

# **CS5242 Neural Networks and Deep Learning**

## **Lecture 03: From Shallow to Deep Neural Networks**

Wei WANG

[cs5242@comp.nus.edu.sg](mailto:cs5242@comp.nus.edu.sg)



**NUS**  
National University  
of Singapore

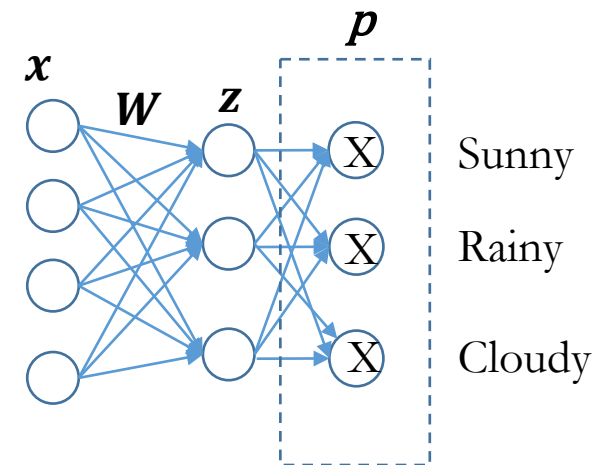
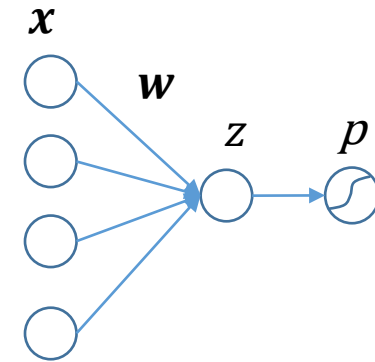
School of  
Computing

# Agenda

- Recap
  - Logistic regression model
  - Softmax regression model
- Multi-layer perceptron model
- Back-propagation algorithm

# Recap

- Binary classification model
  - Logistic function  $\rightarrow$  probability
  - Binary cross-entropy loss
- Multi-class classification model
  - Softmax/multinomial regression
    - Softmax function  $\rightarrow$  a vector of probabilities
  - Cross-entropy loss

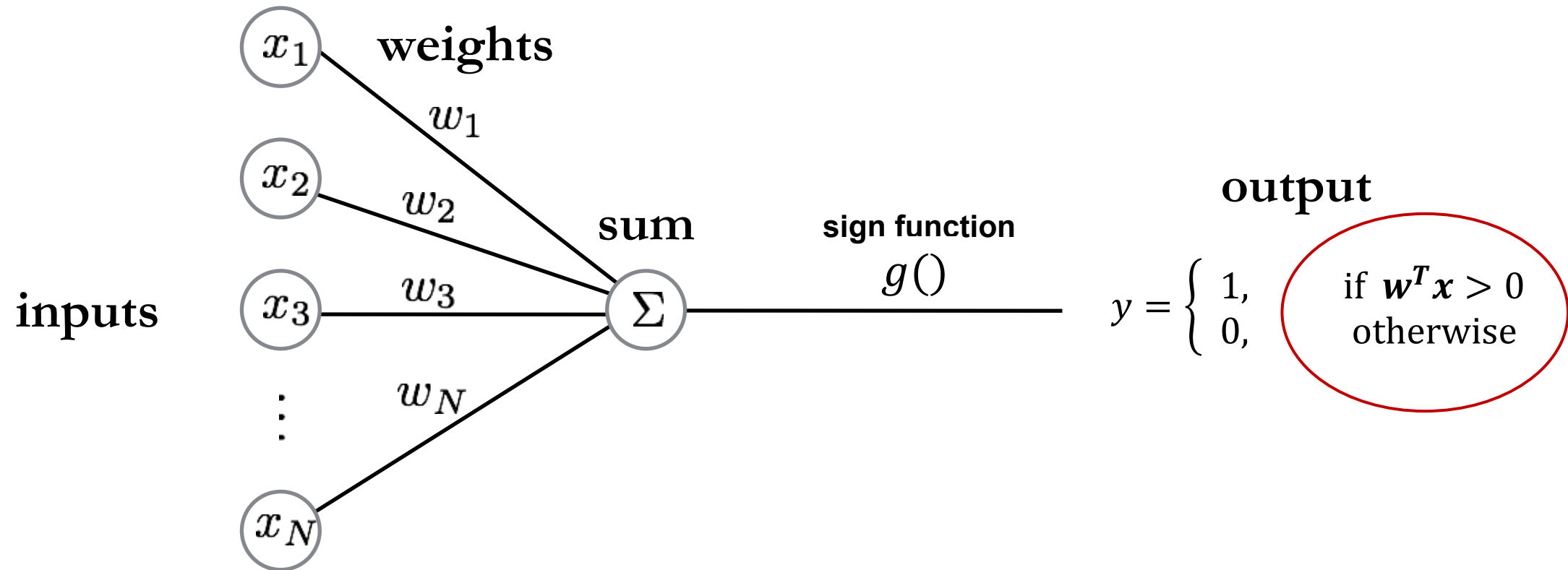


# Multilayer Perceptron (MLP)

Non-linear Activation Functions

Definition of MLPs

# The Perceptron



# Linear functions are limited

- Simple binary example:

$$y = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ 0, & \text{otherwise} \end{cases}$$

Possible weights for AND function:

$$w_1 = 1, w_2 = 1, b = -1.5$$

AND function

x1	x2	y
1	1	1
1	0	0
0	0	0
0	1	0

XOR function

x1	x2	y
0	0	0
0	1	1
1	0	1
1	1	0

Find the weights for the XOR function?

HWQ: can logistic regression model fit the XOR data?

# Non-Linear feature transformations

- $\mathbf{h} = \max(0, W\mathbf{x} + \mathbf{c})$ ,  $W \in R^{2 \times 2}$ ,  $\mathbf{c} \in R^2$
- $\mathbf{h}^T \mathbf{w} + b$ ,  $\mathbf{w} \in R^2$ ,  $b \in R$

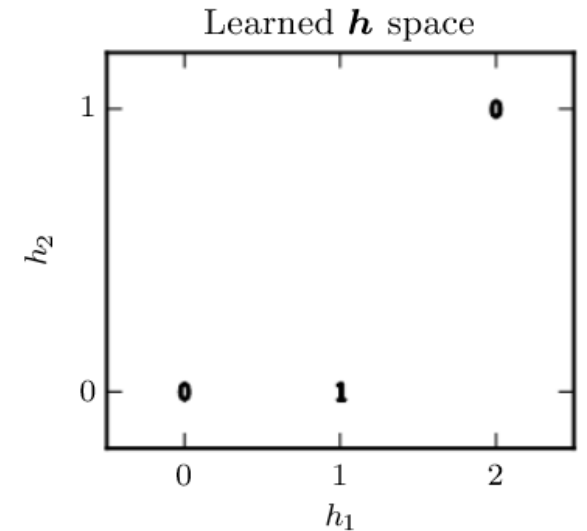
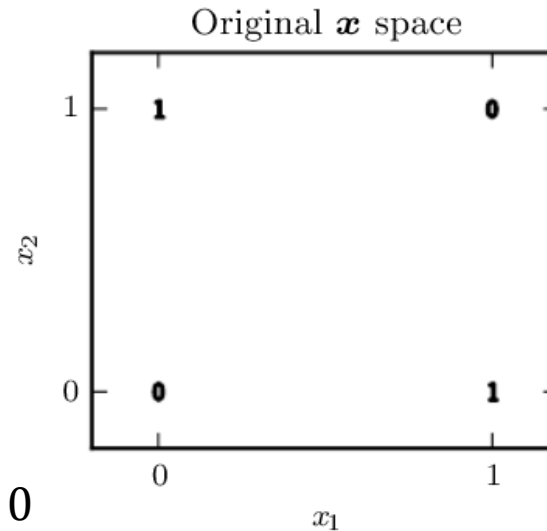
$$X = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}.$$

