

# 344.175 VL: Natural Language Processing

## Neural Networks for NLP – a Walkthrough



Navid Rekab-saz

Email: [navid.rekabsaz@jku.at](mailto:navid.rekabsaz@jku.at)

Office hours: <https://navid-officehours.youcanbook.me>

## Notation – recap

- $a \rightarrow$  scalar
- $\mathbf{b} \rightarrow$  vector
  - $i^{\text{th}}$  element of  $\mathbf{b}$  is the scalar  $b_i$
- $\mathbf{C} \rightarrow$  matrix
  - $i^{\text{th}}$  vector of  $\mathbf{C}$  is  $\mathbf{c}_i$
  - $j^{\text{th}}$  element of the  $i^{\text{th}}$  vector of  $\mathbf{C}$  is the scalar  $c_{i,j}$
- Tensor: generalization of scalar, vector, matrix to any arbitrary dimension

# Probability

- Conditional probability, given two random variables  $X$  and  $Y$ :

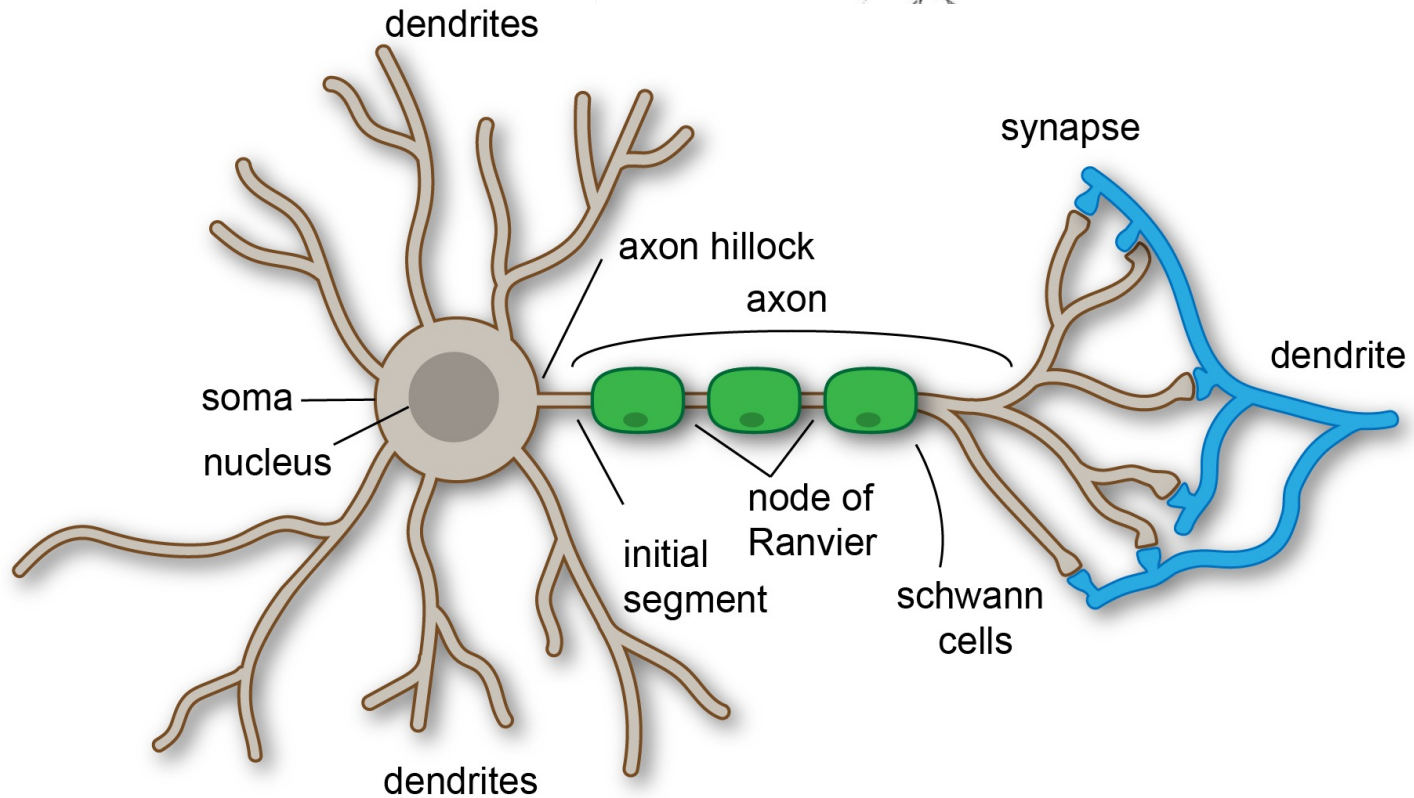
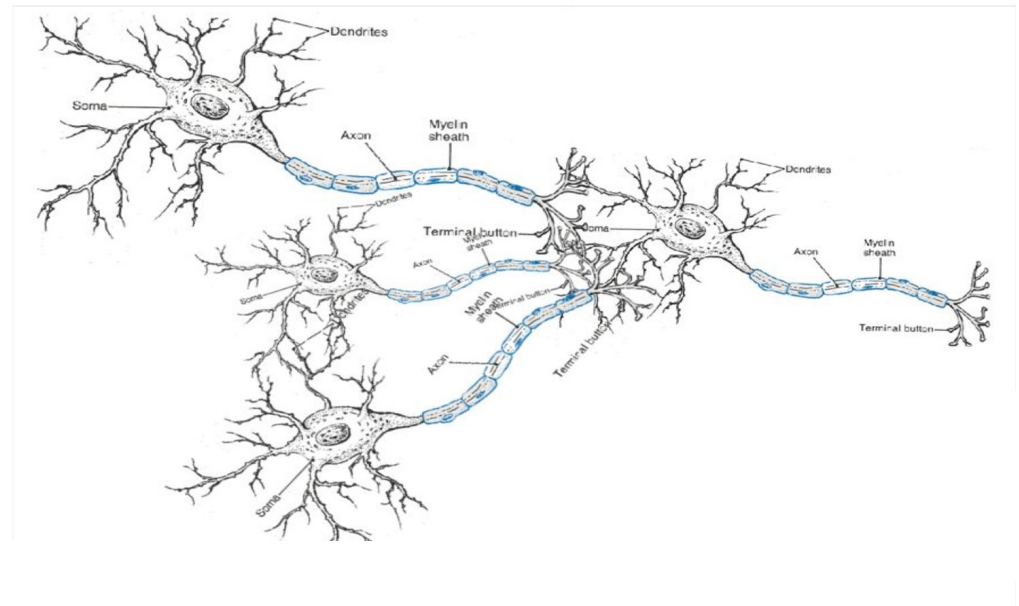
$$P(Y|X)$$

- Probability distribution
  - For a **discrete** random variable  $Y$  with  $K$  states (classes)
    - $0 \leq P(Y_i) \leq 1$
    - $\sum_{i=1}^K P(Y_i) = 1$
  - E.g. with  $K = 4$  states: [0.2 0.3 0.45 0.05]
- Expected value over a set  $\mathcal{D}$

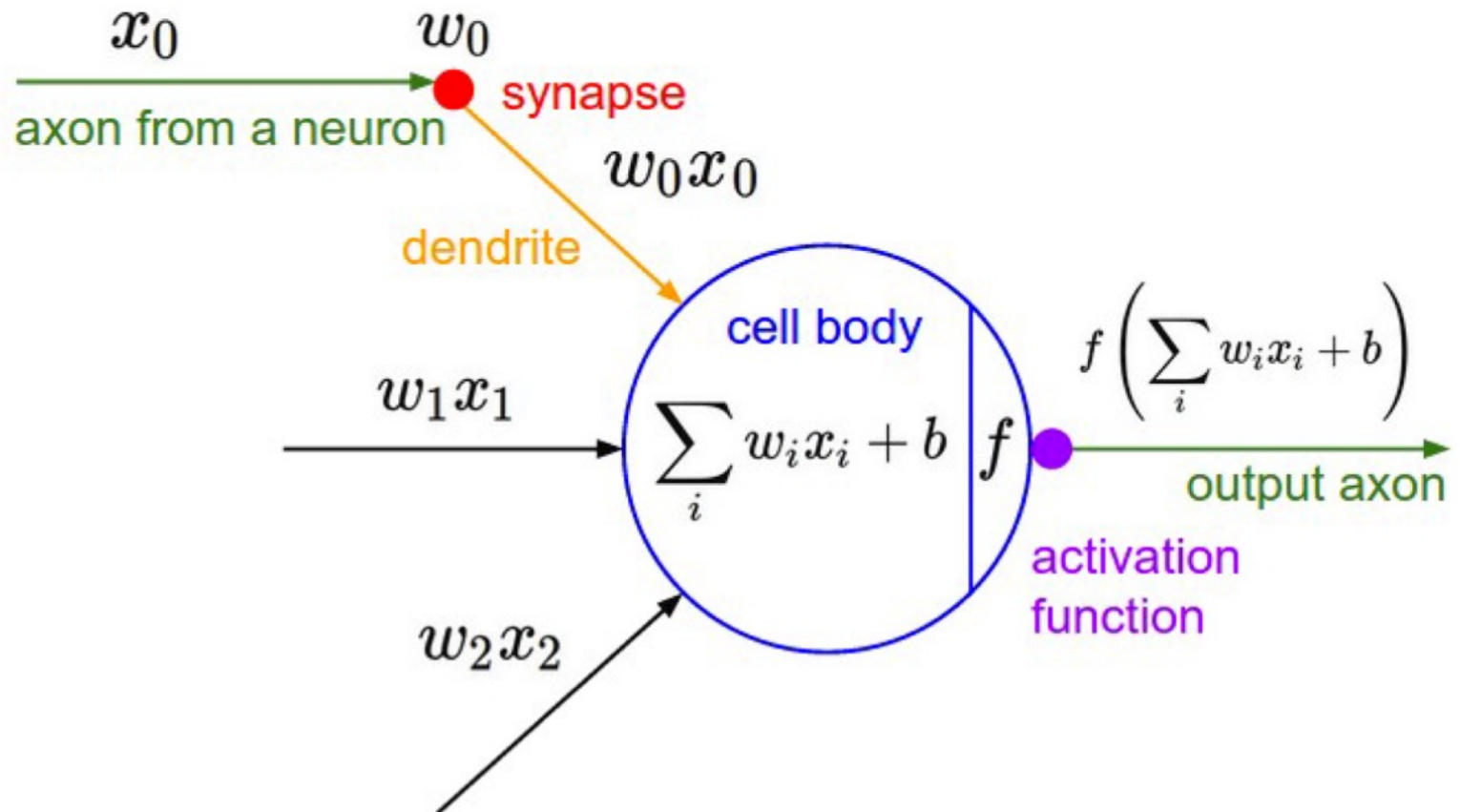
$$\mathbb{E}_{\mathcal{D}}[f] = \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} f(x)$$

Note: The definition of expected value is not completely precise. Though, it suffices for our use in this lecture

# Neural Computation



# An Artificial Neuron



# Learning with Neural Networks

- Design the network's architecture
- Loop until some **exit criteria** are met
  - Sample a **(mini)batch** from training data  $\mathcal{D}$
  - Execute **forward pass**: predict the output tensor of each given input tensor
  - Calculate **loss**
  - **Optimize** the network to reduce loss
    - Calculate the **gradient** of each parameter regarding the loss function using the **backpropagation** algorithm
    - **Update** parameters using their gradients

# Learning with Neural Networks

- Design the network's architecture
- Loop until some **exit criteria** are met
  - Sample a **(mini)batch** from training data  $\mathcal{D}$
  - Execute **forward pass**: predict the output tensor of each given input tensor
  - Calculate **loss**
  - **Optimize** the network to reduce loss
    - Calculate the **gradient** of each parameter regarding the loss function using the **backpropagation** algorithm
    - **Update** parameters using their gradients

