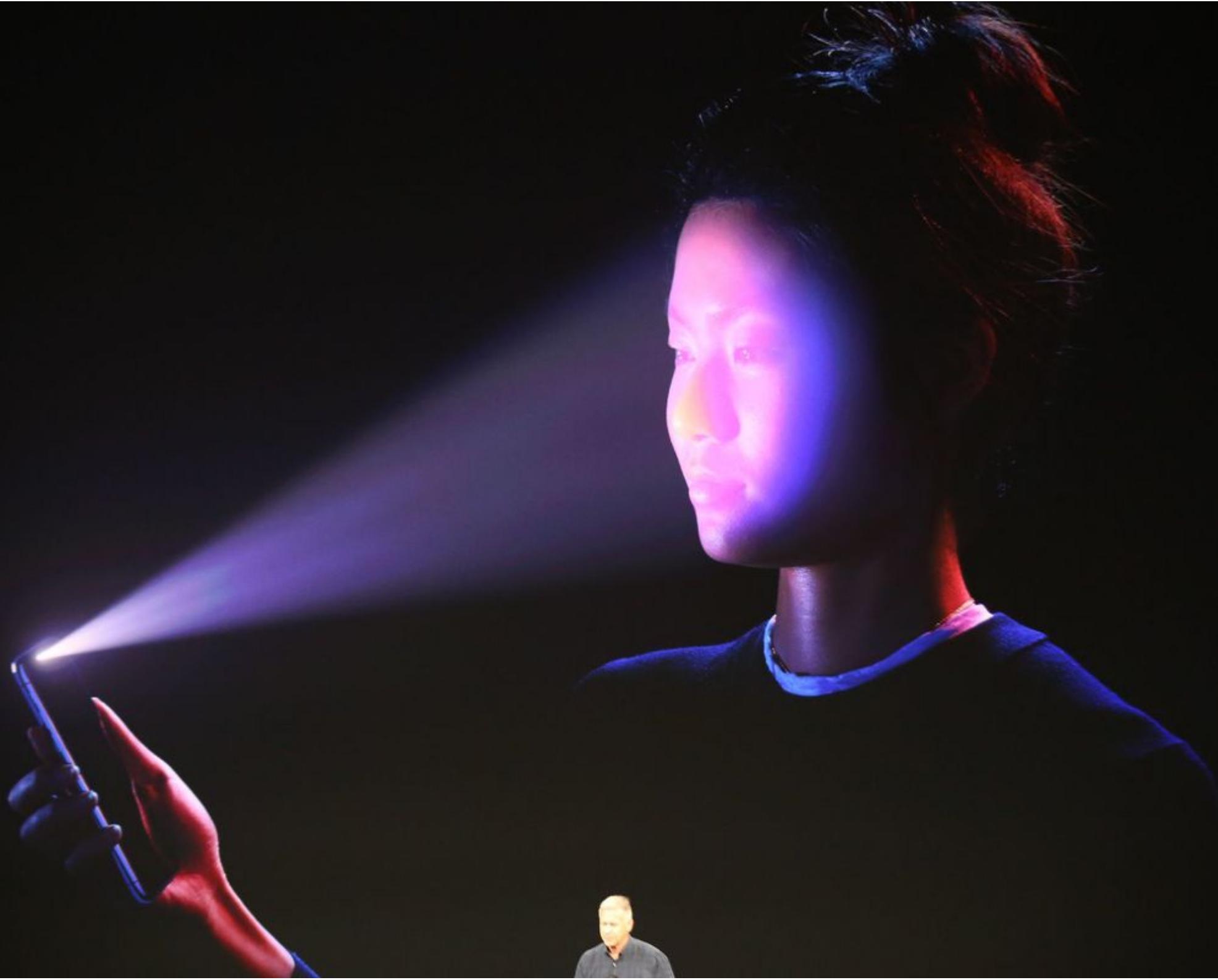
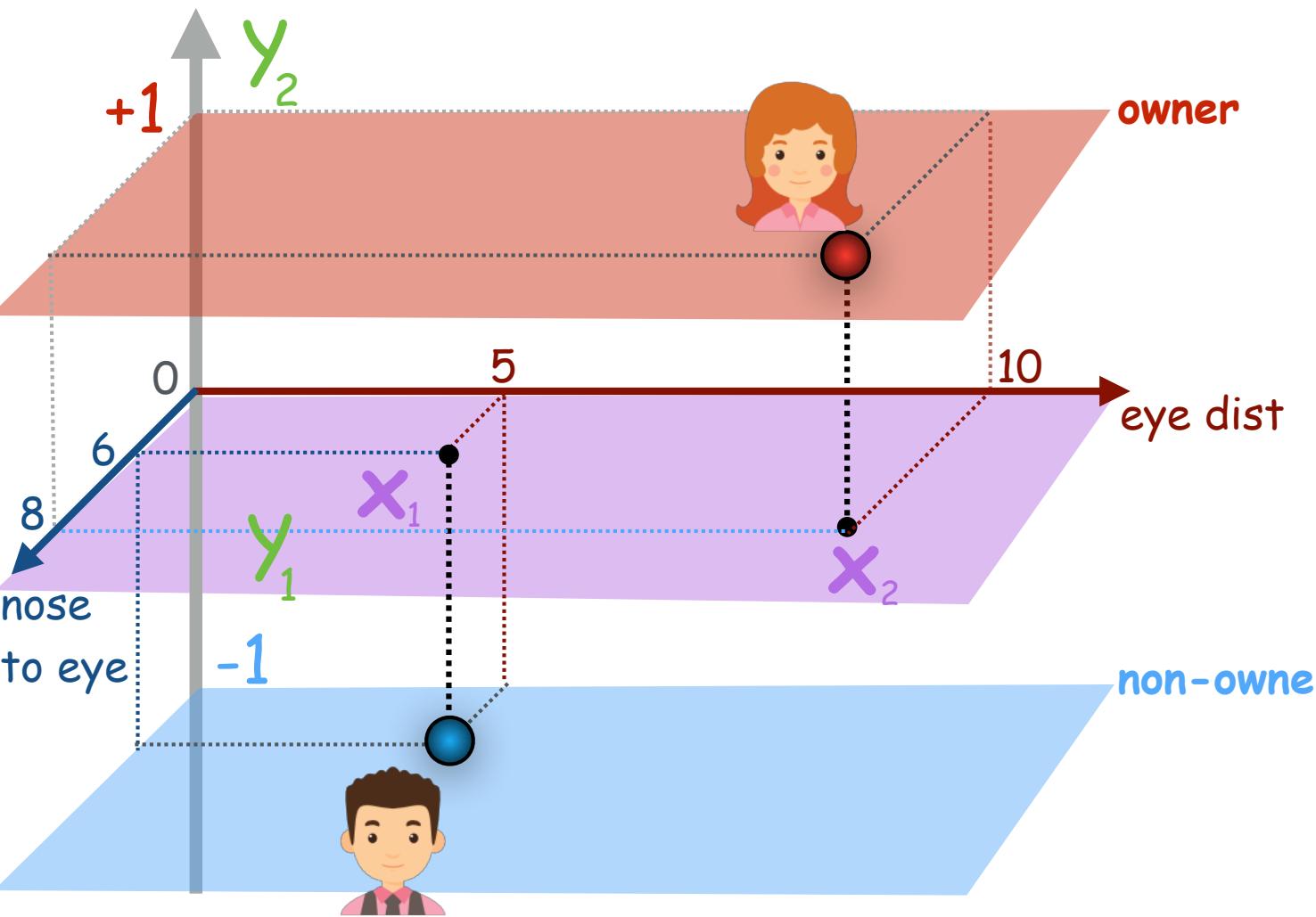


Neural Networks



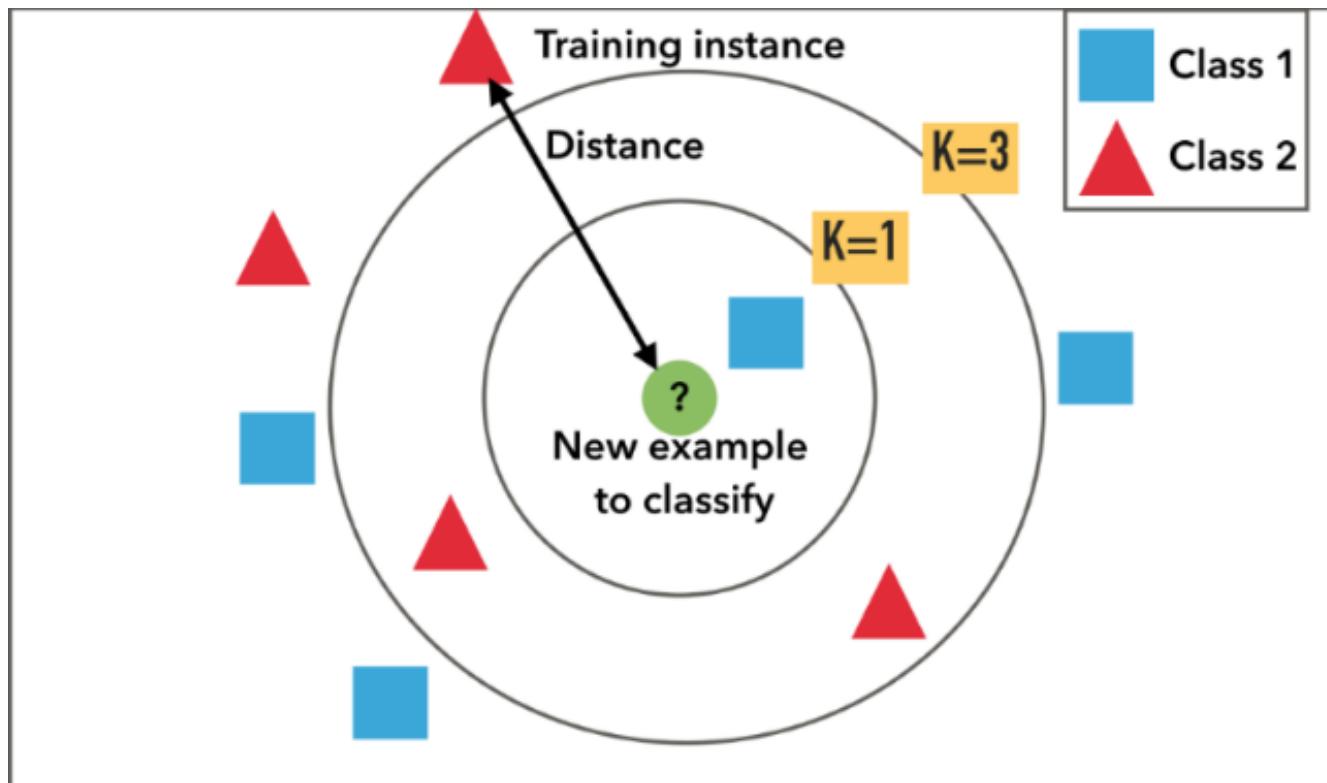
Classification Task

Owner?

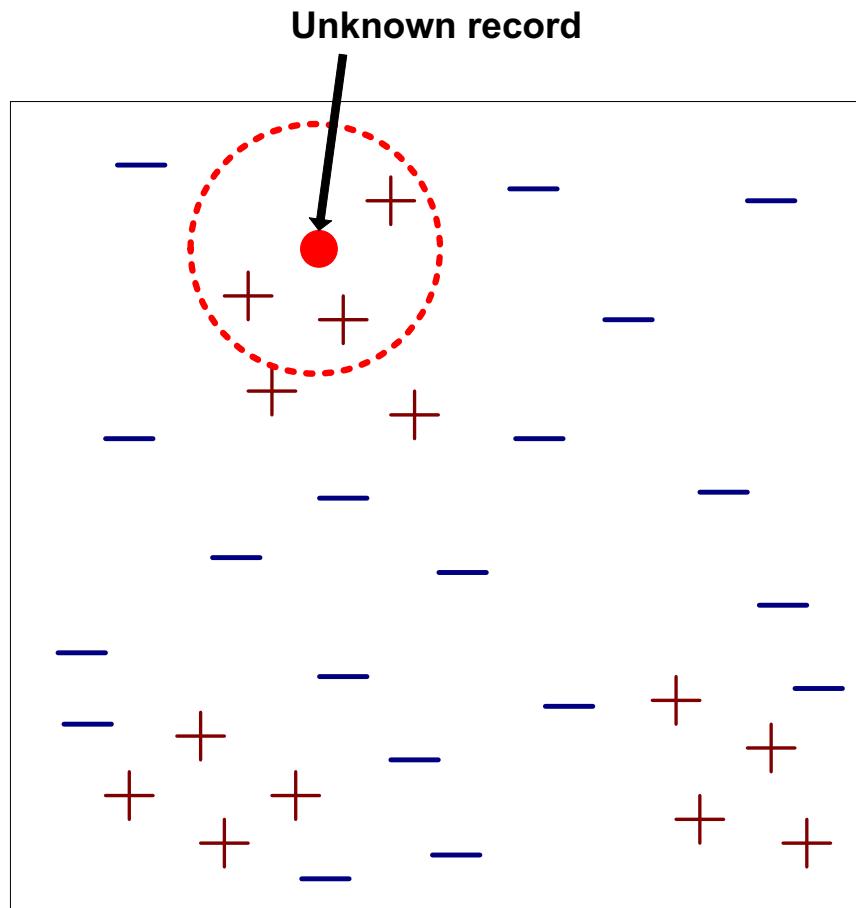


$\mathbf{x} \in \mathbb{R}^2$ $y \in \{-1, 1\}$ find a function $y = f(\mathbf{x})$

Nearest Neighbor Classifier



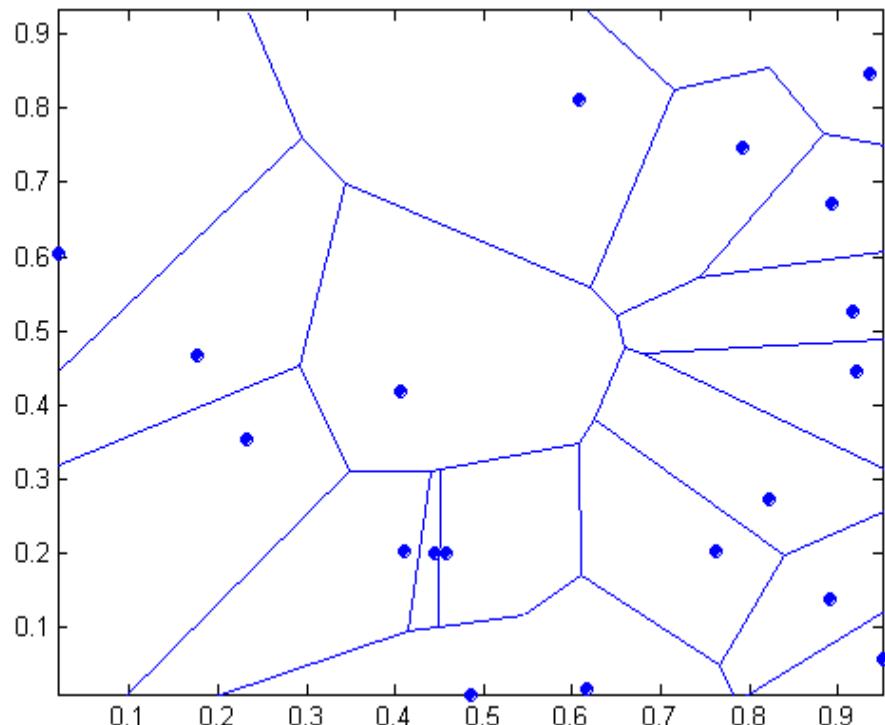
Nearest Neighbor Classifier



- Requires three things
 - The set of stored records
 - Distance Metric to compute distance between records
 - The value of k , the number of nearest neighbors to retrieve
- To classify an unknown record:
 - Compute distance to other training records
 - Identify k nearest neighbors
 - Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

KNN

model: training instances



fast training

large memory cost (GB of data)

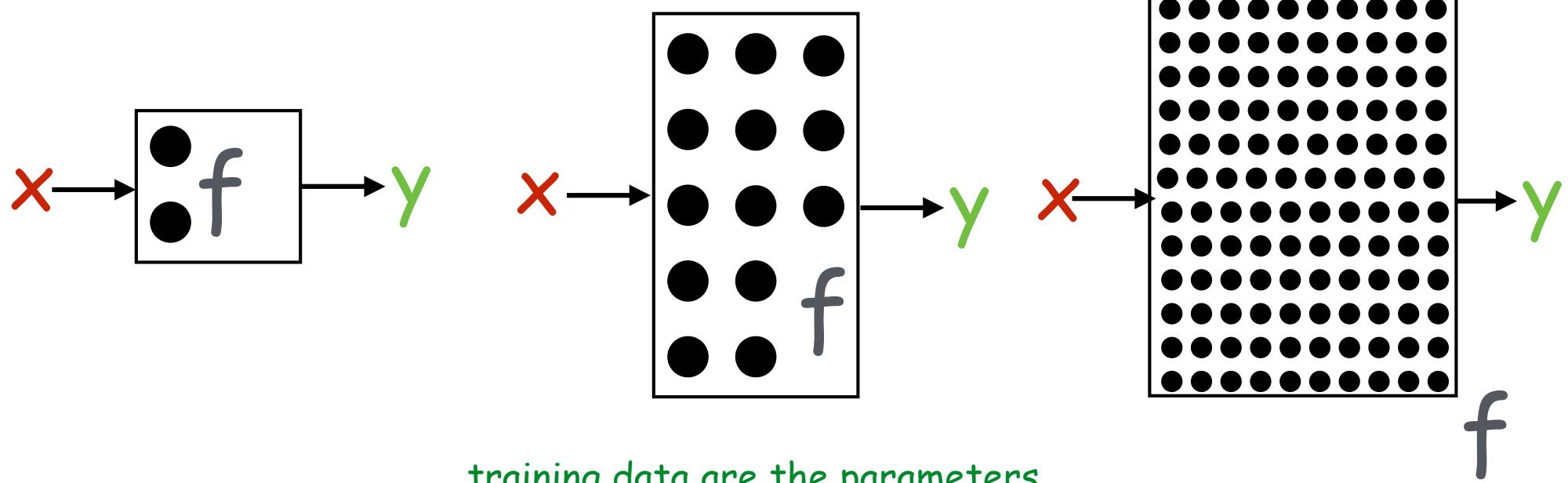
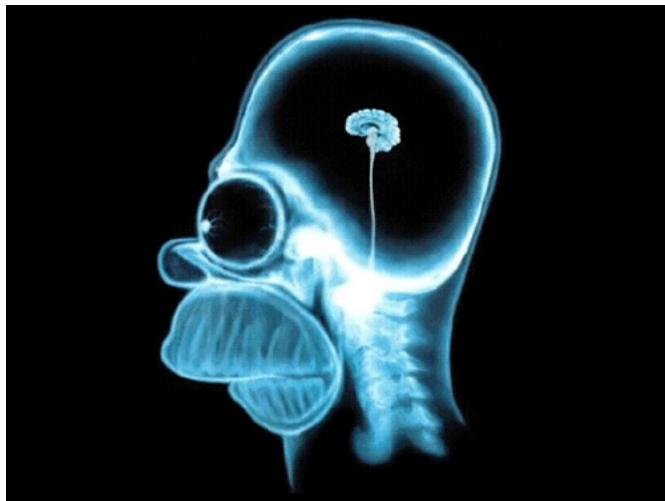
slow inference with large dataset

$O(p n)$ for each test data

n : # training instances

Also called Non-Parametric Models

"Non-parametric"



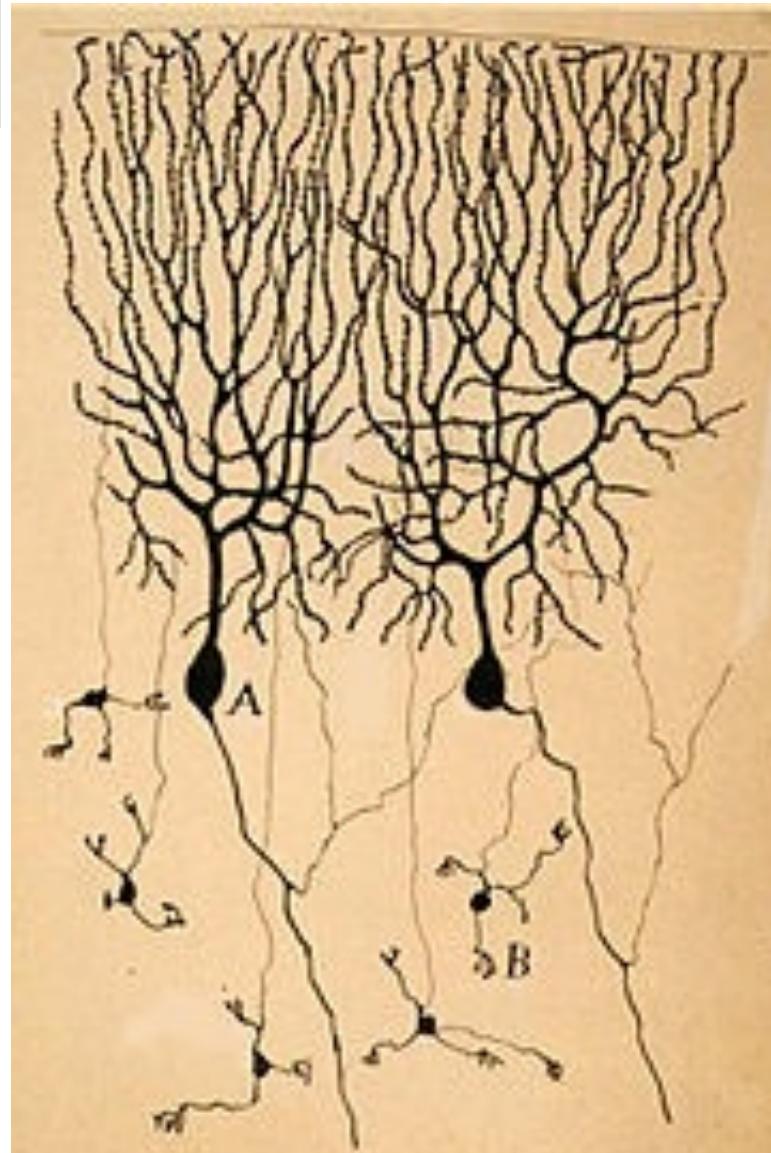
Neuroscience



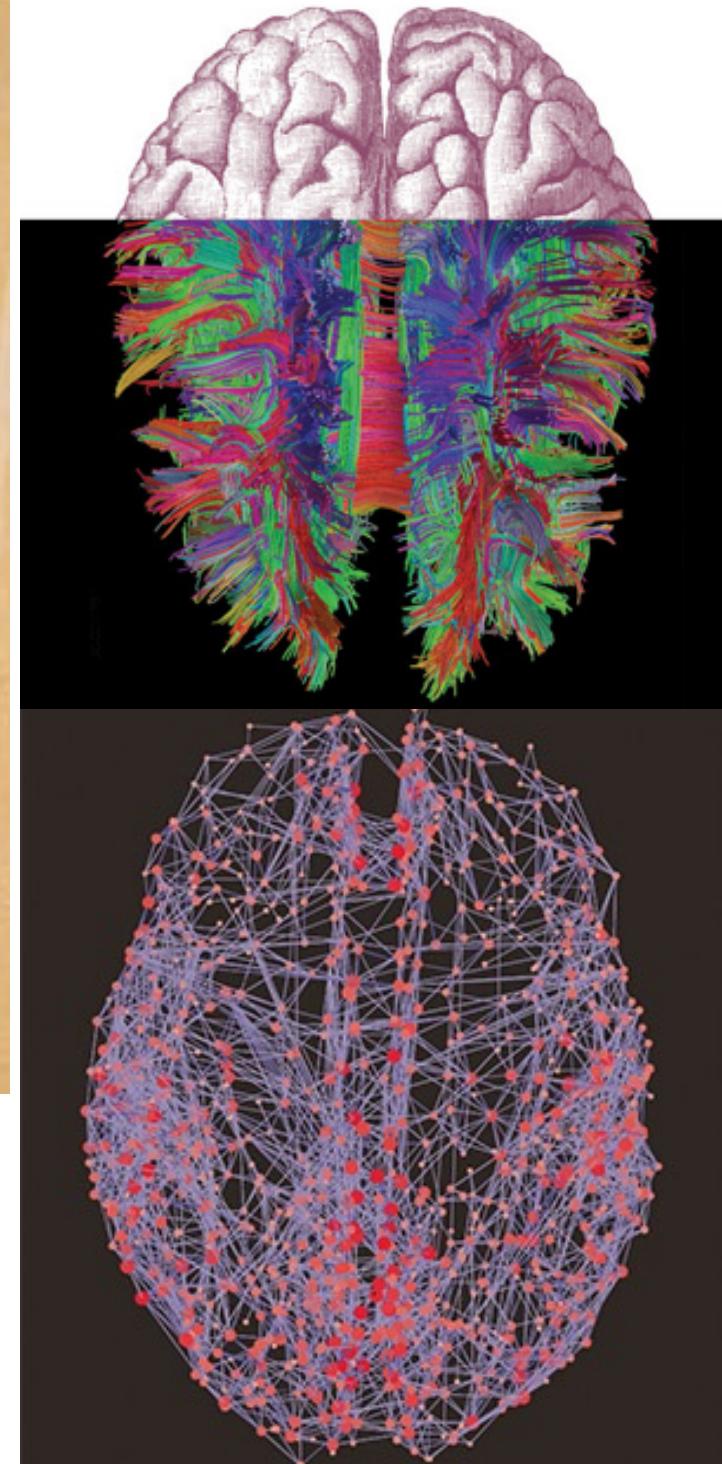
Santiago Ramón y Cajal

Nobel Prize in Medicine

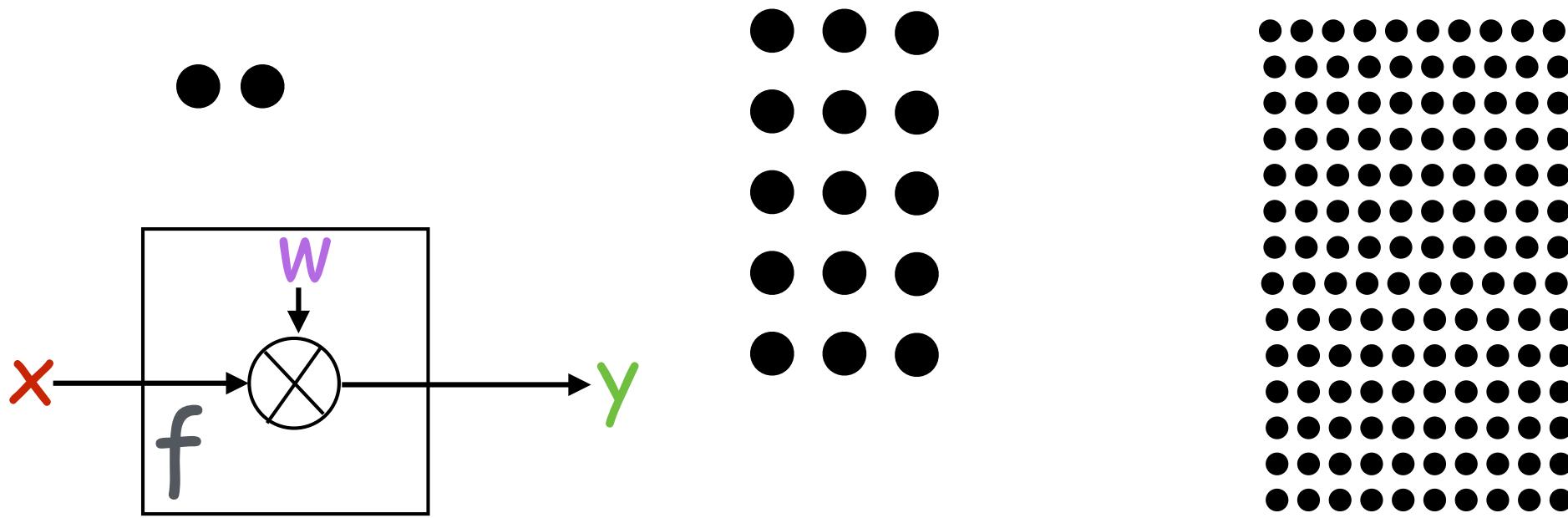
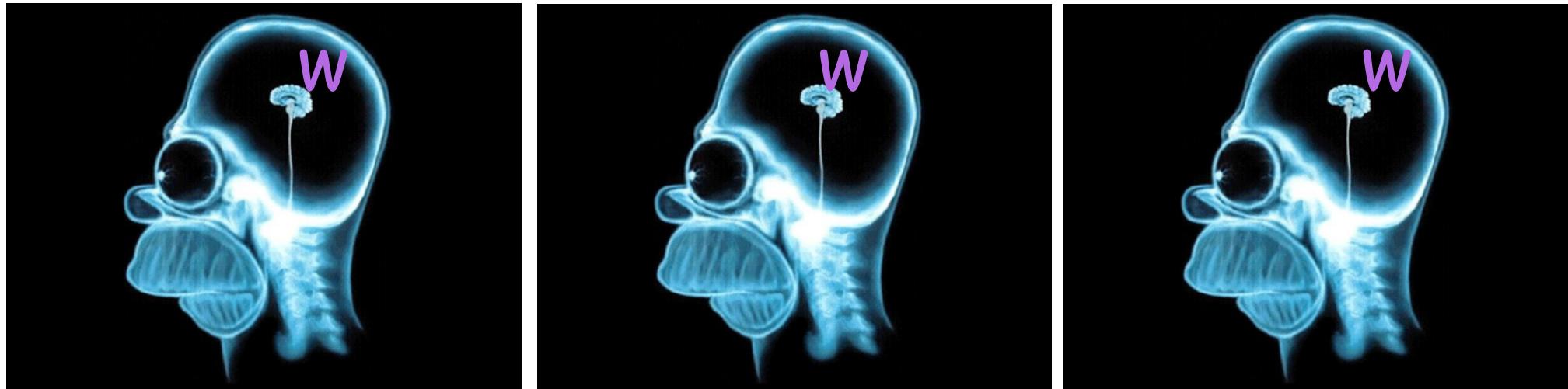
1899



Networks of the Brain



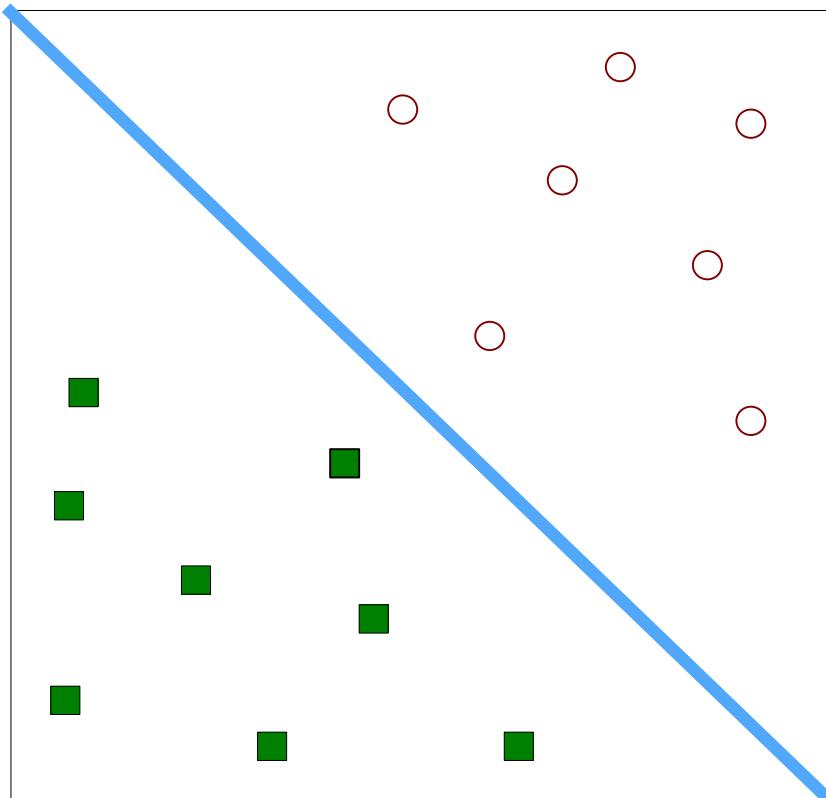
Parametric Approach



" Parametric approach"