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$

$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix},$$

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

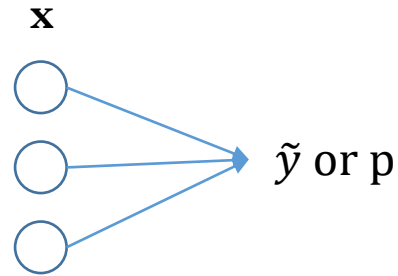
$$b = 0$$



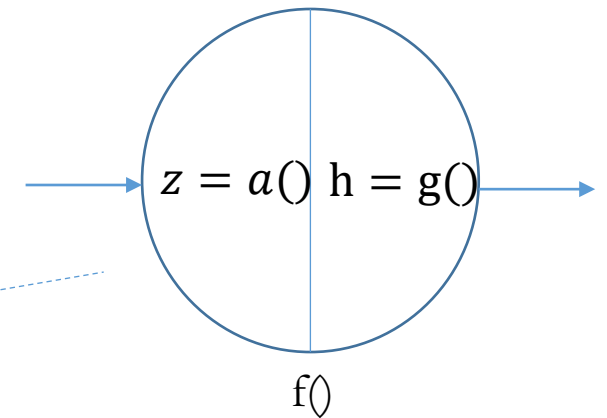
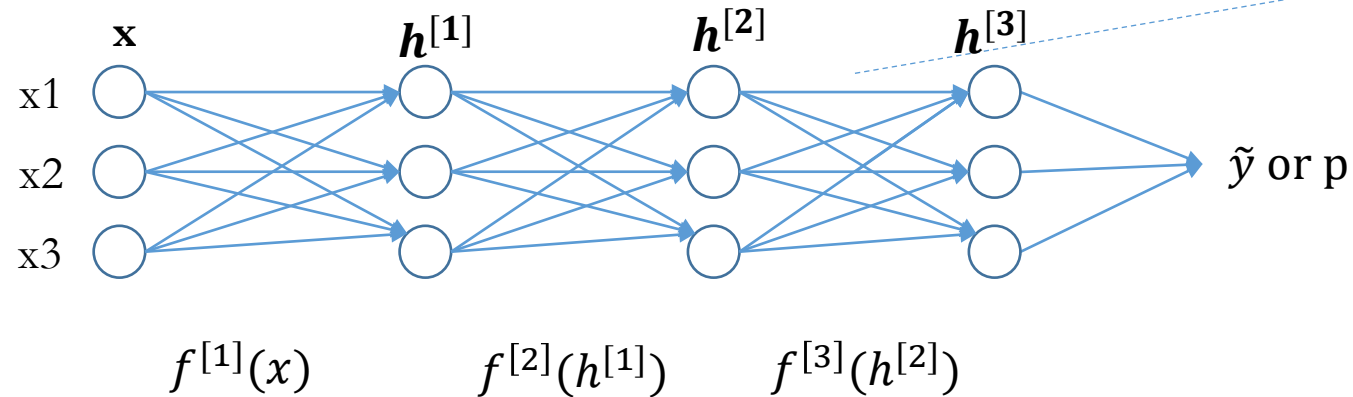
Try it! Compute  $\mathbf{h}$  and  $\tilde{\mathbf{y}}$  to see the classification results  
More examples at [google playground](#)

# Multilayer Perceptron

Single layer  
perceptron  
(no hidden layers)



Multi-layer  
perceptron





# Multilayer Perceptron

A net with multiple layers that transform input features into hidden features and then make predictions

- At least one **(non-linear) hidden layer**
- $i^{\text{th}}$  layer consists of a **linear/affine** transformation function

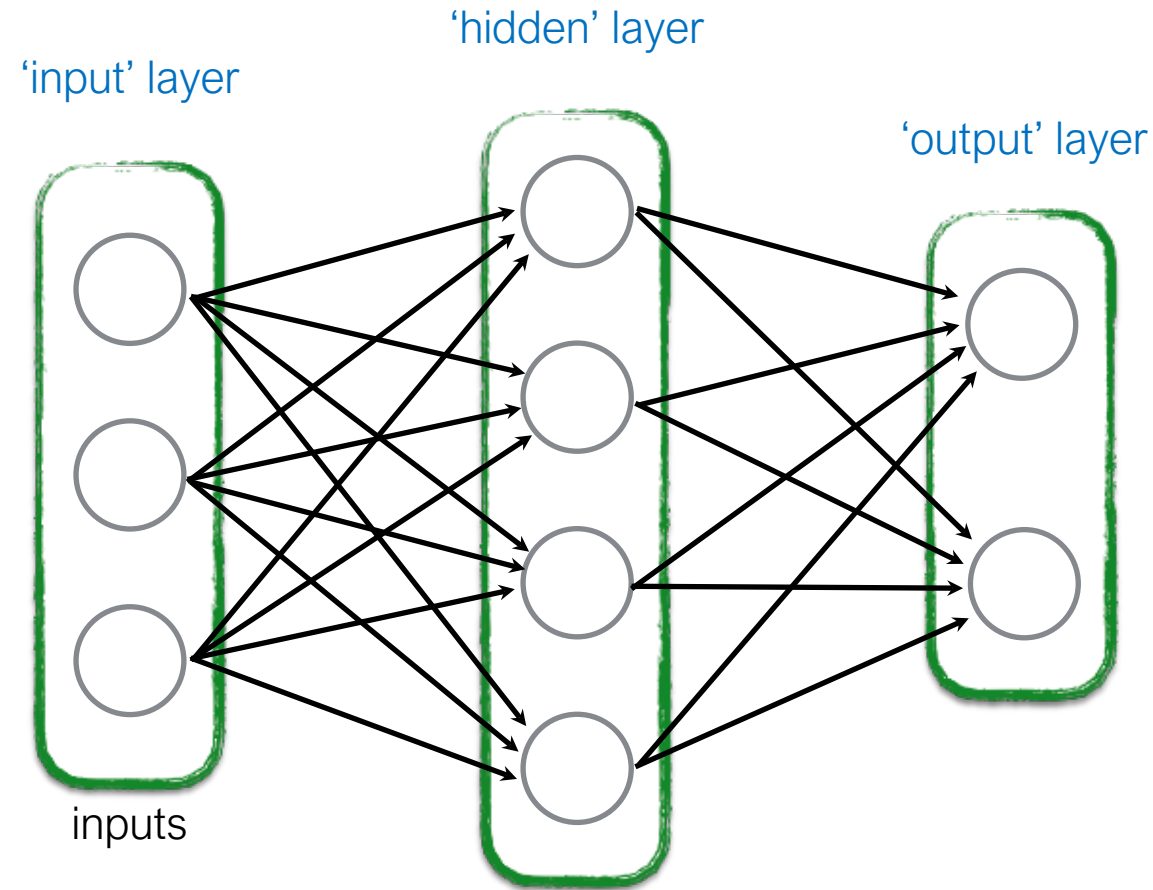
$$\mathbf{z}^{[i]} = a^{[i]}(\mathbf{h}^{[i-1]}) = W^{[i]}\mathbf{h}^{[i-1]} + \mathbf{b}^{[i]}$$

$$W^{[i]} \in R^{n_i \times n_{i-1}}, \mathbf{b}^i \in R^{n_i}$$

- $m_i$  is the number of hidden units at the  $i^{\text{th}}$  layer and is a hyper-parameter
- followed by a **non-linear activation** function