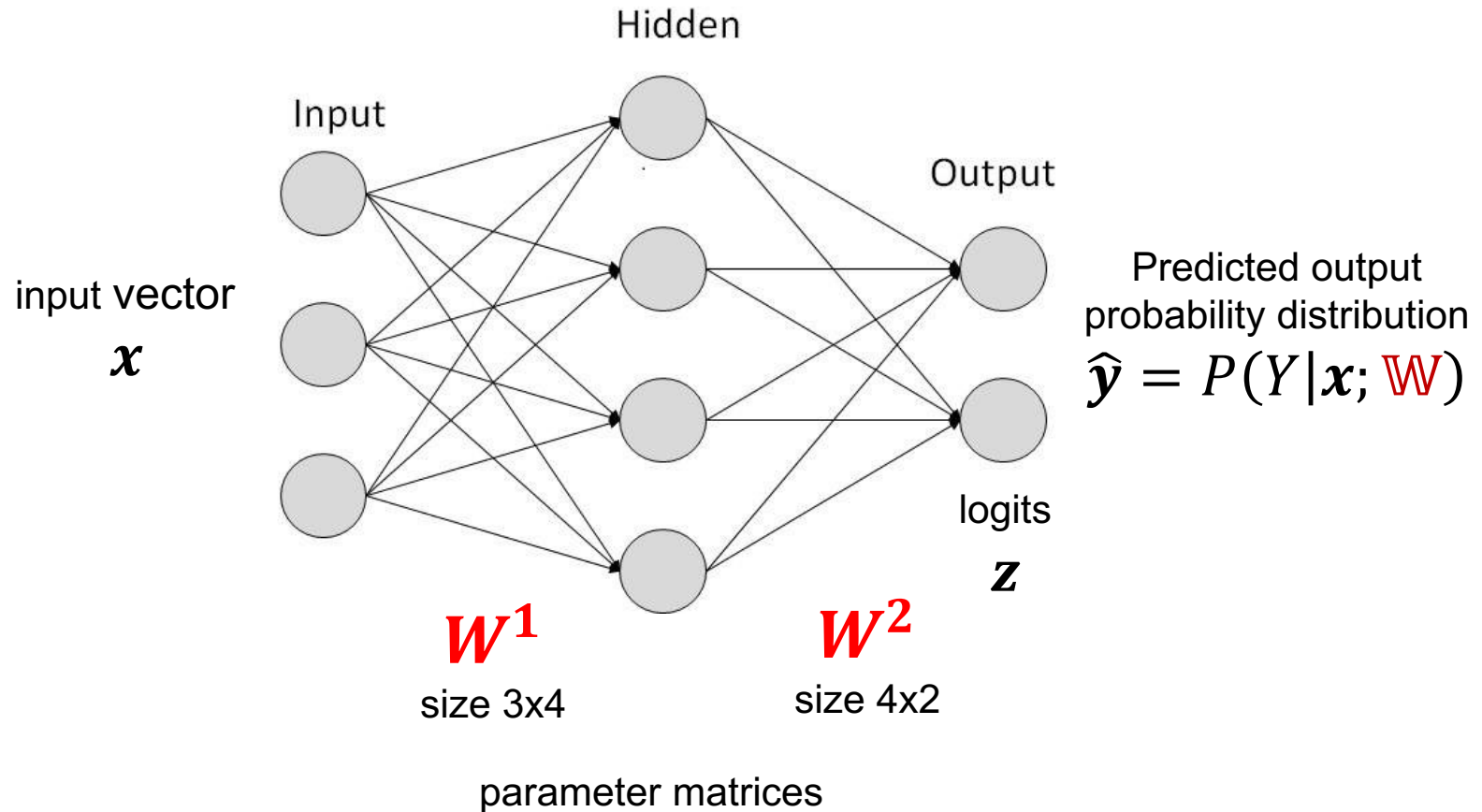
# Artificial Neural Networks

- Neural Networks are **non-linear functions** and **universal approximators**
- Neural networks can readily be defined as **probabilistic models** which estimate  $P(Y|X)$
- Considering model parameter,  $P(Y|X)$  can be written as  $P(Y|\mathbf{x}; \mathbb{W})$ 
  - $\mathbf{x}$  is an input vector and  $\mathbb{W}$  is the set of model parameters
  - The model's predicted probability distribution is:

$$\hat{\mathbf{y}} = P(Y|\mathbf{x}; \mathbb{W})$$



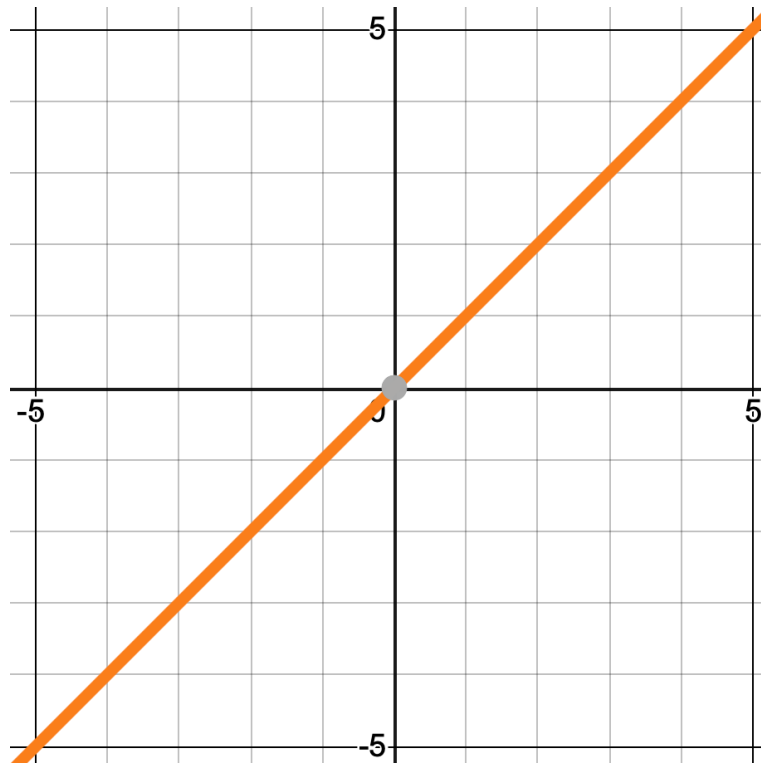
# A sample neural network (Multi Layer Perceptron)



Hidden nodes/layers apply non-linear functions to their inputs

# Linear

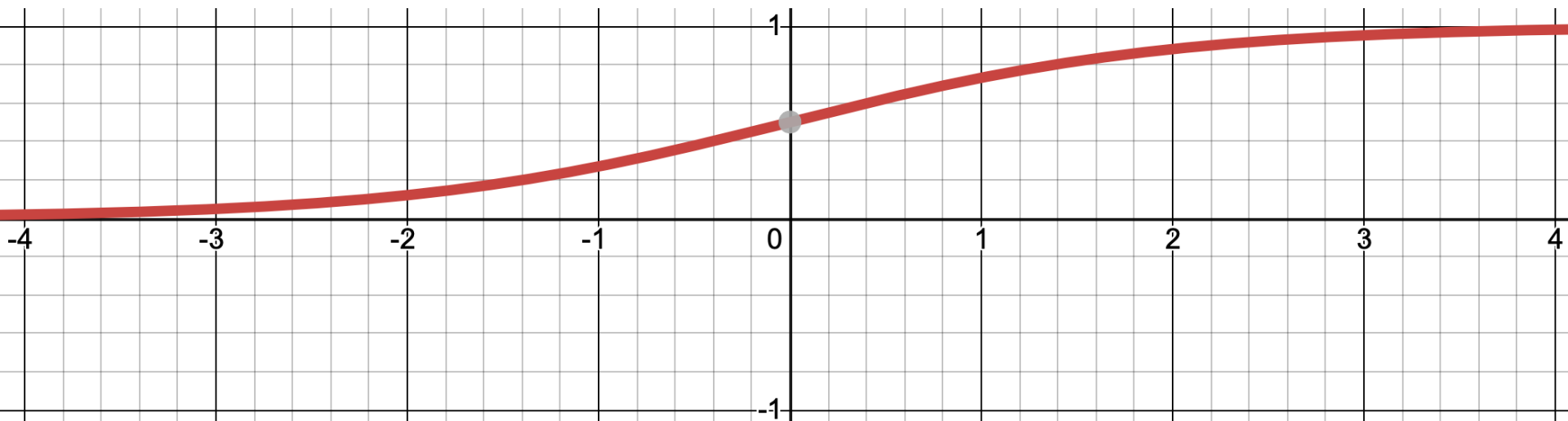
$$f(x) = x$$



## Non-linearities – Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

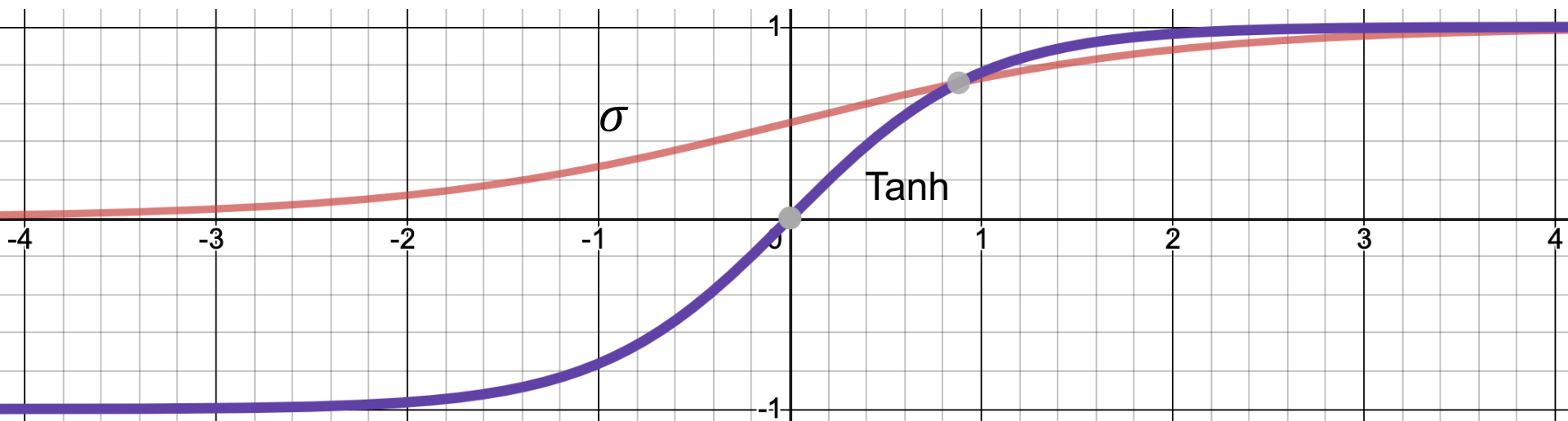
- squashes input between 0 and 1
- Output becomes like a probability value



