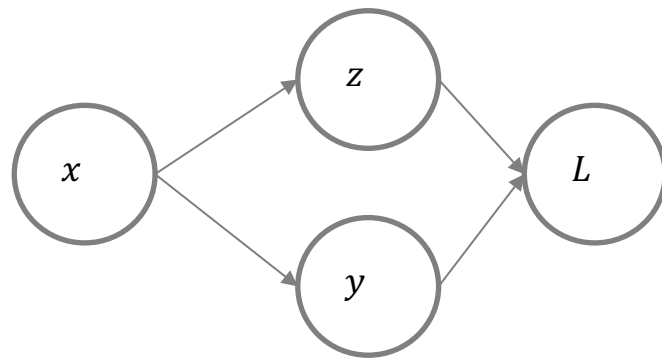
$$y = h(x)$$

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

Target: calculate $\nabla_x L$



$$\nabla_x L = \nabla_z L \nabla_x z + \nabla_y L \nabla_x y$$

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

Computing Derivatives – Multivariable Matrix Multiplication

$$L = f(z)$$

$$z = Wx$$

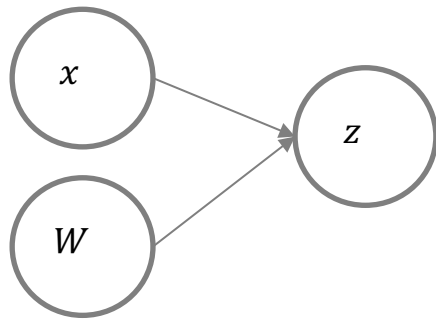
x is a $D \times 1$ vector

z is a $K \times 1$ vector

W is a $K \times D$ matrix

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

Target: calculate $\nabla_x L$, $\nabla_W L$



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

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

Computing Derivatives – Multivariable Generalized Chain Rule

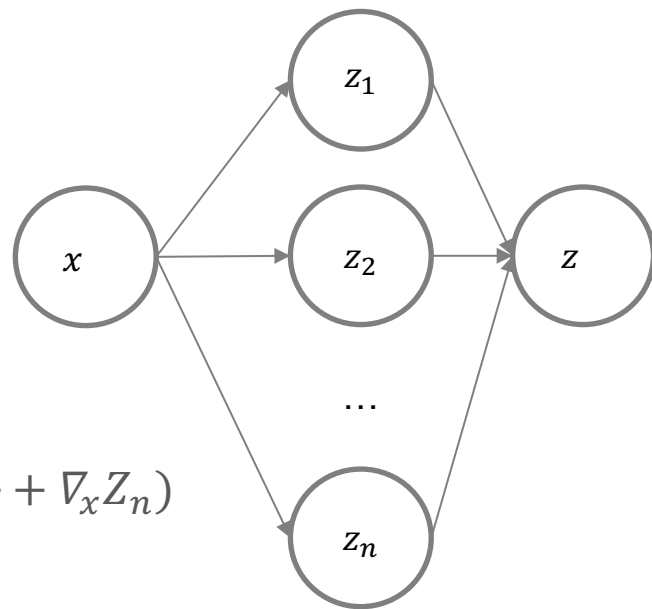
$$z = z_1 + z_2 + \cdots + z_n = g_1(x) + g_2(x) + \cdots + g_n(x)$$

x is a $D \times 1$ vector

z is a $K \times 1$ vector

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

Target: calculate $\nabla_x L$



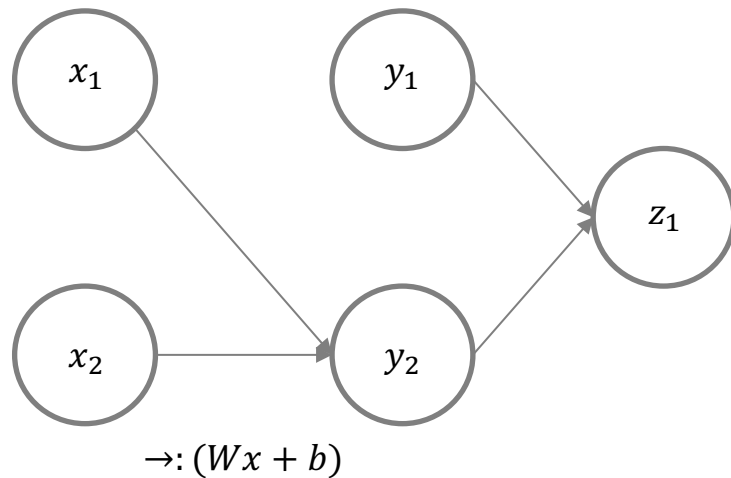
$$\nabla_x L = \nabla_z L \nabla_x Z = \nabla_z L (\nabla_x Z_1 + \nabla_x Z_2 + \cdots + \nabla_x Z_n)$$

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 $z = xy$
 - Work your way backward through the derivatives of the simple relations