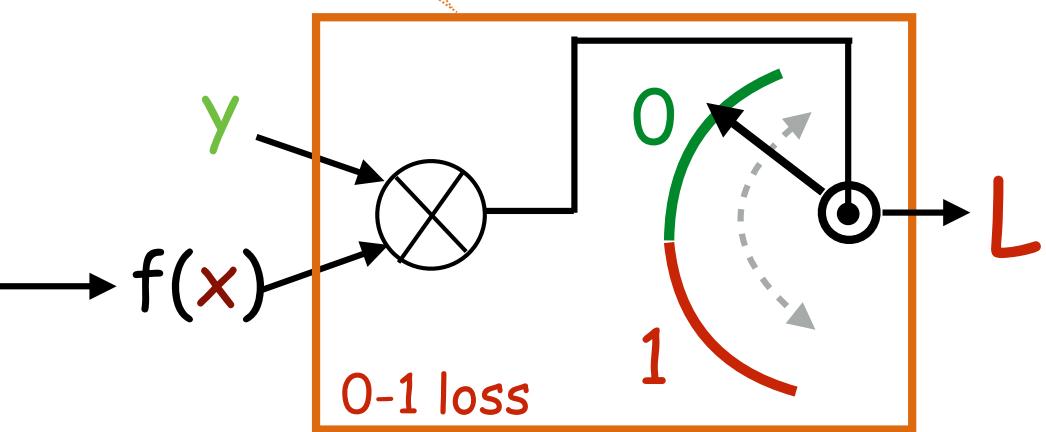
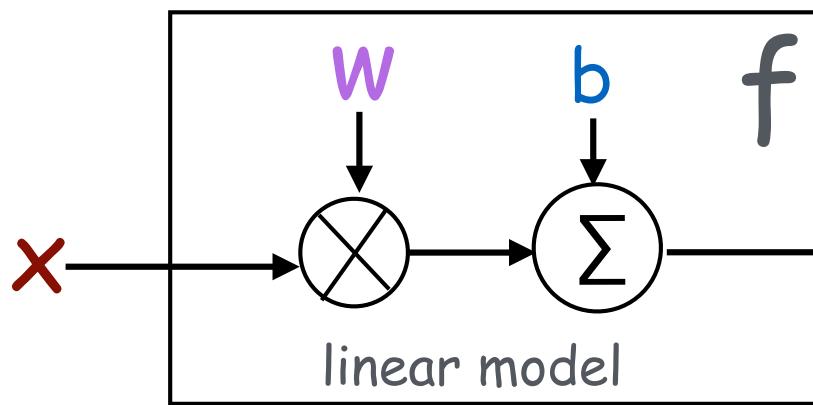
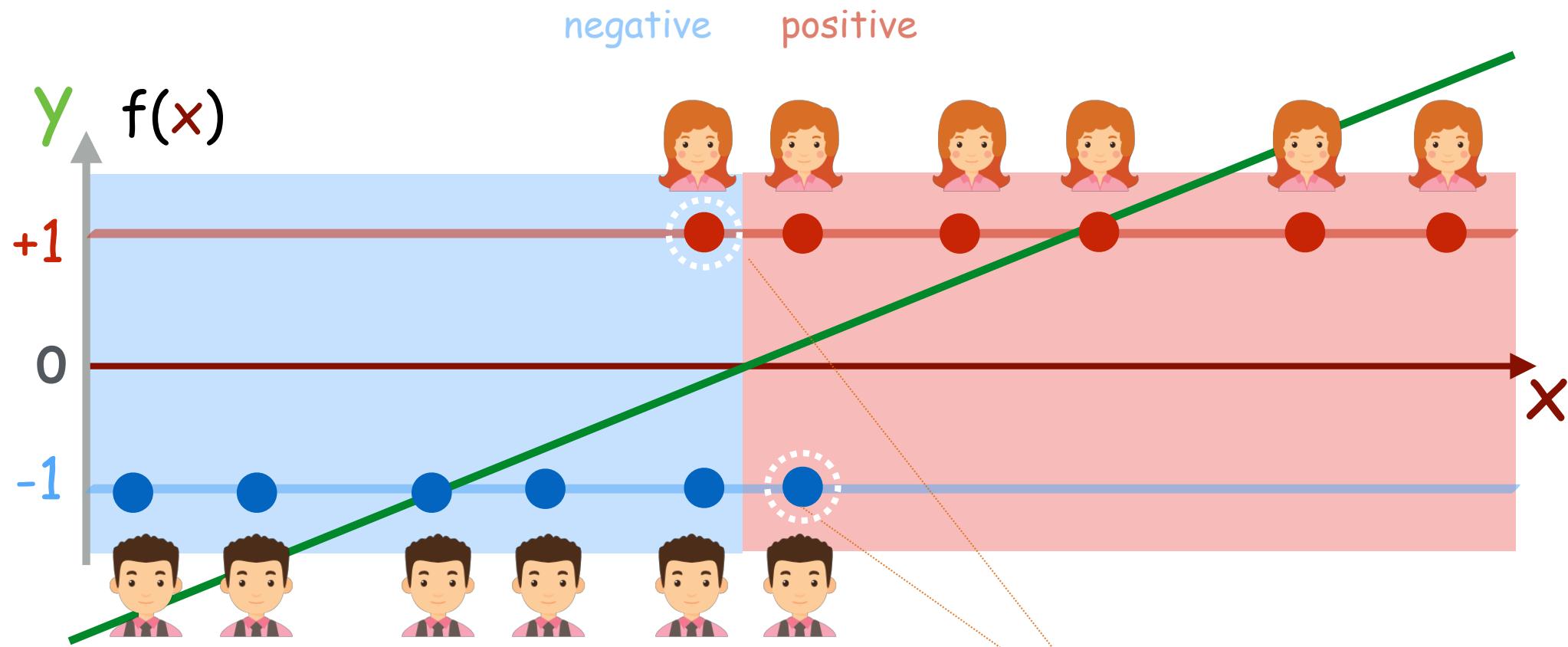
Parametric Models

model: a set of parameters
(the line)

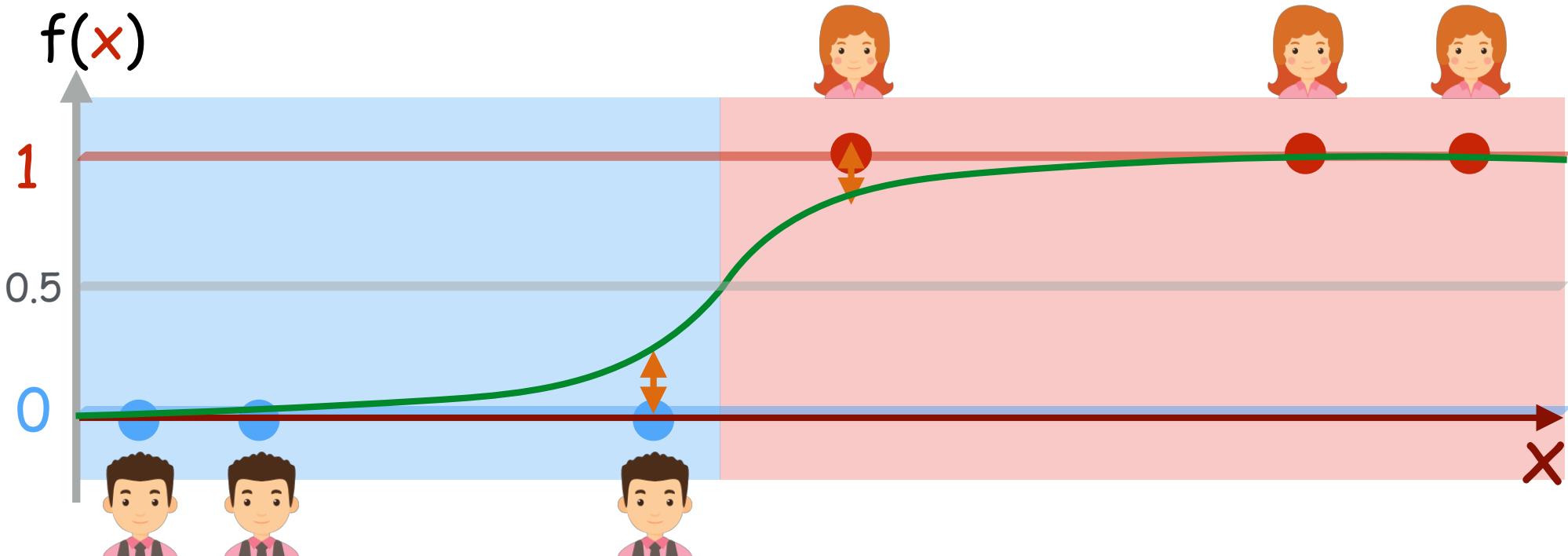


slow in training
low memory cost
fast inference with large dataset
 $O(p)$ for each test data

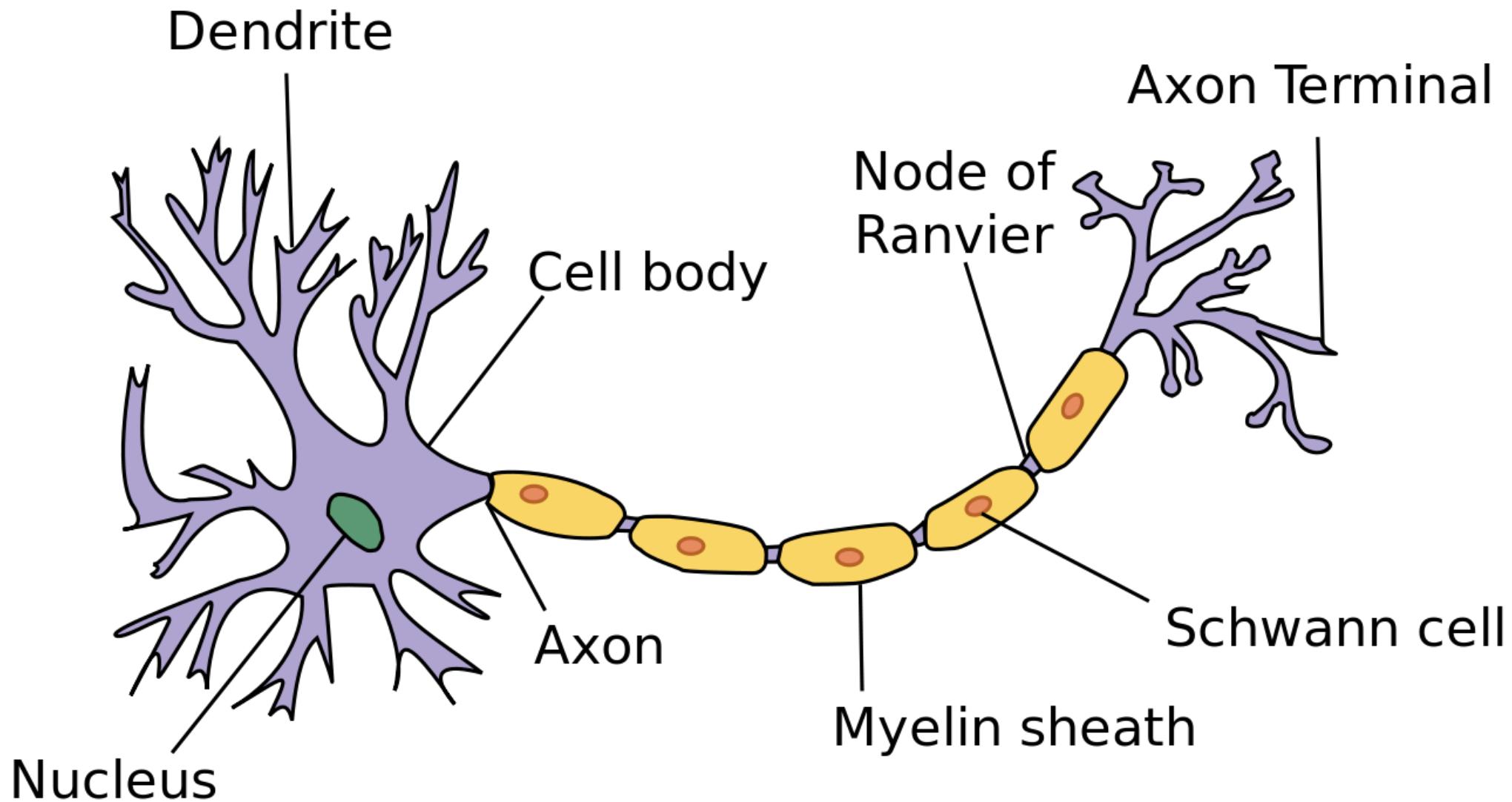
Ideally minimize 0-1 Loss



Logistic Regression

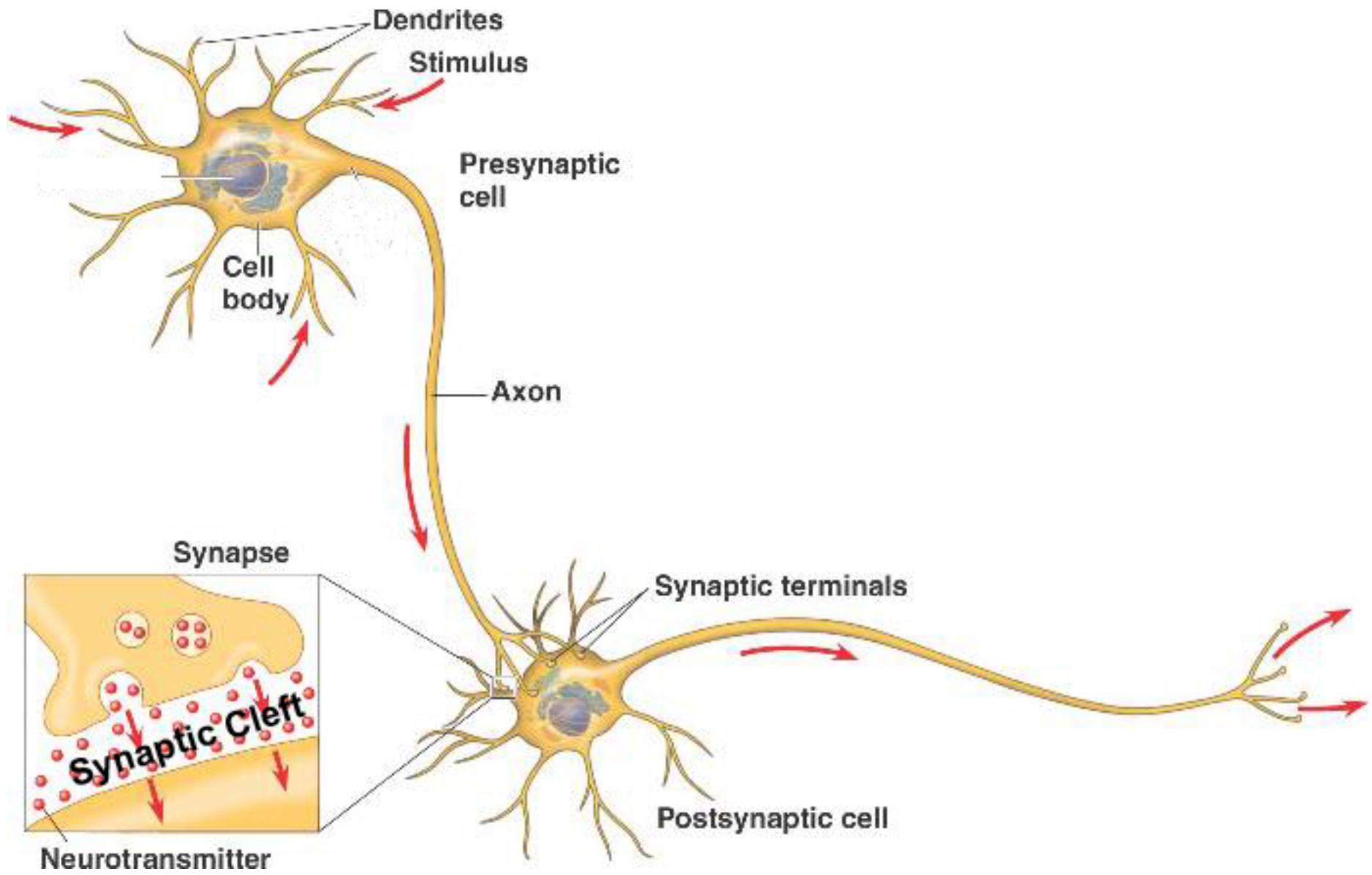


Neuron



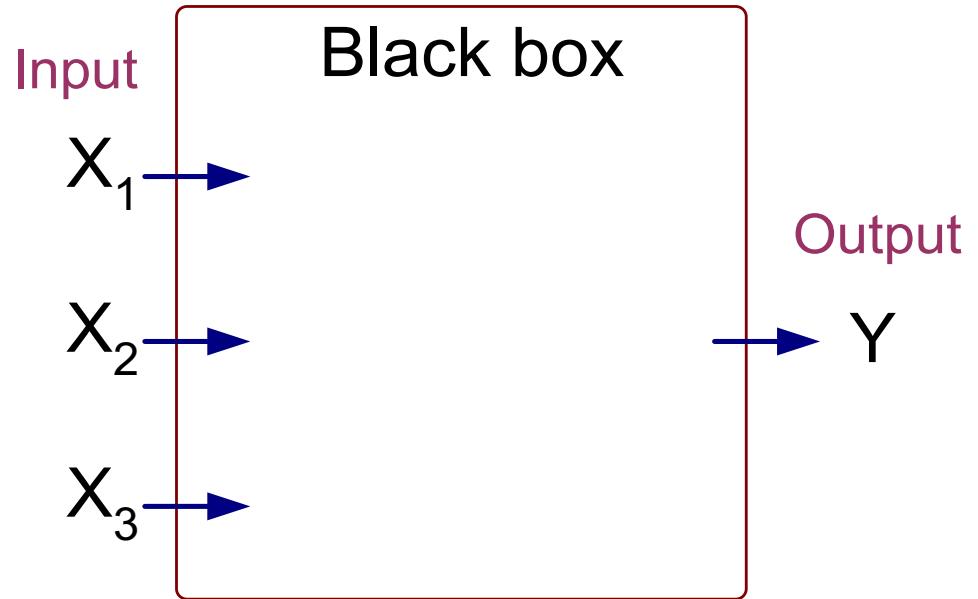
<https://www.youtube.com/watch?v=fHRC8SILcH0>

Neuron



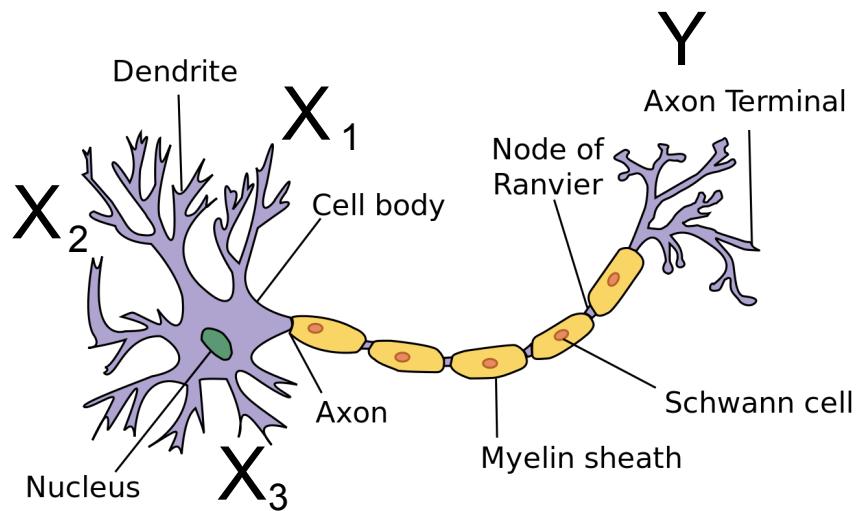
Neuron

| X_1 | X_2 | X_3 | Y |
|-------|-------|-------|-----|
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 |

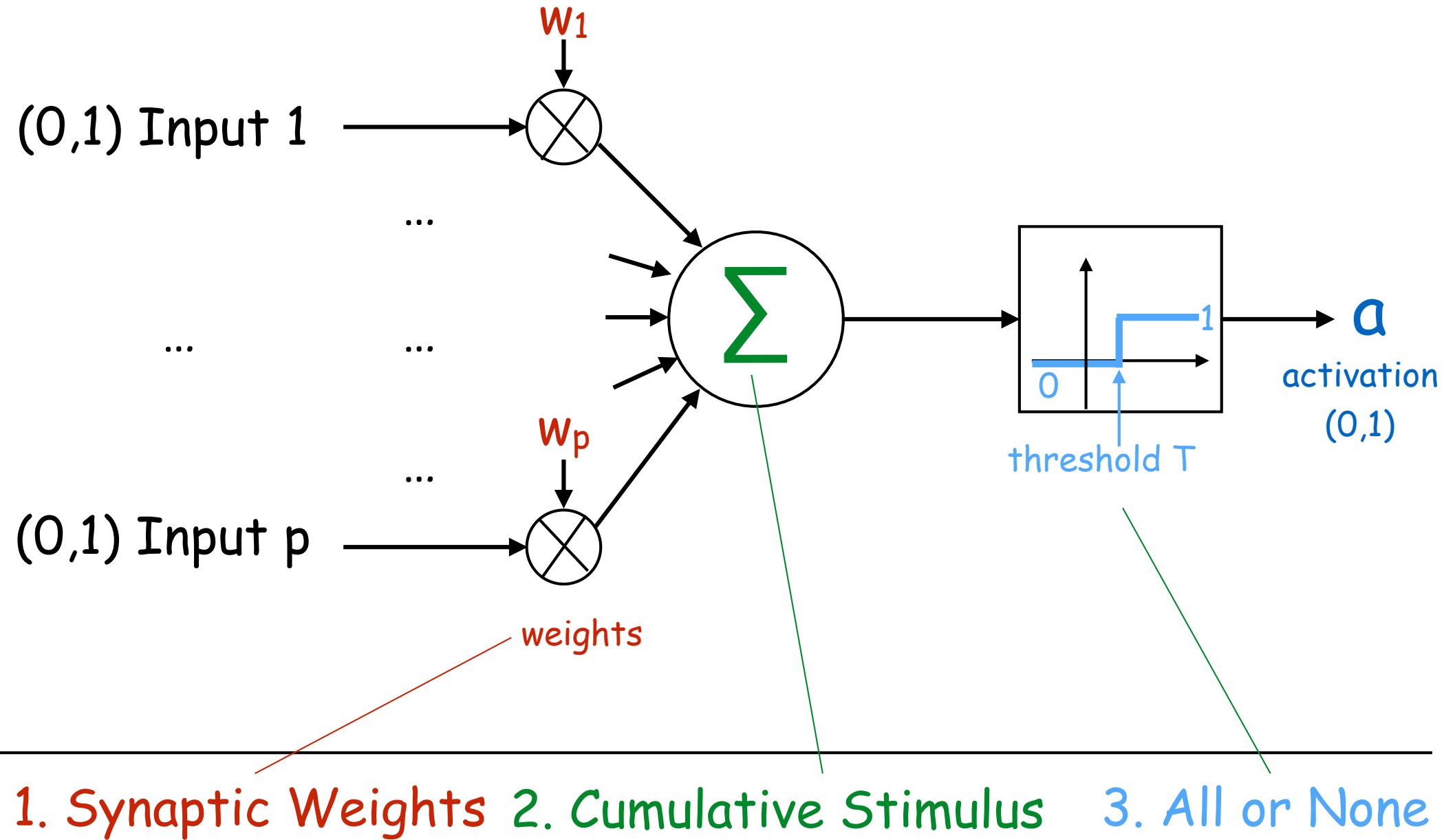


Properties:

1. Synaptic Weights
2. Cumulative Stimulus
3. All or None

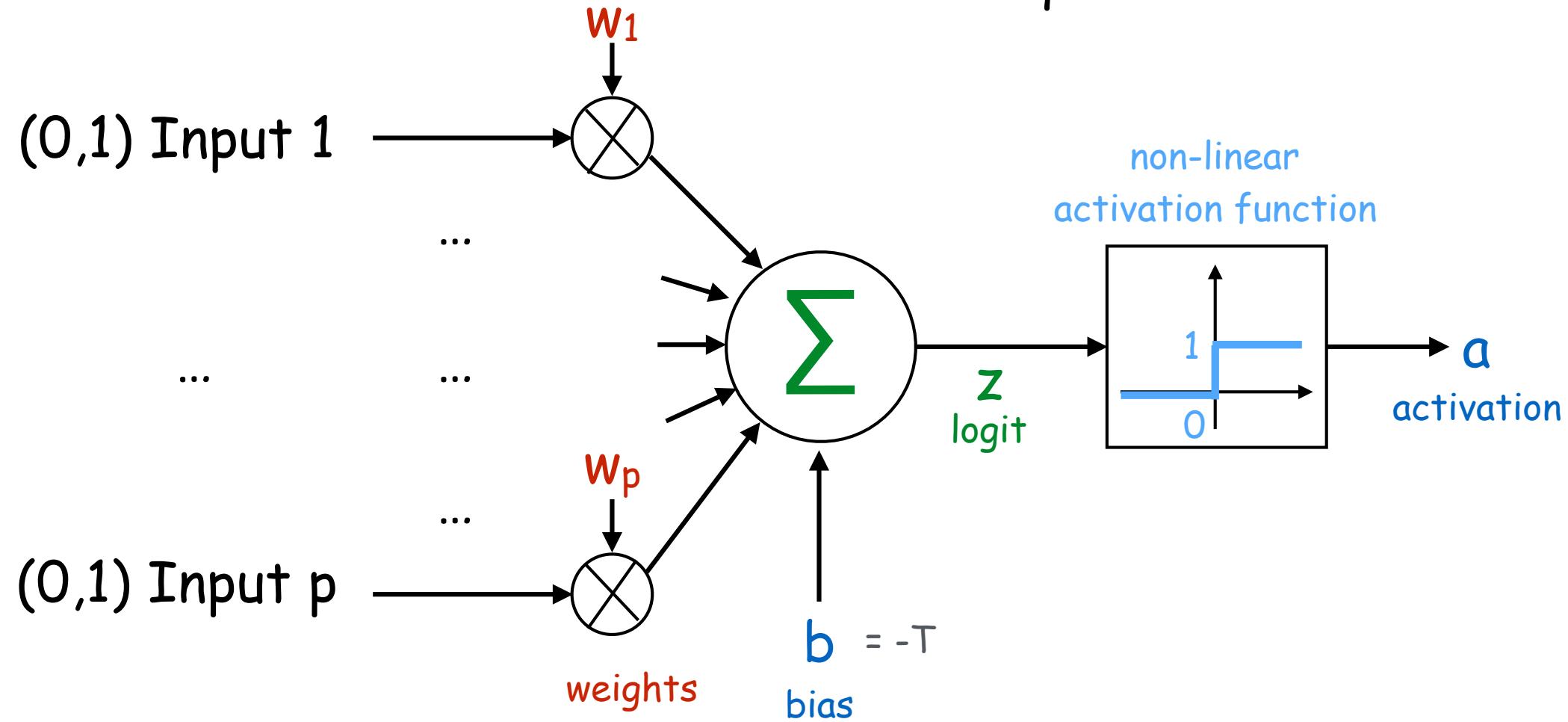


Neuron



Neuron

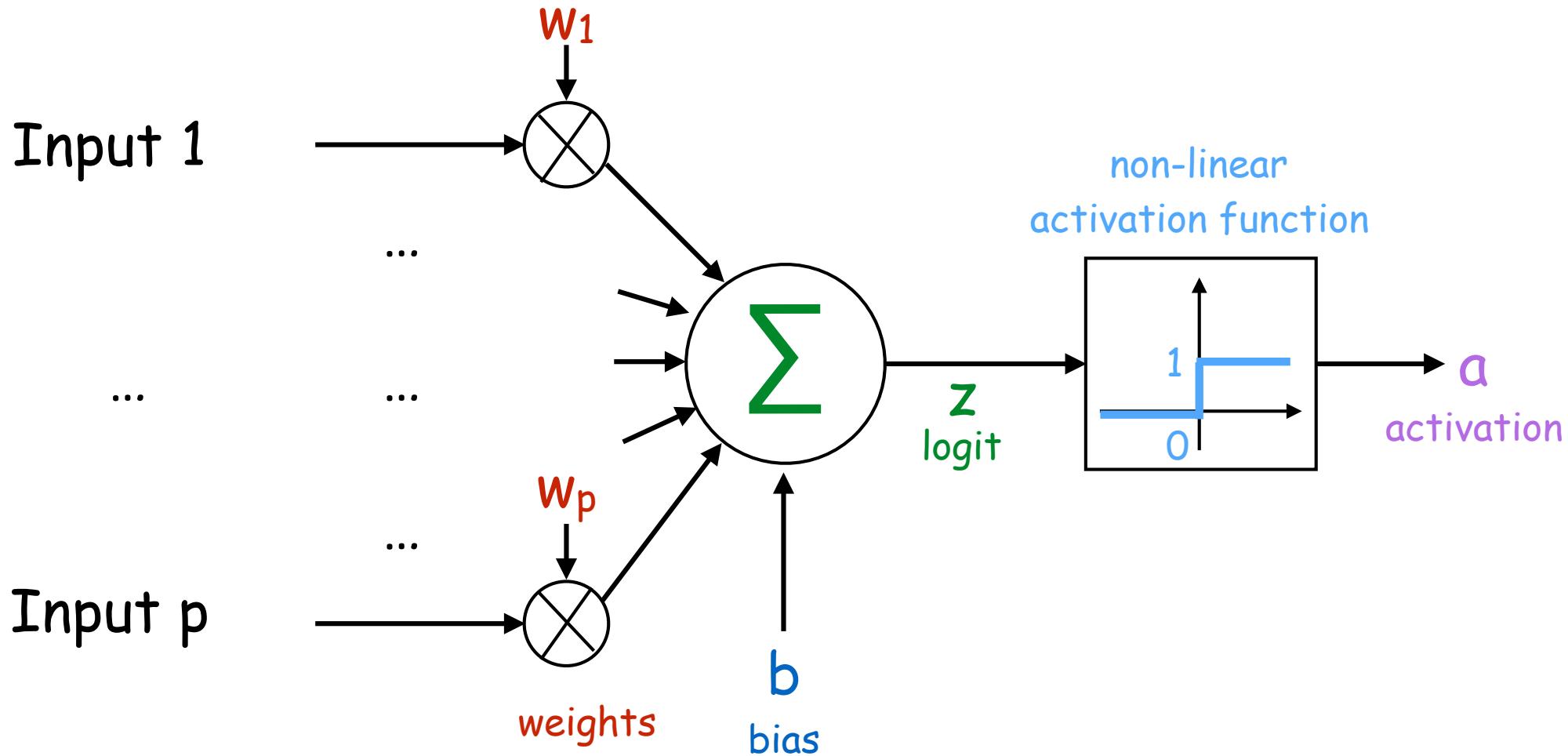
An equivalent model



$$z = w_1 x_1 + \dots + w_p x_p + b \longrightarrow \begin{array}{ll} \text{if } z \geq 0 & a = 1 \\ \text{if } z < 0 & a = 0 \end{array}$$

Neuron

Generalize to continuous input



$$z = w_1 x_1 + \dots + w_p x_p + b \longrightarrow \begin{array}{ll} \text{if } z \geq 0 & a = 1 \\ \text{if } z < 0 & a = 0 \end{array}$$