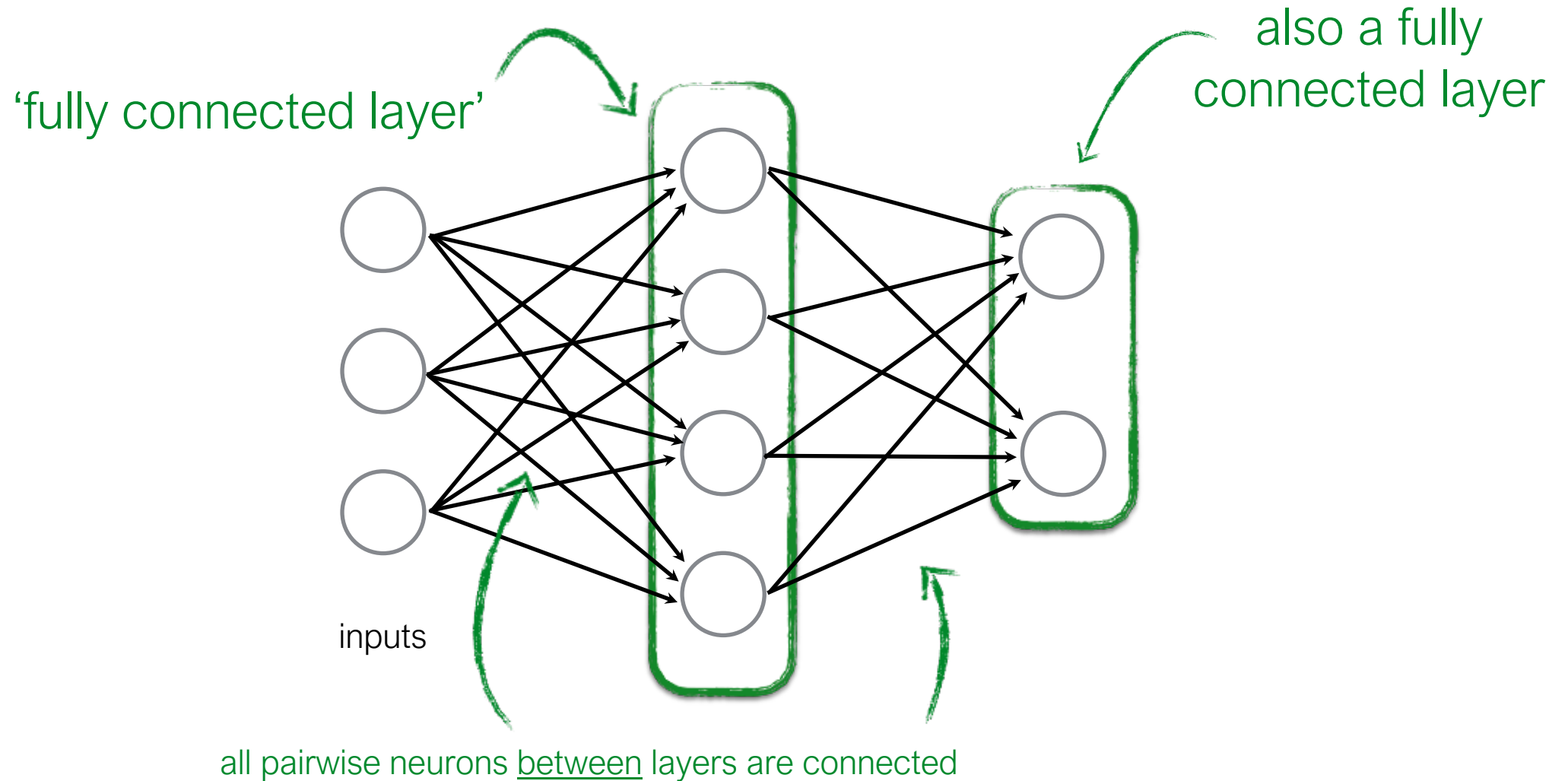
$$\mathbf{h}^{[i]} = g^{[i]}(\mathbf{z}^{[i]}), \in R^{n_i}$$

# Multilayer Perceptron



...also called a **Multi-layer Perceptron (MLP)**

# Multilayer Perceptron

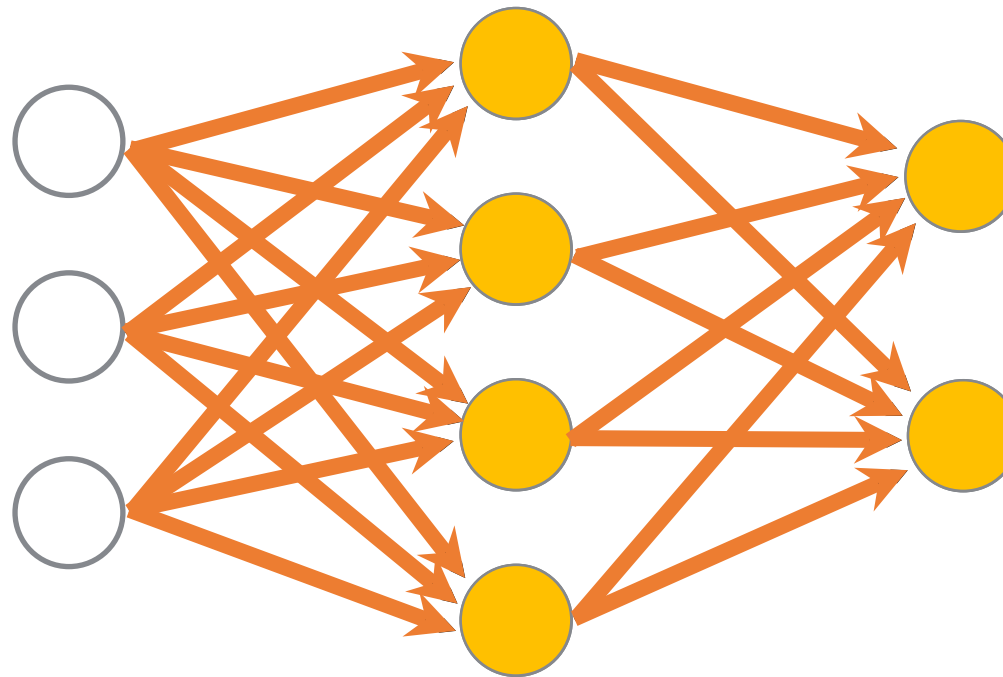


How many neurons (perceptrons)?

$$4 + 2 = 6$$

How many weights (edges)?

$$(3 \times 4) + (4 \times 2) = 20$$



How many learnable parameters total?

$$20 + (4 + 2) = 26$$

6 bias terms  
1 per perceptrons



# Why non-linear activation?

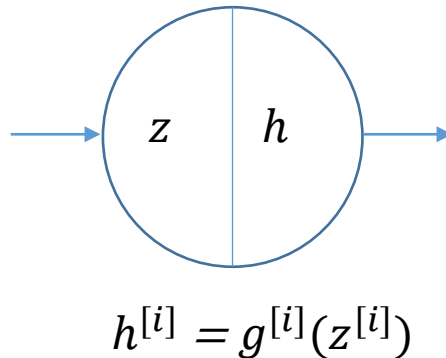
This result says that if all activation functions were linear, we can always define an equivalent collapsed network.

Because successive linear transformations together form yet another linear transformation, a multi-layered linear network is not so interesting. Therefore, we need non-linear activations.

# Activation function $g()$

- Logistic (Sigmoid,  $\sigma$ ) VS ReLU

Logistic (a.k.a. Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$	$(0, 1)$
Rectified linear unit (ReLU) <sup>[9]</sup>		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$[0, \infty)$



## Logistic activation

If  $z_k$  is large, e.g.  $>10$ , then  $h_k$  is near 1

If  $z_k$  is small, e.g.  $<-10$ , Then  $h$  is near 0

For both cases,  $\frac{\partial h_k}{\partial z_k} \approx 0 \rightarrow$  gradient vanishing  
 $\rightarrow$  gradient vanishing

## ReLU activation

If  $z_k$  is positive,  $\frac{\partial h_k}{\partial z_k} = 1$ , no gradient vanishing

If  $z_k$  is negative,  $\frac{\partial h_k}{\partial z_k} = 0$  no gradients ☹

better than sigmoid, since  $z_k$  has a larger working domain  
Leaky ReLU (see previous slides) resolves the gradient vanishing problem for negative  $z_k$

# Activation function $g()$

- Activation functions perform non-linear feature transformations

name	plot	equation	derivative	range
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$	$\{0, 1\}$
Logistic (a.k.a. Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$	$(0, 1)$
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$	$(-1, 1)$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$	$(-\frac{\pi}{2}, \frac{\pi}{2})$
Softsign [7][8]		$f(x) = \frac{x}{1 +  x }$	$f'(x) = \frac{1}{(1 +  x )^2}$	$(-1, 1)$
Rectified linear unit (ReLU) <sup>[9]</sup>		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$[0, \infty)$
Leaky rectified linear unit (Leaky ReLU) <sup>[10]</sup>		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$
Parametric rectified linear unit (PReLU) <sup>[11]</sup>		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$
Randomized leaky rectified linear unit (RRReLU) <sup>[12]</sup>		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ <sup>[1]</sup>	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$

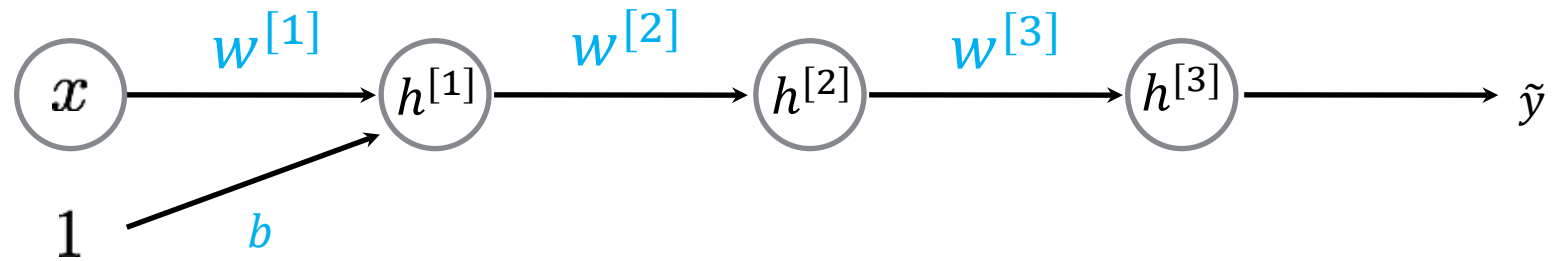
Source from: [https://en.wikipedia.org/wiki/Activation\\_function](https://en.wikipedia.org/wiki/Activation_function)