# Hyperbolic Tangent (Tanh)

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

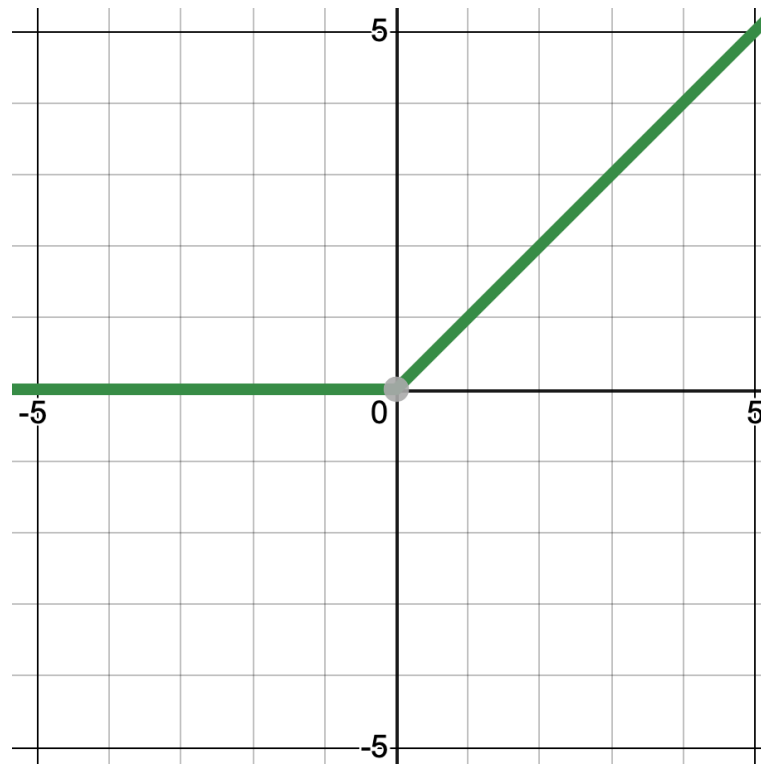
- squashes input between -1 and 1



# Rectified Linear Unit (ReLU)

$$\text{ReLU}(x) = \max(0, x)$$

- fits to deep architectures, as it prevents vanishing gradient



## Examples

$$\mathbf{x} = [1 \quad 3] \quad \mathbf{W} = \begin{bmatrix} 0.5 & -0.5 & 2 & 0 & 0 \\ 0 & 0 & 0 & 4 & -1 \end{bmatrix}$$

- Linear transformation  $\mathbf{xW}$ :

$$\mathbf{xW} = [1 \quad 3] \begin{bmatrix} 0.5 & -0.5 & 2 & 0 & -1 \\ 0 & 0 & 0 & 4 & -1 \end{bmatrix} = [0.5 \quad -0.5 \quad 2 \quad 12 \quad -4]$$

- Non-linear transformation  $\text{ReLU}(\mathbf{xW})$ :

$$\text{ReLU}([0.5 \quad -0.5 \quad 2 \quad 12 \quad -4]) = [0.5 \quad 0.0 \quad 2 \quad 12 \quad 0.0]$$

- Non-linear transformation  $\sigma(\mathbf{xW})$ :

$$\sigma([0.5 \quad -0.5 \quad 2 \quad 12 \quad -4]) = [0.62 \quad 0.37 \quad 0.88 \quad 0.99 \quad 0.11]$$

- Non-linear transformation  $\tanh(\mathbf{xW})$ :

$$\tanh([0.5 \quad -0.5 \quad 2 \quad 12 \quad -4]) = [0.46 \quad -0.46 \quad 0.96 \quad 0.99 \quad -0.99]$$

# Learning with Neural Networks

- Design the network's architecture
- Loop until some **exit criteria** are met
  - Sample a **(mini)batch** from training data  $\mathcal{D}$
  - Execute **forward pass**: predict the output tensor of each given input tensor
  - Calculate **loss**
  - **Optimize** the network to reduce loss
    - Calculate the **gradient** of each parameter regarding the loss function using the **backpropagation** algorithm
    - **Update** parameters using their gradients

# Early Stopping

- Run the model for several steps (epochs), and in each step evaluate the model on the validation set
- Store the model if the evaluation results improve
- At the end, take the stored model with the best validation results as the final model



# Learning with Neural Networks

- Design the network's architecture
- Loop until some **exit criteria** are met
  - Sample a **(mini)batch** from training data  $\mathcal{D}$
  - Execute **forward pass**: predict the output tensor of each given input tensor
  - Calculate **loss**
  - **Optimize** the network to reduce loss
    - Calculate the **gradient** of each parameter regarding the loss function using the **backpropagation** algorithm
    - **Update** parameters using their gradients

# Toy neural network

- A sample neural network is going to calculate the following function:

$$z(x; \mathbb{W}) = w_2^2 * (2 * x * w_1 + w_0)$$

- $x$  is input and  $\mathbb{W}$  is the tensor of parameters
- Parameters are initialized with

$$w_0 = 1 \quad w_1 = 3 \quad w_2 = 2$$

- A neural network first redefines this function as subfunctions of basic/atomic operations with new intermediary variables:\*

$$a = 2 * x * w_1$$

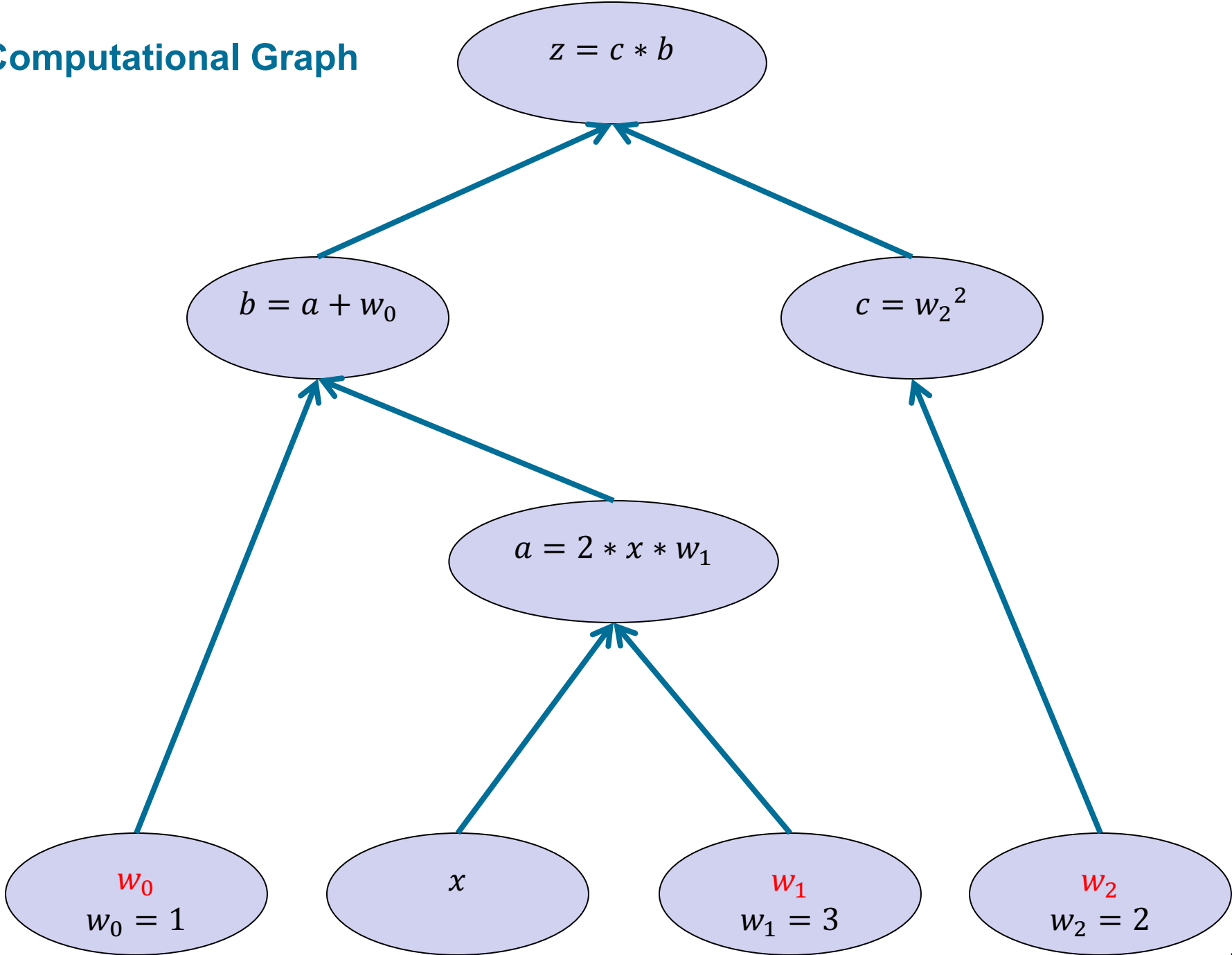
$$b = a + w_0$$

$$c = w_2^2$$

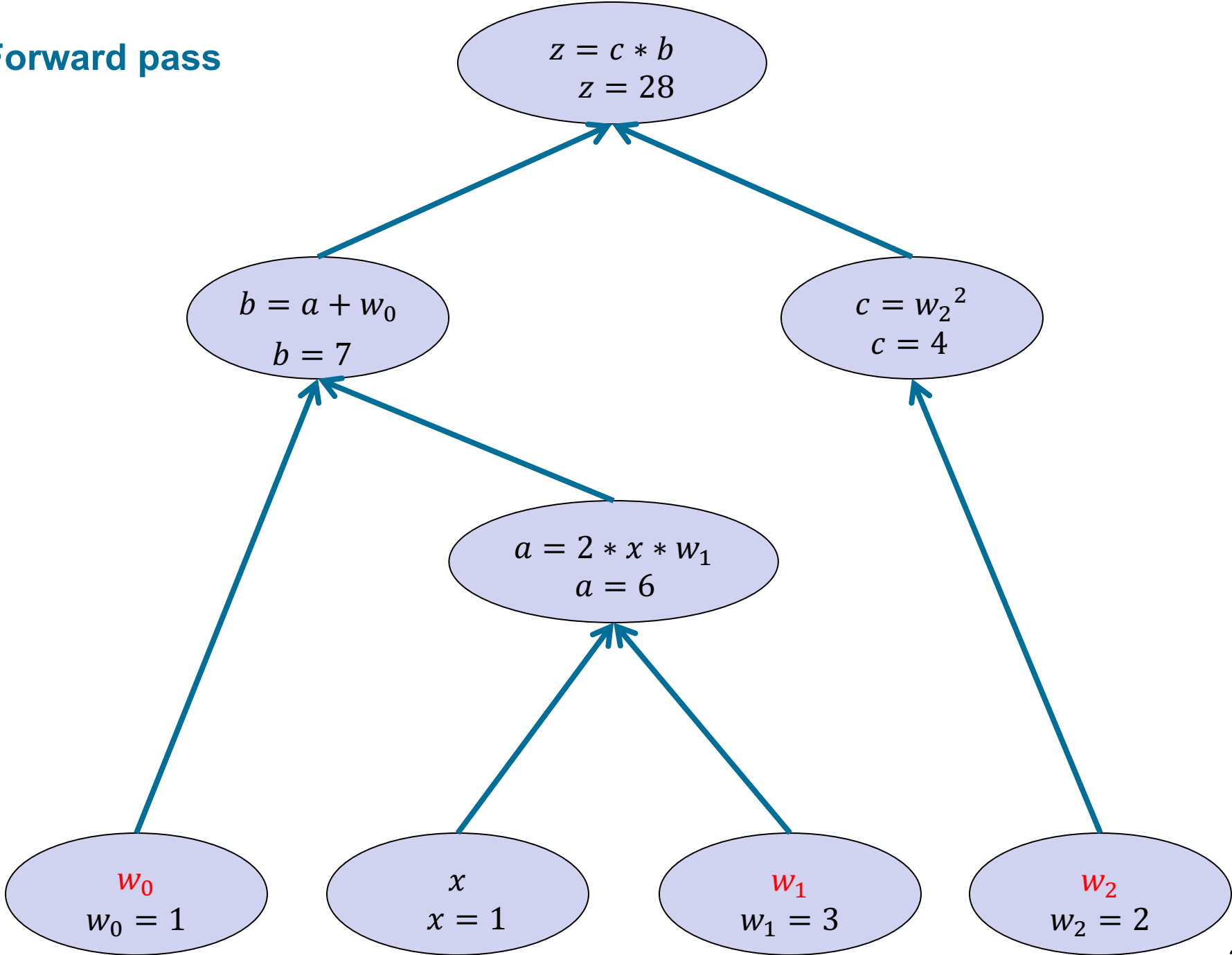
$$z = c * b$$

\* To keep the example simple, the splitting is not applied to all basic operation