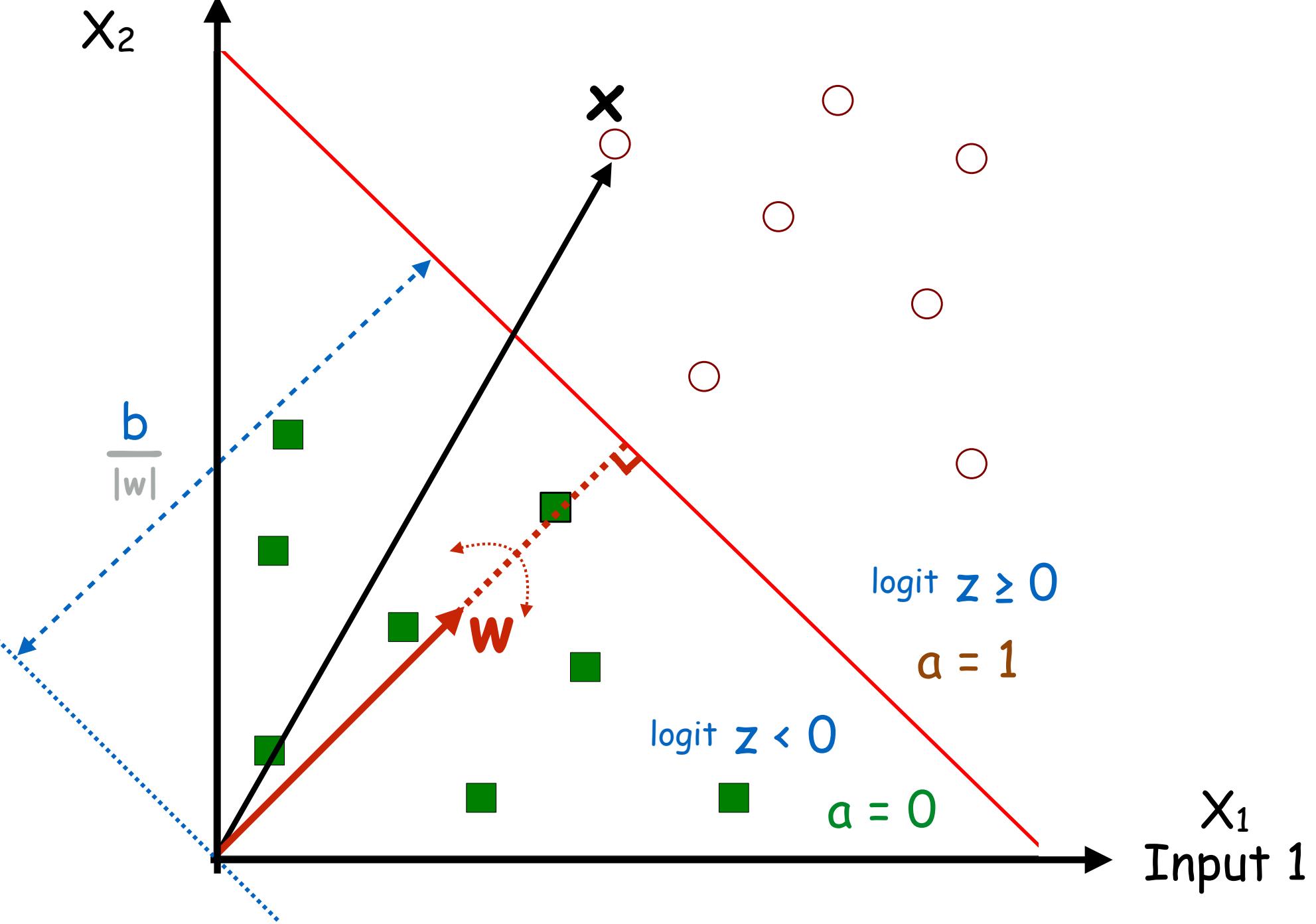
Draw a Line

Input 2
 x_2

x_1
Input 1

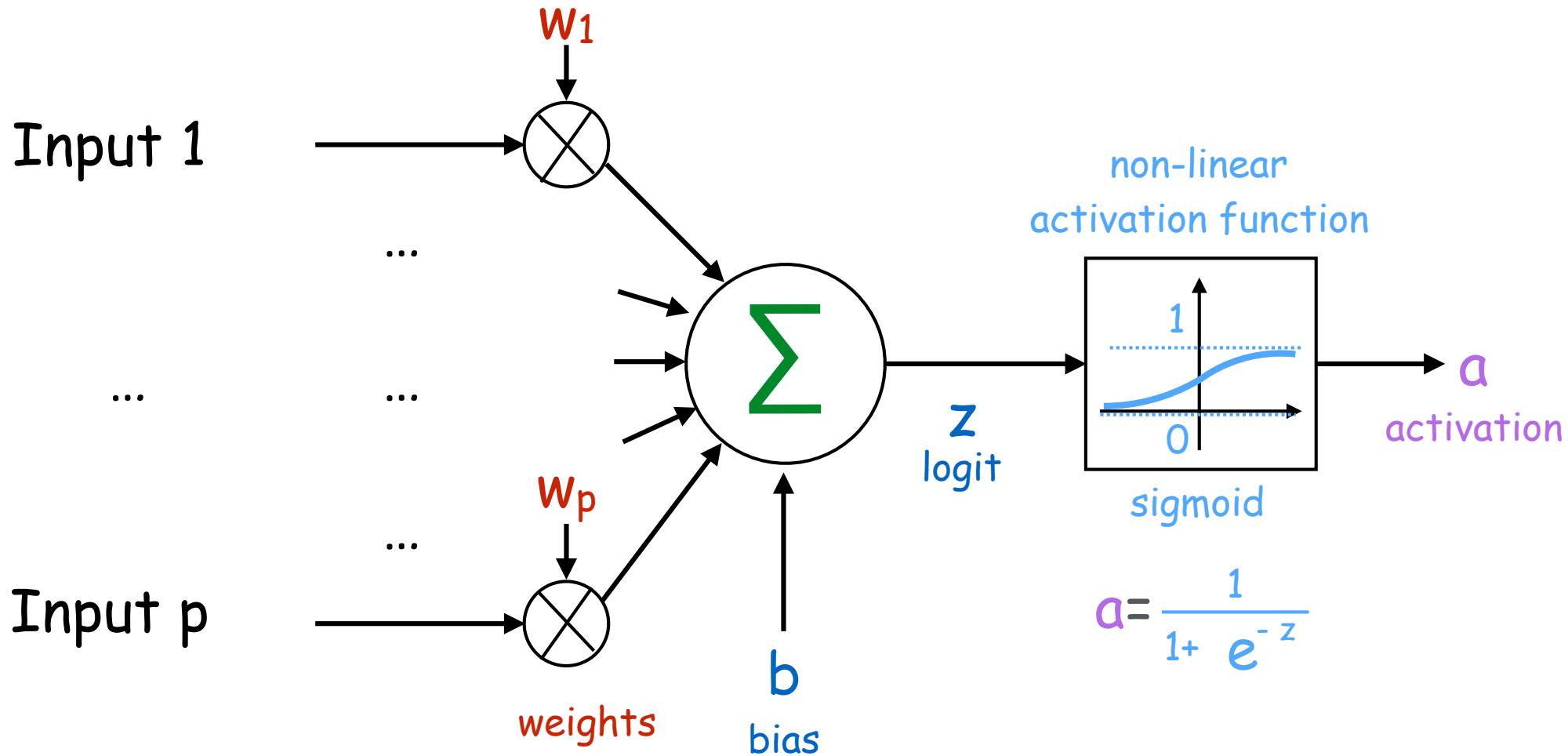
x_2

$$\frac{b}{|w|}$$



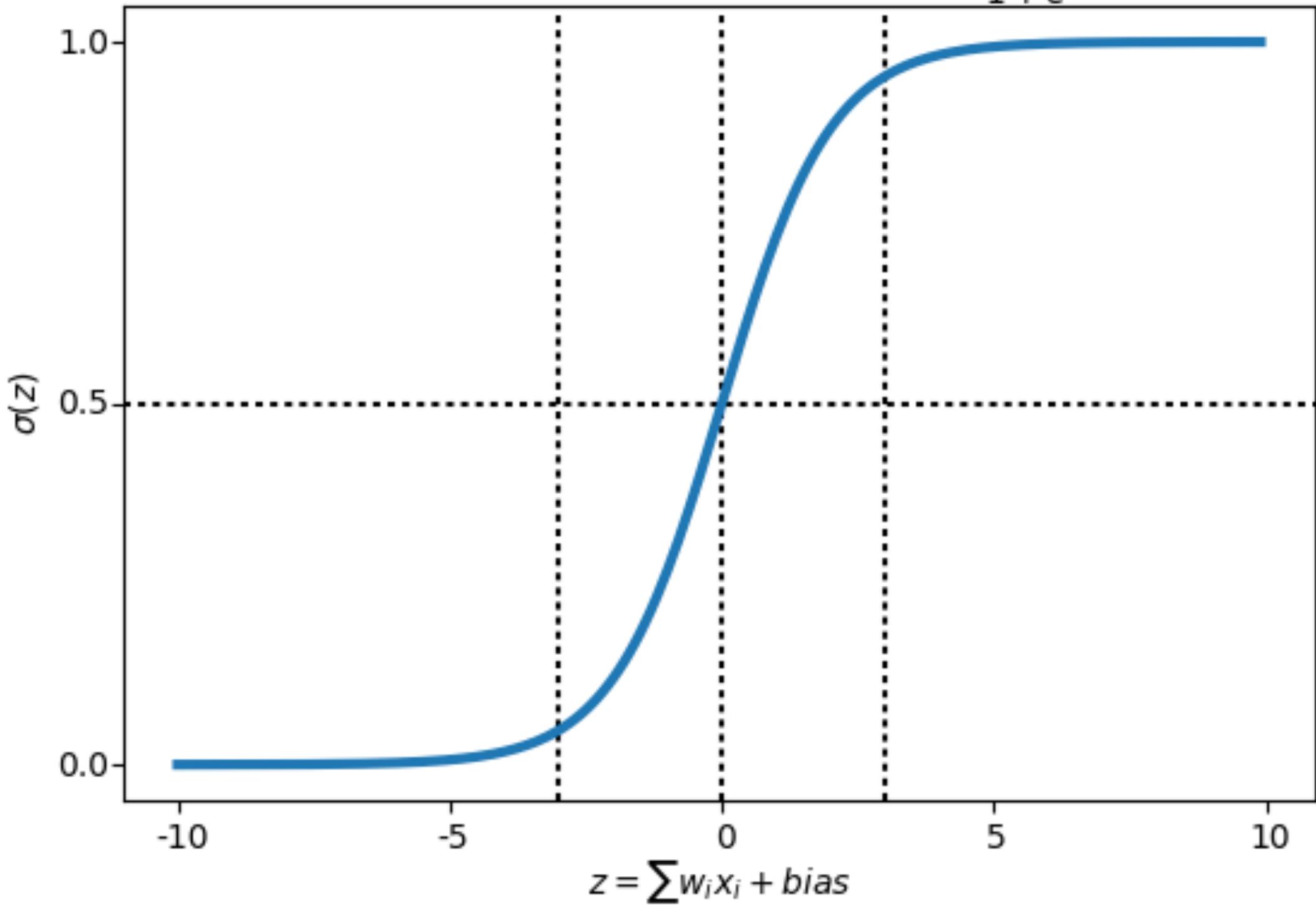
Neuron

Generalize to continuous output

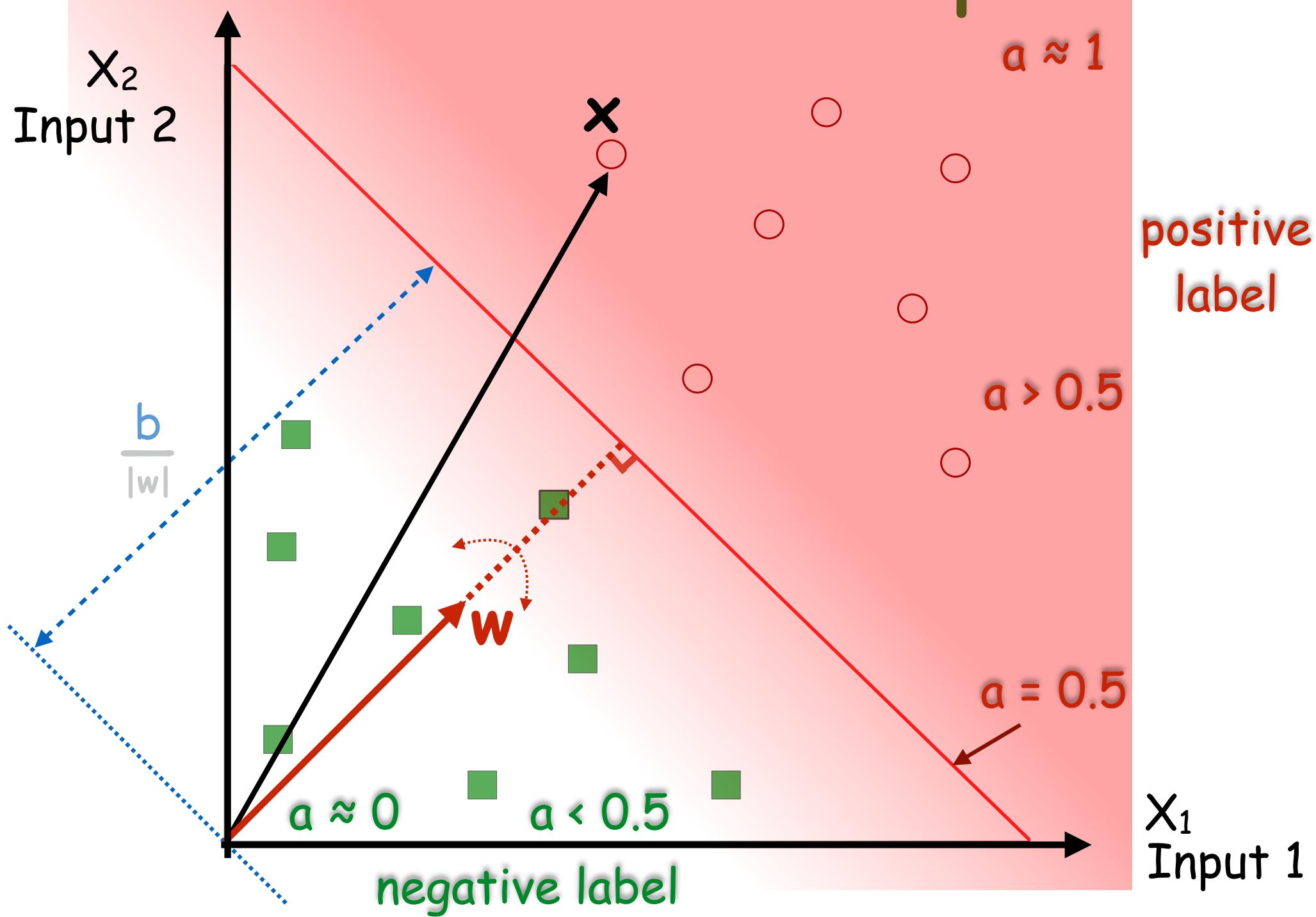


$$w_1 x_1 + \dots + w_p x_p + b = z \longrightarrow a = f(z)$$

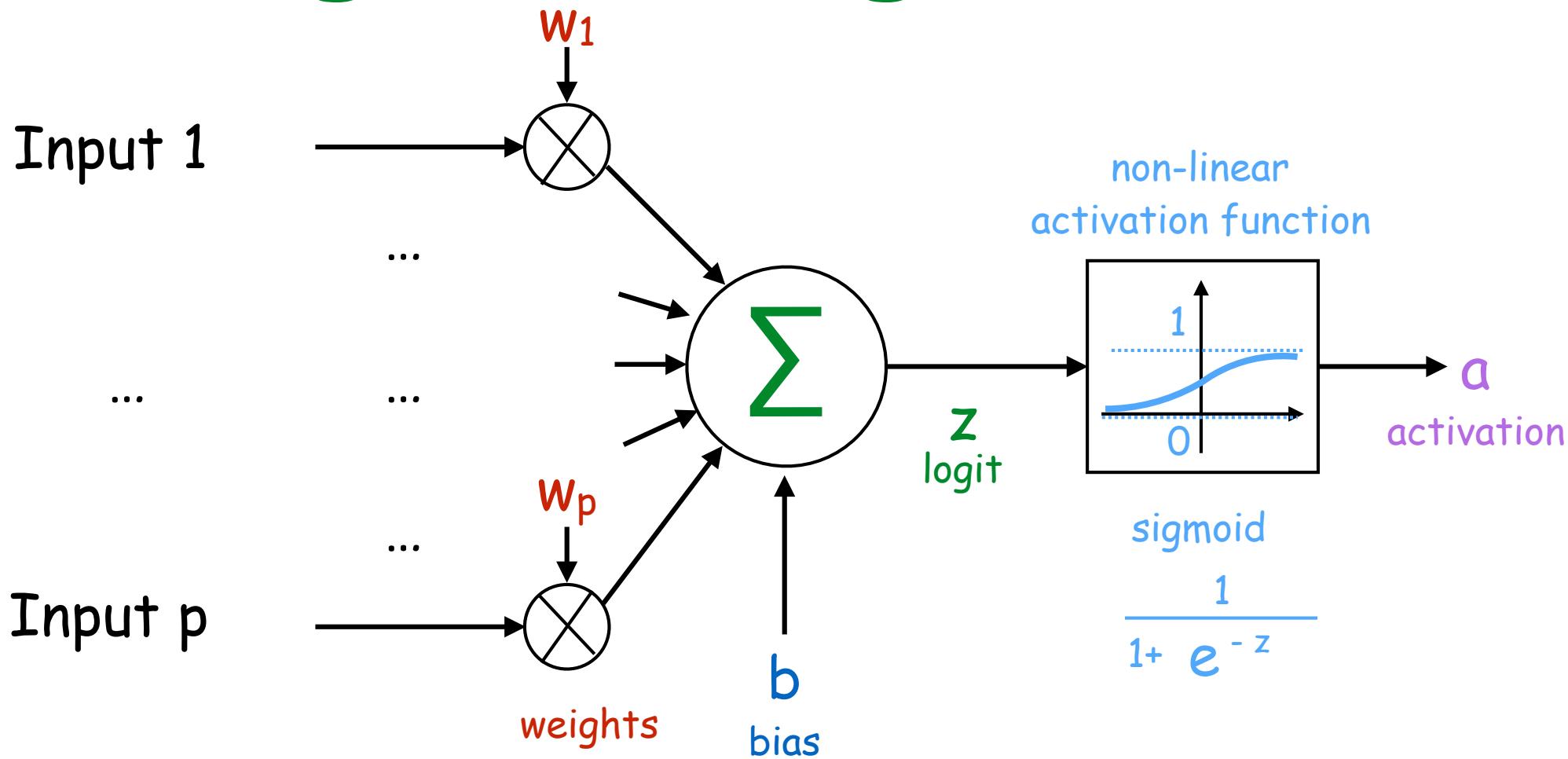
$$\text{Sigmoid Function } \sigma(z) = \frac{1}{1 + e^{-z}}$$



Continuous Output



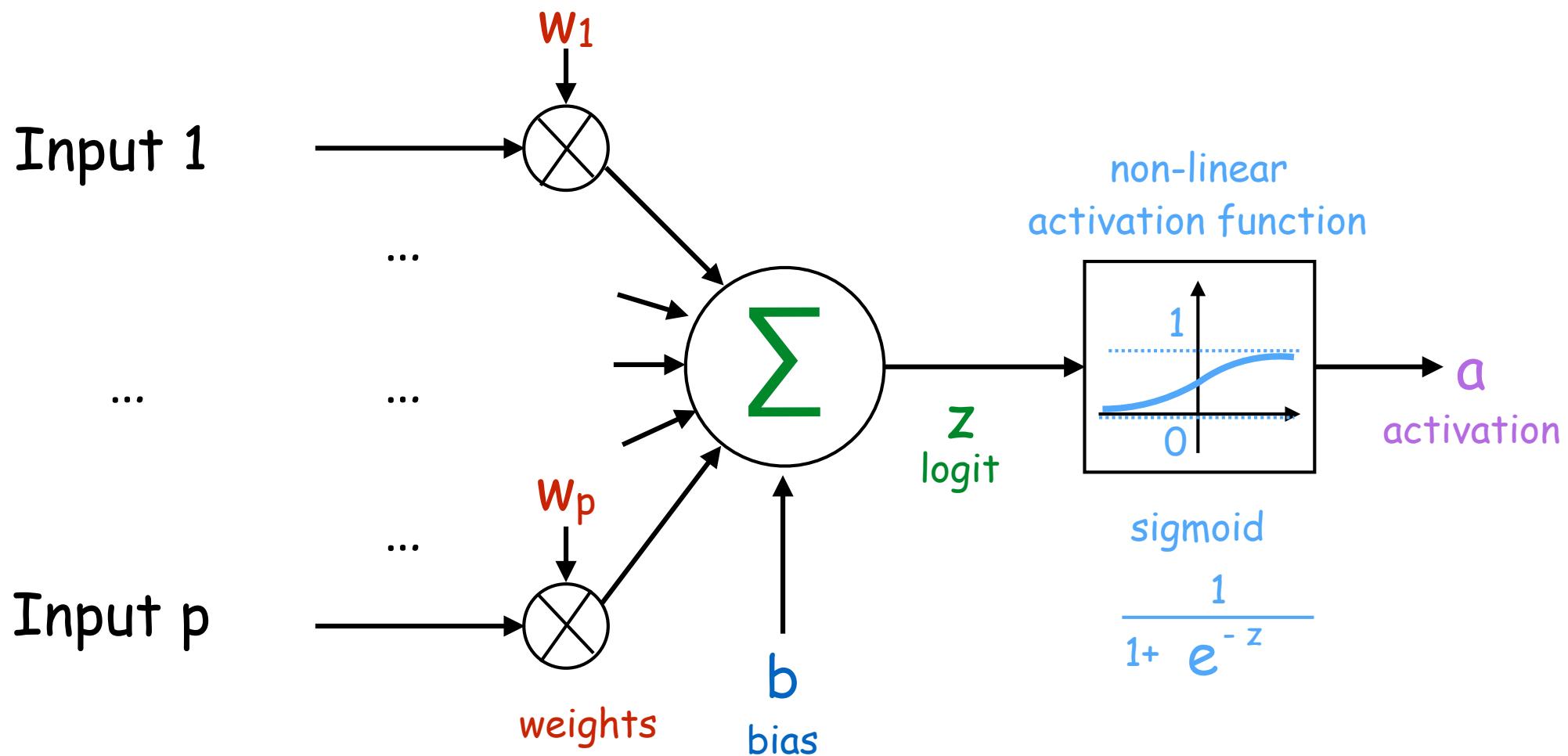
Logistic Regression



$$a = \Pr (\text{label} = + \mid \text{inputs} = X)$$

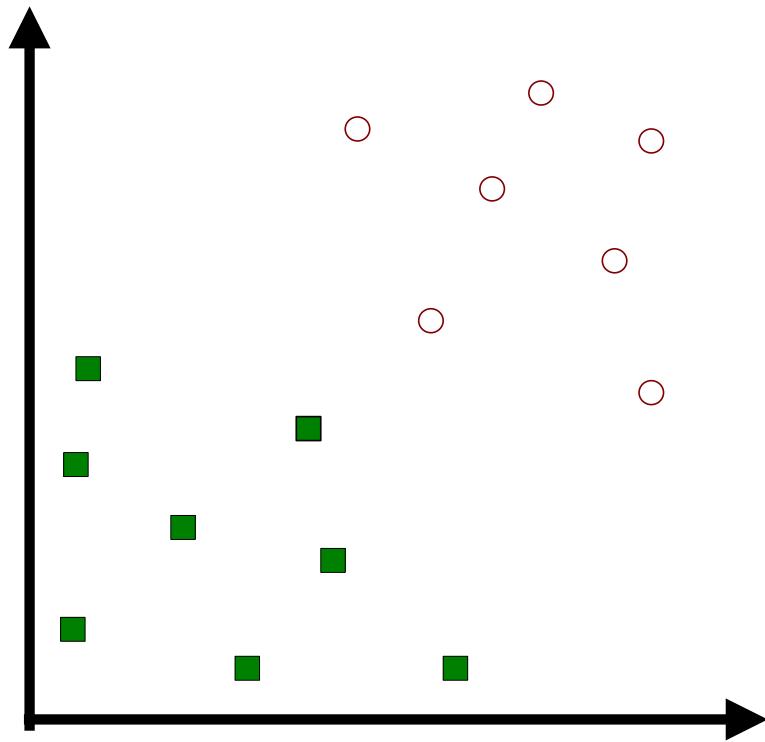
output value represents the probability of the instance having positive label

Parameters w, b



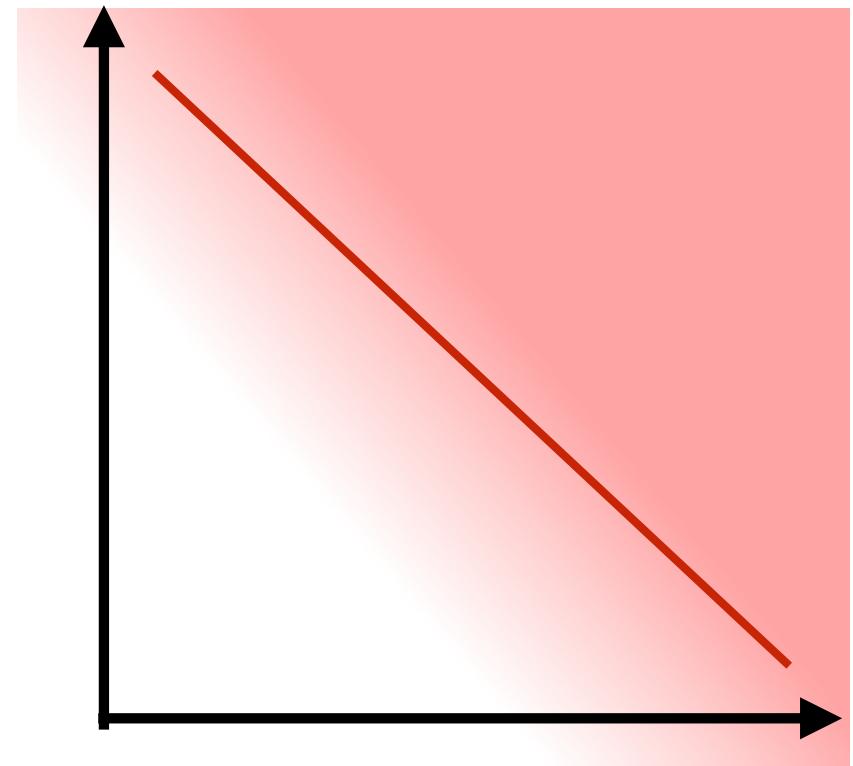
Training Model

training set



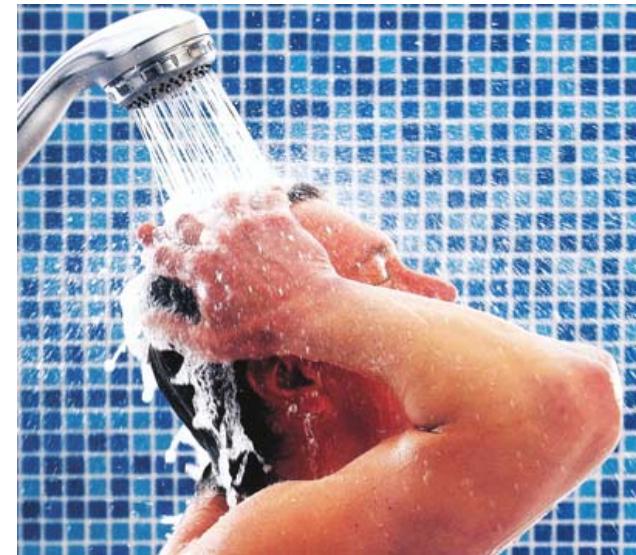
Feature Matrix X
(n by p)
Label Vector y
(length n) (0 or 1 value)

good parameters



Weight vector w
(length p)
bias b
(scalar)

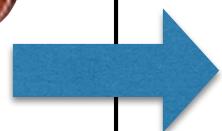
What are good parameters? define "good"



| X_1 | X_2 | good? |
|-------|-------|---------------------------|
| 0 | 0 | bad (no water) |
| 0 | 5 | ok (some water, but cold) |
| 3 | 0 | ok (some water, but hot) |
| 3 | 5 | better (some water, warm) |
| 6 | 10 | great (more water, warm) |
| 0 | 10 | bad (too hot) |

Good Parameter Value?

training set



$A = ?$
probability of head

Parameter value 1: An unbiased coin



50%

50%



$P(H | A)$



$P(T | A)$

$A = 0.5$

parameter value of Coin 1

Parameter value 2: A biased coin



$$A = 0.8$$

parameter value of Coin 2

Probability of observing the result

– Likelihood

observed result



T



H

Coin 1 ($A=0.5$):

$$\text{Likelihood} = 0.5 \times 0.5 = \underline{\underline{0.25}}$$

Coin 2: ($A= 0.8$)

Likelihood =

$$0.2 \times 0.8 = \underline{\underline{0.16}}$$

Probability of observing the result – Likelihood

observed result



Coin 1 ($A=0.5$):

$$\text{Likelihood} = 0.5 \times 0.5 \times 0.5 = 0.125$$

Coin 2: ($A= 0.8$)

$$\text{Likelihood} = 0.2 \times 0.8 \times 0.8 = \underline{\underline{0.128}}$$

Probability of observing the result

– Likelihood

observed result



Coin 1 ($A=0.5$):

Likelihood =

$$0.5 \times 0.5 \times 0.5 \times 0.5 = 0.0625$$

Coin 2: ($A= 0.8$)

Likelihood =

$$0.2 \times 0.8 \times 0.8 \times 0.8 = \underline{0.1024}$$

Maximum Likelihood

Higher likelihood → better parameter

training set



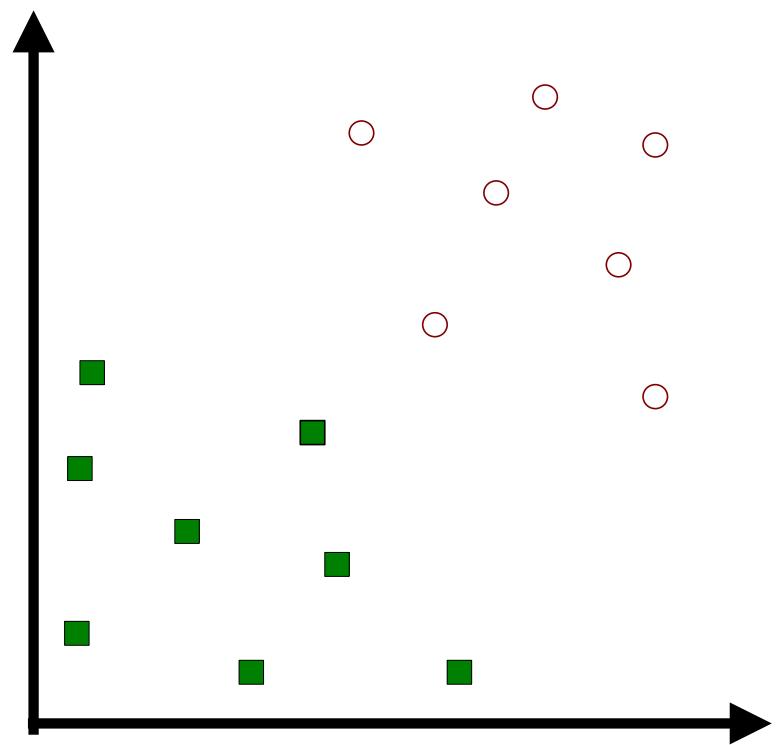
good parameter value



$$A^* = \operatorname{argmax} L(A)$$

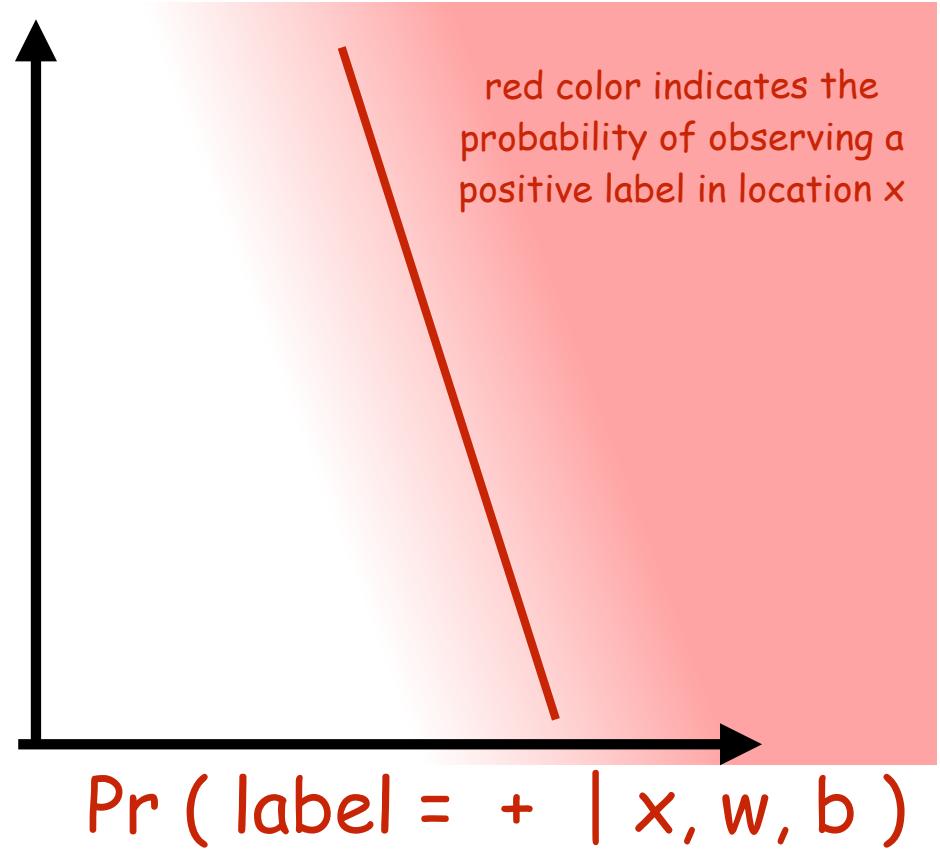
Conditional Likelihood

observed labels



Feature Matrix $X_{(n \text{ by } p)}$
Label Vector y
(length n) (0 or 1 value)