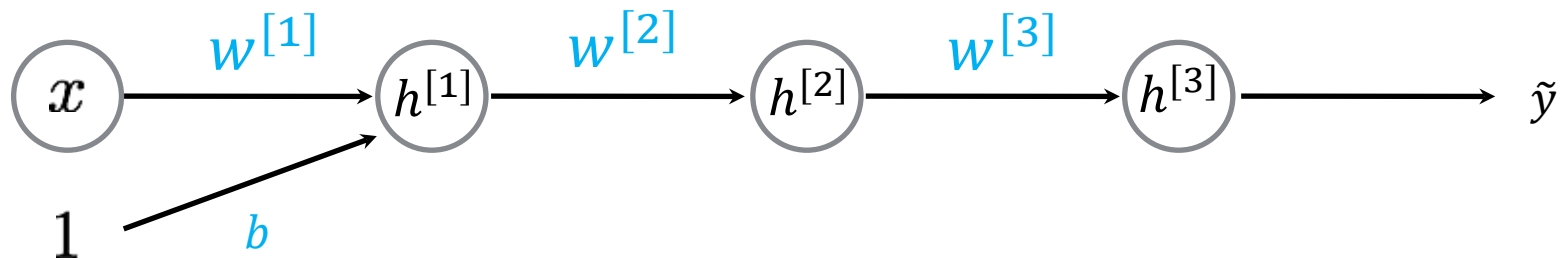
# Training MLP

- Cross-entropy loss for classification
- Squared Euclidean distance for regression
- GD algorithm
  - Compute  $J(\mathbf{X}, \mathbf{Y})$
  - Compute gradient :  $\frac{\partial J}{\partial \mathbf{W}^{[1]}}, \frac{\partial J}{\partial \mathbf{b}^{[1]}}, \frac{\partial J}{\partial \mathbf{W}^{[2]}}, \frac{\partial J}{\partial \mathbf{b}^{[2]}} \dots$
  - Update:  $\mathbf{W}^{[k]} = \mathbf{W}^{[k]} - \alpha \frac{\partial J}{\partial \mathbf{W}^{[k]}}$ ,  $\mathbf{b}^{[k]} = \mathbf{b}^{[k]} - \alpha \frac{\partial J}{\partial \mathbf{b}^{[k]}}$



# A simple MLP example





$$h^{[1]} = g^{[1]}(w^{[1]}x + b)$$

$$h^{[2]} = g^{[2]}(w^{[2]}h^{[1]})$$

$$h^{[3]} = g^{[3]}(w^{[3]}h^{[2]})$$

$$\tilde{y} = h^{[3]}$$

Entire network can be written out as one long equation

The diagram illustrates the equation for the output of a neural network. A blue arrow labeled "known" points from the input  $x$  to the output  $\tilde{y}$ . Four green arrows labeled "unknown" point from the weights  $w^{[1]}$ ,  $w^{[2]}$ ,  $w^{[3]}$ , and the bias  $b$  to the equation.

$$\tilde{y} = h^{[3]} = g^{[3]}(w^{[3]}g^{[2]}(w^{[2]}g^{[1]}(w^{[1]}x + b)))$$

We need to train the network:

What is known? What is unknown?

# Gradient descent algorithm

Random initialize  $w^{[1]}, w^{[2]}$  ...

Repeat

$J = 0$

For each sample  $\mathbf{x}^{(i)}, y^{(i)}$

a. Forward pass

b. Accumulate Loss

c. Accumulate partial gradient (derivative)  $\frac{\partial J}{\partial w^{[1]}}, \frac{\partial J}{\partial w^{[2]}}$  ...

how to compute?

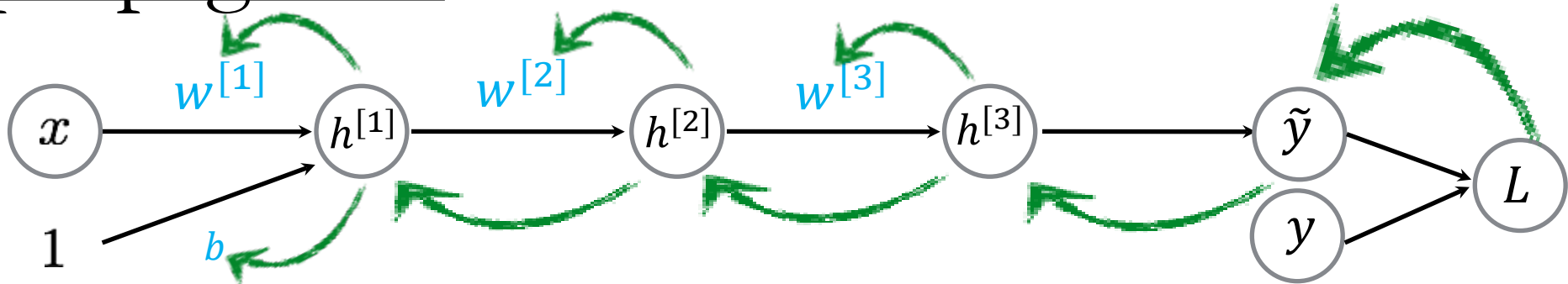
Update

$$w^{[1]} = w^{[1]} - \alpha \frac{\partial J}{\partial w^{[1]}}$$

$$w^{[2]} = w^{[2]} - \alpha \frac{\partial J}{\partial w^{[2]}}$$

move in opposite direction of partial derivatives

# Back-propagation



$$\text{loss} = L(\tilde{y}, y) \quad \tilde{y} = h^{[3]}$$

$$h^{[3]} = g^{[3]}(w^{[3]}h^{[2]})$$

$$h^{[2]} = g^{[2]}(w^{[2]}h^{[1]})$$

$$h^{[1]} = g^{[1]}(w^{[1]}x + b)$$

$$\frac{\partial L}{\partial w^{[3]}} = \frac{\partial L}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial h^{[3]}} \frac{\partial h^{[3]}}{\partial w^{[3]}}$$

$$\frac{\partial L}{\partial w^{[2]}} = \frac{\partial L}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial h^{[3]}} \frac{\partial h^{[3]}}{\partial h^{[2]}} \frac{\partial h^{[2]}}{\partial w^{[2]}}$$

$$\frac{\partial L}{\partial w^{[1]}} = \frac{\partial L}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial h^{[3]}} \frac{\partial h^{[3]}}{\partial h^{[2]}} \frac{\partial h^{[2]}}{\partial h^{[1]}} \frac{\partial h^{[1]}}{\partial w^{[1]}}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial h^{[3]}} \frac{\partial h^{[3]}}{\partial h^{[2]}} \frac{\partial h^{[2]}}{\partial h^{[1]}} \frac{\partial h^{[1]}}{\partial b}$$

The term back-propagation comes from the application of the chain rule, in which the gradients or partial derivatives from down-stream are used upstream.

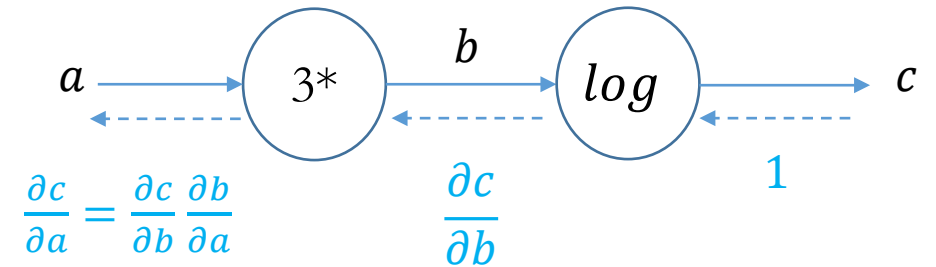
# Backpropagation

# How to compute the gradient?

- GD algorithm needs the gradient of the loss w.r.t each parameter
- Linear regression model
  - $\frac{\partial L}{\partial \mathbf{w}} = (\tilde{y} - y)\mathbf{x}$
- Logistic regression model
  - $\frac{\partial L}{\partial \mathbf{w}} = (p - y)\mathbf{x}$
- Softmax regression model
  - $\frac{\partial L}{\partial \mathbf{w}} = (\mathbf{p} - \mathbf{y})\mathbf{x}^T$
- How about MLP model with 10 hidden layers?
  - We need a **modular** way to compute the gradients

# Chain rule I

- $c = \log 3a$  Given  $a = 1$ , what is  $\frac{\partial c}{\partial a}$ ?