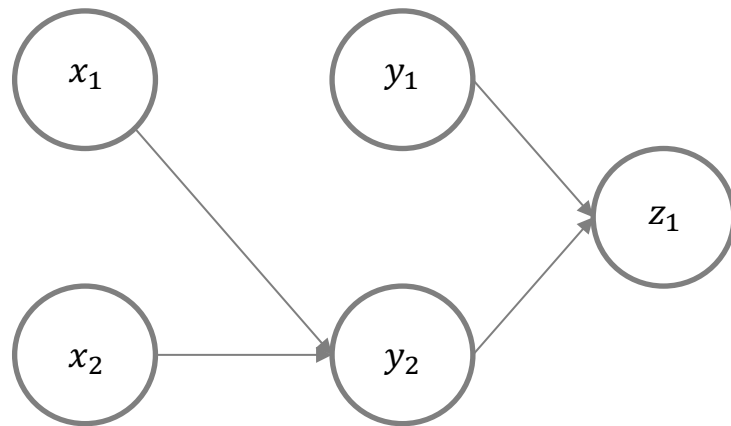
Example:

- $y_2 = \tanh(W_{x_1}x_1 + b_{x_1} + W_{x_2}x_2 + b_{x_2})$
- $z_1 = \tanh(W_{y_1}y_1 + b_{y_1} + W_{y_2}y_2 + b_{y_2})$
- $\nabla_{z_1} L$ is given
- Target:
 - $\nabla_{x_1} L, \nabla_{x_2} L, \nabla_{y_1} L, \nabla_{y_2} L$
 - $\nabla_{W_{x_1}} L, \nabla_{b_{x_1}} L, \nabla_{W_{x_2}} L, \nabla_{b_{x_2}} L$
 - $\nabla_{W_{y_1}} L, \nabla_{b_{y_1}} L, \nabla_{W_{y_2}} L, \nabla_{b_{y_2}} L$



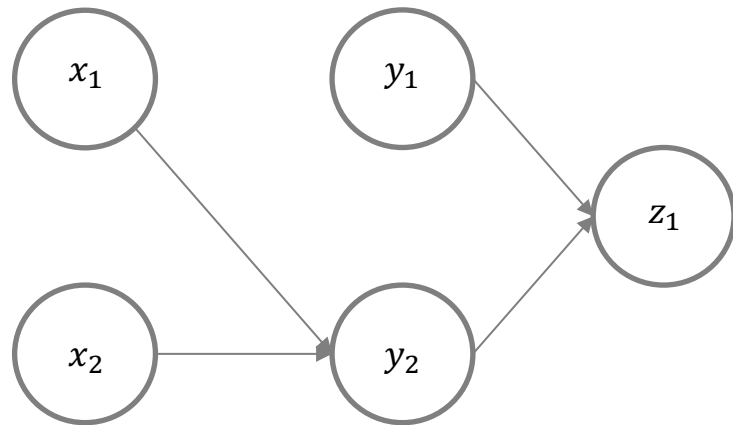
Example:

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



Example:

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

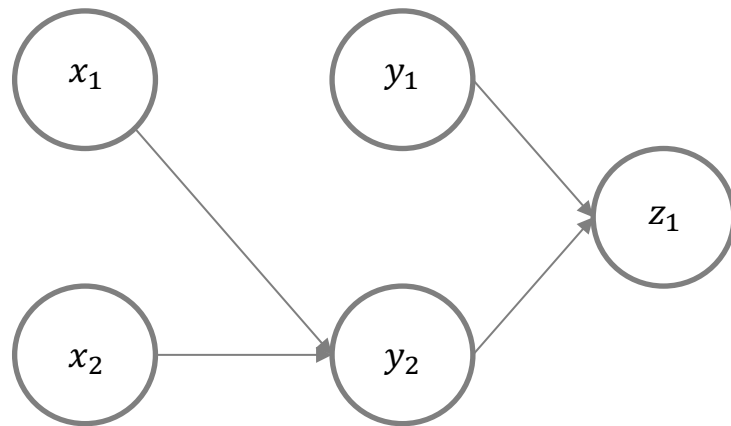


Example:

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

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

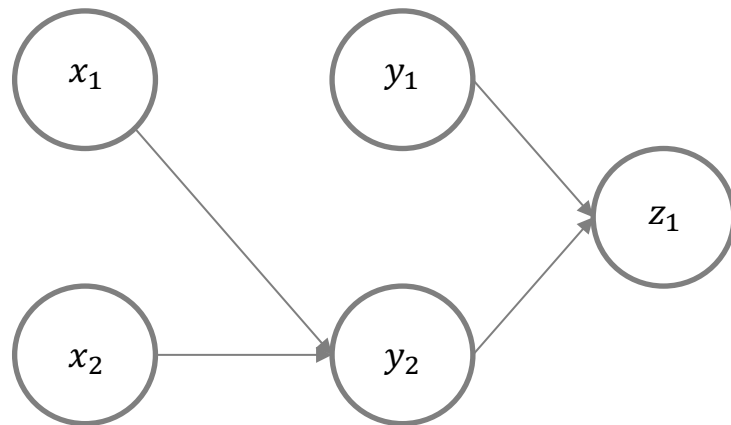
- $i_4 = W_{y_1}y_1$
- $i_5 = W_{y_2}y_2$
- $i_6 = i_4 + b_{y_1} + i_5 + b_{y_2}$
- $z_1 = \tanh(i_6)$