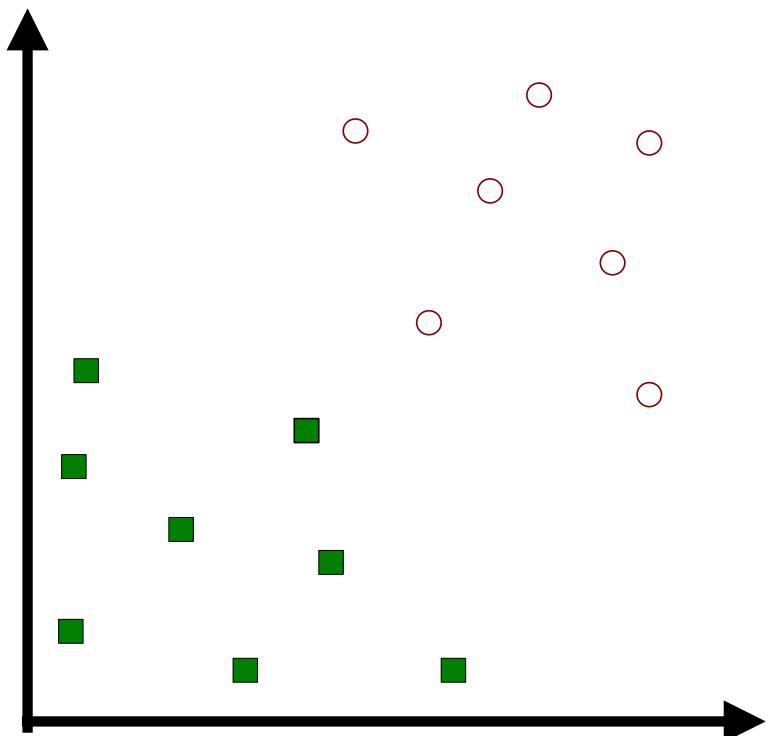
parameter value



Weight vector w
bias b
(length p)
(scalar)

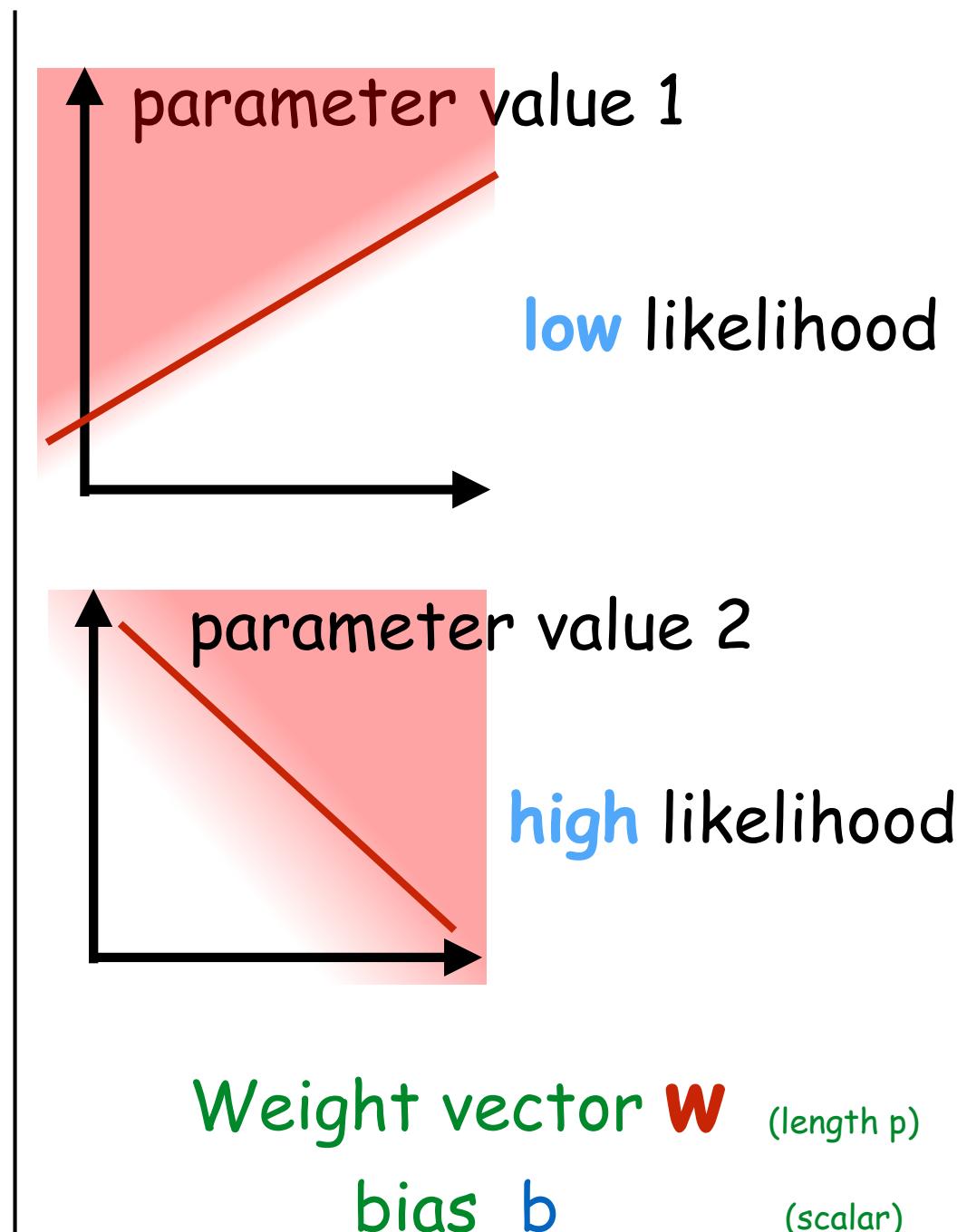
Conditional Likelihood

observed labels



Feature Matrix X
(n by p)

Label Vector y
(length n) (0 or 1 value)



Minimize negative log Likelihood



Coin 1 ($A=0.5$):

- log Likelihood =

$$-\log(0.5 \times 0.5 \times \dots \times 0.5) = -(\log 0.5 + \log 0.5 + \dots + \log 0.5)$$

$$= \underline{\underline{13.8629436}}$$

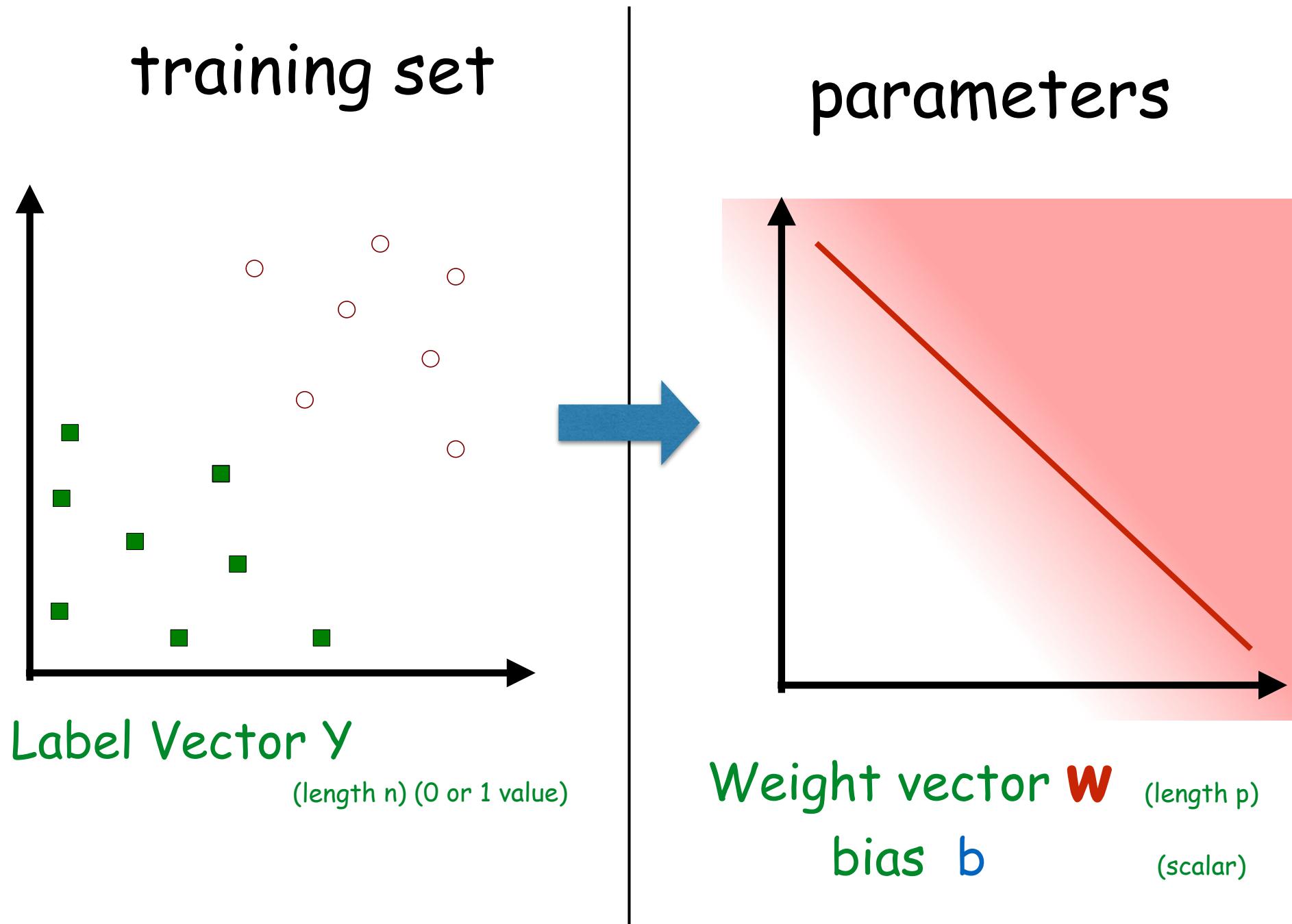
Coin 2: ($A= 0.8$)

- log Likelihood =

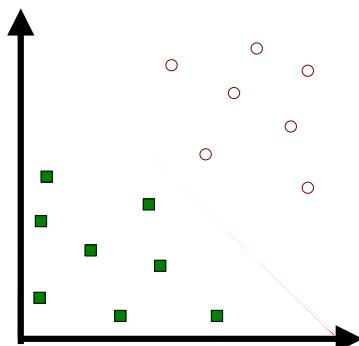
$$-\log(0.2 \times 0.8 \times \dots \times 0.8) = -(\log 0.2 + \log 0.8 + \dots + \log 0.8)$$

$$= \underline{\underline{4.462871}}$$

Minimize negative log Likelihood

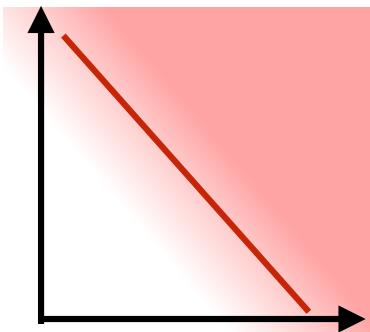


Minimize negative log Likelihood



Labels:
observed data

| y_1 | y_2 | y_3 | ... | y_n |
|-------|-------|-------|-----|-------|
| 1 | 0 | 1 | ... | 0 |



Weight vector w
bias b

| Outputs: probabilities | a_1 | a_2 | a_3 | a_n |
|---------------------------|-------|-------|-------|-------|
| | .8 | .3 | .6 | .2 |

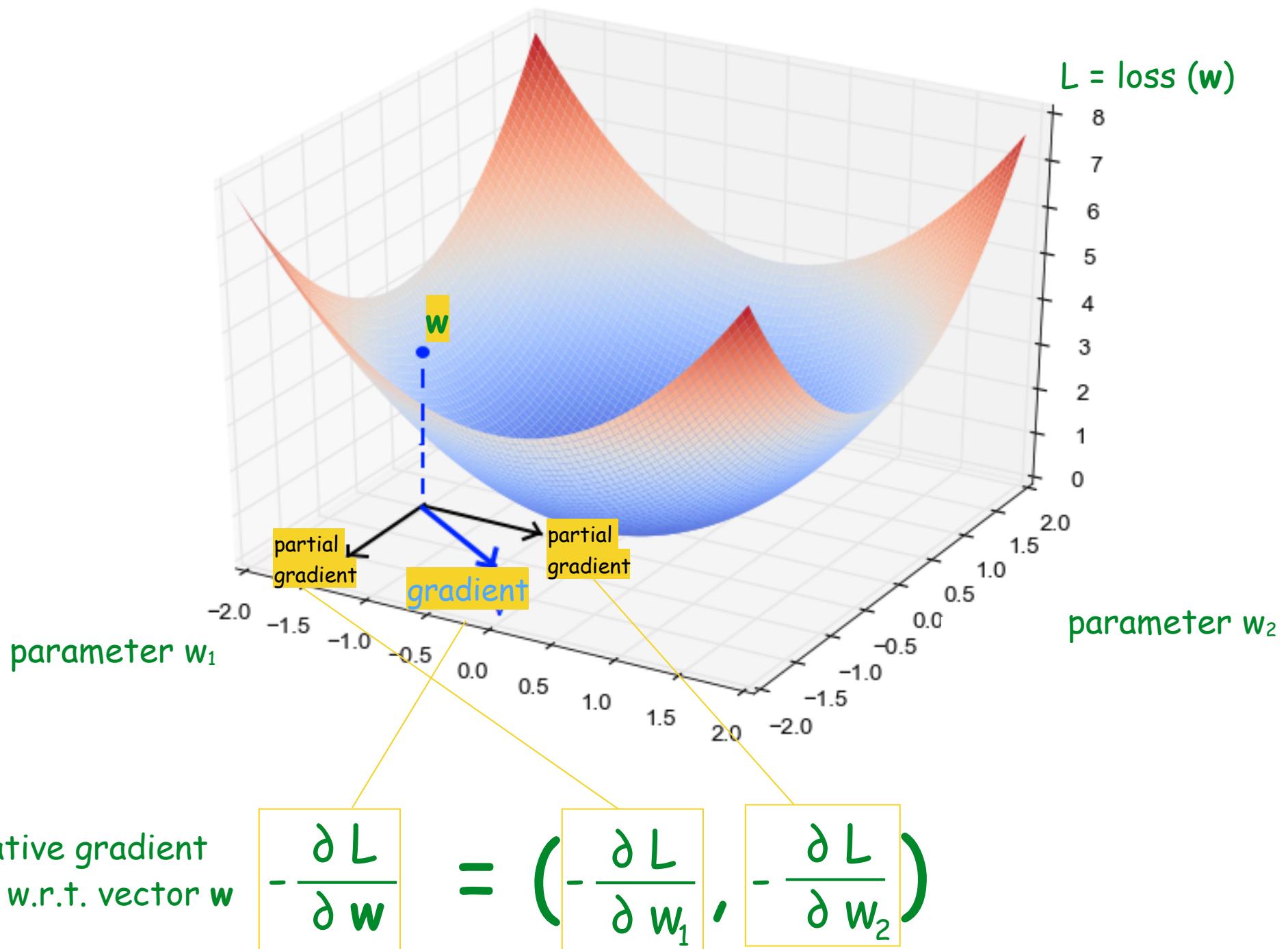
- log Likelihood =

$$\sum_{i=1}^n -y_i \log a_i - (1-y_i) \log(1-a_i)$$

which is also called cross entropy loss

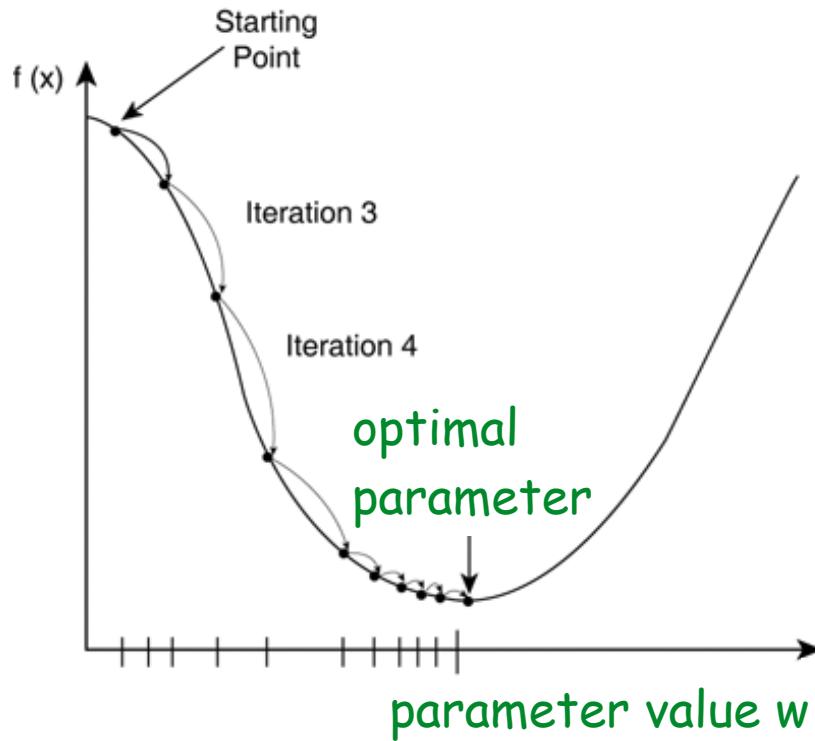
multiplied with a constant factor n

Minimize Cross Entropy Loss



Gradient Descent (update model)

$L = \text{loss}(w)$



$$w \leftarrow w - a \frac{\partial L}{\partial w}$$

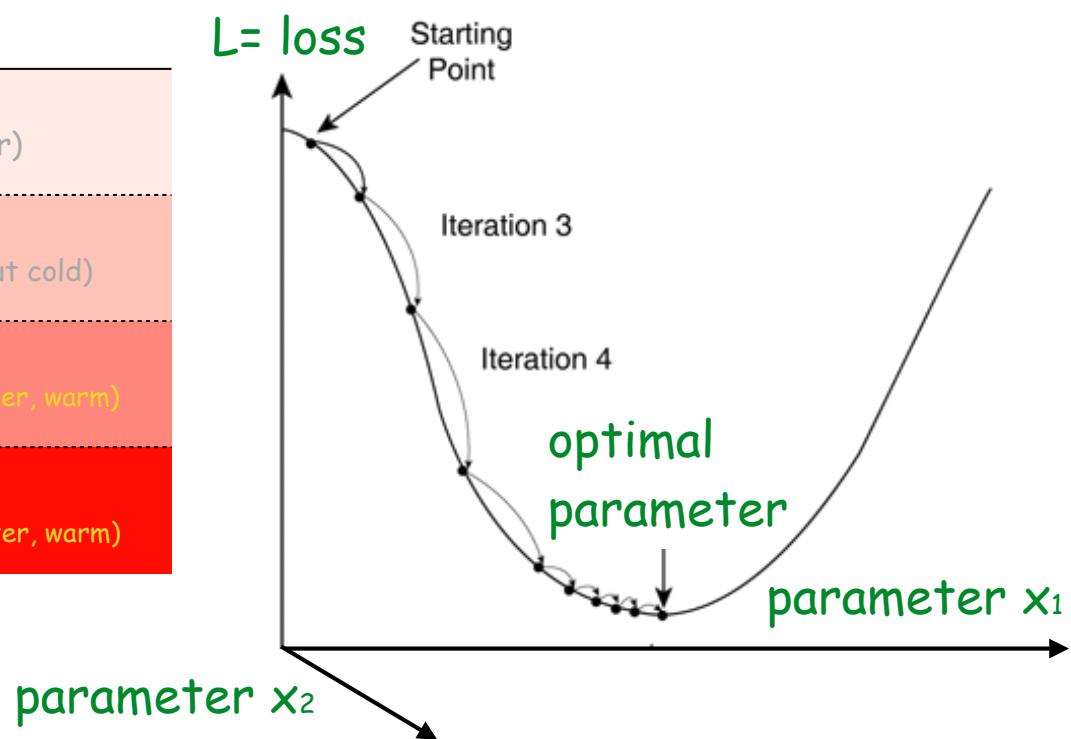
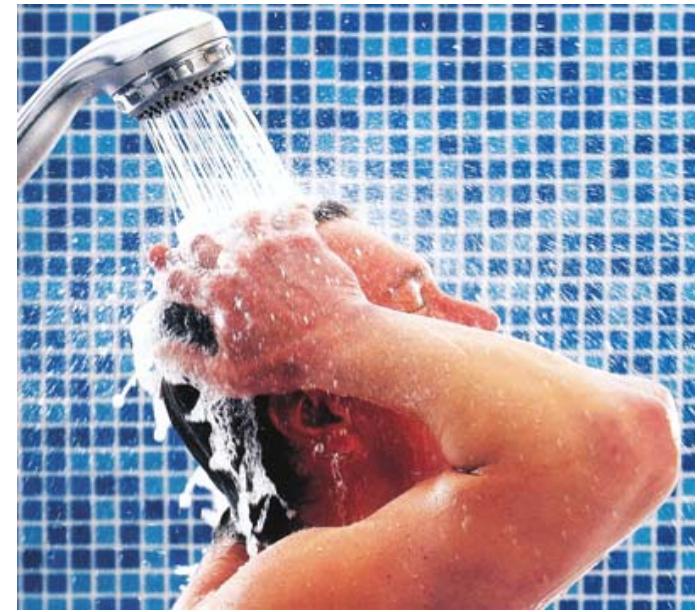
$$b \leftarrow b - a \frac{\partial L}{\partial b}$$

a step size (a constant scalar)

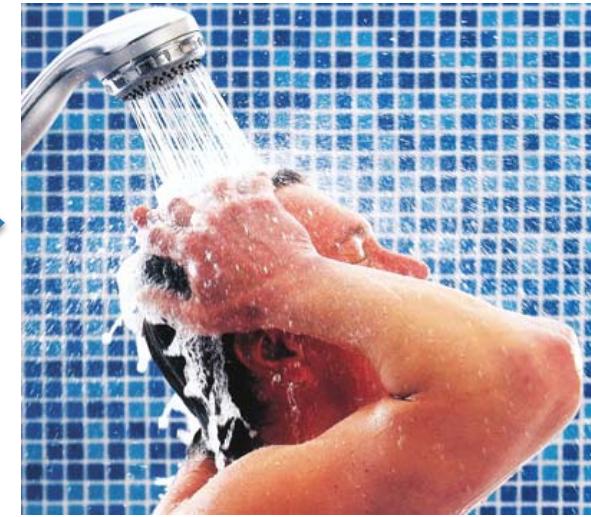
Example



| | x_1 | x_2 | L |
|---|-------|-------|------------------------------|
| 1 | 0 | 0 | 100 bad (no water) |
| 2 | 0 | 5 | 80 ok (some water, but cold) |
| 3 | 3 | 5 | 60 better (some water, warm) |
| 4 | 6 | 10 | 20 great (more water, warm) |



Gradient



x_1

$$\frac{\partial L}{\partial x_1}$$

0

-3

0

-6

3

-3

6

0

8

3

9

6

10

x_2

$$\frac{\partial L}{\partial x_2}$$

0

-10

5

-6

6

-3

10

0

6

-3

3

-6

0

-10

Loss

100 bad
(no water)

80 ok
(some water, but cold)

60 better
(some water, warm)

20 great
(more water, warm)

60 good
(more water, hot)

80 ok
(more water, too hot)

200 bad
(more water, way too hot)

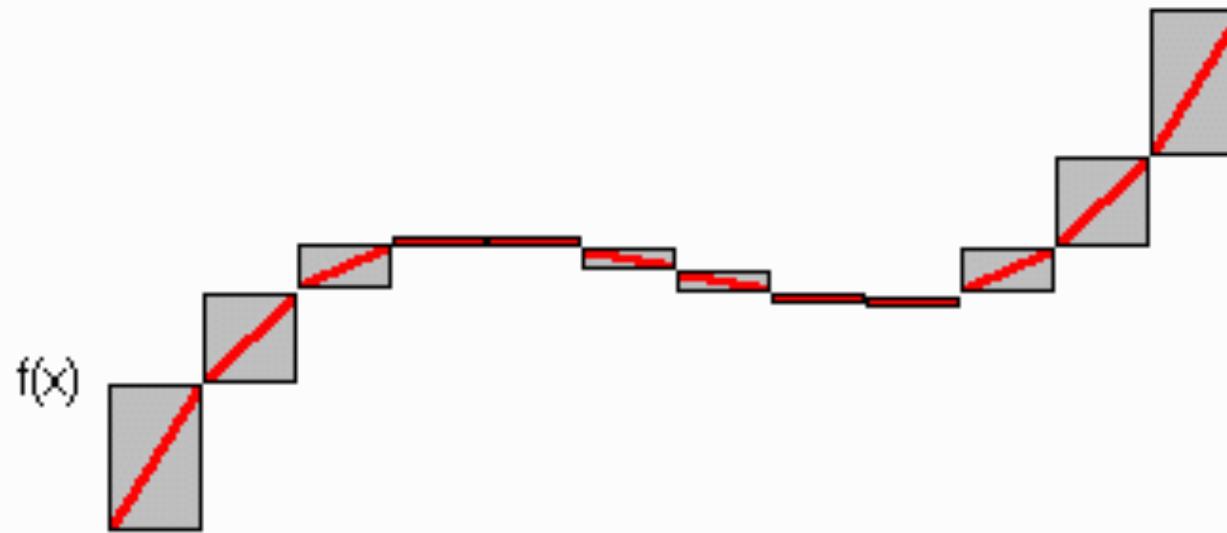
Gradient of L w.r.t. a **scalar**?

$$\frac{\partial L}{\partial b}$$

Gradient of L w.r.t. a **vector**?

$$\frac{\partial L}{\partial w}$$

Gradient of a function f with respect to (w.r.t.) a scalar x



differentiate



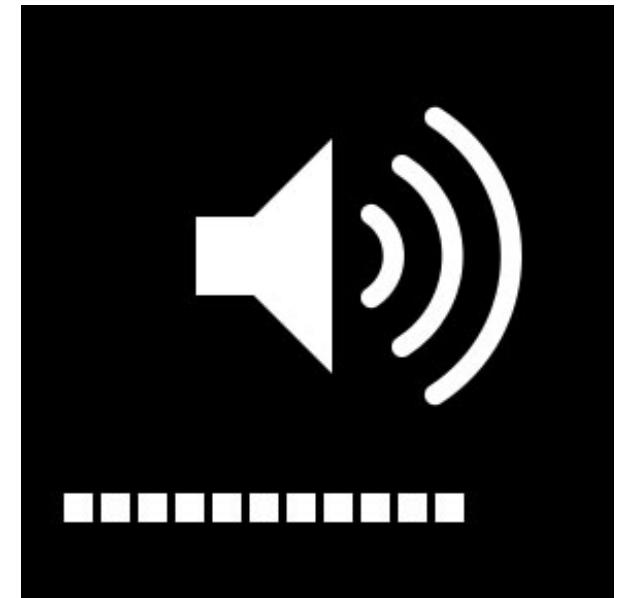
$$\frac{\partial f}{\partial x}$$

df/dx

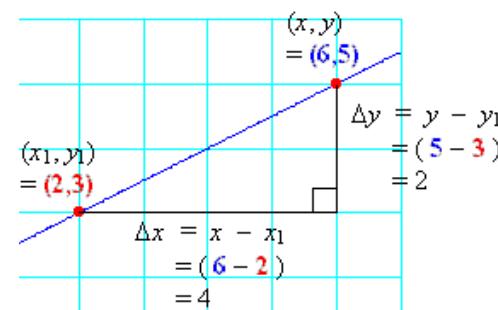
x axis

Gradient of f function w.r.t. input x

Gradient of L w.r.t. a scalar x



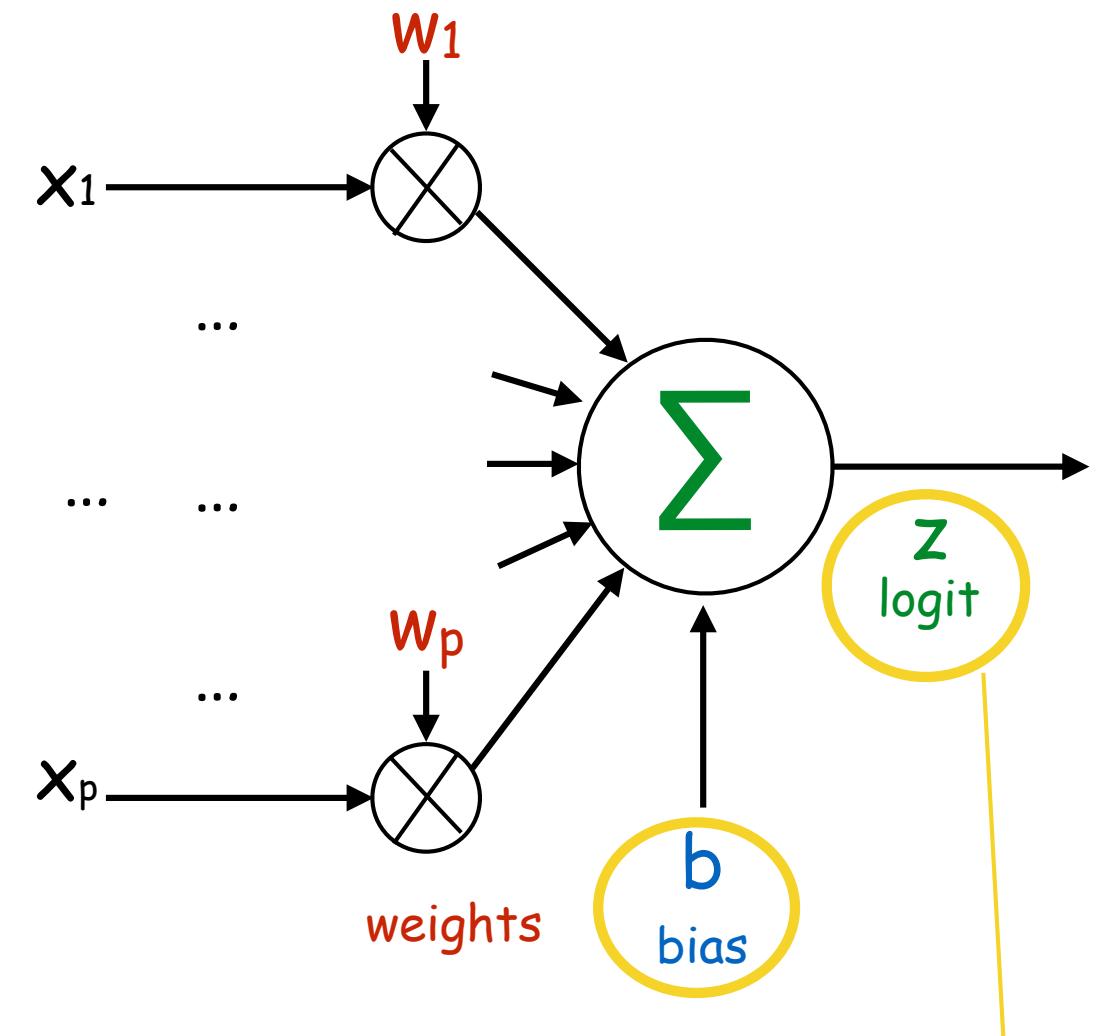
\times
control



L
response

$\frac{\partial L}{\partial x}$ a constance number

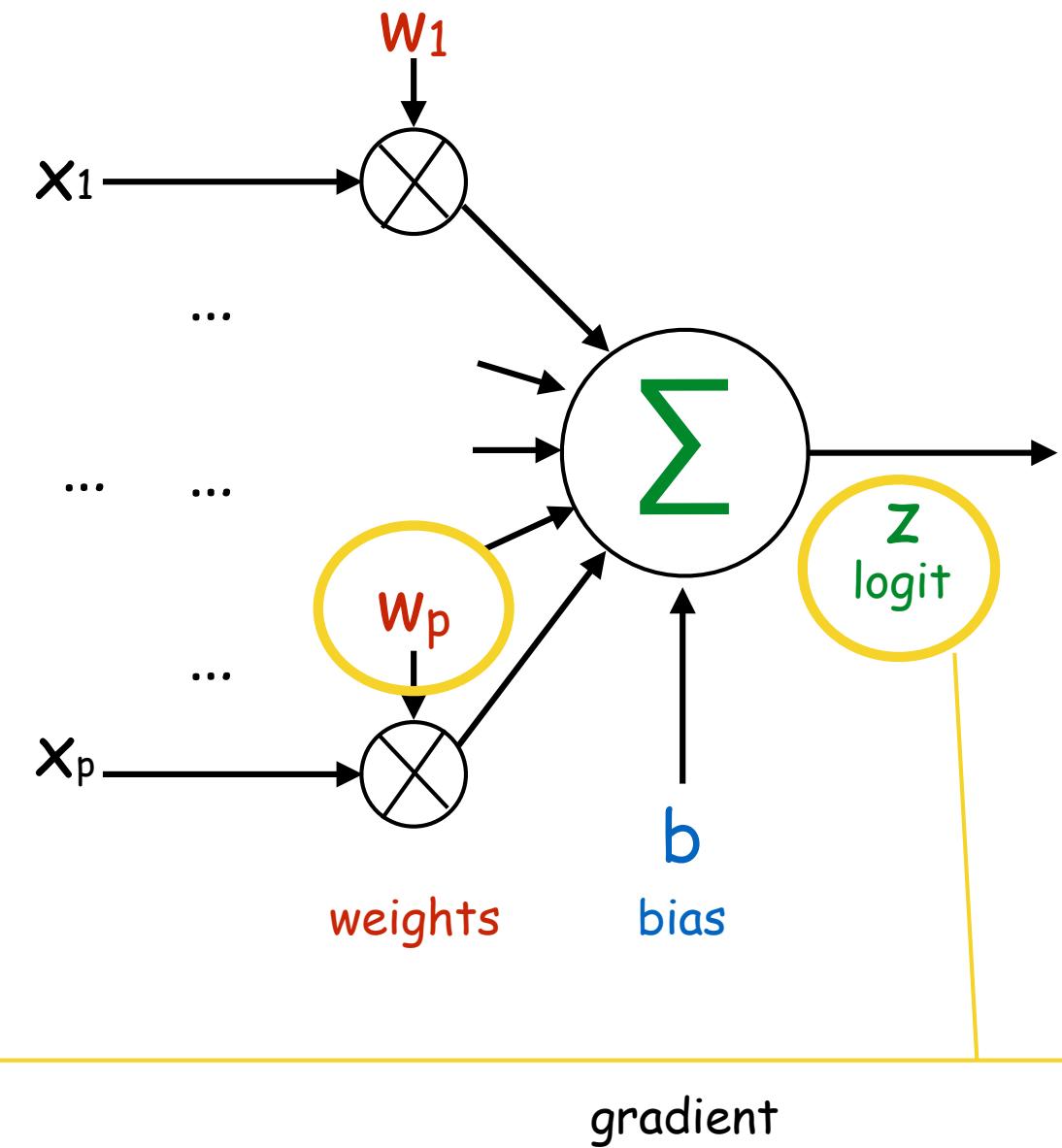
Example



gradient

$$\frac{\partial z}{\partial b} = 1$$

Example



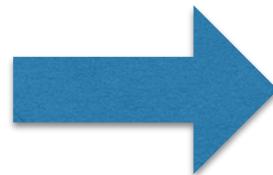
$$\frac{\partial z}{\partial w_p} = x_p$$