- $b = f_1(a) = 3a, \quad c = f_2(b) = \log b \rightarrow c = \log 3a$
- $\frac{\partial c}{\partial a} = \frac{\partial c}{\partial b} \frac{\partial b}{\partial a} = f_2'(b) f_1'(a) = \frac{1}{b} * 3 = \frac{1}{3a} * 3 = \frac{3}{3} = 1$



# Chain rule I

- Basic chain rule

- $v_2 = f_1(v_1), v_3 = f_2(v_2) \quad \dots \quad v_k = f_{k-1}(v_{k-1})$

- $v_i$  could be a scalar, vector, matrix, tensor

- $\frac{\partial v_k}{\partial v_i} = \text{mul}\left(\frac{\partial v_k}{\partial v_{i+1}}, \frac{\partial v_{i+1}}{\partial v_i}\right)$

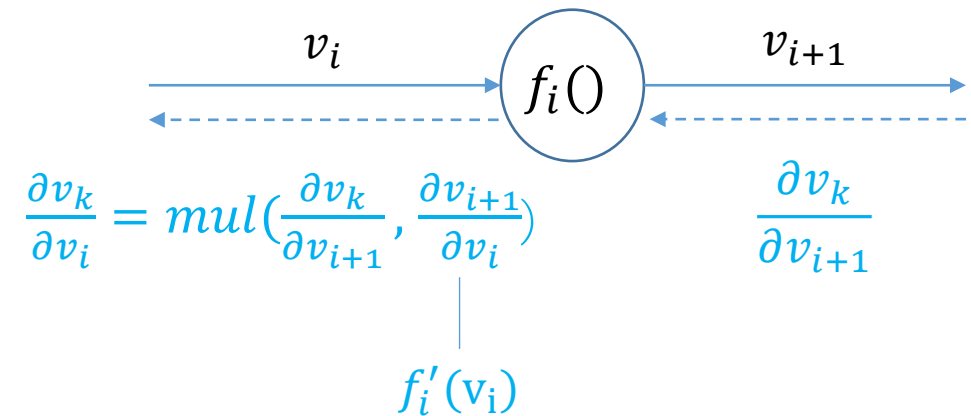
General multiplication operation;  
to avoid the scalar/matrix/tensor  
transpose/ordering messiness ...

- Reorder and transpose

- $\frac{\partial v_k}{\partial v_{i+1}} \frac{\partial v_{i+1}}{\partial v_i}$  if all are scalars

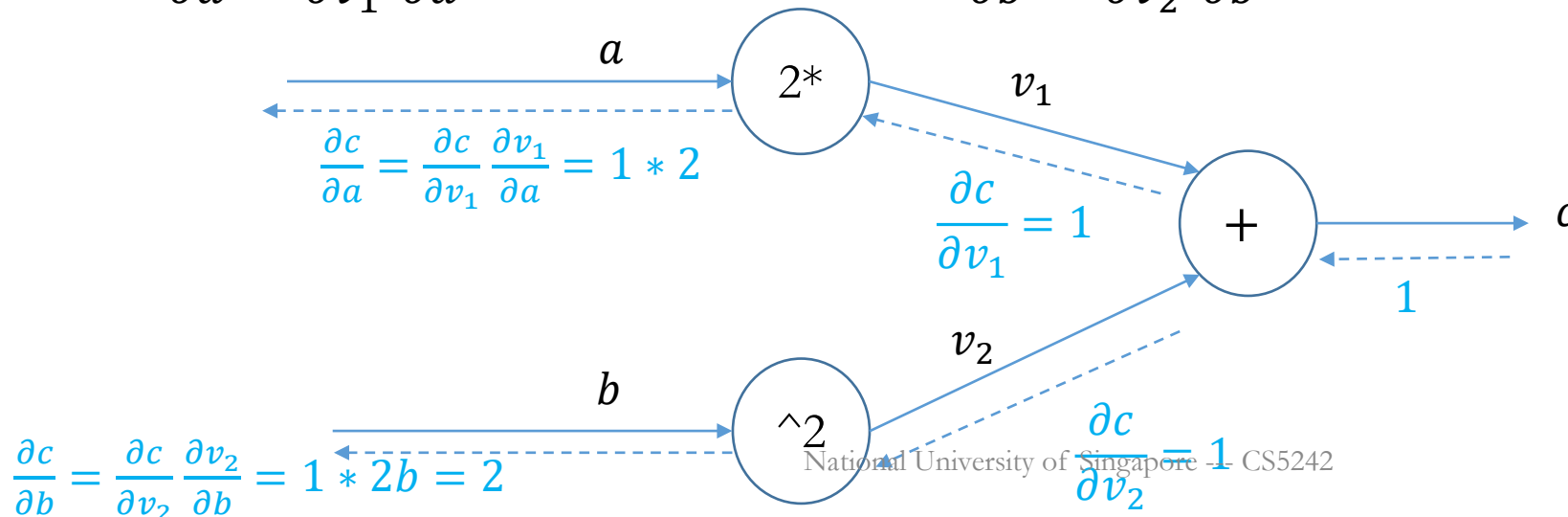
- $\frac{\partial v_{i+1}}{\partial v_i} \frac{\partial v_k}{\partial v_{i+1}}$  if  $v_k$  is a scalar,  $v_{i+1}, v_i$  are vectors

Some  
special  
cases;  
shape-  
check to  
confirm!



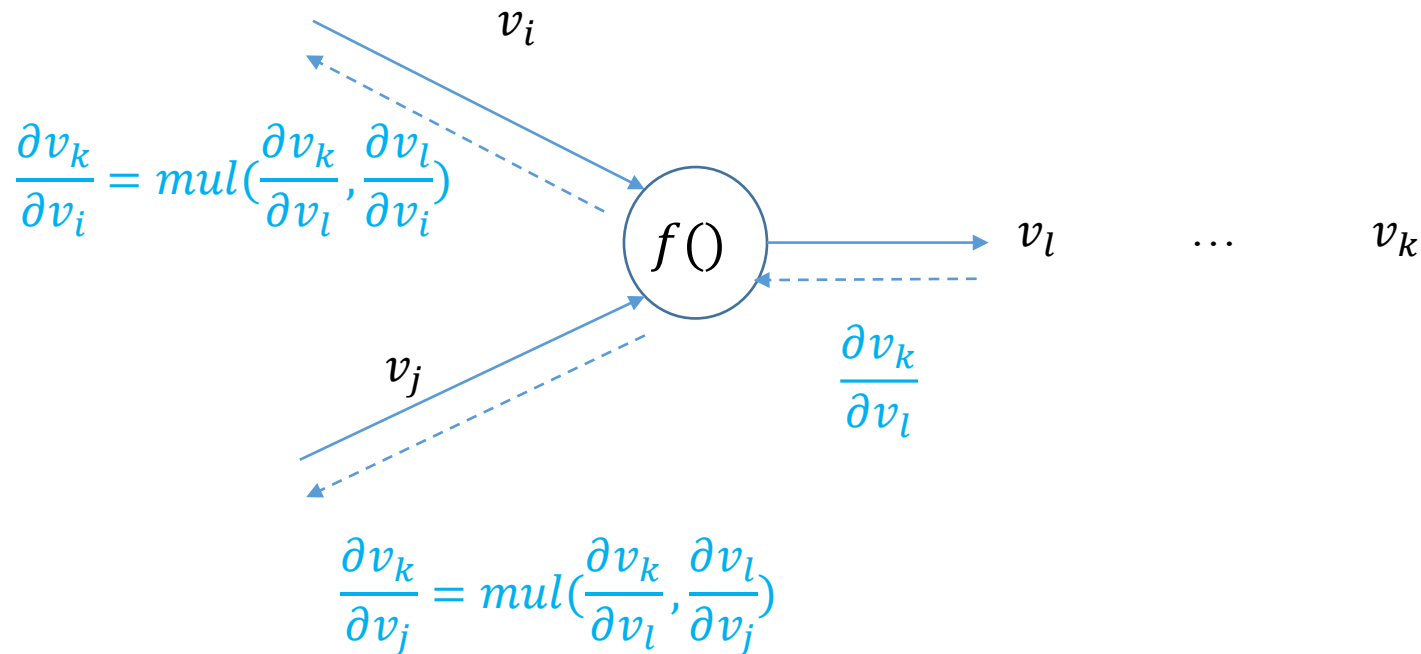
# Chain rule II

- $c = f(a, b) = 2a + b^2$ 
  - Given,  $a = 1, b = 1$
  - $\frac{\partial c}{\partial a} = 2, \frac{\partial c}{\partial b} = 2b = 2$
- $v_1 = f_1(a) = 2a, v_2 = f_2(b) = b^2, c = f_3(v_1, v_2) = v_1 + v_2$ 
  - $\frac{\partial c}{\partial a} = \frac{\partial c}{\partial v_1} \frac{\partial v_1}{\partial a} = 1 * 2 = 2$        $\frac{\partial c}{\partial b} = \frac{\partial c}{\partial v_2} \frac{\partial v_2}{\partial b} = 1 * 2b = 2$