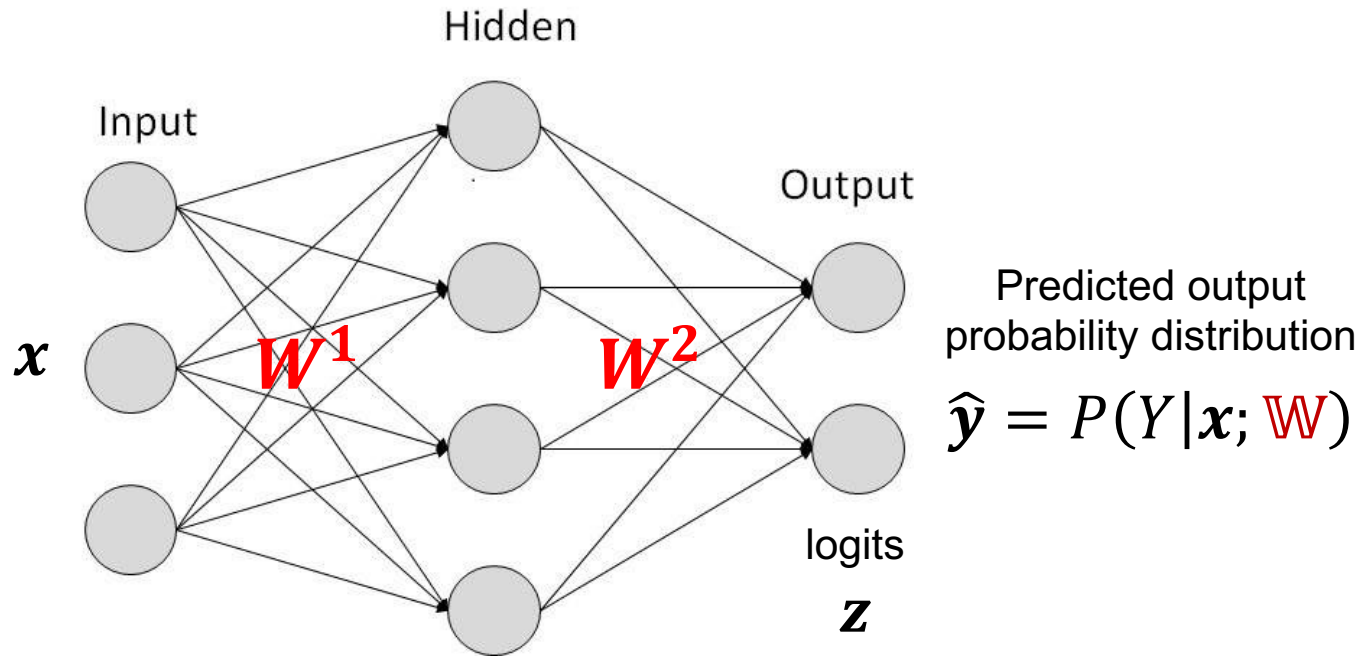
# Computational Graph



# Forward pass




# Output probability distribution



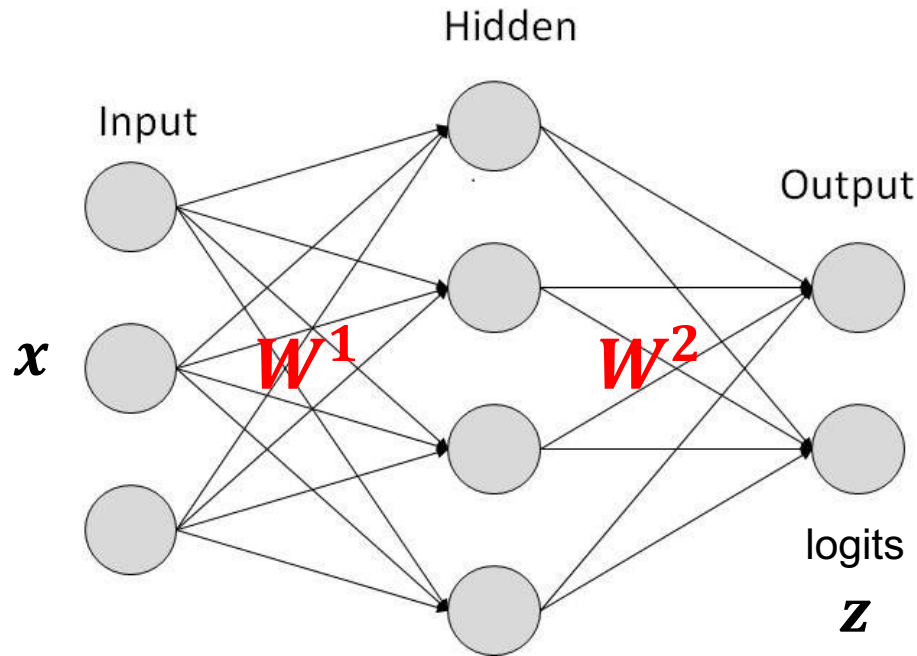
# Softmax

- As discussed, neural networks can readily turn to probabilistic models
- To do it, we need to transform the **output vector  $\mathbf{z}$**  of a neural network with  **$K$  output classes** to a probability distribution
  - In the context of neural networks,  $\mathbf{z}$  is usually called **logits**
- softmax turns a vector to a probability distribution
  - $\mathbf{z}$  could be the output vector of a neural network

$$\text{softmax}(\mathbf{z})_l = \frac{e^{z_l}}{\sum_{i=1}^K e^{z_i}}$$


normalization term

# Output probability distribution



Predicted output  
probability distribution

$$\hat{\mathbf{y}} = P(Y|\mathbf{x}; \mathbb{W})$$
$$= \text{softmax}(\mathbf{z})$$

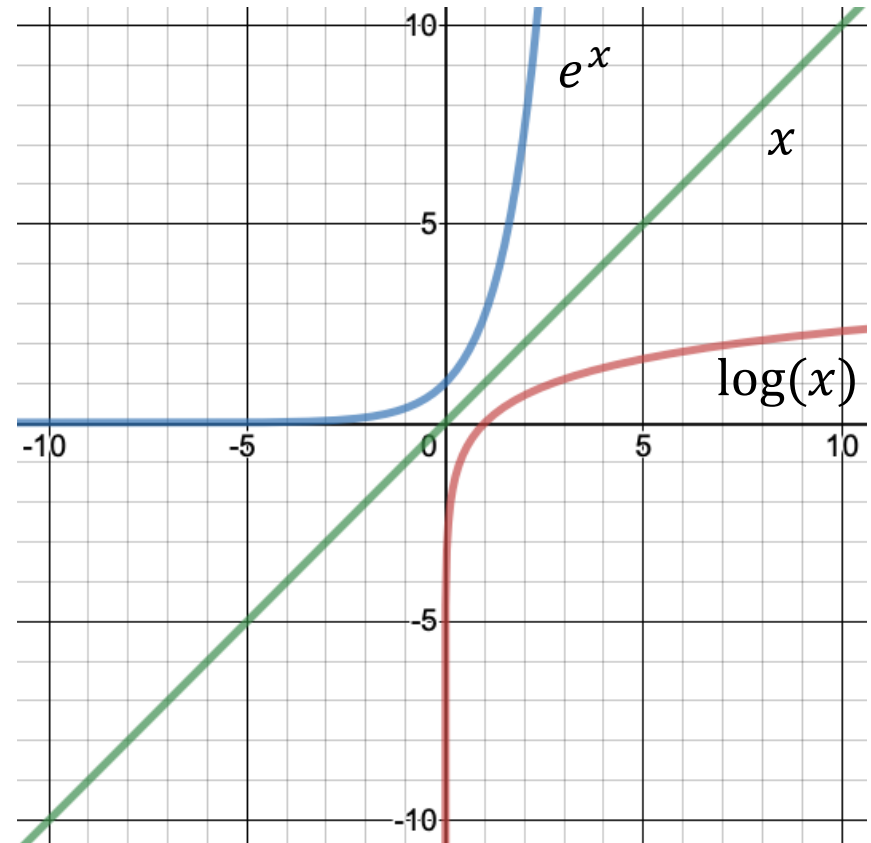
# Softmax – example

$K = 4$  classes

$$\text{softmax}(\mathbf{z})_l = \frac{e^{z_l}}{\sum_{i=1}^K e^{z_i}}$$

$$\mathbf{z} = \begin{bmatrix} 1 \\ 2 \\ 5 \\ 6 \end{bmatrix}$$

$$\text{softmax}(\mathbf{z}) = \begin{bmatrix} 0.004 \\ 0.013 \\ 0.264 \\ 0.717 \end{bmatrix}$$





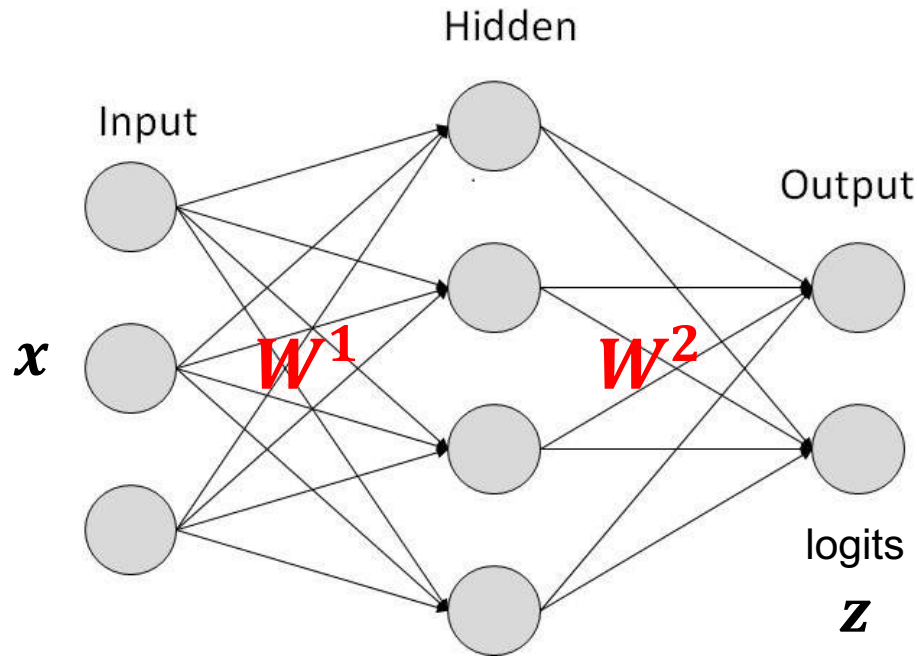
## Softmax characteristics

- The exponential function in softmax makes the **maximum** becomes much higher than the others
- Softmax identifies the “*max*” but in a “*soft*” way!
- Softmax imposes competition between the predicted output values, as in fact “*winner takes (almost) all!*”
  - Winner-takes-all is the case when one value is 1 and the rest are 0
  - Softmax provides a soft distribution of winner-takes-all
  - This resembles the competition between nearby neurons in the cortex

# Learning with Neural Networks

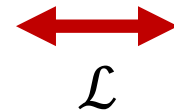
- Design the network's architecture
- Loop until some **exit criteria** are met
  - Sample a **(mini)batch** from training data  $\mathcal{D}$
  - Execute **forward pass**: predict the output tensor of each given input tensor
  - Calculate **loss**
  - **Optimize** the network to reduce loss
    - Calculate the **gradient** of each parameter regarding the loss function using the **backpropagation** algorithm
    - **Update** parameters using their gradients

# Sample neural network



Predicted output  
probability distribution

$$\hat{\mathbf{y}} = P(Y|\mathbf{x}; \mathbb{W}) \\ = \text{softmax}(\mathbf{z})$$