is a function of x_p

Gradient Changing over x



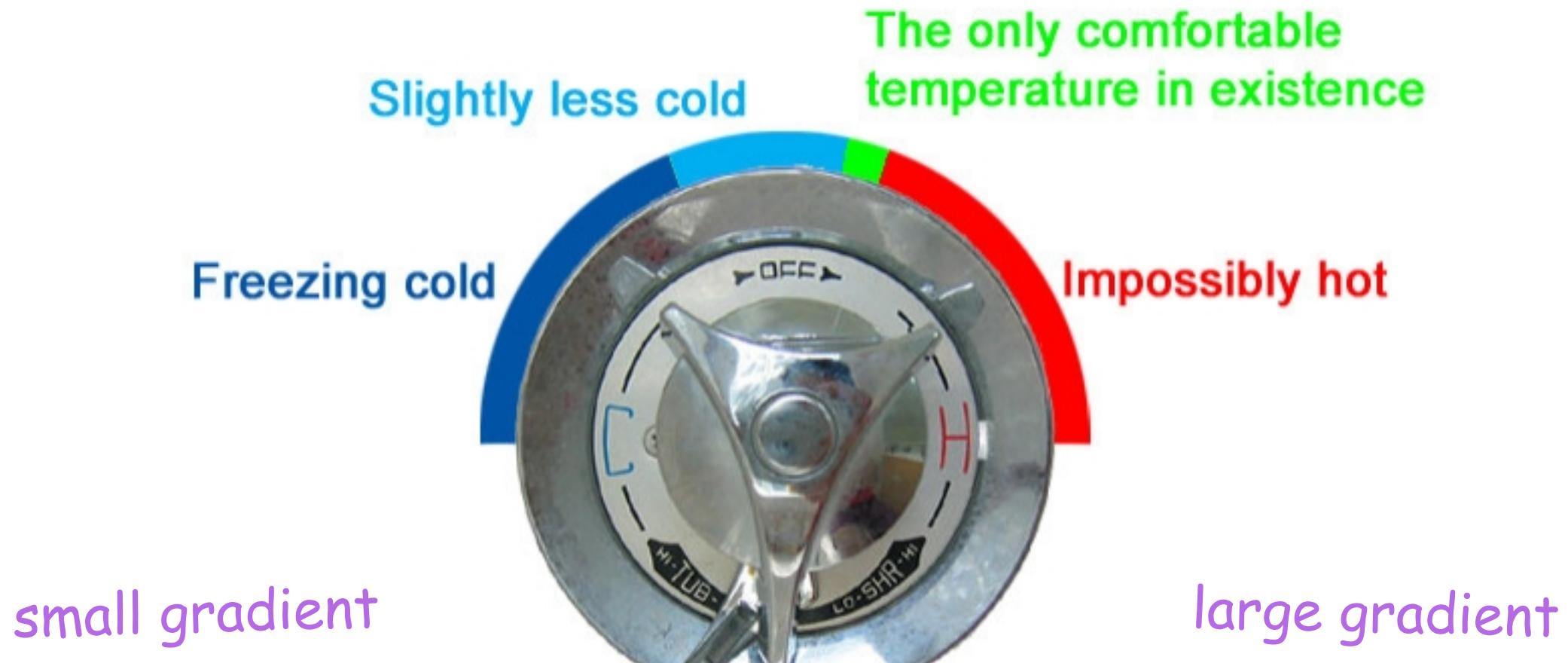
✗



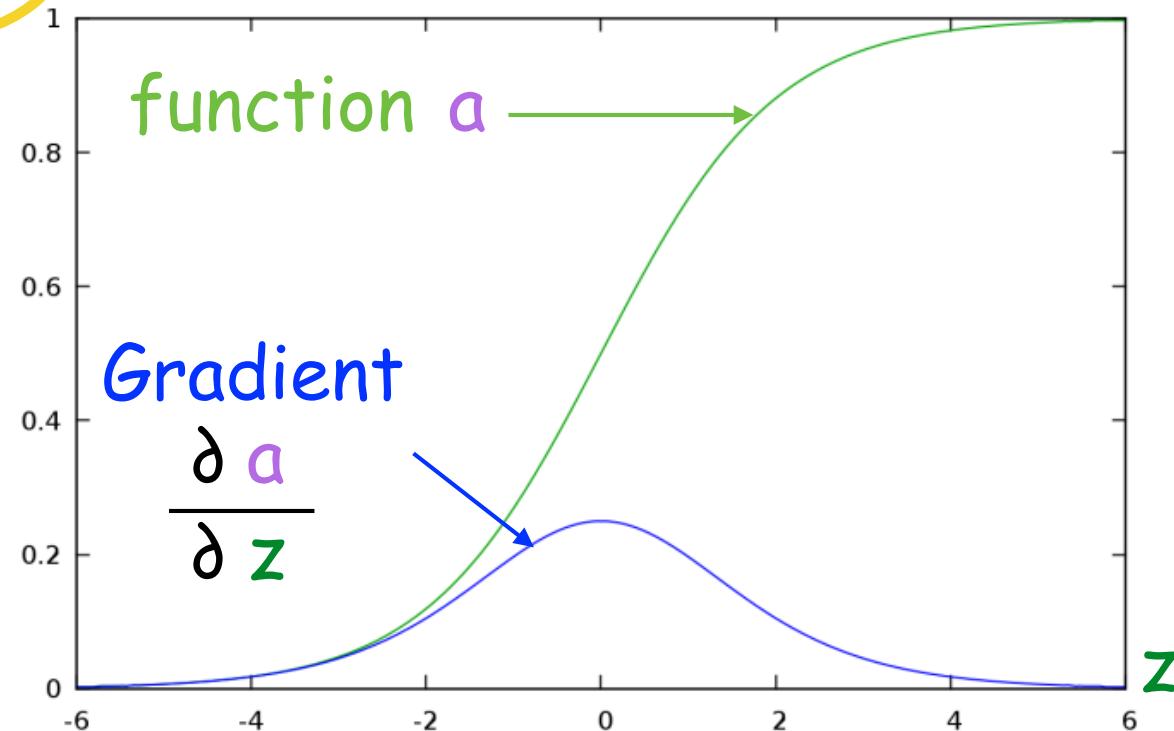
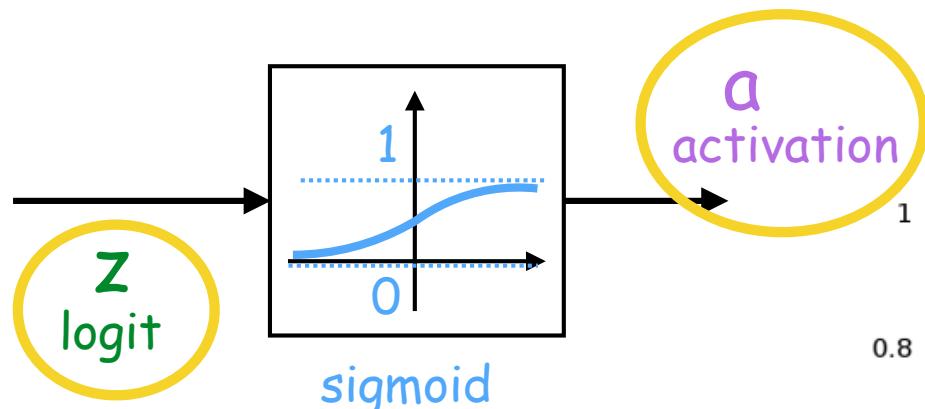
L

water temperature

L is a function of x



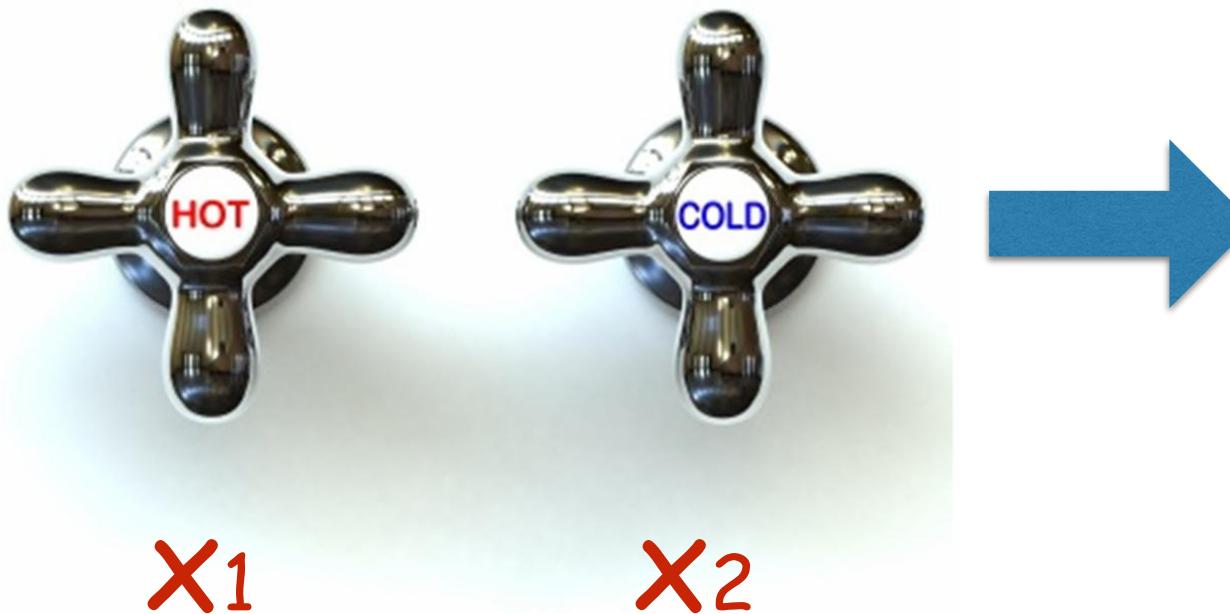
Gradient $\frac{\partial L}{\partial x}$ is also a function of x



$$\begin{aligned}\frac{\partial a}{\partial z} &= a(1-a) \\ &= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}}\right)\end{aligned}$$

is a function of z

Partial Gradients w.r.t. a vector $\frac{\partial L}{\partial x}$



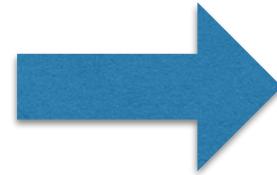
$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2} \right)$$

L
temperature
 $L = f(x_1, x_2)$

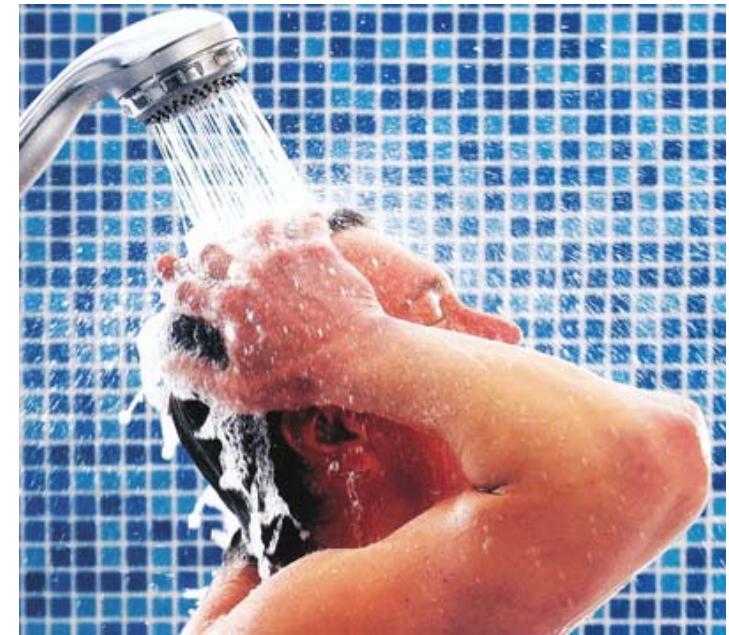
Partial Gradients w.r.t. a vector $\frac{\partial L}{\partial x}$

x_1

x_2



L



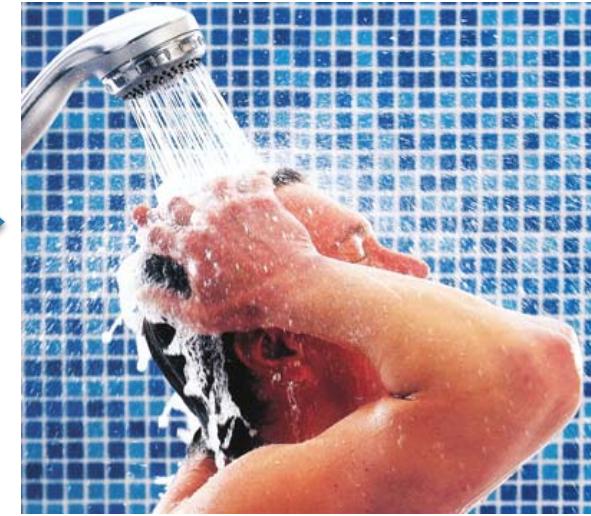
partial gradient – the gradient of f over x_1 fixing x_2

$$\left(\frac{\partial L}{\partial x_1} \right)$$

$$\left. , \quad \frac{\partial L}{\partial x_2} \right)$$

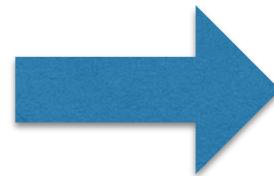
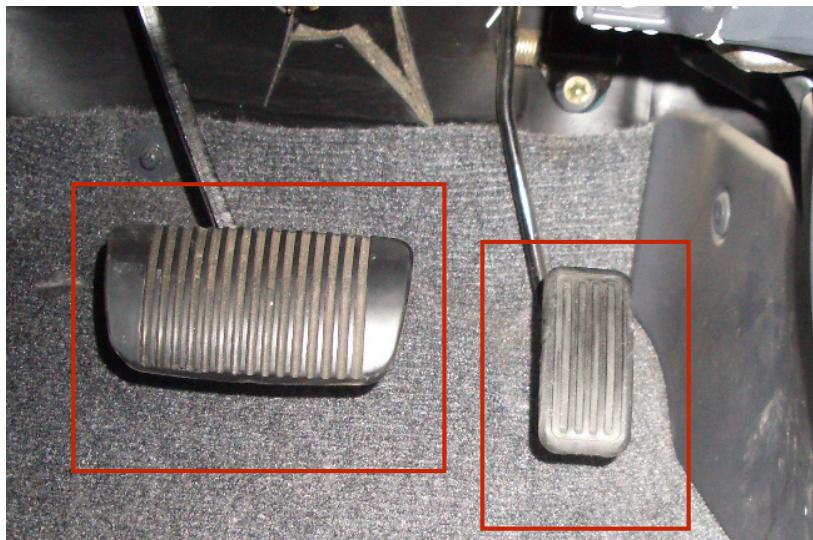
$$L = f(x_1, x_2)$$

Partial Gradients



| x_1 | $\frac{\partial L}{\partial x_1}$ | x_2 | $\frac{\partial L}{\partial x_2}$ | Loss |
|-------|-----------------------------------|-------|-----------------------------------|-----------------------------------|
| 0 | -3 | 0 | -10 | 100 bad (no water) |
| 0 | -6 | 5 | -6 | 80 ok (some water, but cold) |
| 3 | -3 | 6 | -3 | 60 better (some water, warm) |
| 6 | 0 | 10 | 0 | 20 great (more water, warm) |
| 8 | 3 | 6 | -3 | 60 good (more water, hot) |
| 9 | 6 | 3 | -6 | 80 ok (more water, too hot) |
| 10 | 10 | 0 | -10 | 200 bad (more water, way too hot) |

Partial Gradients w.r.t. a vector $\frac{\partial L}{\partial x}$



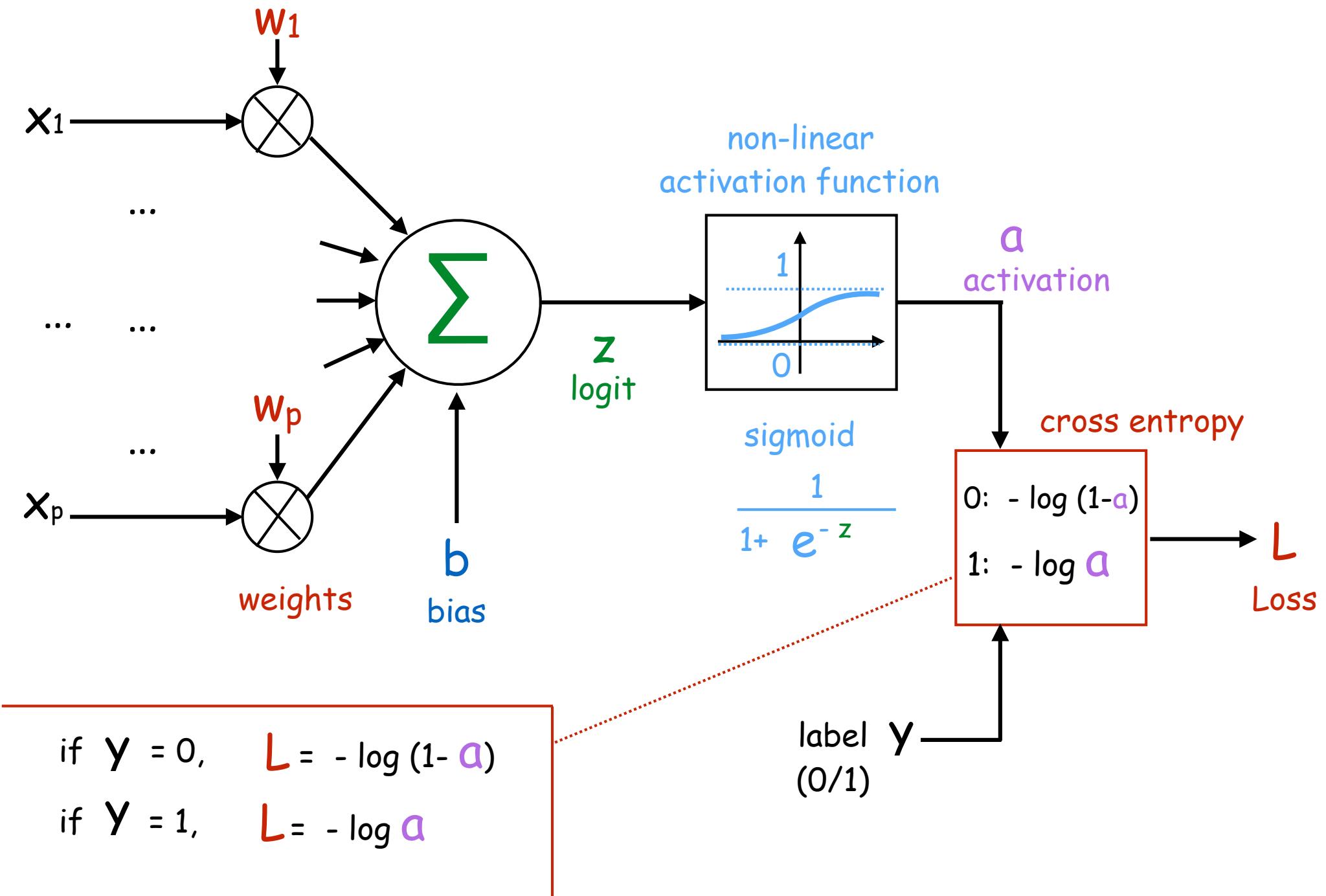
x_1
brake

x_2
gas

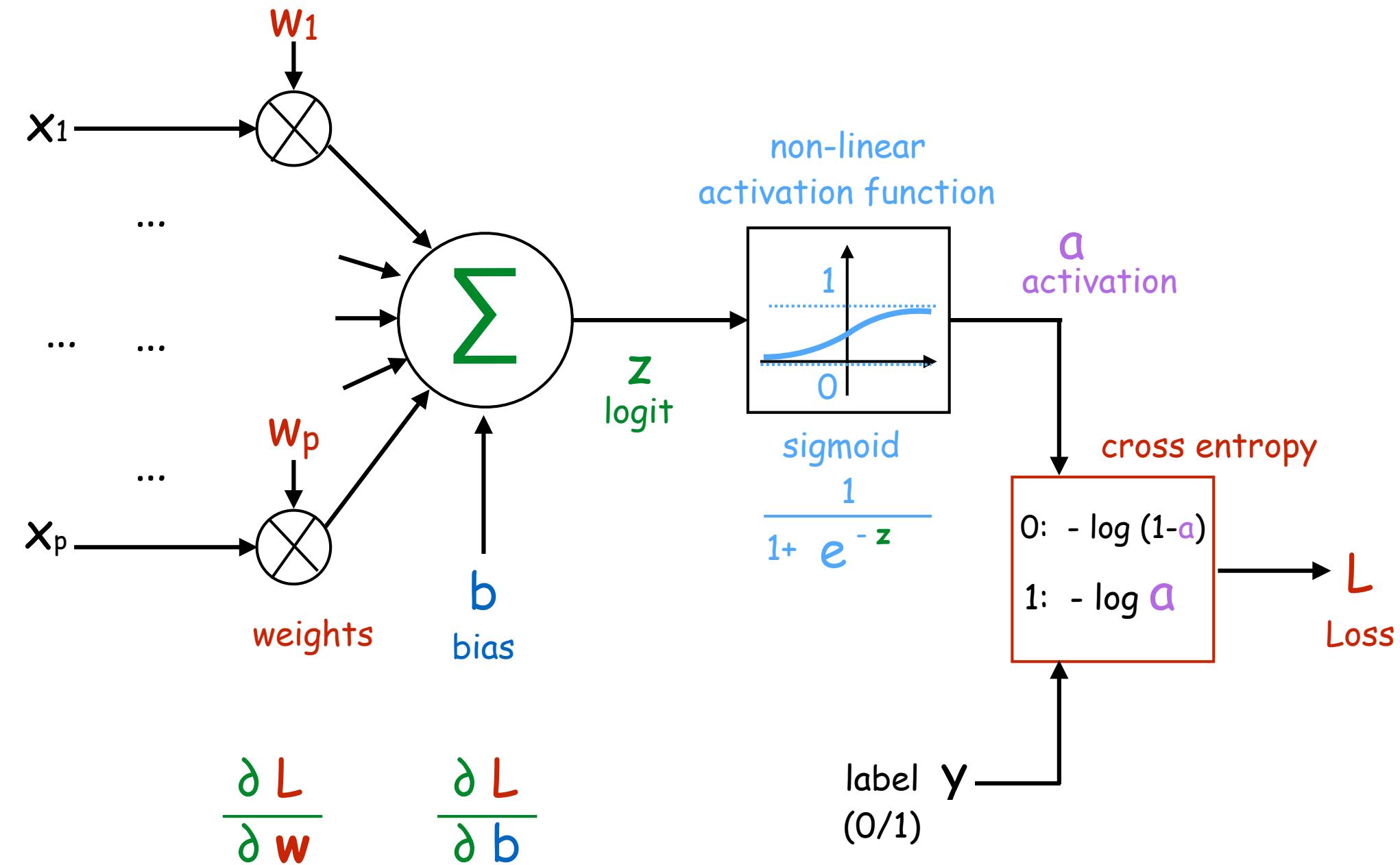
L
speed

$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2} \right)$$

Loss on one data instance



Global Gradient of w, b



Chain Rule

$$y = f(g(x)) \quad y = f(u) \quad u = g(x)$$

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \times \frac{\partial u}{\partial x}$$

global gradient

local gradient

local gradient

$$y = f(g(h(x))) \quad y = f(u) \quad u = g(v) \quad v = h(x)$$

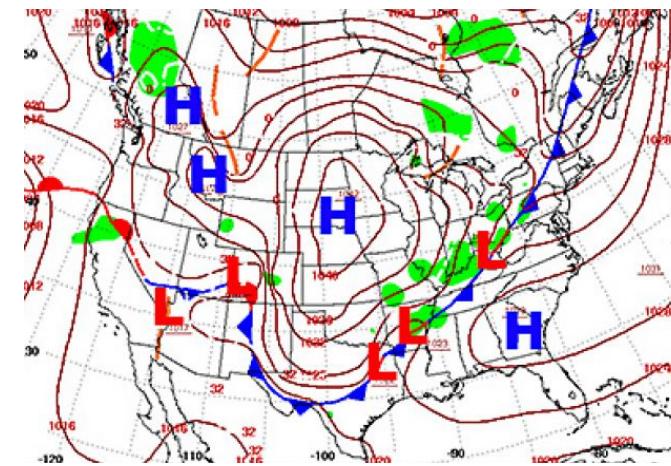
$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \times \frac{\partial u}{\partial v} \times \frac{\partial v}{\partial x}$$

global gradient

local gradient local gradient local gradient

Chain Rule

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



X

low air pressure

Y

rain

L

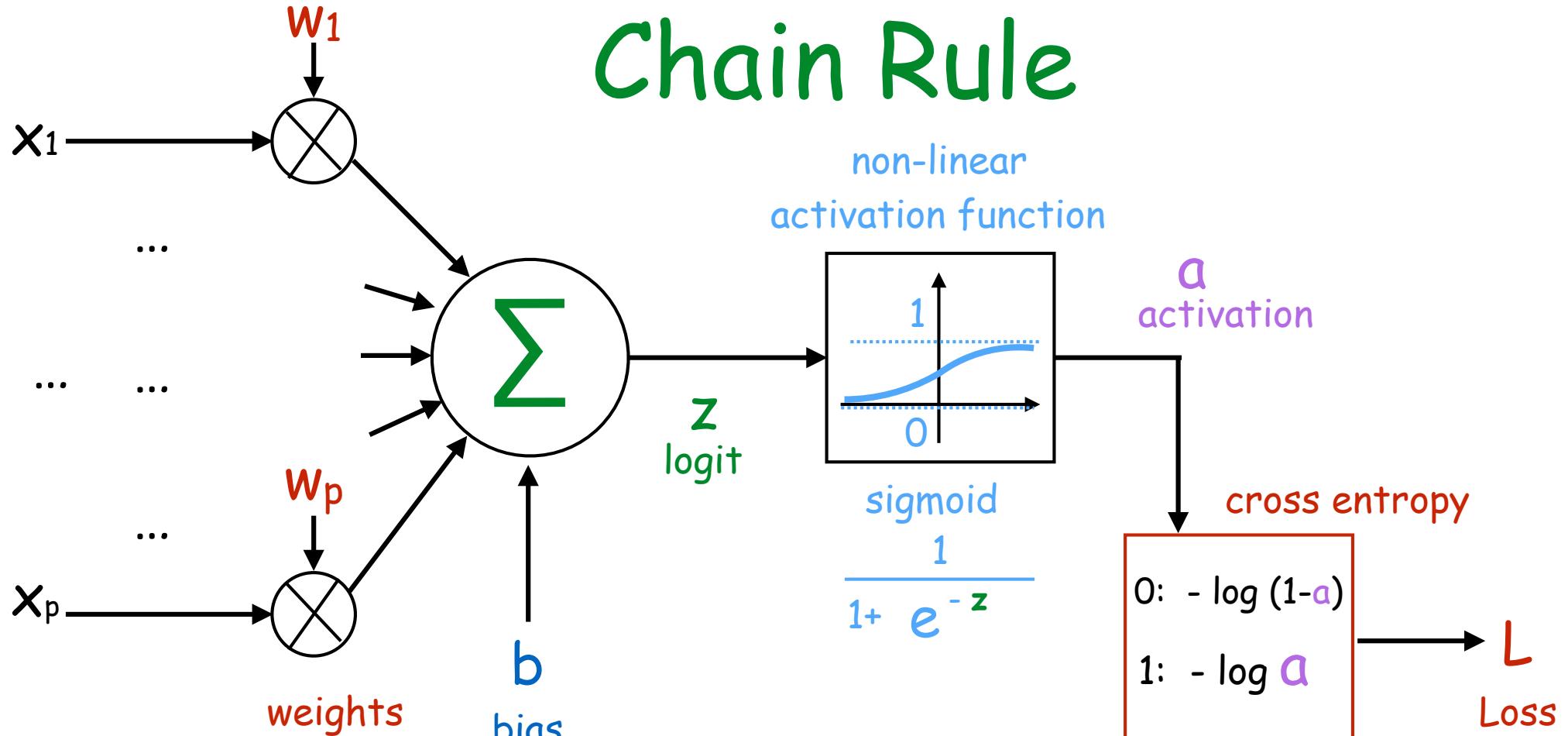
slow traffic

$\frac{\partial L}{\partial y}$: how much the traffic speed L will be changed if we change the rain Y

$\frac{\partial y}{\partial x}$: how much the rain Y will be changed if we change the air pressure X

$\frac{\partial L}{\partial x}$: how much the traffic speed L will be changed if we change the air pressure X

Chain Rule

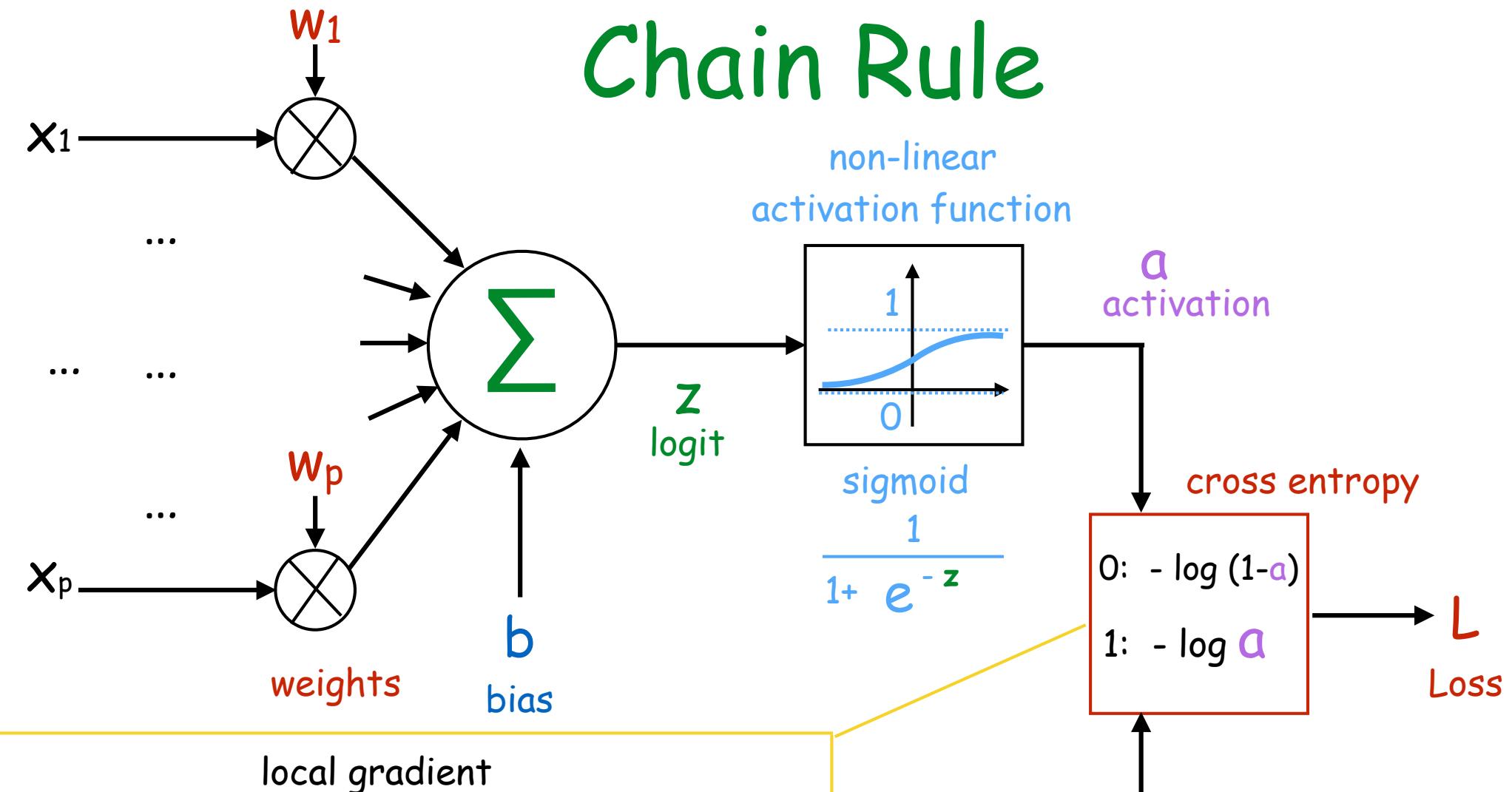


$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w_i}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial b}$$