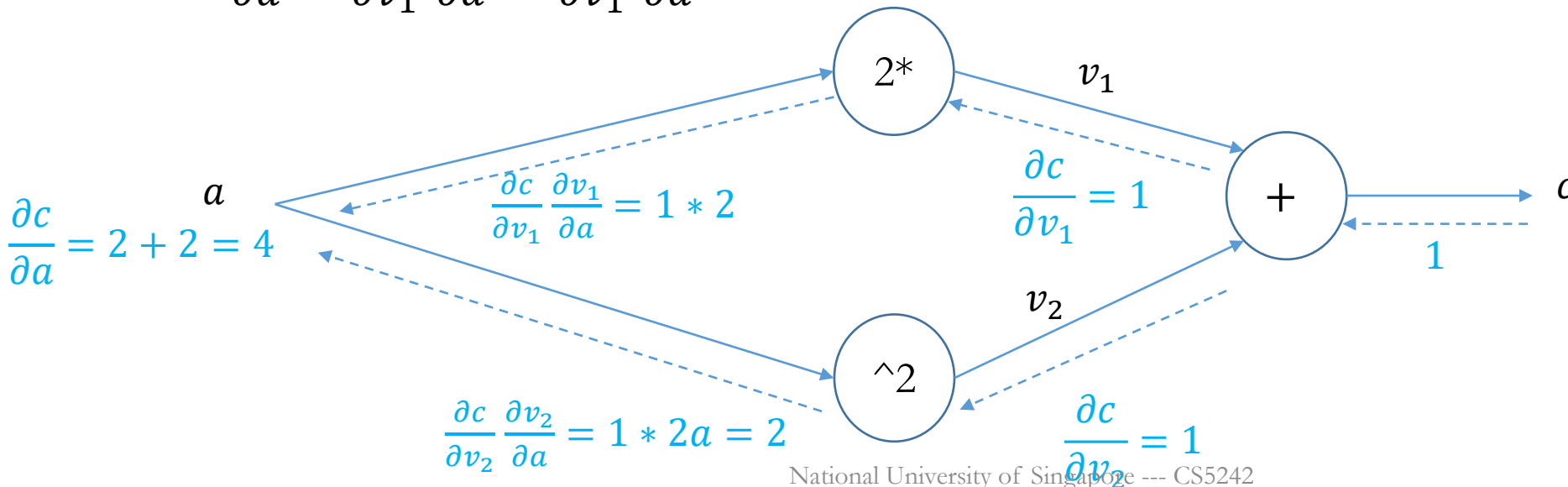
# Chain rule II

- $v_l = f(v_i, v_j)$ ,
- $v_k$  is the final output



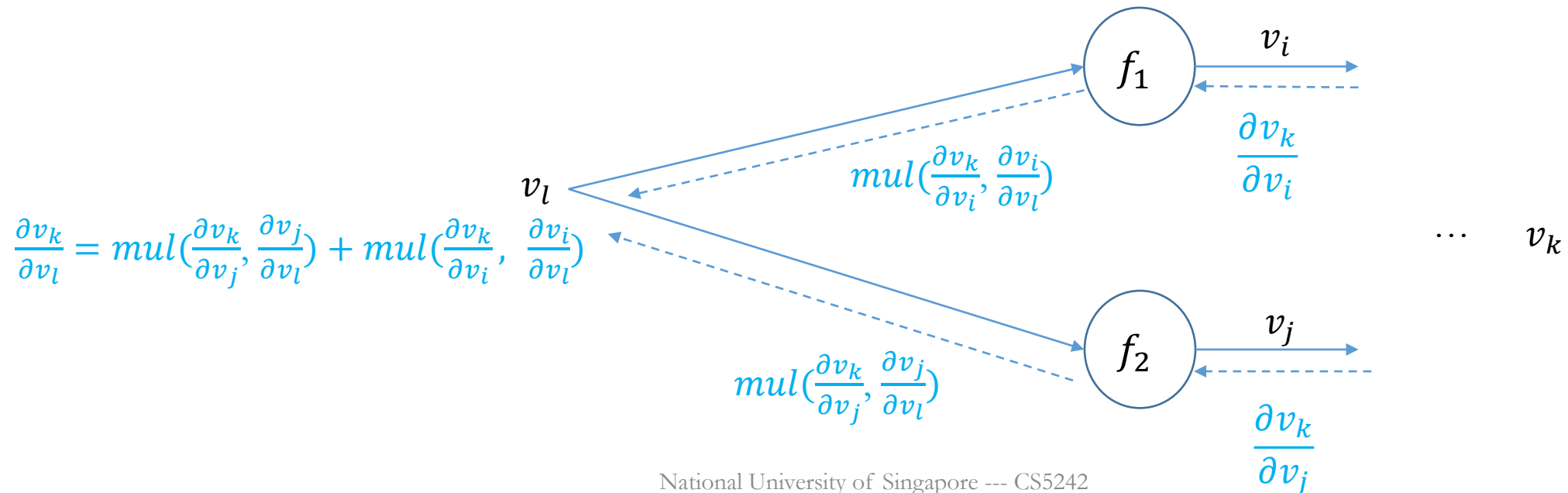
# Chain rule III

- $c = 2a + a^2$ 
  - Given  $a = 1$ ,  $\frac{\partial c}{\partial a} = 2 + 2a = 4$
- $v_1 = f_1(a) = 2a, v_2 = f_2(a) = a^2, c = f_3(v_1, v_2) = v_1 + v_2$ 
  - $\frac{\partial c}{\partial a} = \frac{\partial c}{\partial v_1} \frac{\partial v_1}{\partial a} + \frac{\partial c}{\partial v_2} \frac{\partial v_2}{\partial a} = 1 * 2 + 1 * 2a = 4$