Labels

$\mathbf{y}$

# Cross Entropy Loss

- Given a classification task with  $K$  classes
  - known as **multi-class classification**
- $\hat{\mathbf{y}}$  → predicted probability distribution of the classes
- $\mathbf{y}$  → actual probability distribution of the classes (labels)
- Cross Entropy loss is defined as:

$$\mathcal{L} = -\mathbb{E}_{\mathcal{D}} \sum_{i=1}^K y_i \log \hat{y}_i$$

- $\mathcal{D}$  → the set of training data
- In neural networks, we can write it as:

$$\mathcal{L}(\mathbf{W}) = -\mathbb{E}_{\mathcal{D}} \sum_{i=1}^K y_i \log P(Y_i|\mathbf{x}; \mathbf{W})$$

# Cross Entropy Loss – example 1

- A multi-label scenario:

$$\hat{\mathbf{y}} = \begin{bmatrix} 0.004 \\ 0.013 \\ 0.264 \\ 0.717 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 0 \\ 0.25 \\ 0 \\ 0.75 \end{bmatrix}$$

$$\mathcal{L} = - \sum_{i=1}^K y_i \log \hat{y}_i$$

$$\mathcal{L} = -(0 \times \log 0.004 + 0.25 \times \log 0.013 + 0 \times \log 0.264 + 0.75 \times \log 0.717)$$

$$\mathcal{L} = -(0 - 0.471 + 0 - 0.108)$$

$$\mathcal{L} = 0.579$$

## Cross Entropy Loss – example 2

- A single-label scenario:

$$\hat{\mathbf{y}} = \begin{bmatrix} 0.004 \\ 0.013 \\ 0.264 \\ 0.717 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\mathcal{L} = - \sum_{i=1}^K y_i \log \hat{y}_i$$

$$\mathcal{L} = -(0 \times \log 0.004 + 0 \times \log 0.013 + 0 \times \log 0.264 + 1 \times \log 0.717)$$

$$\mathcal{L} = -(0 + 0 + 0 - 0.144)$$

$$\mathcal{L} = 0.144$$

# Negative Log Likelihood (NLL) Loss

- Single-label classification is the most common scenario
- In this case, we can simplify Cross Entropy formulation to

$$\mathcal{L}(\mathbb{W}) = -\mathbb{E}_{\mathcal{D}} \sum_{i=1}^K y_i \log P(Y_i|\mathbf{x}; \mathbb{W}) = -\mathbb{E}_{\mathcal{D}} \log P(Y_l|\mathbf{x}; \mathbb{W})$$

- where  $l$  is the index of the correct class
- This loss function is known as **Negative Log Likelihood (NLL)**
  - NLL is a special case of Cross Entropy

## NLL + softmax

- What happens when we use NLL and softmax in the output layer of a neural network?

$$\mathcal{L}(\mathbb{W}) = -\mathbb{E}_{\mathcal{D}} \log P(Y_l | \mathbf{x}; \mathbb{W}) = -\mathbb{E}_{\mathcal{D}} \log \text{softmax}(\mathbf{z})_l$$

$\mathbf{z} \rightarrow$  output vector before softmax (logits)

$$\mathcal{L}(\mathbb{W}) = -\mathbb{E}_{\mathcal{D}} \log \frac{e^{z_l}}{\sum_{i=1}^K e^{z_i}} = -\mathbb{E}_{\mathcal{D}} \left[ \log e^{z_l} - \log \sum_{i=1}^K e^{z_i} \right]$$

$$\mathcal{L}(\mathbb{W}) = -\mathbb{E}_{\mathcal{D}} \left[ z_l - \underbrace{\log \sum_{i=1}^K e^{z_i}} \right]$$

This term is (almost)  
equal to  $\max(\mathbf{z})$



## NLL + softmax – example 1

$$\mathcal{L} = - \left[ z_l - \log \sum_{i=1}^K e^{z_i} \right]$$

$$\mathbf{z} = [1 \quad 2 \quad 0.5 \quad 6]$$

- If the correct class is the first one,  $l = 1$ :

$$\mathcal{L} = -[1 - \log(e^1 + e^2 + e^{0.5} + e^6)] = -1 + 6.02 = \mathbf{5.02}$$

- If the correct class is the third one,  $l = 3$ :

$$\mathcal{L} = -[0.5 - \log(e^1 + e^2 + e^{0.5} + e^6)] = -0.5 + 6.02 = \mathbf{5.52}$$

- If the correct class is the fourth one,  $l = 4$ :

$$\mathcal{L} = -[6 - \log(e^1 + e^2 + e^{0.5} + e^6)] = -6 + 6.02 = \mathbf{0.02}$$

## NLL + softmax – example 2

$$\mathcal{L} = - \left[ z_l - \log \sum_{i=1}^K e^{z_i} \right]$$

$$\mathbf{z} = [1 \quad 2 \quad 5 \quad 6]$$

- If the correct class is the first one,  $l = 1$ :

$$\mathcal{L} = -[1 - \log(e^1 + e^2 + e^5 + e^6)] = -1 + 6.33 = \mathbf{5.33}$$

- If the correct class is the third one,  $l = 3$ :

$$\mathcal{L} = -[5 - \log(e^1 + e^2 + e^5 + e^6)] = -5 + 6.33 = \mathbf{1.33}$$

- If the correct class is the fourth one,  $l = 4$ :

$$\mathcal{L} = -[6 - \log(e^1 + e^2 + e^5 + e^6)] = -6 + 6.33 = \mathbf{0.33}$$

# Learning with Neural Networks

- Design the network's architecture
- Loop until some **exit criteria** are met
  - Sample a **(mini)batch** from training data  $\mathcal{D}$
  - Execute **forward pass**: predict the output tensor of each given input tensor
  - Calculate **loss** function of the (mini)batch
  - **Optimize** the network to reduce loss
    - Calculate the **gradient** of each parameter regarding the loss function using the **backpropagation** algorithm
    - **Update** parameters using their gradients

# Toy neural network

$$z(x; \mathbb{W}) = w_2^2 * (2 * x * w_1 + w_0)$$

- Initialization:  $w_0 = 1$   $w_1 = 3$   $w_2 = 2$
- Intermediary variables:

$$a = 2 * x * w_1$$

$$b = a + w_0$$

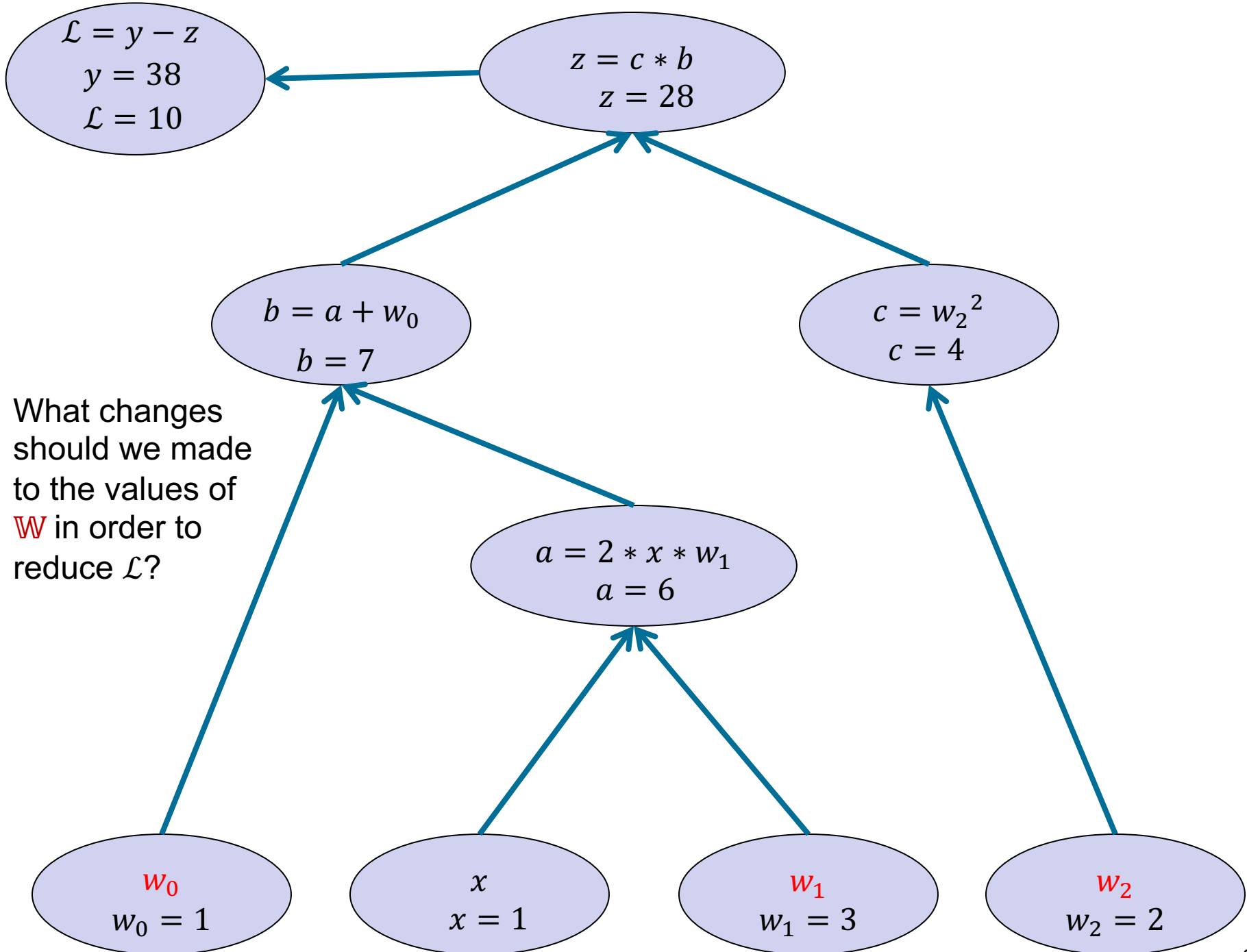
$$c = w_2^2$$

$$z = c * b$$

- An “imaginary” loss:

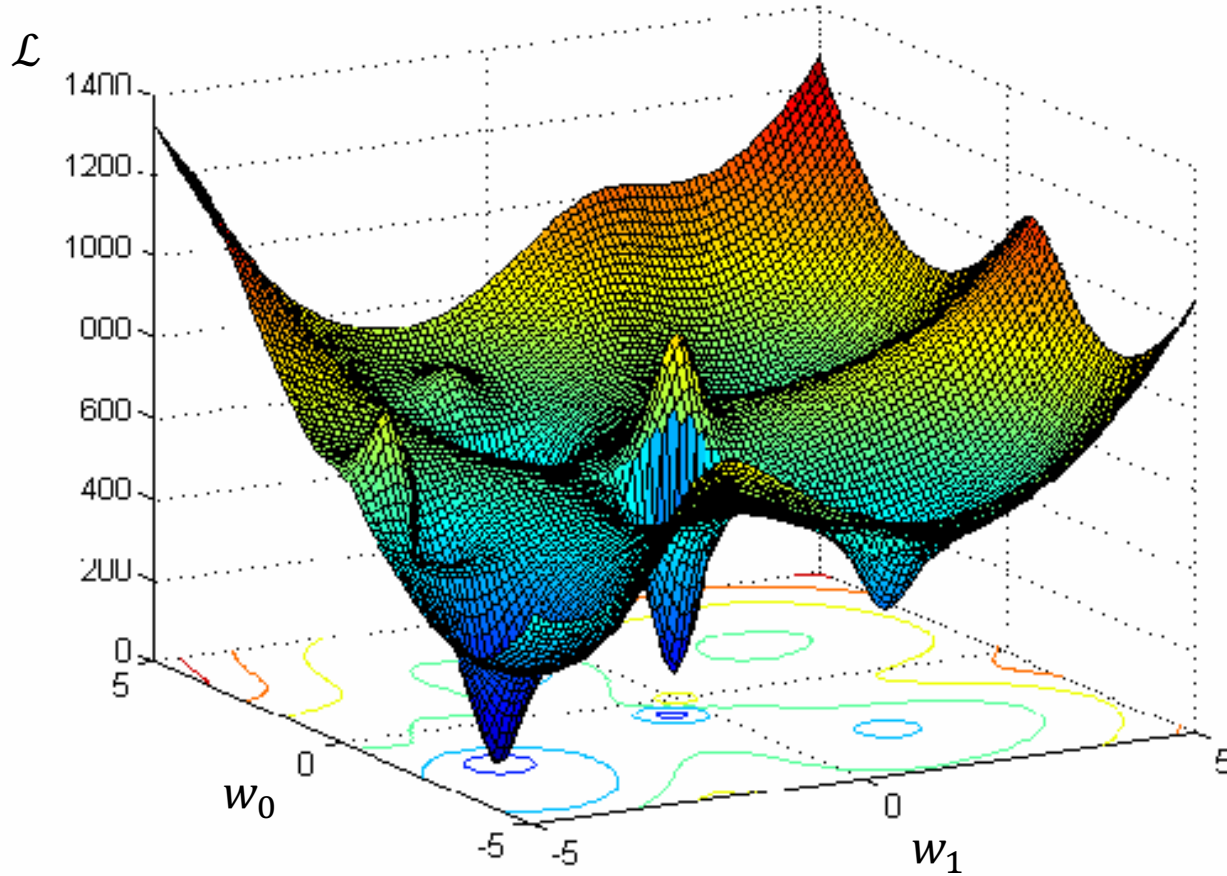
$$\mathcal{L} = y - z$$

For the current datapoint  $x$  we have  $y = 38$



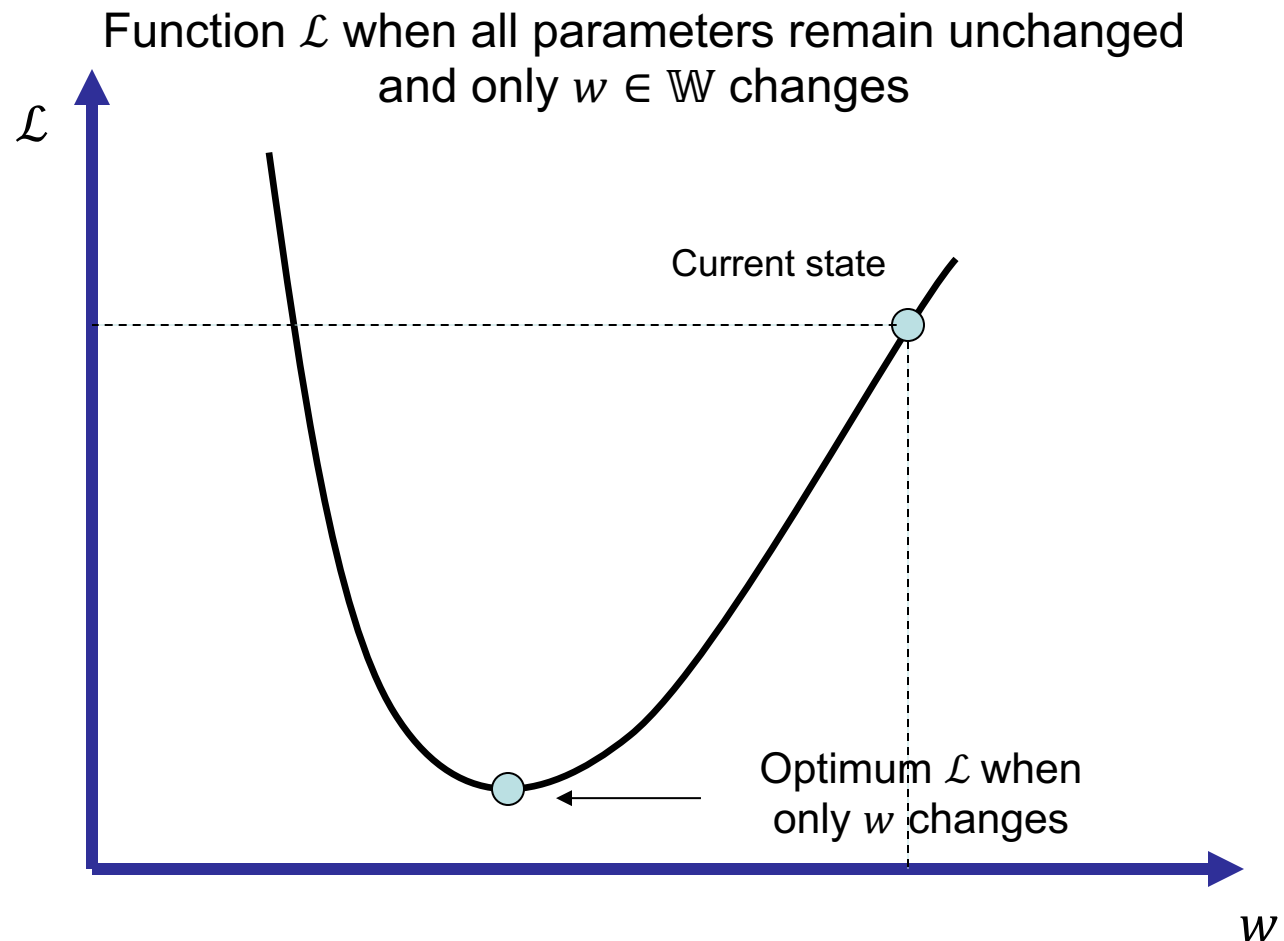
What changes should we made to the values of  $W$  in order to reduce  $\mathcal{L}$ ?

# Optimization



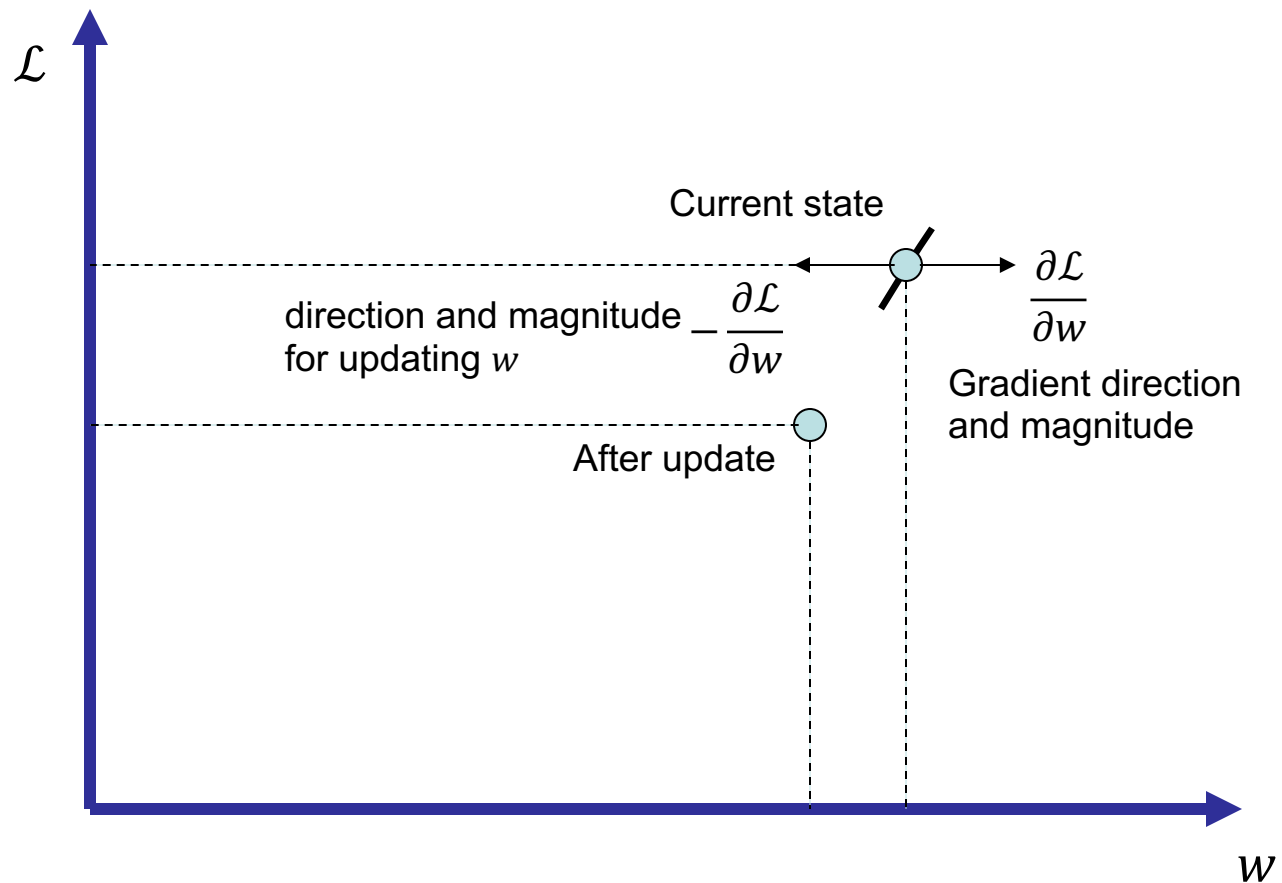
# Gradient-based optimization

- Assumption 1: optimize  $\mathcal{L}$  in respect to each parameter  $w \in \mathbb{W}$  independently regardless of other parameters



# Gradient Descent optimization

- Assumption 2: decide about your course of change for  $w \in \mathbb{W}$  according to the local changes in  $\mathcal{L}$





# Gradient Descent optimization

- We hence need the derivatives of  $\mathcal{L}$  in respect to each  $w \in \mathbb{W}$ :

$$\nabla_{\mathbb{W}}\mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_0} & \frac{\partial \mathcal{L}}{\partial w_1} & \frac{\partial \mathcal{L}}{\partial w_2} & \dots \end{bmatrix}$$

- $\nabla_{\mathbb{W}}\mathcal{L}$  is often called **gradient tensor**, whose elements are the **partial derivatives** of  $\mathcal{L}$  in respect to each parameter:

# Gradient Descent algorithm

- A model with parameters  $\mathbb{W}$  at time step  $t \rightarrow \mathbb{W}^{(t)}$ , **learning rate**  $\eta$ , and set of datapoints  $\mathcal{D}$
- Loop for some epochs
  - Compute gradient tensor  $\mathbb{G}$  of parameters  $\mathbb{W}$  averaged over datapoints  $\mathcal{D}$ :

$$\mathbb{G} \leftarrow \frac{1}{|\mathcal{D}|} \nabla_{\mathbb{W}} \sum_{(x,y) \in \mathcal{D}} \mathcal{L}(x, y; \mathbb{W})$$