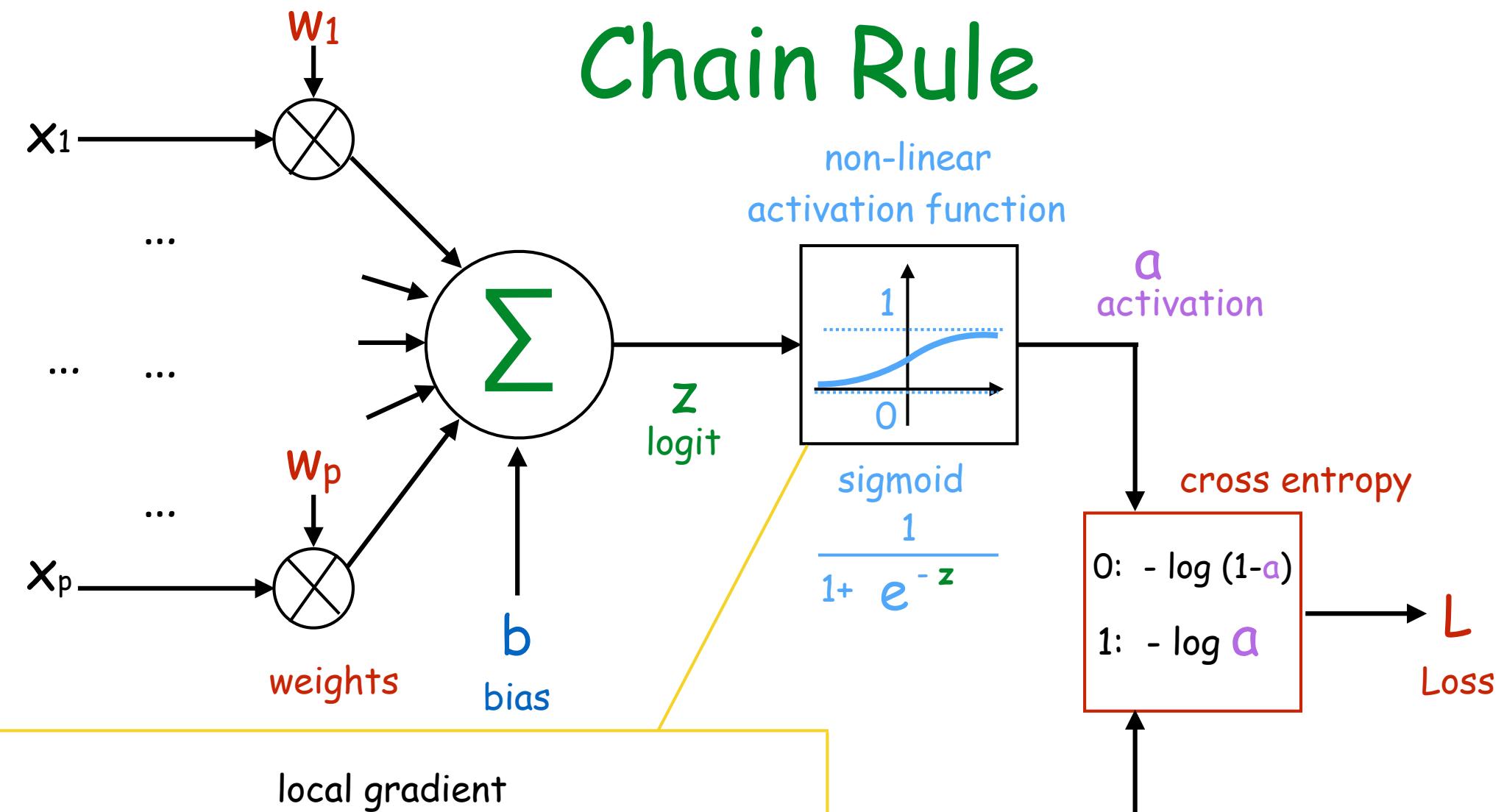
products of local gradients

Chain Rule



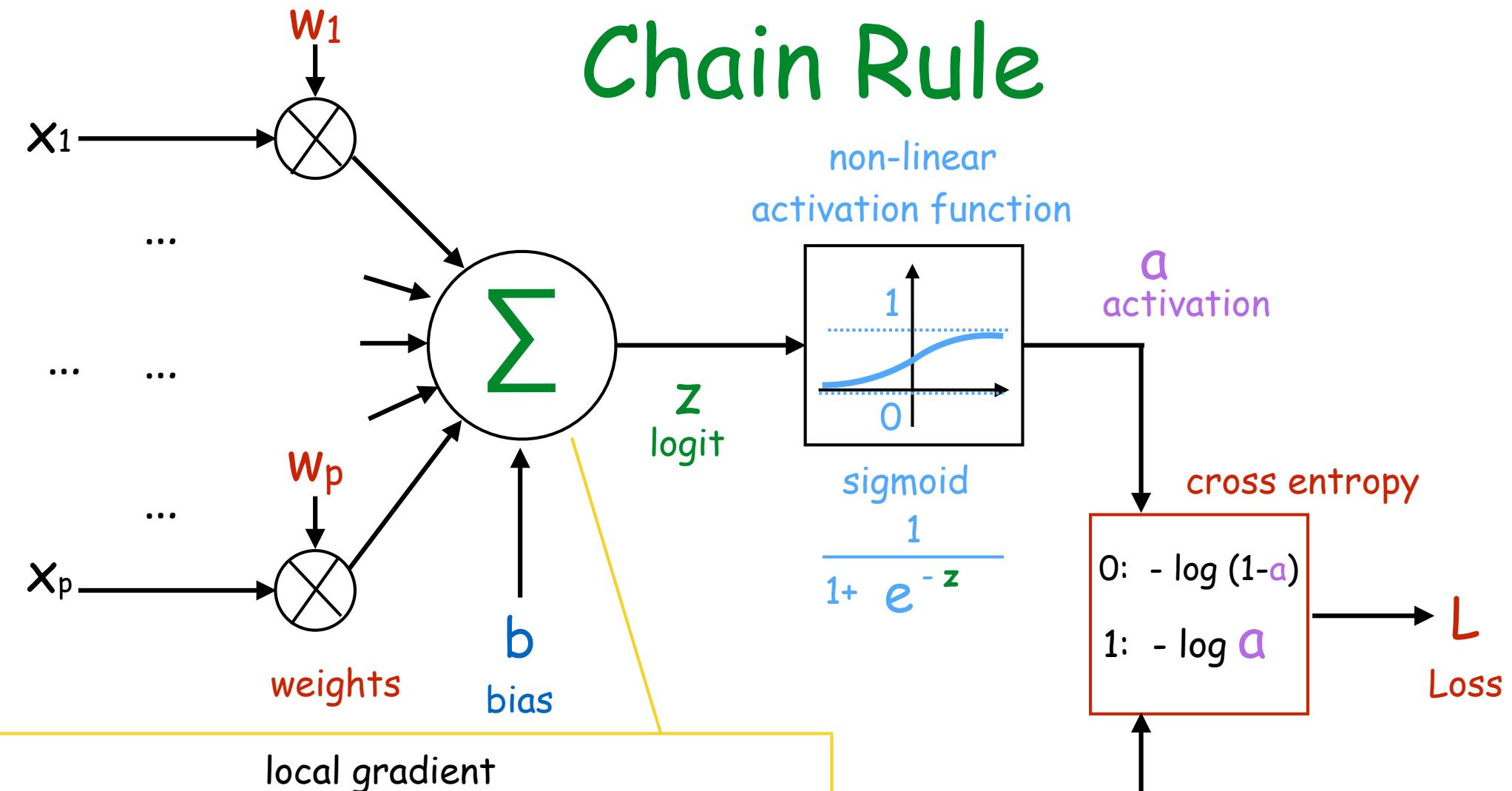
$$\frac{\partial L}{\partial a} = \begin{cases} 1/(1-a) & \text{if } y = 0 \\ -1/a & \text{if } y = 1 \end{cases}$$

Chain Rule



$$\frac{\partial a}{\partial z} = a(1-a)$$

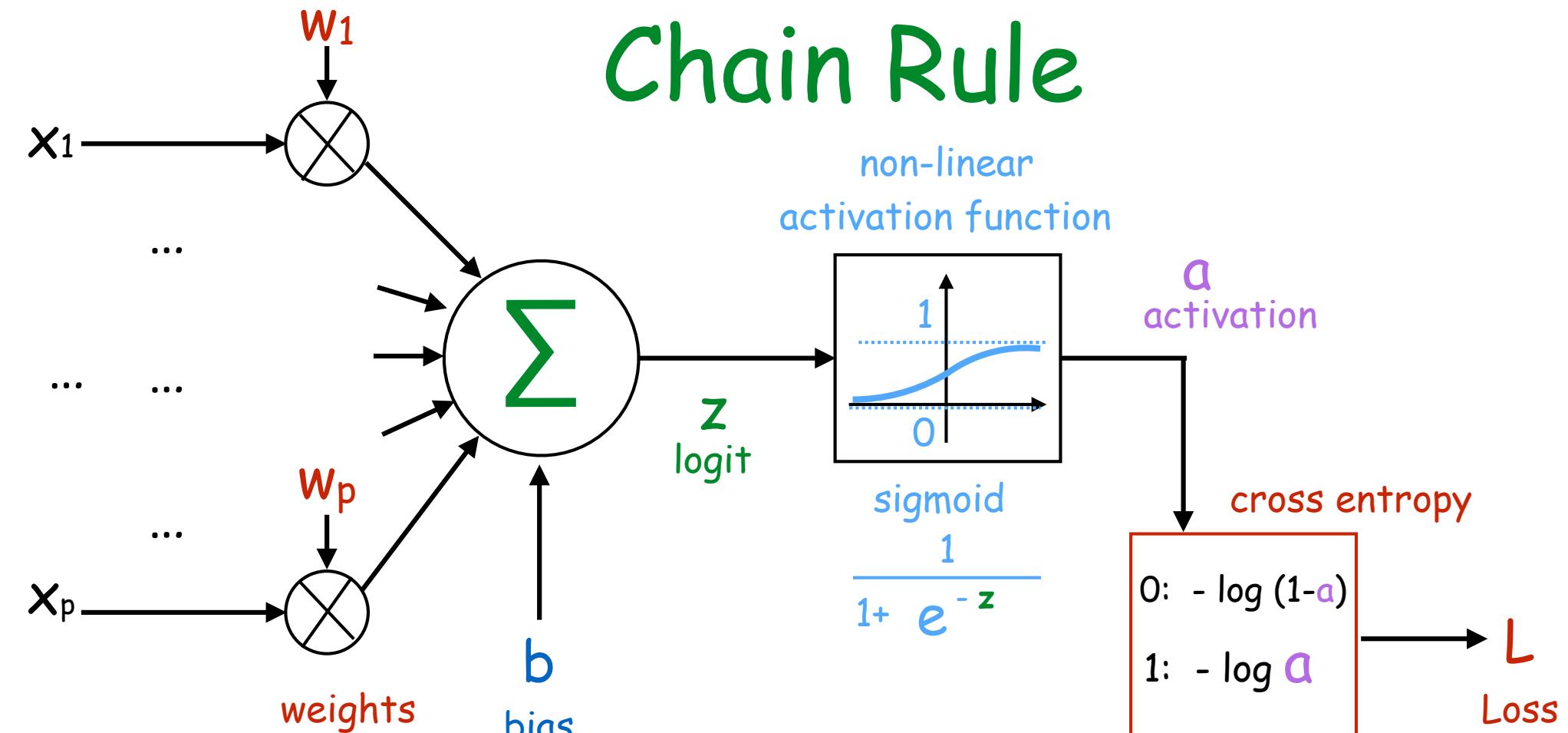
Chain Rule



$$\frac{\partial z}{\partial w_i} = x_i$$

$$\frac{\partial z}{\partial b} = 1$$

Chain Rule

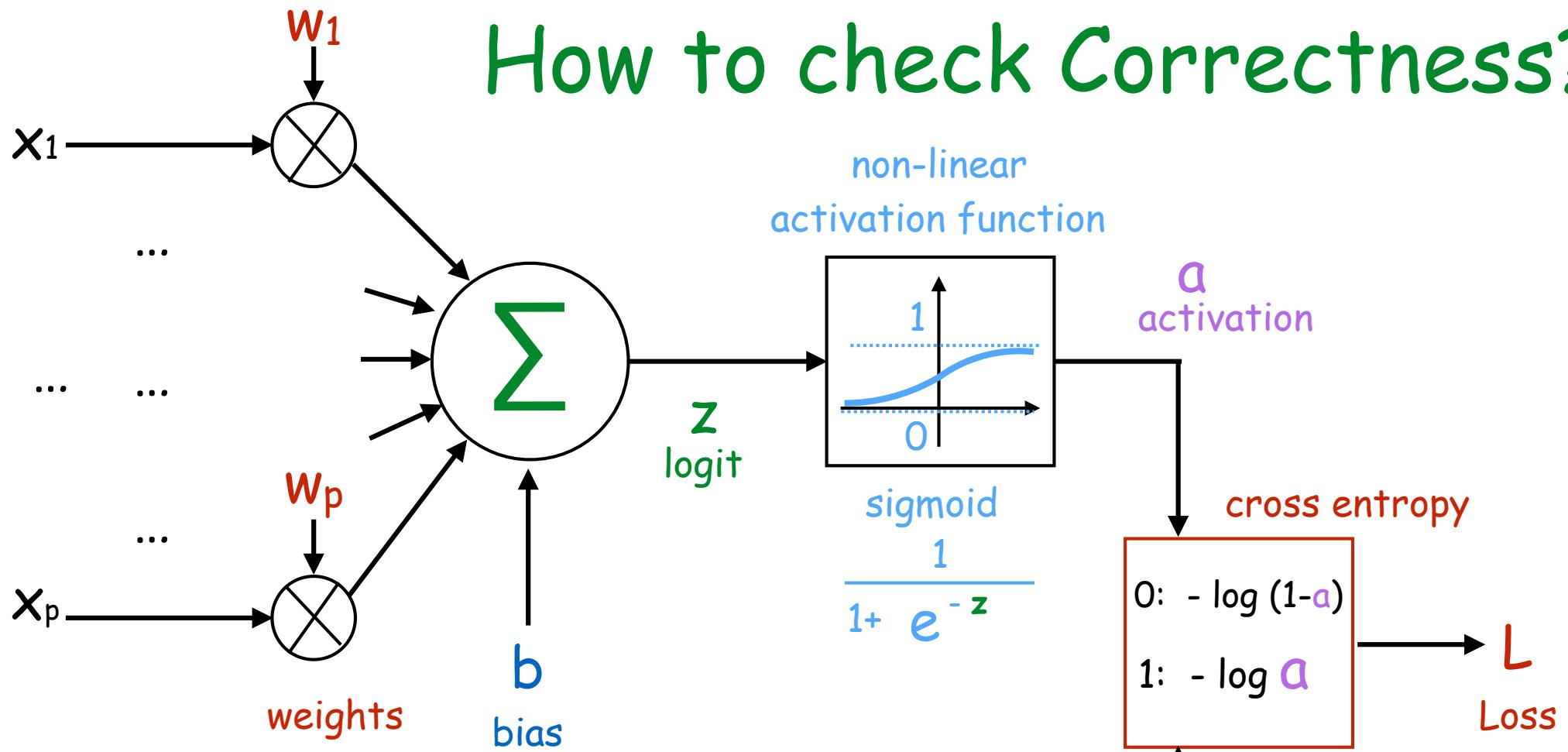


$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w_i}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial b}$$

products of local gradients

How to check Correctness?



$$\frac{\partial L}{\partial w_i}$$

very complex

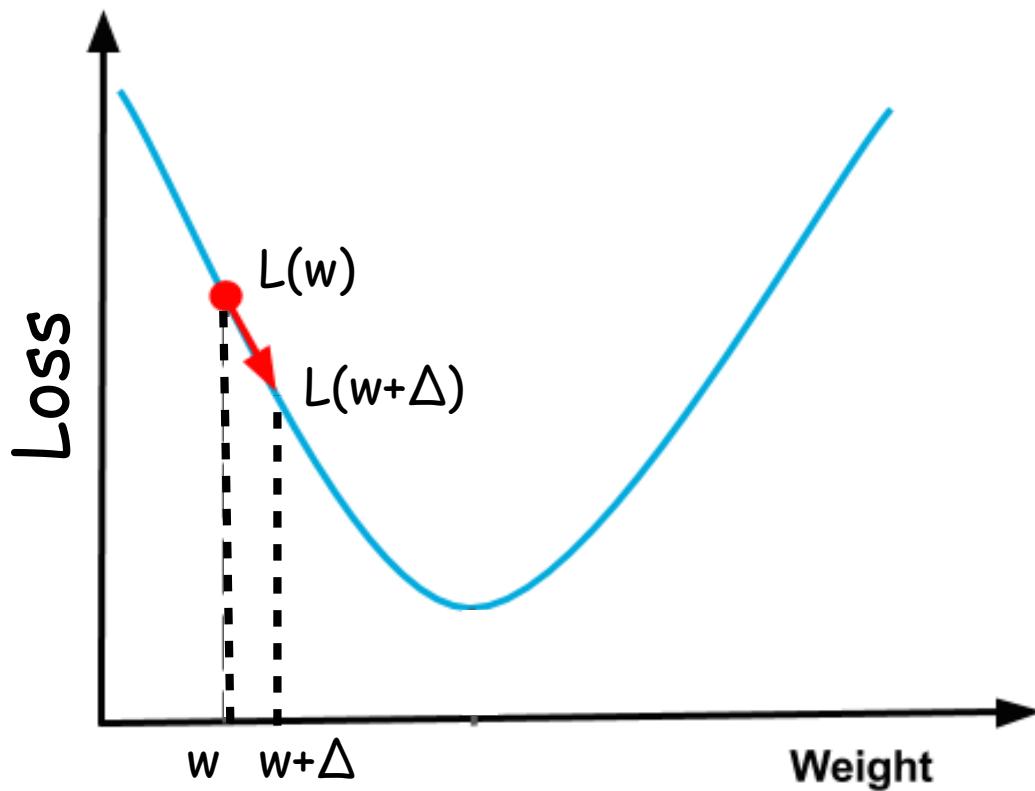
$\frac{\partial L}{\partial b}$ how to check for correctness?

Gradient Checking

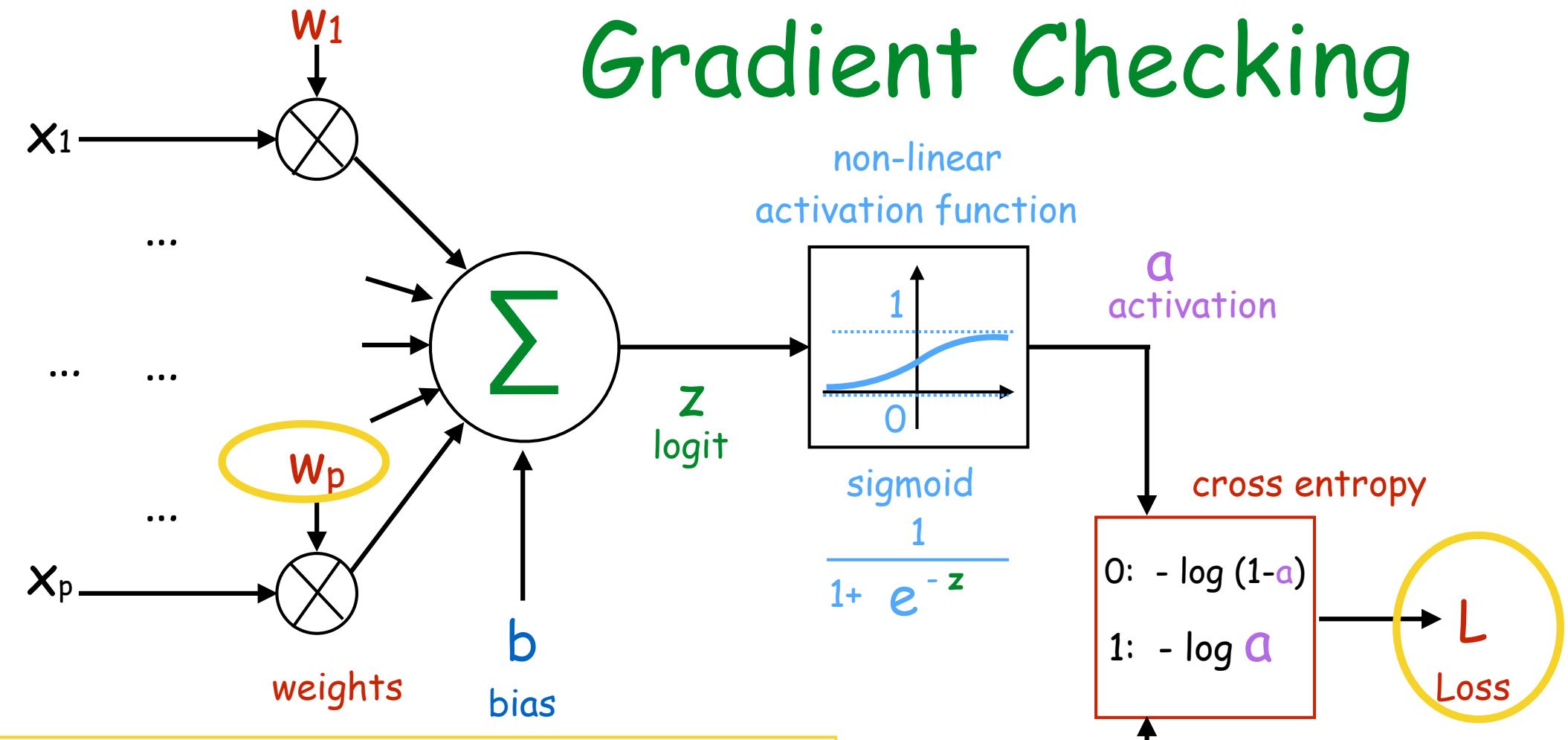
(when Δ is small)

$$\frac{\partial L}{\partial w} \approx \frac{L(w + \Delta) - L(w)}{\Delta}$$

Numerical gradient

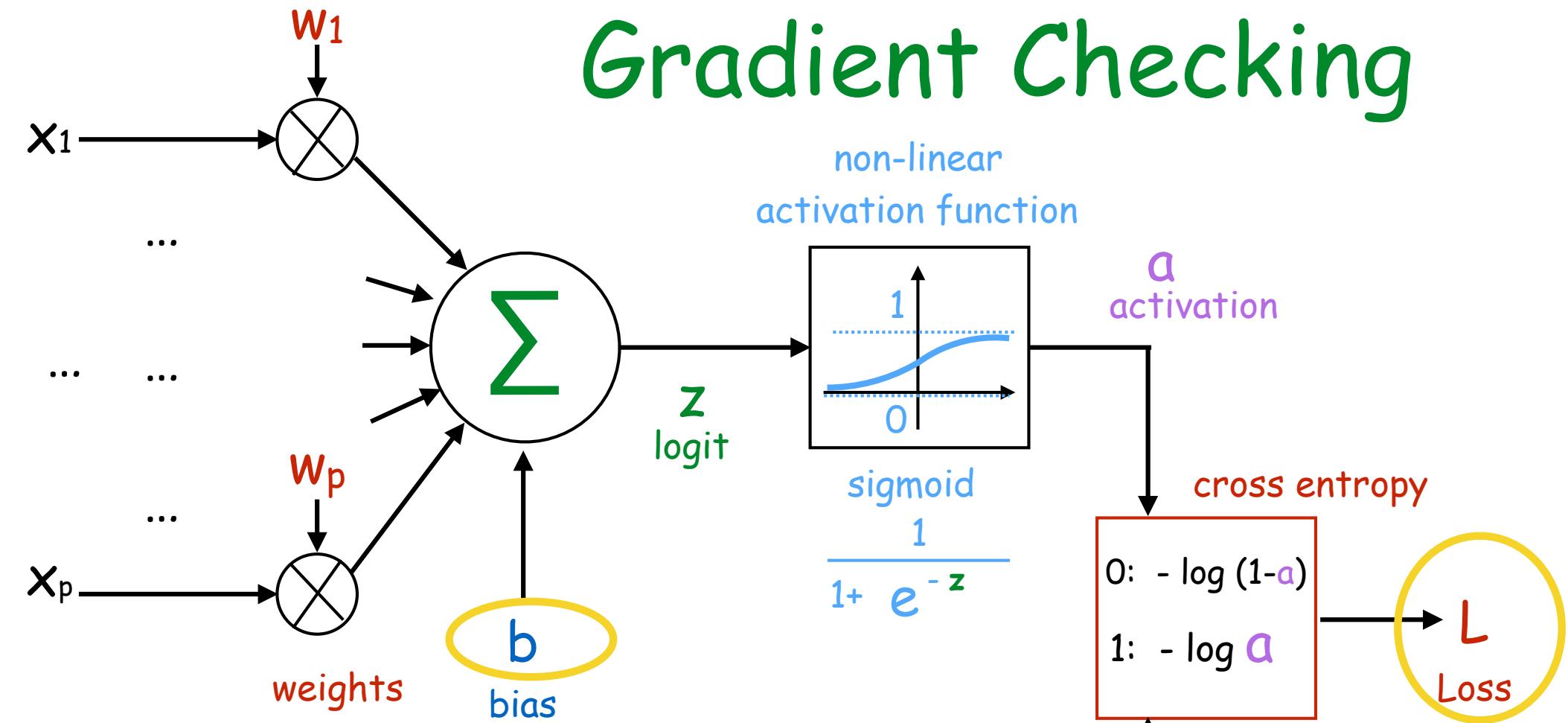


Gradient Checking



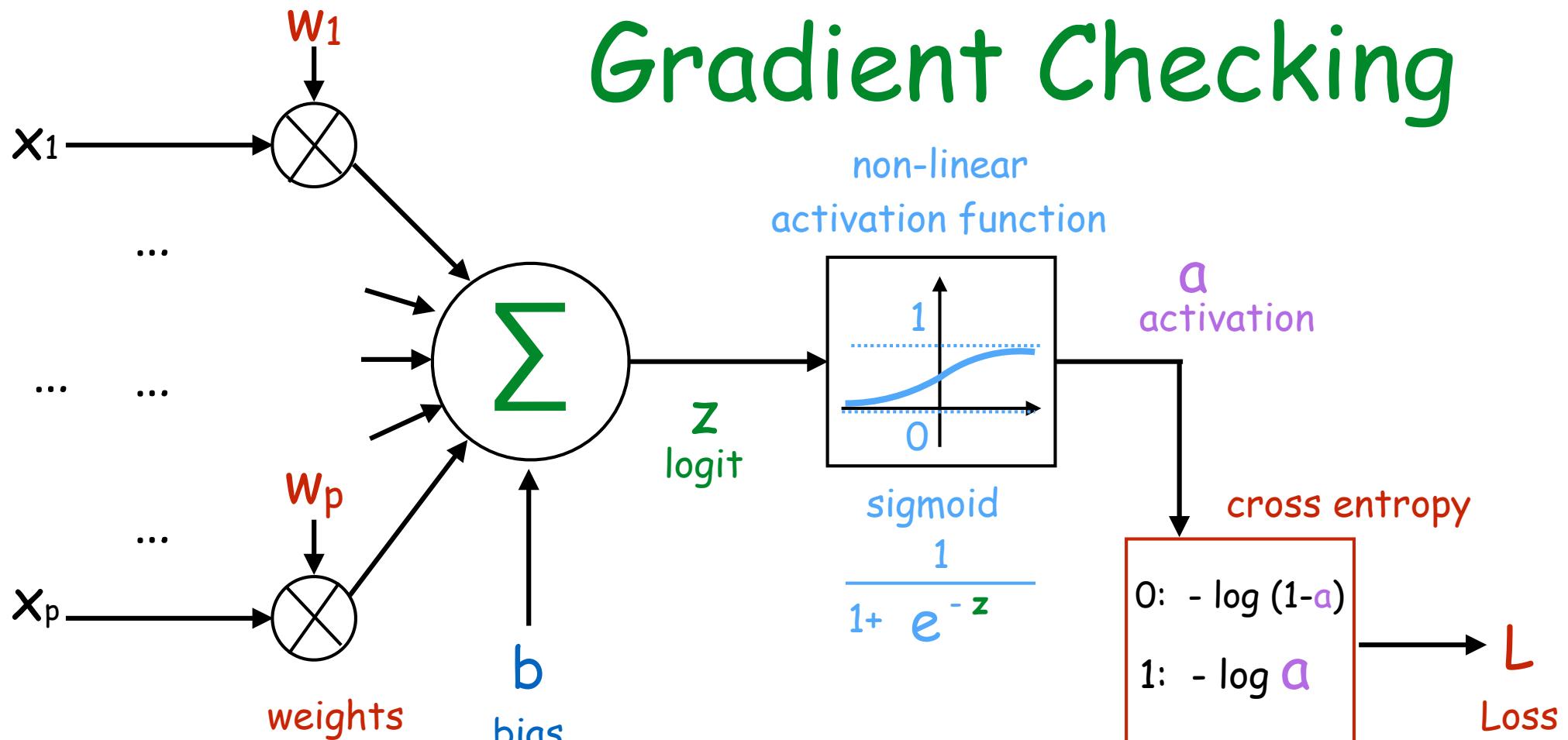
$$\frac{\partial L}{\partial w_p} \approx \frac{L(w_p + \Delta) - L(w_p)}{\Delta}$$