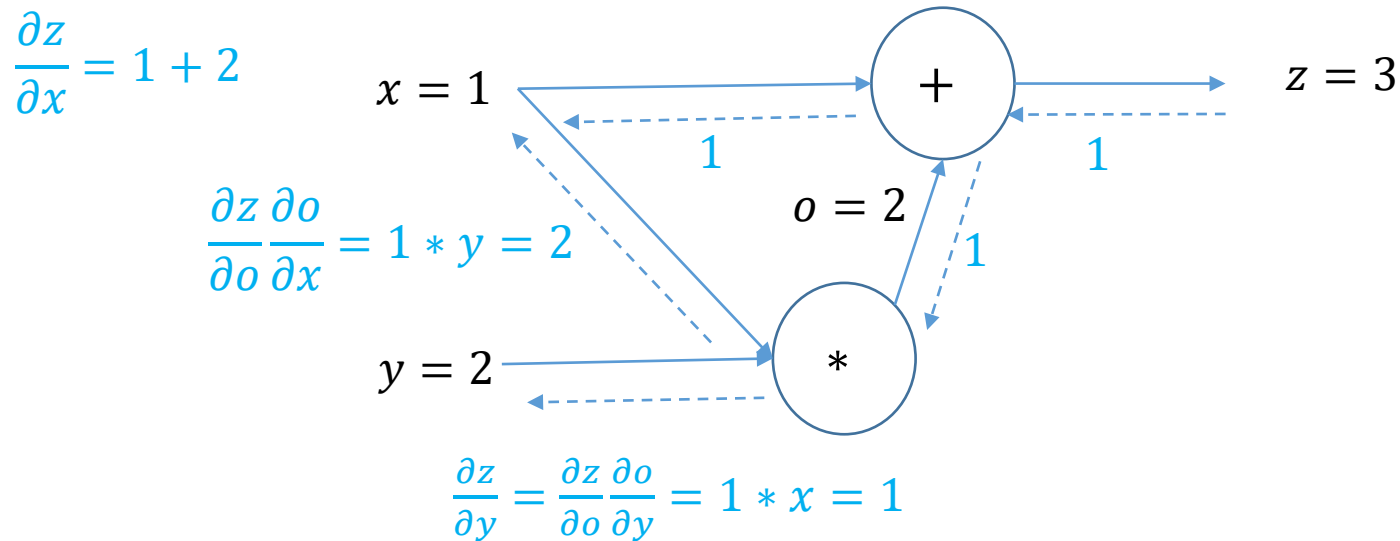
# Chain rule III

- $v_i = f_1(v_l)$
- $v_j = f_2(v_l)$
- $v_k$  is the final output



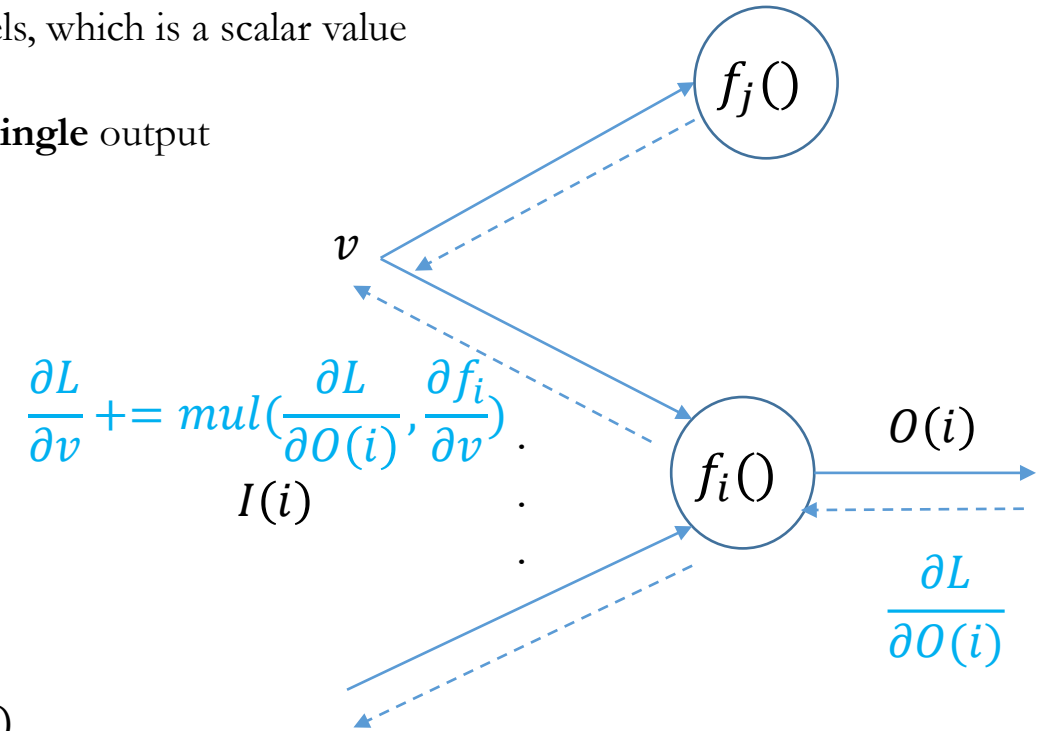
# Practice

- $z = f(x, y) = x * y + x$
- Given  $x = 1, y = 2$ , compute  $\frac{\partial z}{\partial x}, \frac{\partial z}{\partial y}$ ?



# Backpropagation (BP)

- Computation graph
  - Denote all variables as  $v_1, v_2, \dots$  sorted in topological order
    - The last variable is the loss (L) for neural network models, which is a scalar value
  - Denote all operations as  $f_1, f_2, \dots, f_k$ 
    - Each operation accepts multiple inputs and generate **a single** output
      - which may be used by multiple other operations
    - $I(i)$  denotes all the input variables to  $f_i$
    - $O(i)$  denotes the single output variable of  $f_i$
- Forward pass
  - For  $i = 1, 2, \dots, k$ 
    - Run operation  $f_i: I(i) \rightarrow O(i)$
- Backward pass
  - Initialize  $\frac{\partial L}{\partial v} = 0$  for all  $v$
  - For  $i = k, \dots, 1$  *Note the reverse order!*
    - For each  $v \in I(i)$  Compute  $\frac{\partial L}{\partial v} += \text{mul}(\frac{\partial L}{\partial O(i)}, \frac{\partial f_i(I(i))}{\partial v})$



# Logistic regression

- $z = \mathbf{w}^T \mathbf{x}, \mathbf{x} \in R^{n \times 1}, \mathbf{w} \in R^{n \times 1}$
- $p = \sigma(z)$
- $L(p, y) = -y \log p - (1 - y) \log(1 - p)$

- Dot

- Forward( $\mathbf{x}, \mathbf{w}$ ):  $\mathbf{w}^T \mathbf{x}$
- Backward( $dz, \mathbf{x}, \mathbf{w}$ ):  $dz \mathbf{x}$

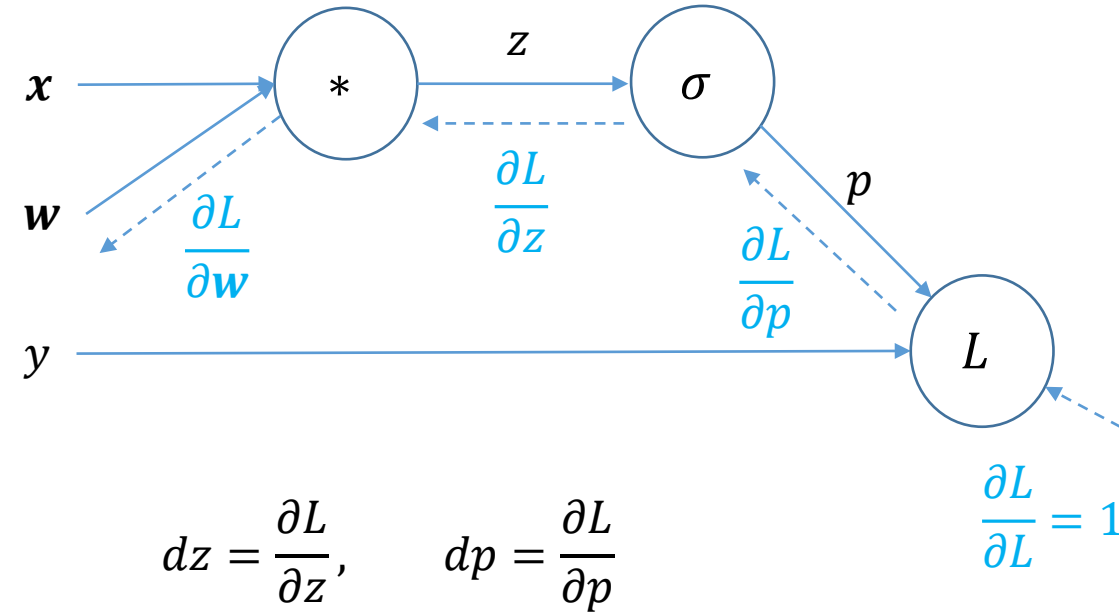
- Logistic

- Forward( $z$ ):  $\sigma(z)$
- Backward( $dp, z$ ):  $p = \sigma(z); dp * p * (1 - p) = (-\frac{y}{p} + \frac{1-y}{1-p}) * p * (1 - p) = p - y$

- Binary-Cross-Entropy

- Forward( $p, y$ ):  $-y \log p - (1 - y) \log(1 - p)$
- Backward( $[dL], p, y$ ):  $-\frac{y}{p} + \frac{1-y}{1-p}$

$dL$  is redundant since it is 1



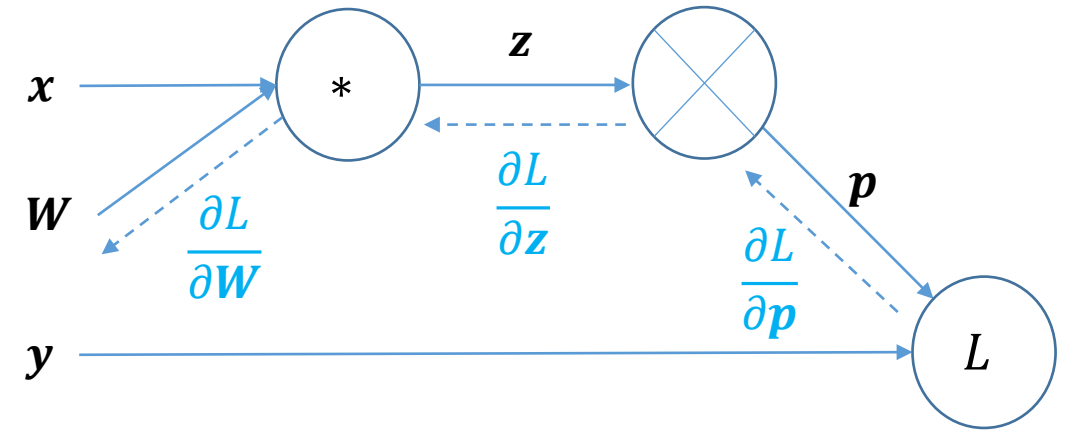
This operation can also be further decomposed; but since binary cross-entropy is used so commonly, we treat it as a whole stand-alone unit.