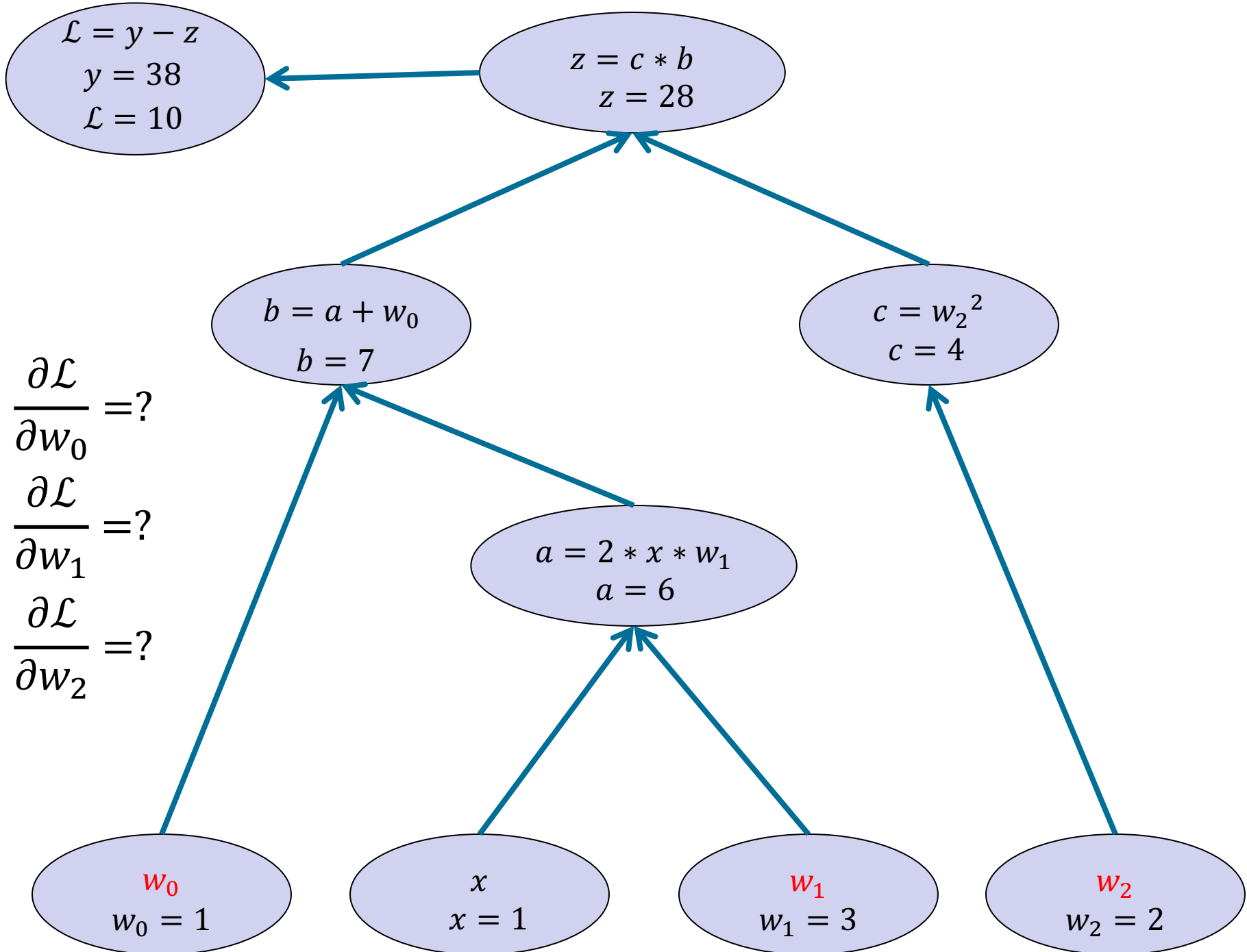
- Update the parameters by taking steps in the opposite direction of the gradient tensor multiplied by  $\eta$ :

$$\mathbb{W}^{(t+1)} \leftarrow \mathbb{W}^{(t)} - \eta \mathbb{G}$$

- Reduce learning rate (**annealing**) if some criteria are met or according to a scheduler

# Learning with Neural Networks

- Design the network's architecture
- Loop until some **exit criteria** are met
  - Sample a **(mini)batch** from training data  $\mathcal{D}$
  - Execute **forward pass**: predict the output tensor of each given input tensor
  - Calculate **loss**
  - **Optimize** the network to reduce loss
    - Calculate the **gradient** of each parameter regarding the loss function using the **backpropagation** algorithm
    - **Update** parameters using their gradients



## Chain rule

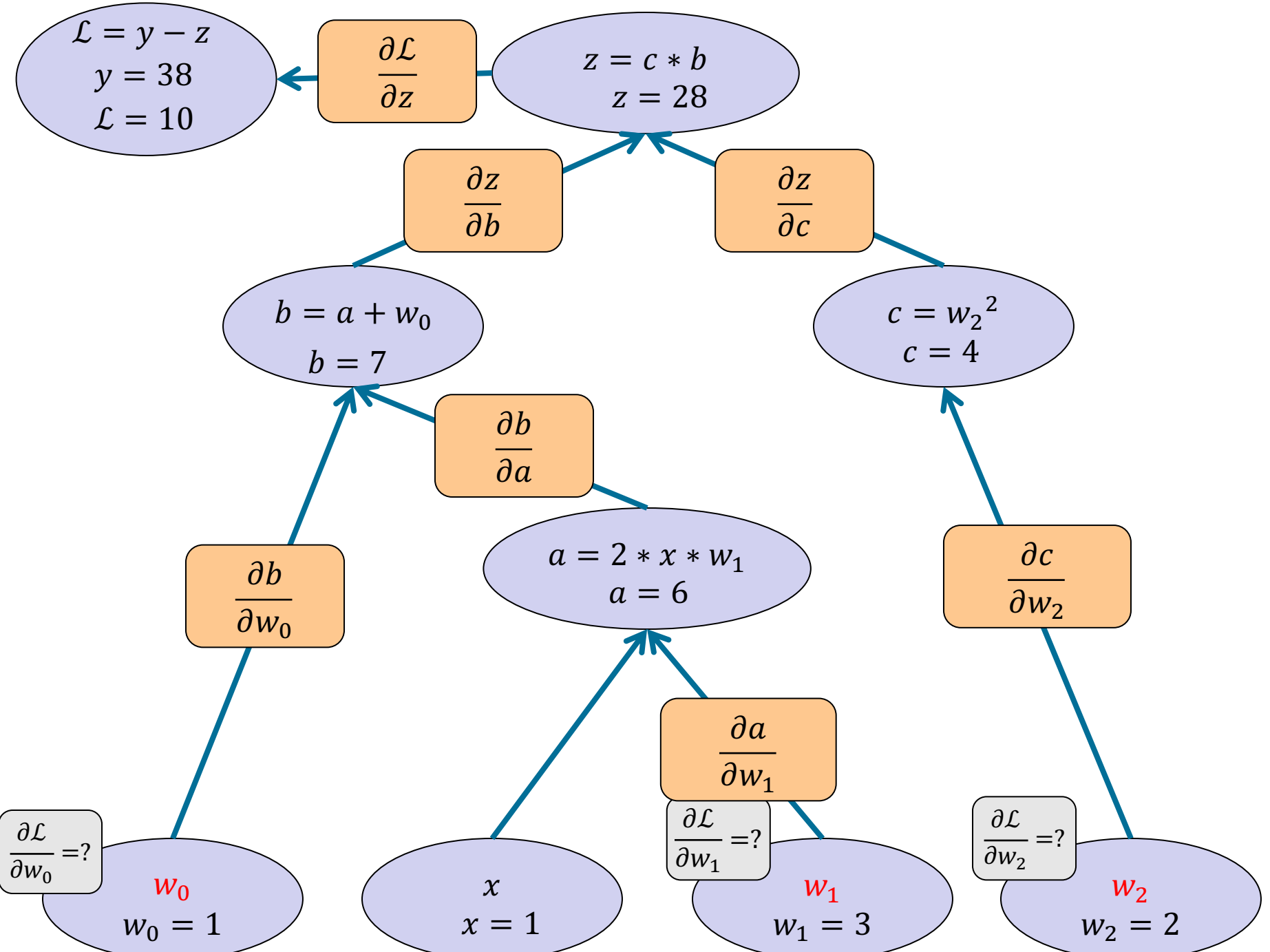
- Gradient tensor:  $\nabla_{\mathbb{W}}\mathcal{L} = \left[ \frac{\partial \mathcal{L}}{\partial w_0} =? \quad \frac{\partial \mathcal{L}}{\partial w_1} =? \quad \frac{\partial \mathcal{L}}{\partial w_2} =? \right]$
- Partial derivatives can be calculated using **local derivatives** and the **chain rule**:

$$\frac{\partial \mathcal{L}}{\partial w_0} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} \frac{\partial b}{\partial w_0}$$

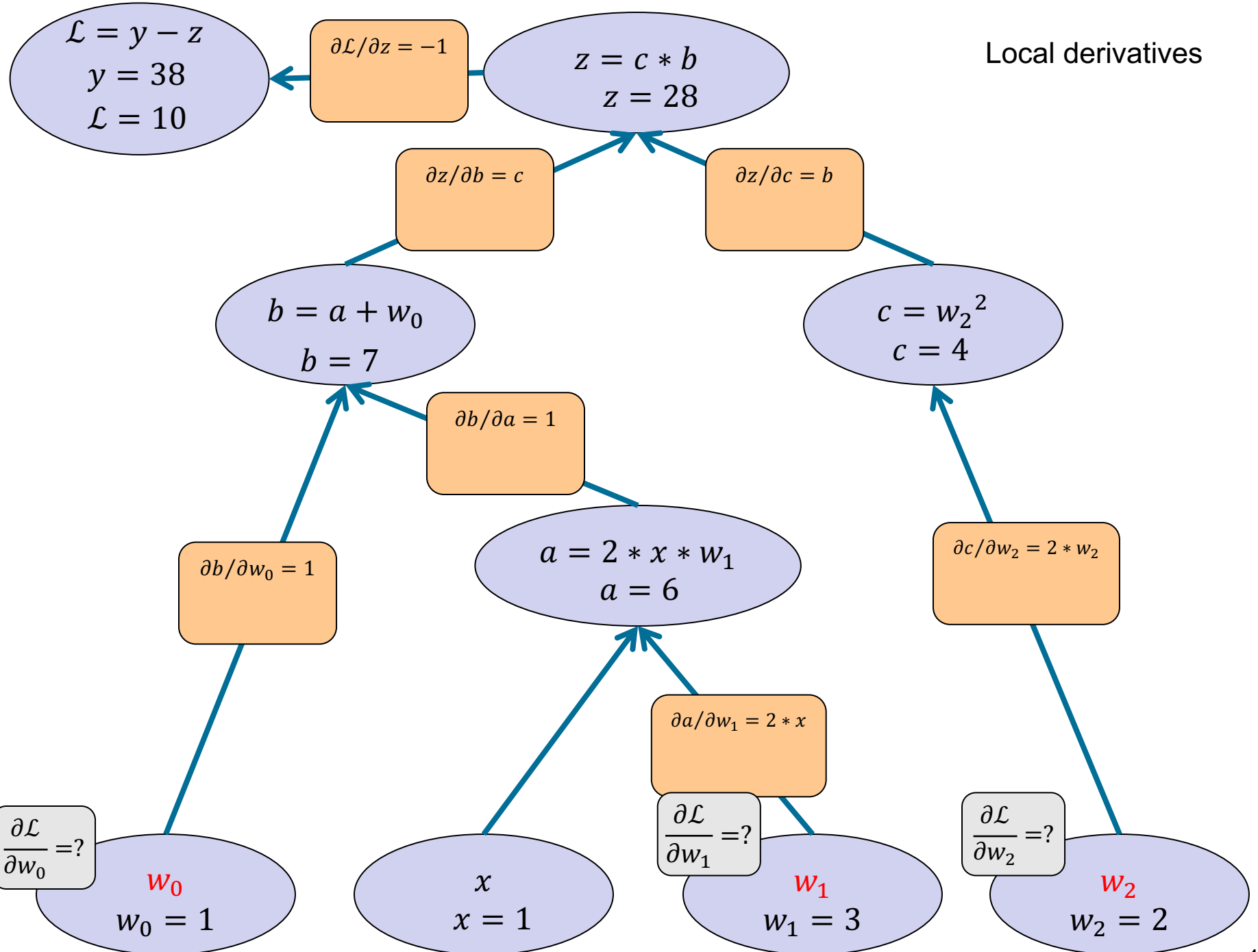
$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial w_1}$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial c} \frac{\partial c}{\partial w_2}$$