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

Example:

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

Example:

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

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:

- 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:

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

Example:

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

Example:

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

Example:

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

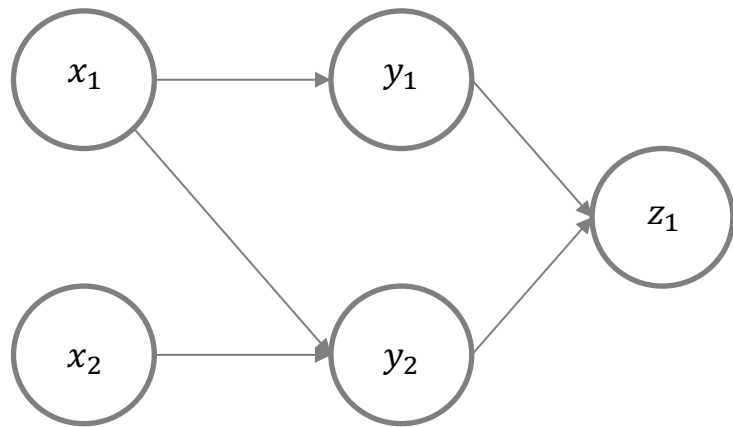
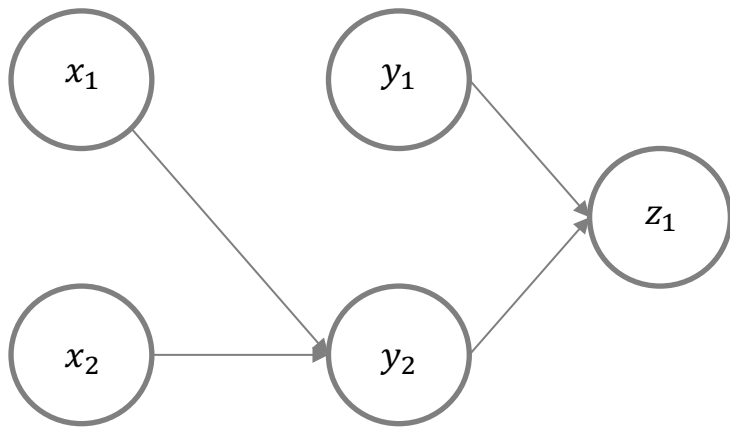
Example:

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

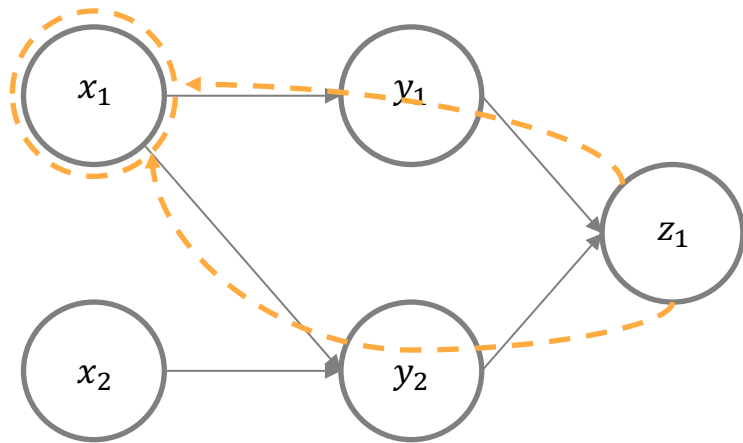
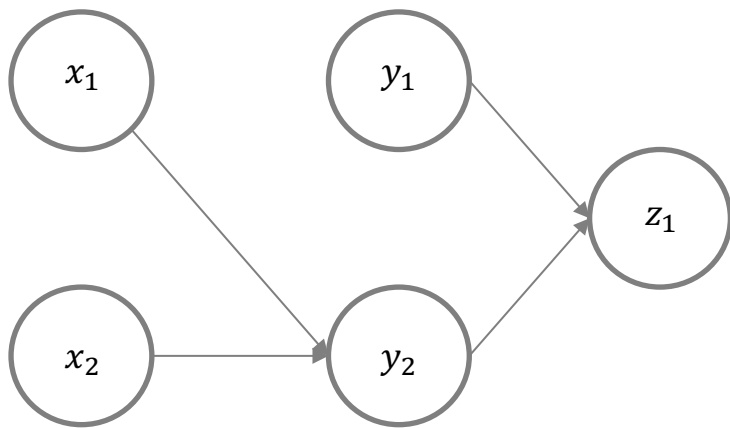
When to use “=” vs “+=”

- 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, the 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

- 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 “+=”?



Please read Prof. Raj's notes about the derivatives and influence diagrams

- <https://piazza.com/class/knsmz2b3z131mn?cid=574>

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>