# Softmax regression

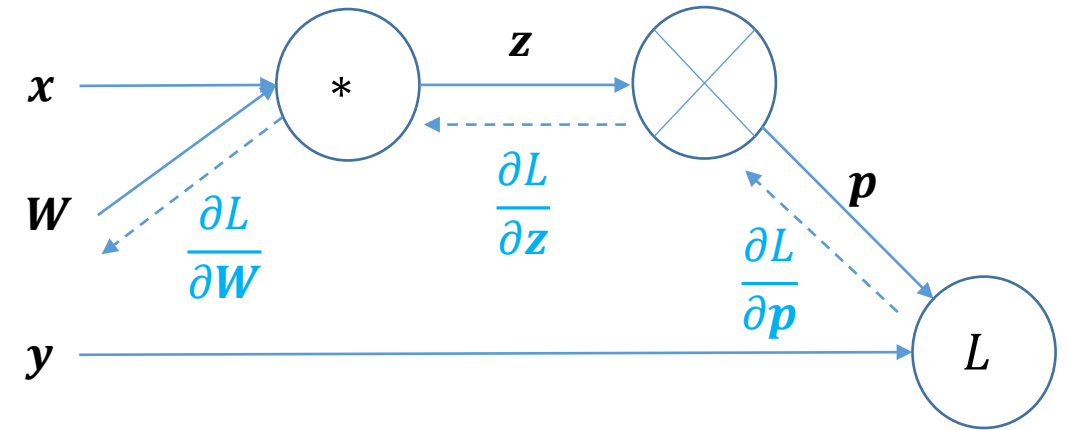
- $\mathbf{z} = \mathbf{W}\mathbf{x}$ ,  $\mathbf{x} \in \mathbb{R}^n, \mathbf{W} \in \mathbb{R}^{k \times n}$
- $\mathbf{p} = \text{softmax}(\mathbf{z})$   $\mathbf{p} \in \mathbb{R}^k$
- $L(\mathbf{p}, \mathbf{y}) = \sum_{i=1} -y_i \log p_i$

- Matmul *matrix multiplication*
  - Forward( $\mathbf{x}, \mathbf{W}$ ):  $\mathbf{W}\mathbf{x}$
  - Backward( $d\mathbf{z}, \mathbf{x}, \mathbf{W}$ )
- Softmax
  - Forward( $\mathbf{z}$ ):  $\mathbf{p}$
  - Backward( $d\mathbf{p}, \mathbf{z}$ )
- Cross-entropy
  - Forward( $\mathbf{p}, \mathbf{y}$ ):  $\sum_{i=1} -y_i \log p_i$
  - Backward( $\mathbf{p}, \mathbf{y}$ )



# Softmax regression

- $\mathbf{z} = \mathbf{W}\mathbf{x}$ ,  $\mathbf{x} \in R^n, \mathbf{W} \in R^{k \times n}$
- $\mathbf{p} = \text{softmax}(\mathbf{z})$   $\mathbf{p} \in R^k$
- $L(\mathbf{p}, \mathbf{y}) = \sum_{i=1} -y_i \log p_i$



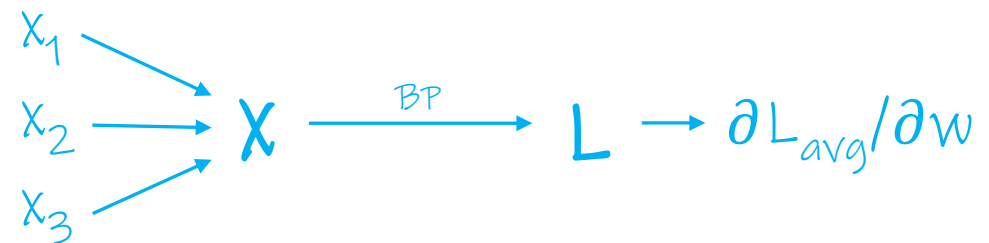
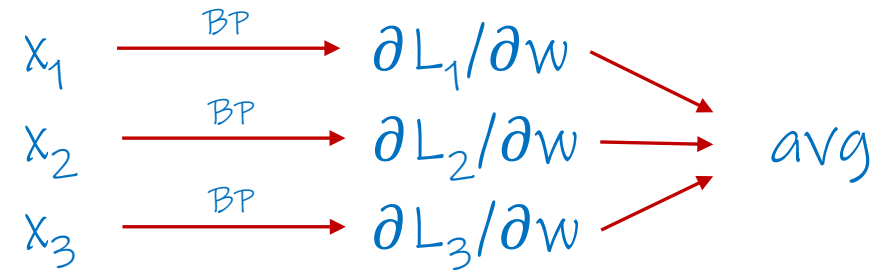
L is a complex node

- Matmul
  - Forward( $\mathbf{x}, \mathbf{W}$ ):  $\mathbf{W}\mathbf{x}$
  - Backward( $d\mathbf{z}, \mathbf{x}, \mathbf{W}$ ):  $d\mathbf{z} \mathbf{x}^T$
- Softmax-Cross-Entropy
  - Forward( $\mathbf{z}$ ):  $\mathbf{p}$  ;  $\sum_{i=1} -y_i \log p_i$
  - Backward( $\mathbf{z}, \mathbf{y}$ ):  $\mathbf{p} - \mathbf{y}$

# BP for multiple examples

So far, we've considered backpropagation for a single example.  
How should we handle all of our training samples?

- Approach I (individual)
  - Backpropagation for each example separately
  - Average the gradients across all examples
- Approach II (vectorised)
  - Use matrix (one row per example) in the BP
  - Average the loss
  - Compute gradient based on averaged loss



which approach is better and why?

# BP operations

These are all derived  
via chain rule

- Add\_bias e.g.  $z = xw + b$ 
  - $\mathbf{A}$  is a matrix;  $b$  is a row vector
  - Forward( $\mathbf{A}, b$ ):  $\mathbf{C} = \mathbf{A} + b$
  - Backward( $d\mathbf{C}, \mathbf{A}, b$ ):  $d\mathbf{A} = d\mathbf{C}, db = \mathbf{1}^T d\mathbf{C}$
- Array and scalar multiplication
  - $\mathbf{v}$  is an array,  $k$  is a scalar
  - Forward( $\mathbf{v}, k$ ):  $\mathbf{c} = k \mathbf{v}$
  - Backward( $d\mathbf{c}, \mathbf{v}, k$ ):  $d\mathbf{v} = k d\mathbf{c}$

# BP operations

- Matmul matrix multiplication operation
  - $\mathbf{A} \in R^{m \times k}, \mathbf{B} \in R^{k \times n}$  (including matrix with a single column or row)
  - Forward( $\mathbf{A}, \mathbf{B}$ ):  $\mathbf{C} = \mathbf{AB} \in R^{m \times n}$
  - Backward( $d\mathbf{C}, \mathbf{A}, \mathbf{B}$ ):  $d\mathbf{A} = ?, d\mathbf{B} = ?$
- Logistic operation
  - $a$  is an array of any shape
  - Forward( $a$ ):  $b = \sigma(a)$
  - Backward( $db, a$ ):  $da = ?$
- Softmax-Cross-entropy operation
  - $\mathbf{Z}$  and  $\mathbf{Y} \in R^{m \times k}$  are matrix of the same shape; each row of  $\mathbf{a}$  (or  $\mathbf{b}$ ) sums to 1
  - Forward ( $\mathbf{Z}, \mathbf{Y}$ ):  $\mathbf{P} = \text{softmax}(\mathbf{Z}); L = \text{sum}(-\mathbf{Y} \log \mathbf{P})/m$
  - Backward( $\mathbf{Z}, \mathbf{Y}$ ):  $d\mathbf{Z} = ?$

# Summary

- Extending simple Perceptron to Multilayer Perceptron model
- Backpropagation algorithm
  - A modular way to compute the gradient of the loss w.r.t parameters
  - Based on chain rules and matrix calculus

# Homework & Practice

- Homework
  - slide 6 & 38,
- Practice
  - Colab notebook (implement the BP operations in Python)