- Local derivatives are pre-defined on each atomic operation in the neural computation graph



Local derivatives



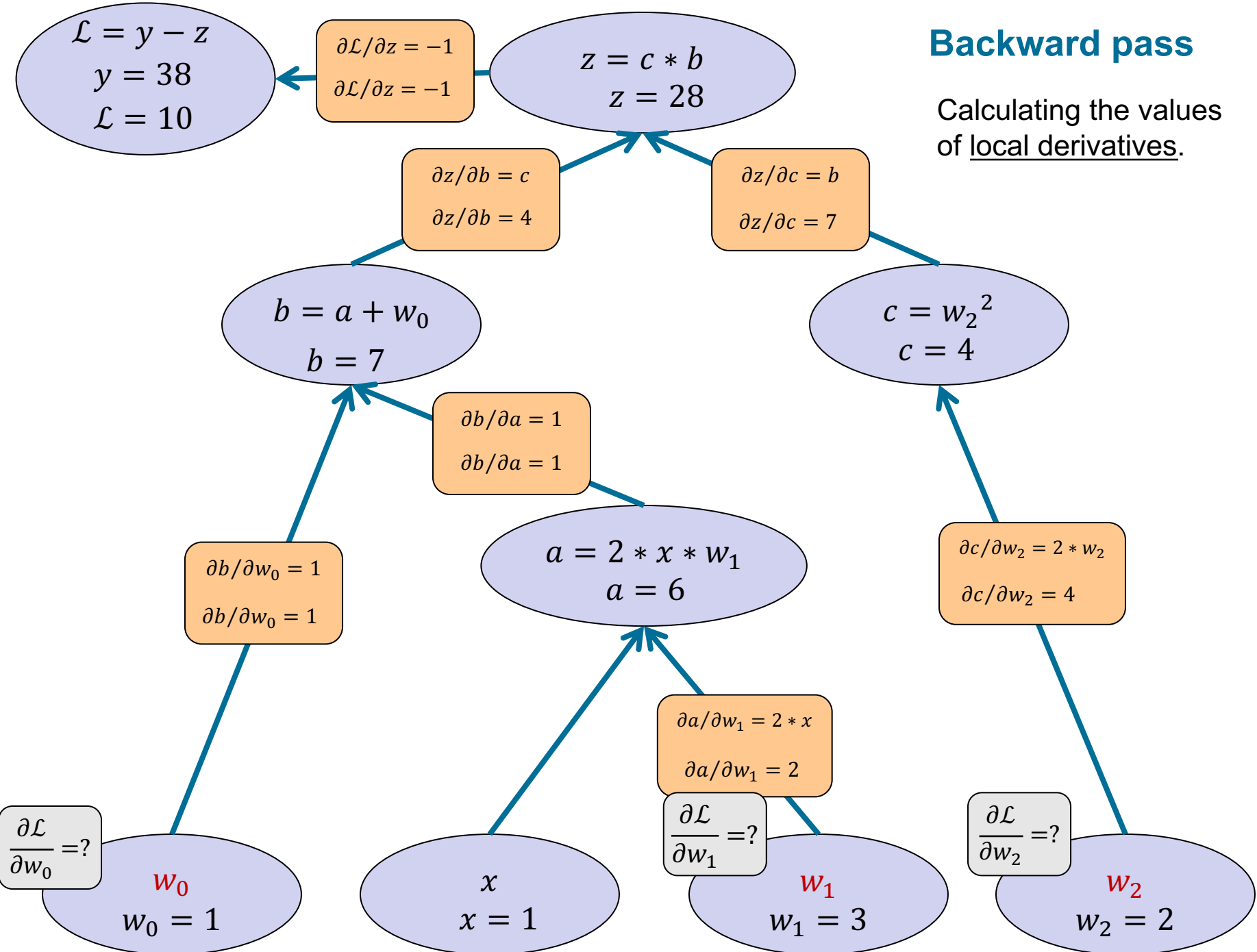
# Backward pass

- Tracing the computation graph from top to bottom and calculating the values of local derivatives
- It means that:
  - We need to keep the values of all intermediate variables after forward pass
  - For the local derivative of every atomic operation, we now have a new stored value



# Backward pass

Calculating the values of local derivatives.



# Backpropagation

Calculating partial derivatives:

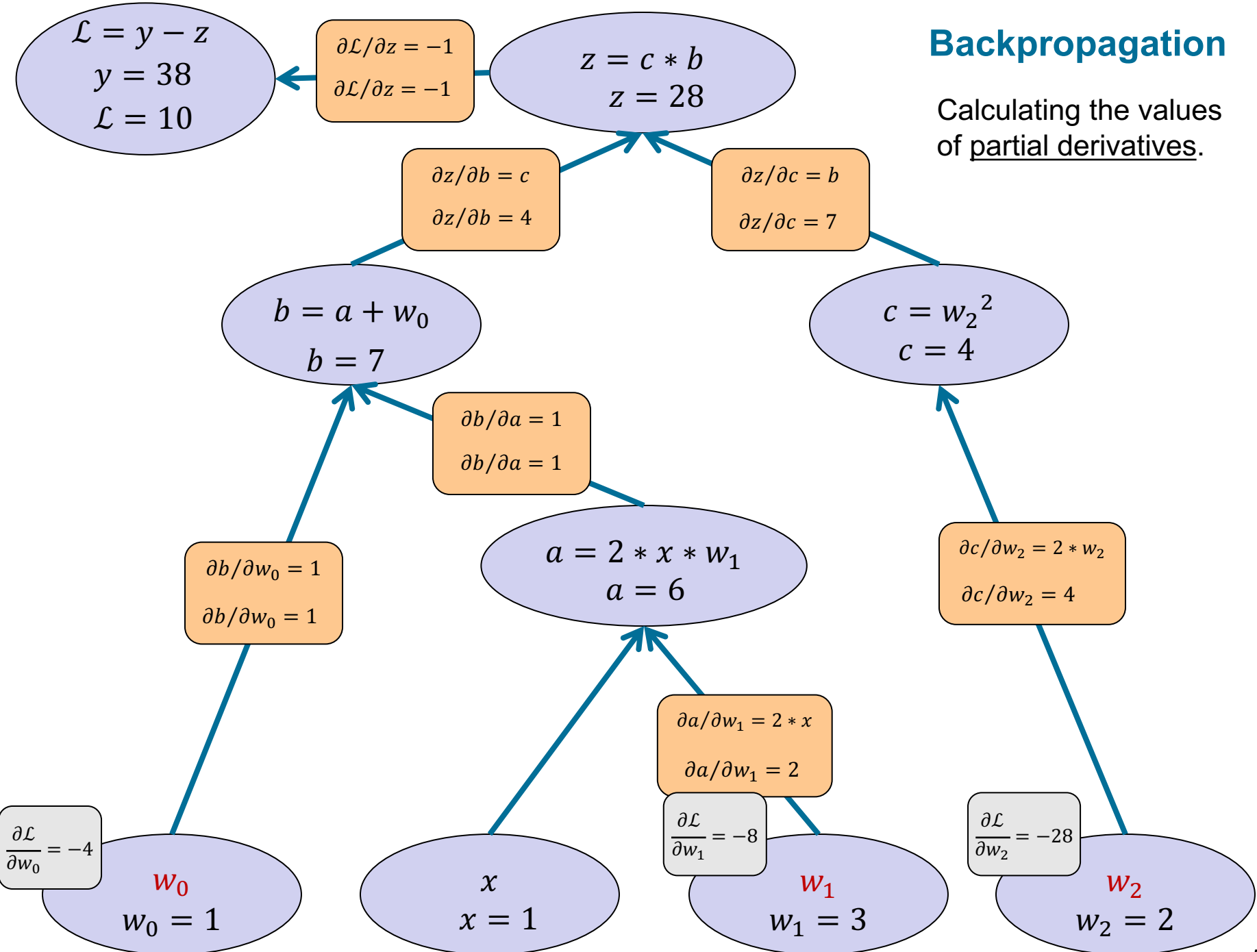
$$\frac{\partial \mathcal{L}}{\partial w_0} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} \frac{\partial b}{\partial w_0} = -1 \times 4 \times 1 = -4$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial b} \frac{\partial b}{\partial a} \frac{\partial a}{\partial w_1} = -1 \times 4 \times 1 \times 2 = -8$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial c} \frac{\partial c}{\partial w_2} = -1 \times 7 \times 4 = -28$$

# Backpropagation

Calculating the values of partial derivatives.



# Learning with Neural Networks

- Design the network's architecture
- Loop until some **exit criteria** are met
  - Sample a **(mini)batch** from training data  $\mathcal{D}$
  - Execute **forward pass**: predict the output tensor of each given input tensor
  - Calculate **loss** function of the (mini)batch
  - **Optimize** the network to reduce loss
    - Calculate the **gradient** of each parameter regarding the loss function using the **backpropagation** algorithm
    - **Update** parameters using their gradients

# Gradient Descent algorithm – recap

- A model with parameters  $\mathbb{W}$  at time step  $t \rightarrow \mathbb{W}^{(t)}$ , **learning rate**  $\eta$ , and set of datapoints  $\mathcal{D}$
- Loop for some epochs
  - Compute gradient tensor  $\mathbb{G}$  of parameters  $\mathbb{W}$  averaged over datapoints  $\mathcal{D}$ :

$$\mathbb{G} \leftarrow \frac{1}{|\mathcal{D}|} \nabla_{\mathbb{W}} \sum_{(x,y) \in \mathcal{D}} \mathcal{L}(x, y; \mathbb{W})$$

- Update the parameters by taking steps in the opposite direction of the gradient tensor multiplied by  $\eta$ :

$$\mathbb{W}^{(t+1)} \leftarrow \mathbb{W}^{(t)} - \eta \mathbb{G}$$

- Reduce learning rate (**annealing**) if some criteria are met or according to a scheduler

# Batch

- In (vanilla) Gradient Descent, first all data points are processed, and their gradients are aggregated, and then a small parameter update is made
  - Training can take very long time
  - Training is not stochastic
- Batch/Mini-batch
  - A (small) set of data to be processed together
  - Suitable for multi-processing capabilities of GPUs
- **Stochastic Gradient Descent**
  - In each step, we process a (mini-)batch of data, calculate their gradients, and update parameters
  - Typical setting for training deep learning models

# (Mini-batch) Stochastic Gradient Descent algorithm

- A model with parameters  $\mathbb{W}$  at time step  $t \rightarrow \mathbb{W}^{(t)}$ , **learning rate**  $\eta$ , and set of datapoints  $\mathcal{D}$
- Loop until some exit criteria are met

- $\hat{\mathcal{D}}$  is the set of datapoints in the **minibatch**
- Compute gradient tensor  $\mathbb{G}$  of parameters  $\mathbb{W}$  averaged over batch datapoints  $\hat{\mathcal{D}}$ :

$$\mathbb{G} \leftarrow \frac{1}{|\hat{\mathcal{D}}|} \nabla_{\mathbb{W}} \sum_{(x,y) \in \hat{\mathcal{D}}} \mathcal{L}(x, y; \mathbb{W})$$

- Update the parameters by taking steps in the opposite direction of the gradient tensor multiplied by  $\eta$ :

$$\mathbb{W}^{(t+1)} \leftarrow \mathbb{W}^{(t)} - \eta \mathbb{G}$$

- Reduce learning rate (**annealing**) if some criteria are met or according to a scheduler

# Other gradient-based optimizations

- Some limitations of the mentioned SGD algorithms
  - Choosing learning rate is hard
  - Choosing annealing method/rate is hard
  - Same learning rate is applied to all parameters
  - Can get trapped in non-optimal local minima and saddle points
  
- Some other commonly used algorithms:
  - Nesterov accelerated gradient
  - Adagrad
  - Adam