Gradient Checking



$$\frac{\partial L}{\partial b} \approx \frac{L(b + \Delta) - L(b)}{\Delta}$$

Gradient Checking

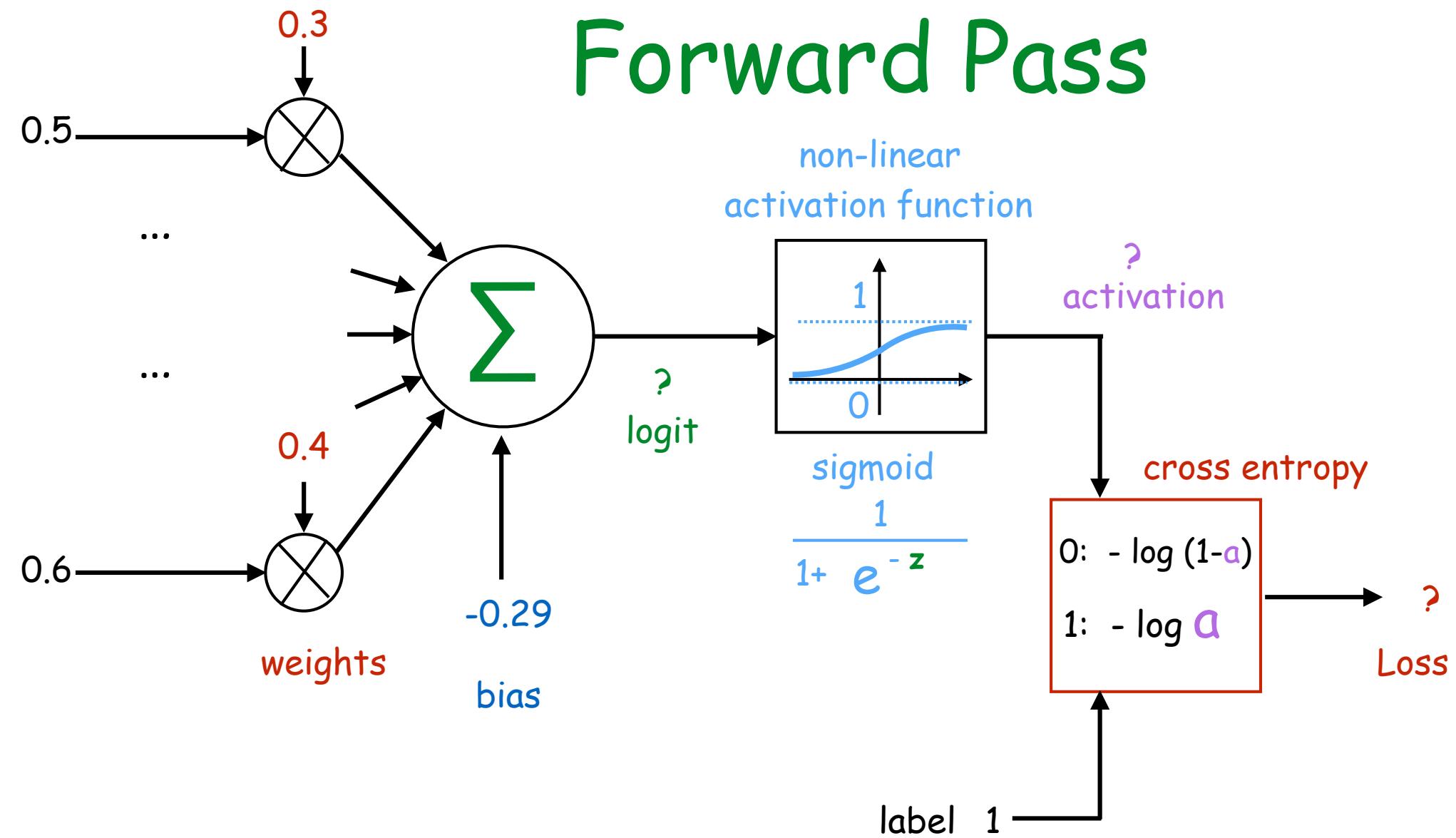


$$\frac{\delta L}{\delta b} \approx \frac{L(b + \Delta) - L(b)}{\Delta}$$

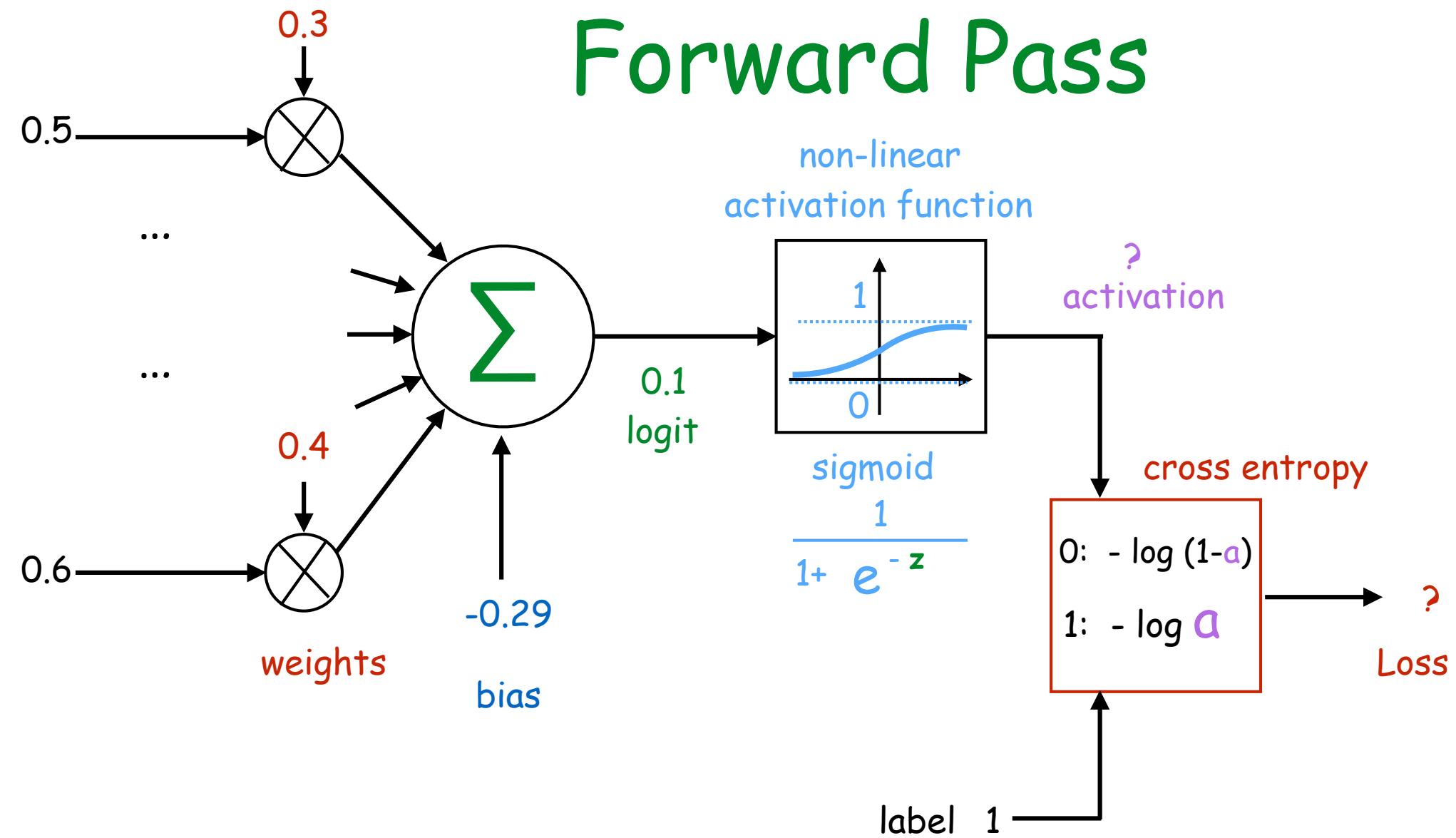
slow (require two forward passes)
 not accurate (require Δ to be very small)
 only used for checking the correctness
 of analytical gradients

Example

Forward Pass



Forward Pass

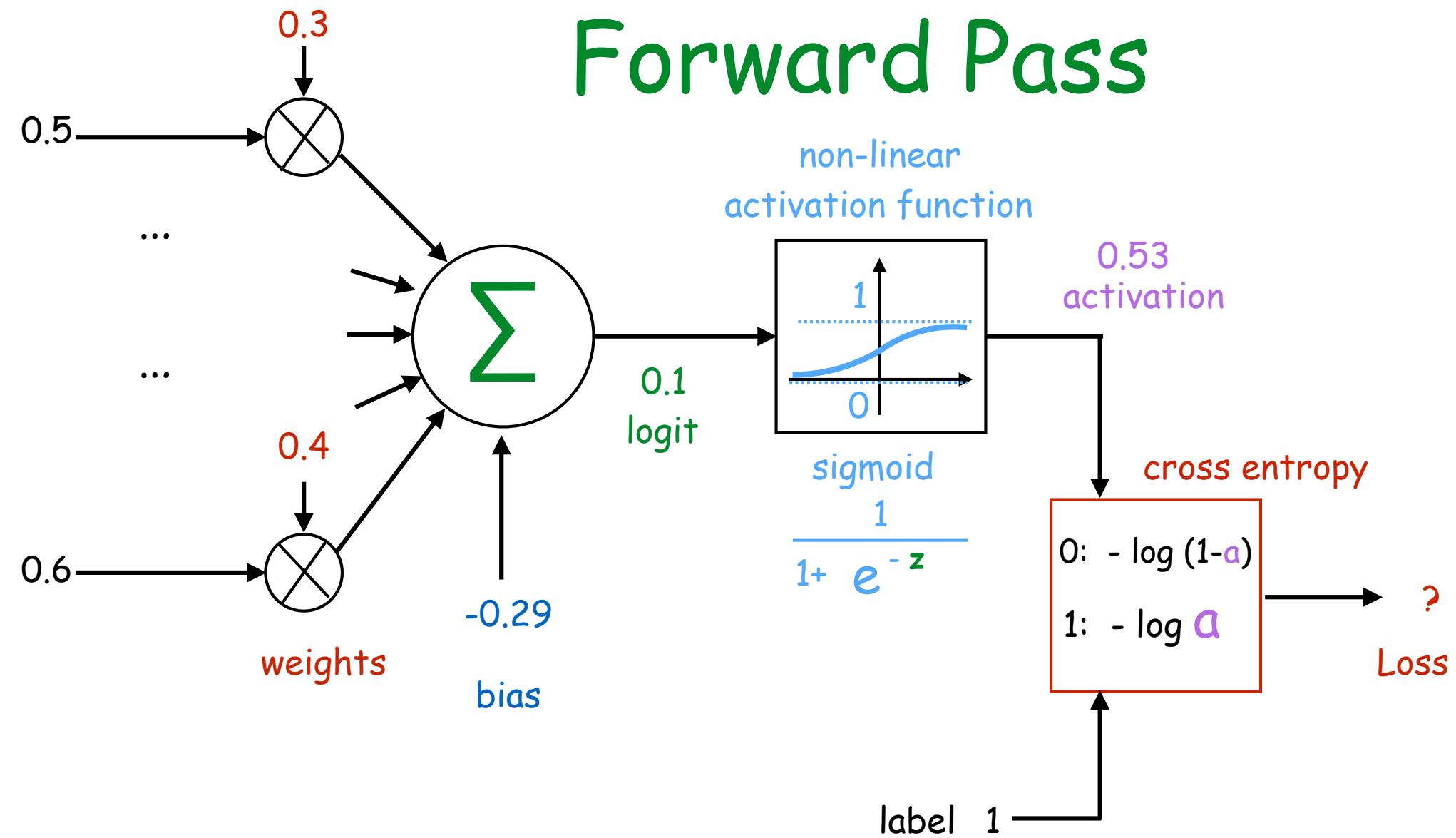


Given a training instance x, y

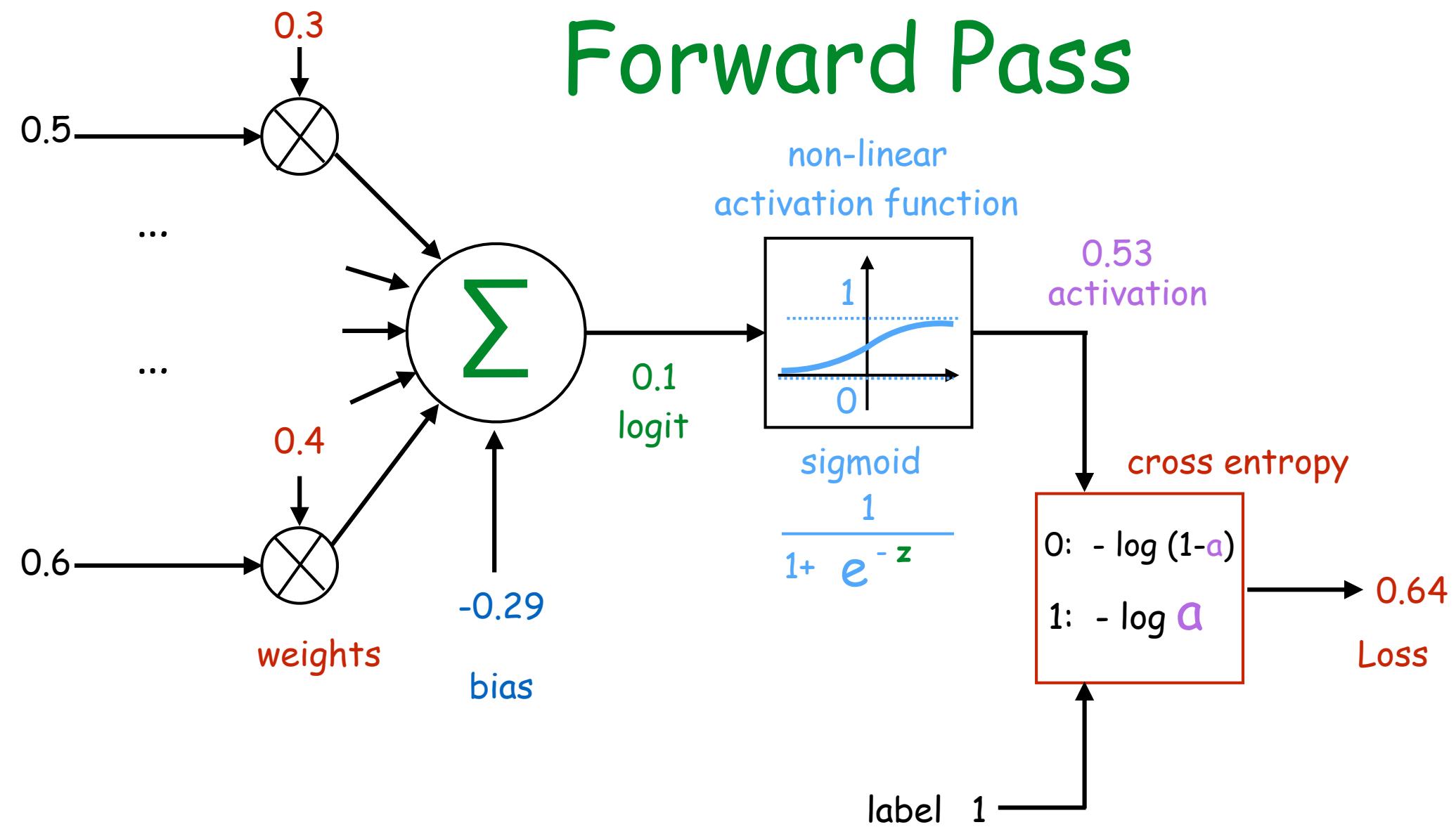
$$x = (0.5, \dots, 0.6)$$

$$y = 1$$

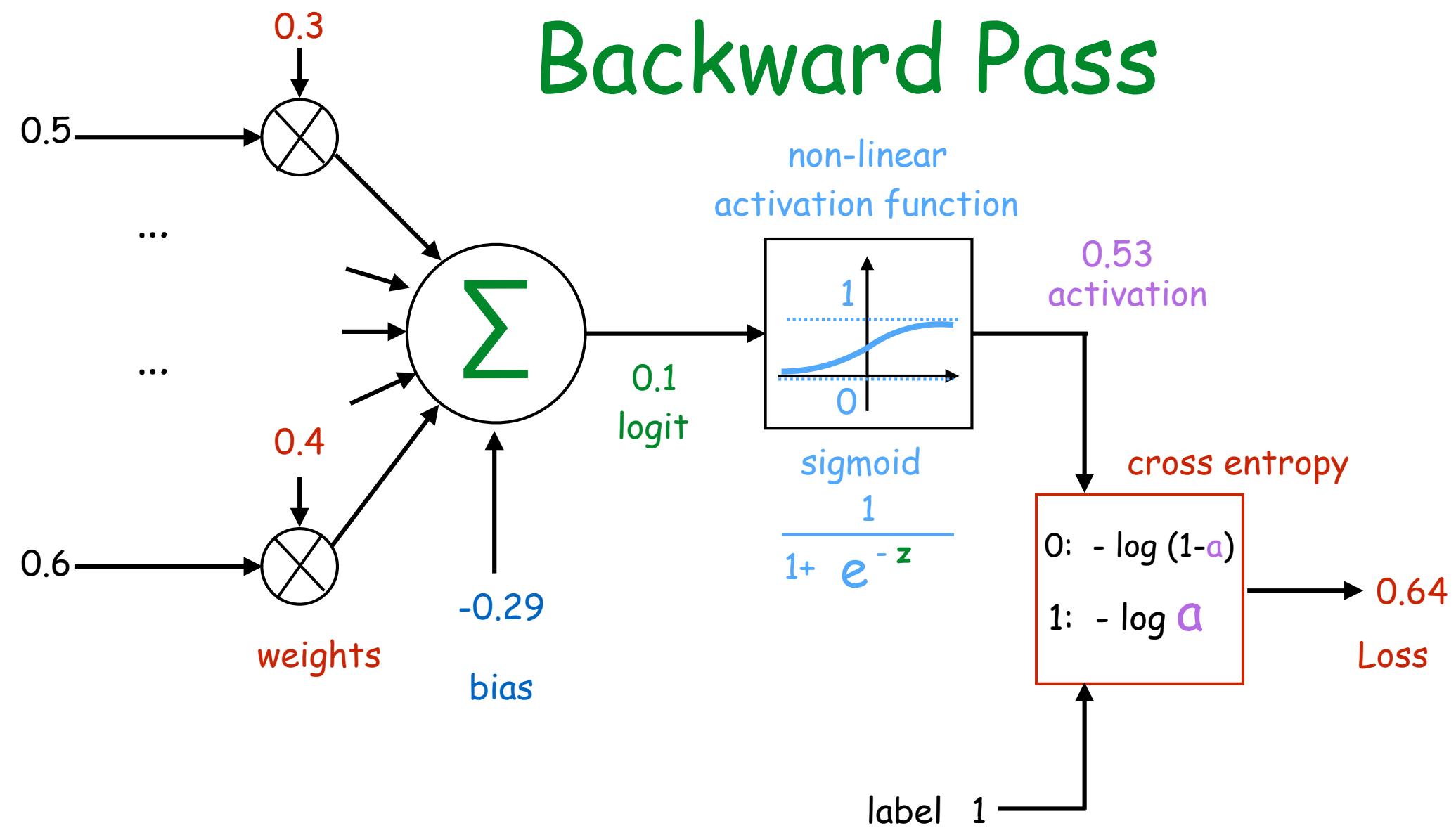
Forward Pass



Forward Pass

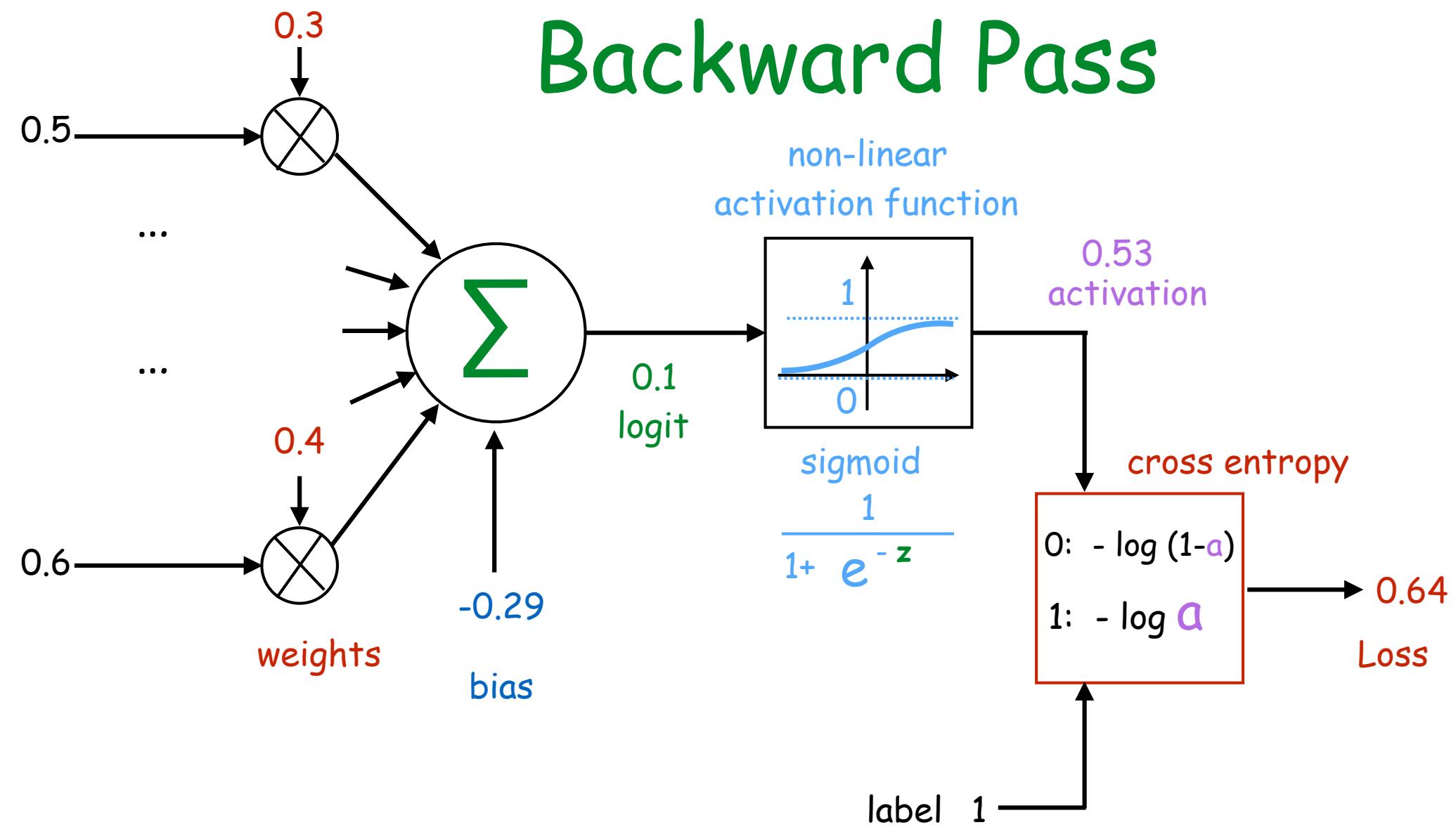


Backward Pass



$$\begin{aligned}\frac{\partial L}{\partial a} &= -1/a \\ &= -1.9\end{aligned}$$

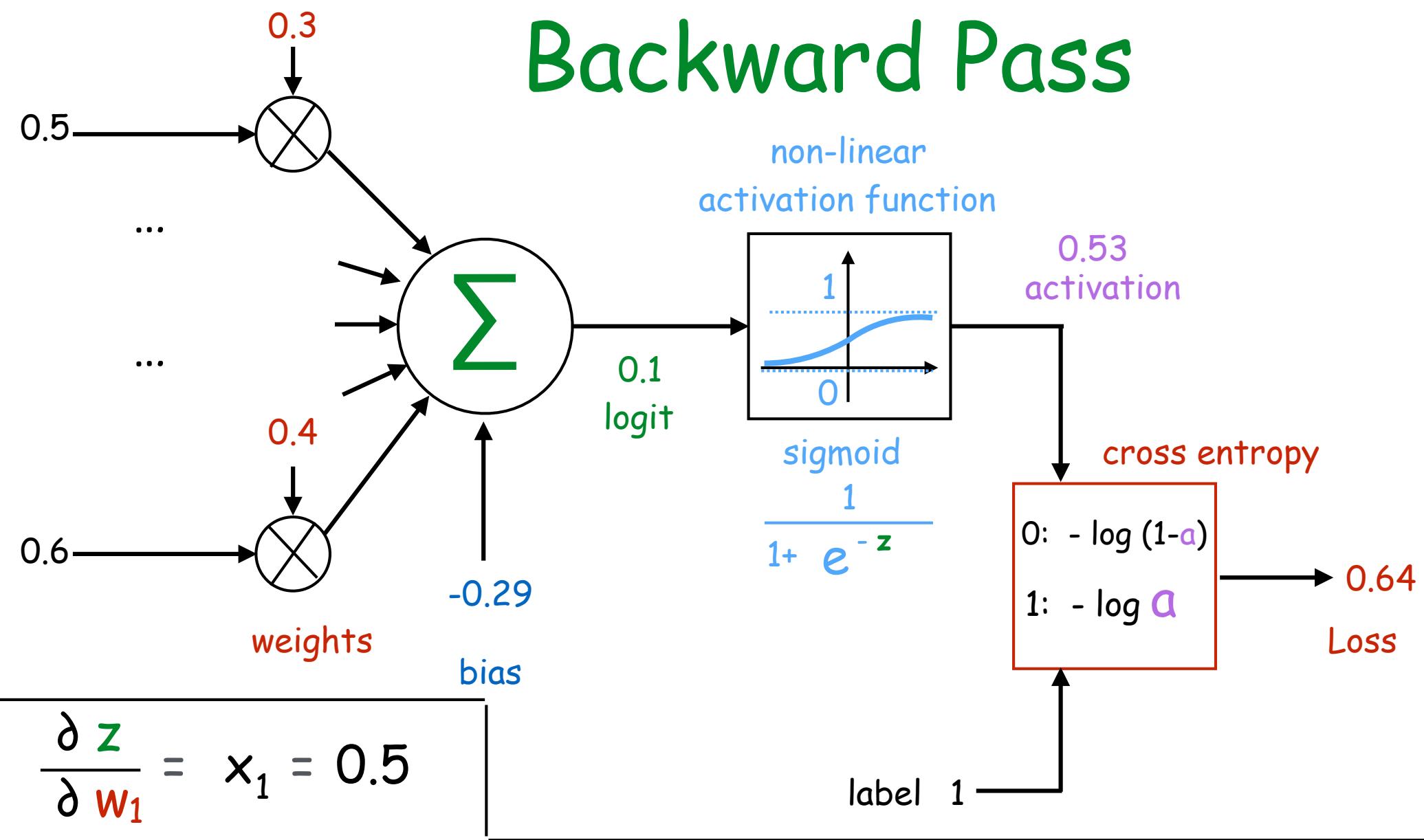
Backward Pass



$$\begin{aligned}\frac{\partial a}{\partial z} &= a(1-a) \\ &= 0.24\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial a} &\\ &= -1.9\end{aligned}$$

Backward Pass



$$\frac{\partial z}{\partial w_1} = x_1 = 0.5$$

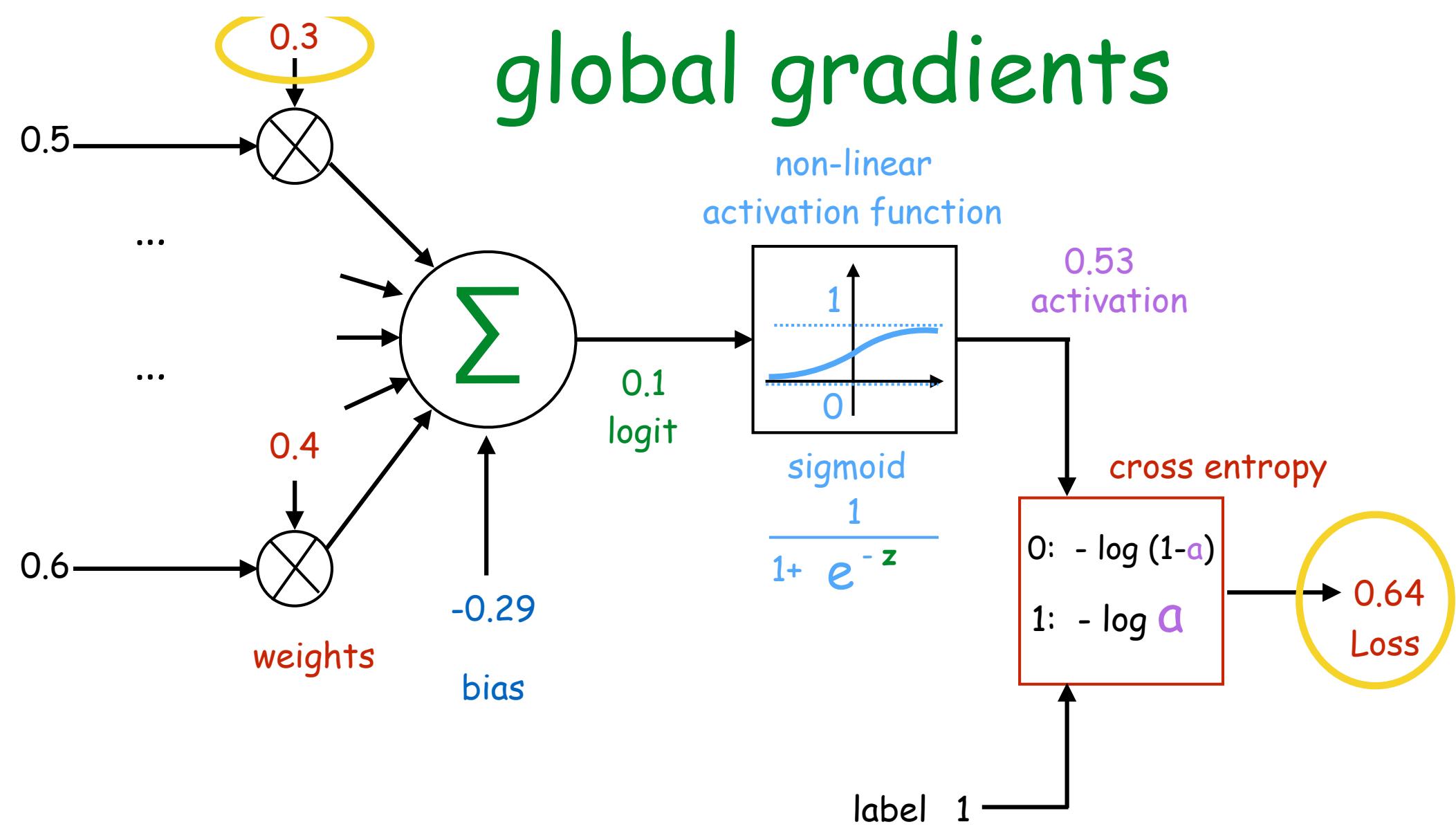
$$\frac{\partial z}{\partial w_p} = x_p = 0.6$$

$$\frac{\partial z}{\partial b} = 1$$

$$\frac{\partial a}{\partial z} = 0.24$$

$$\frac{\partial L}{\partial a} = -1.9$$

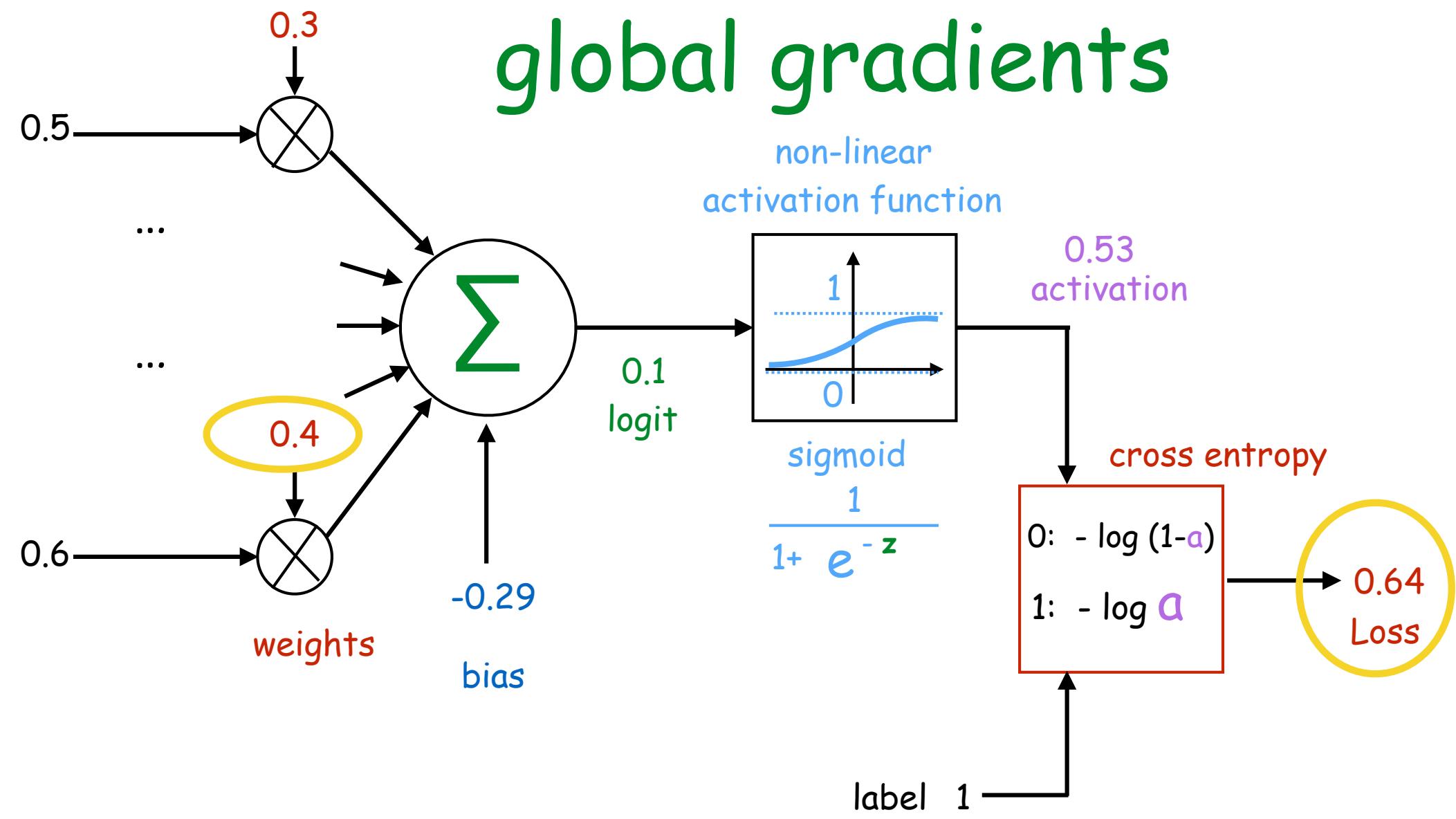
global gradients



$$\frac{\partial L}{\partial w_1} = \frac{\partial z}{\partial w_1} \times \frac{\partial a}{\partial z} \times \frac{\partial L}{\partial a}$$

-0.22 0.5 0.24 -1.9

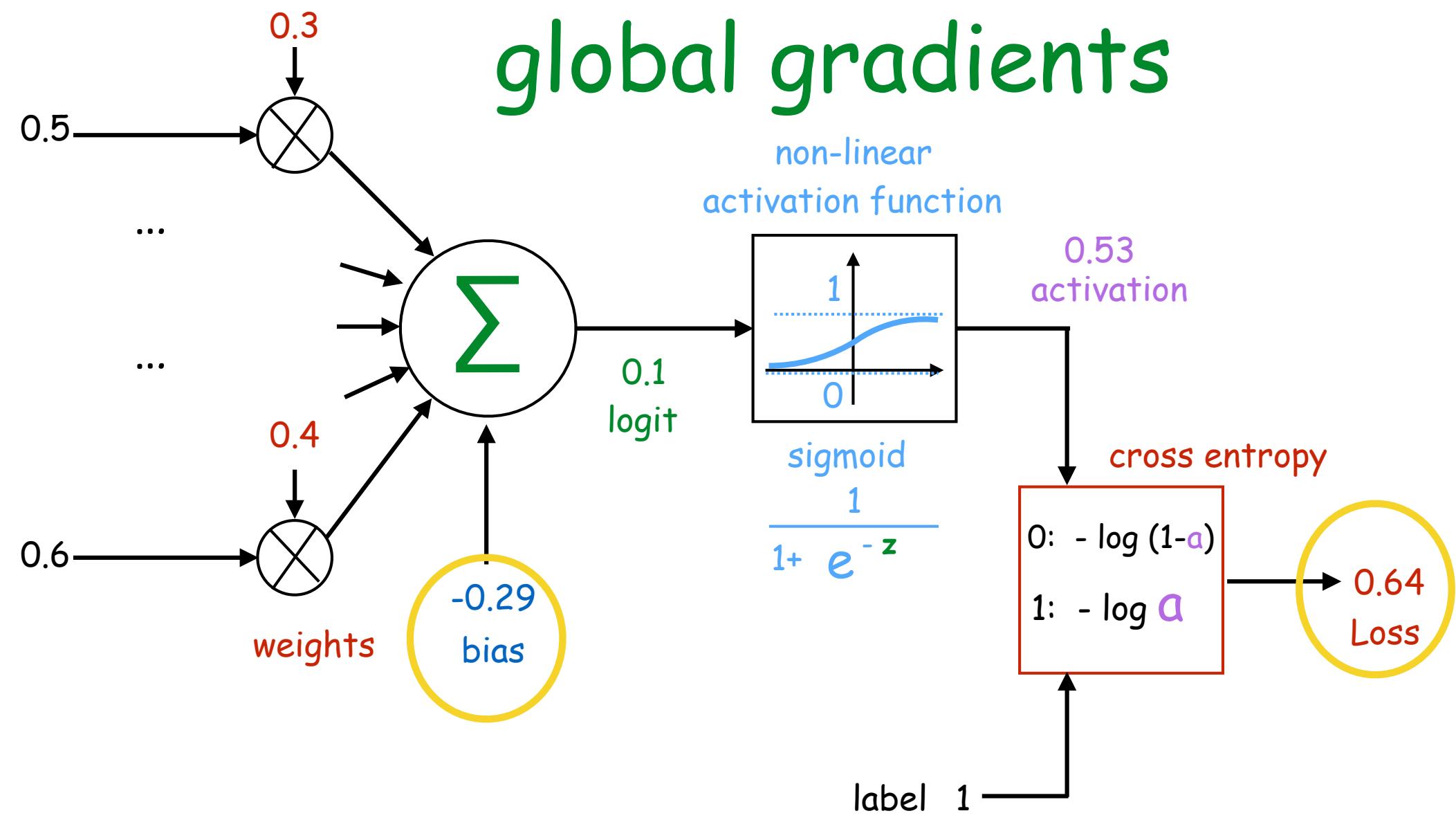
global gradients



$$\frac{\partial L}{\partial w_p} = \frac{\partial z}{\partial w_p} \times \frac{\partial a}{\partial z} \times \frac{\partial L}{\partial a}$$

-0.27 0.6 0.24 -1.9

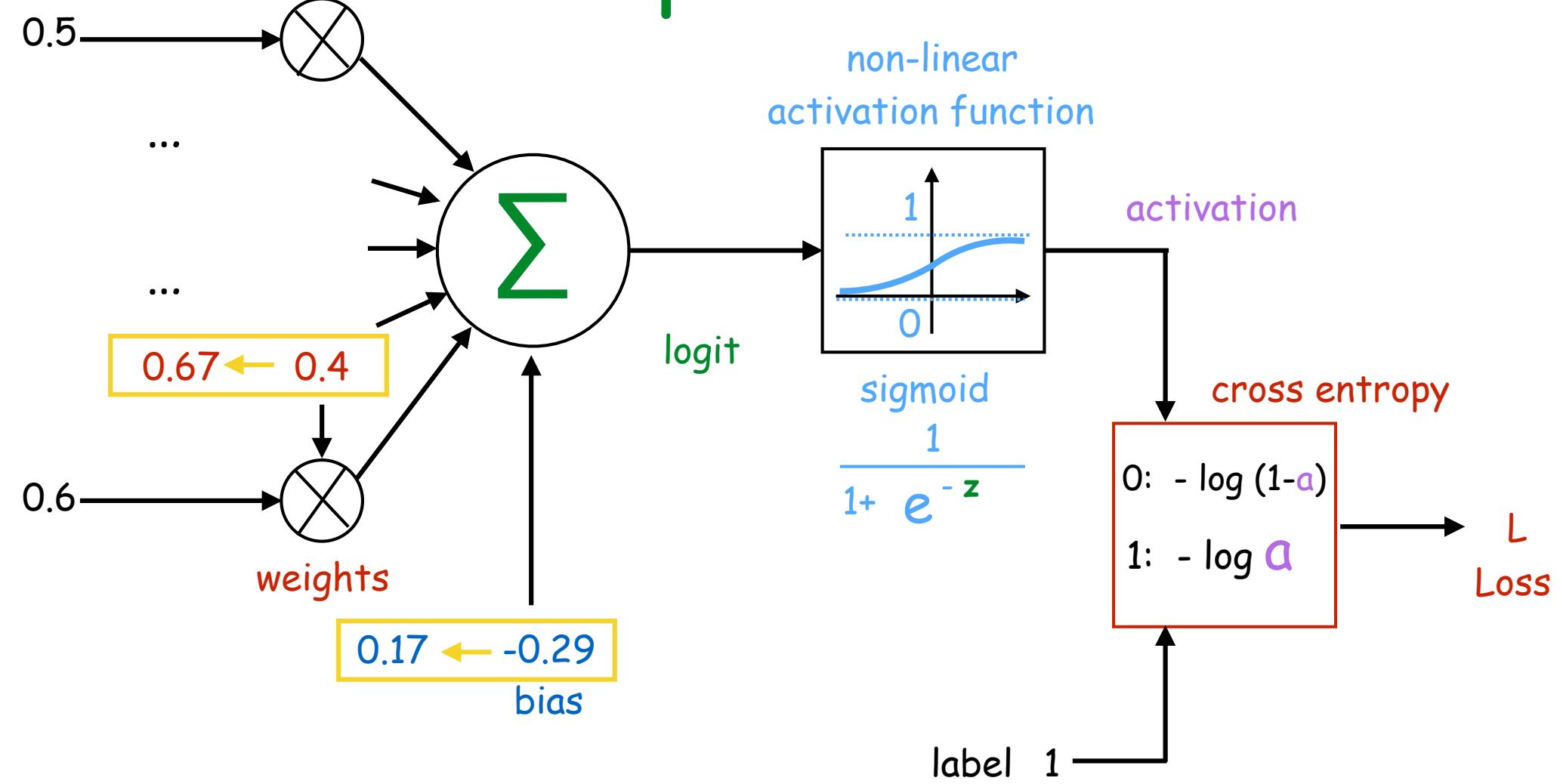
global gradients



$$\frac{\partial L}{\partial b} = \frac{\partial z}{\partial b} \times \frac{\partial a}{\partial z} \times \frac{\partial L}{\partial a}$$

| | | | | |
|---------------------------------|---|---------------------------------|--|--|
| $\frac{\partial L}{\partial b}$ | = | $\frac{\partial z}{\partial b}$ | $\times \frac{\partial a}{\partial z}$ | $\times \frac{\partial L}{\partial a}$ |
| -0.46 | | 1 | 0.24 | -1.9 |

update model



$$w \leftarrow w - a \frac{\partial L}{\partial w}$$

$$b \leftarrow b - a \frac{\partial L}{\partial b}$$

for example

$$a = 1$$

Logistic Regression (train)

initialize w and b

Loop for n_{epoch} iterations:

Loop for each training instance (x, y) in training set

forward pass to compute z , a and L for the instance

backward pass to compute local gradients

$$\frac{\partial L}{\partial a} \quad \frac{\partial a}{\partial z} \quad \frac{\partial z}{\partial b} \quad \frac{\partial z}{\partial w}$$

compute global gradients using chain rule

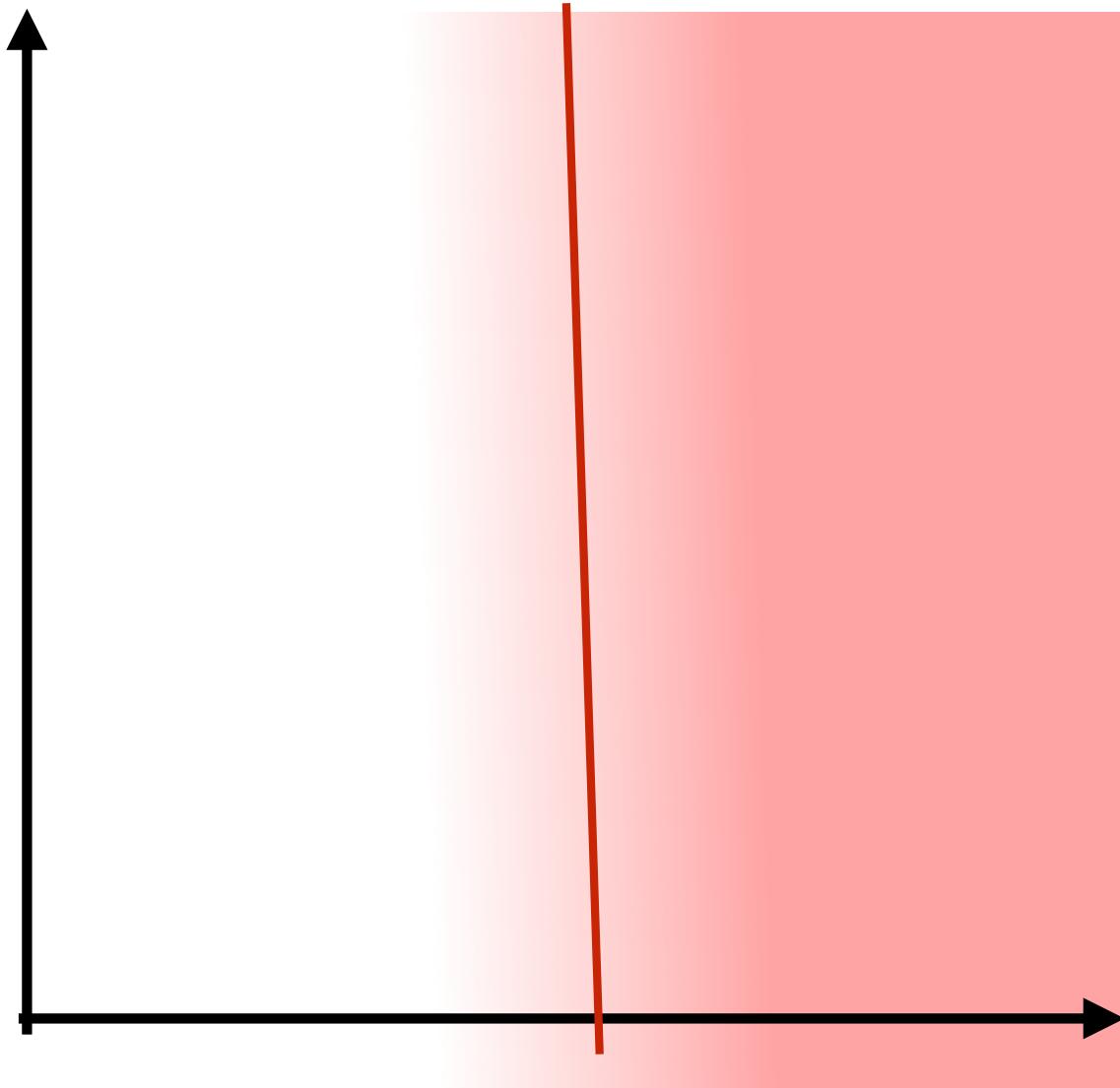
$$\frac{\partial L}{\partial w} \quad \frac{\partial L}{\partial b}$$

update the parameters w and b

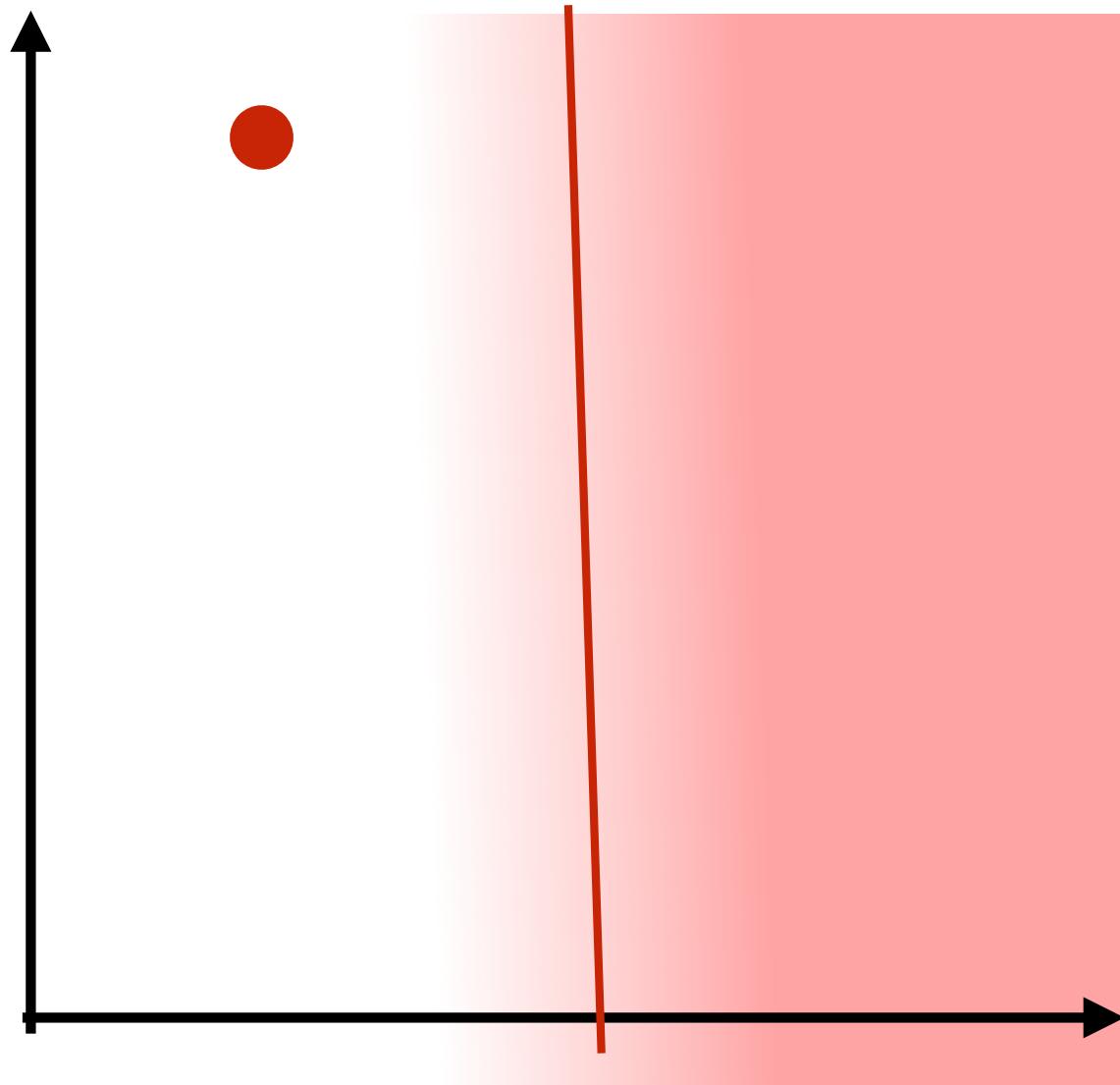
$$w \leftarrow w - a \frac{\partial L}{\partial w}$$

$$b \leftarrow b - a \frac{\partial L}{\partial b}$$

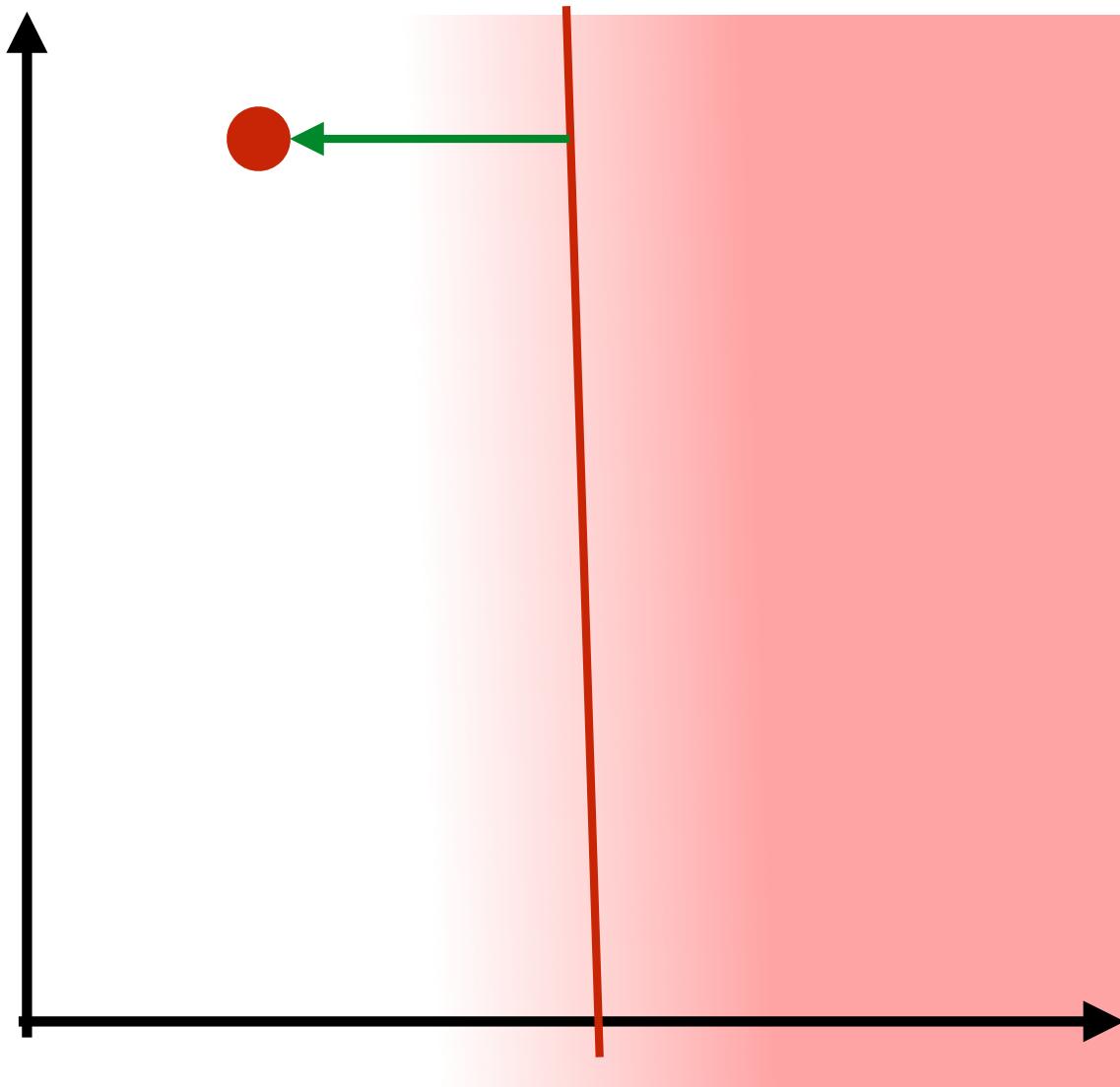
initialize w and b



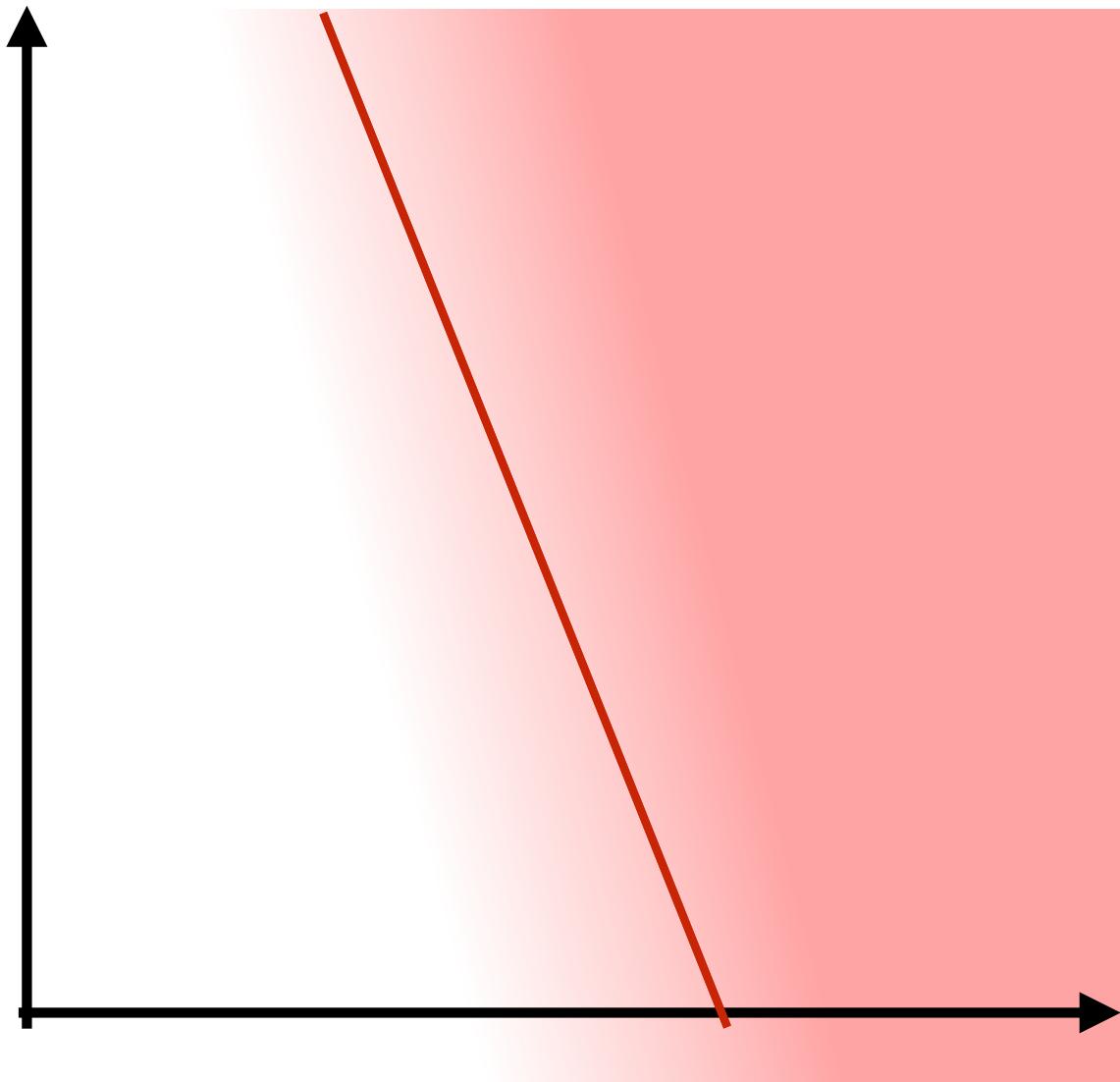
One training instance x and y



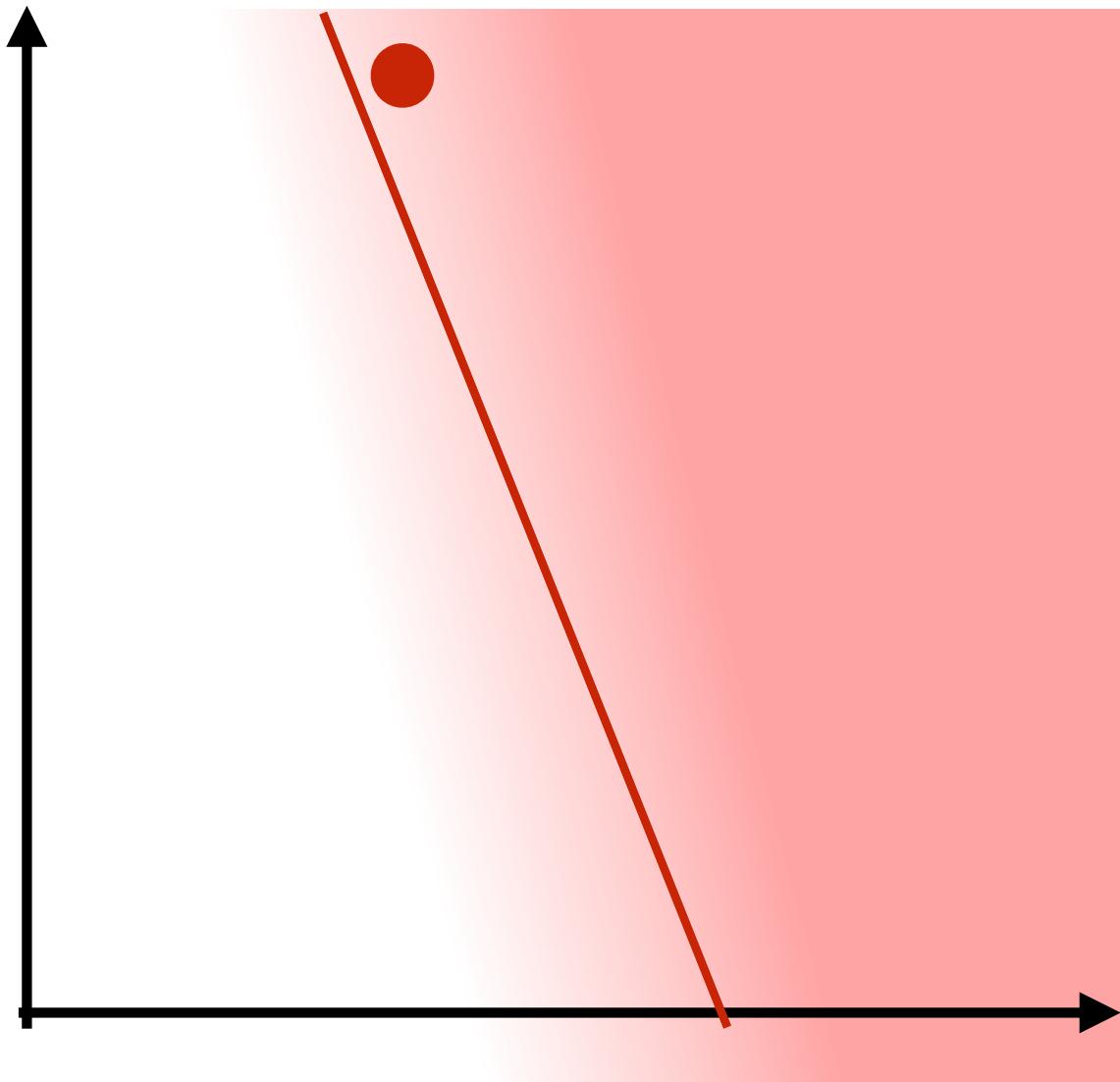
gradient of the loss w.r.t. parameters (w, b)



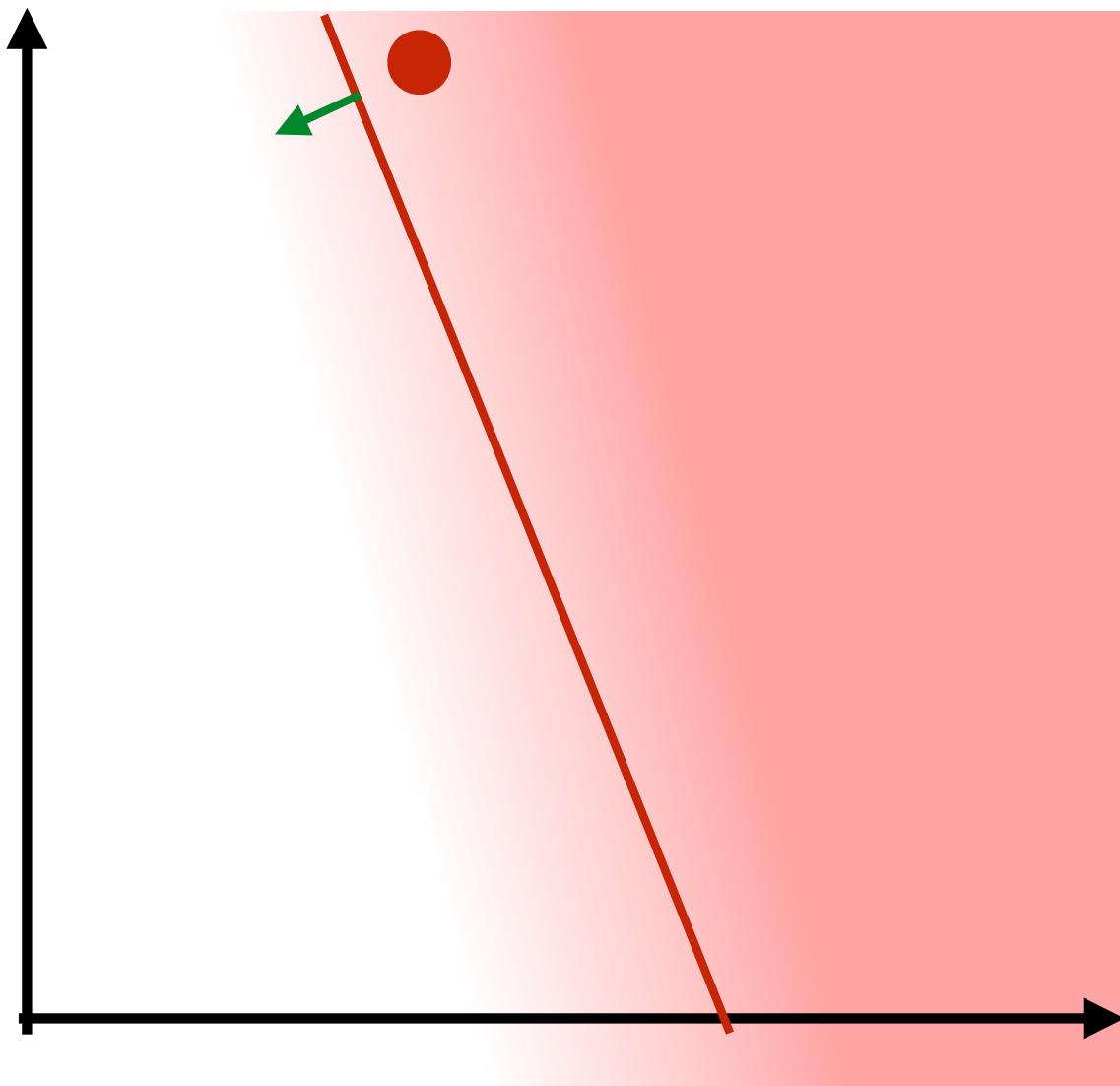
update model



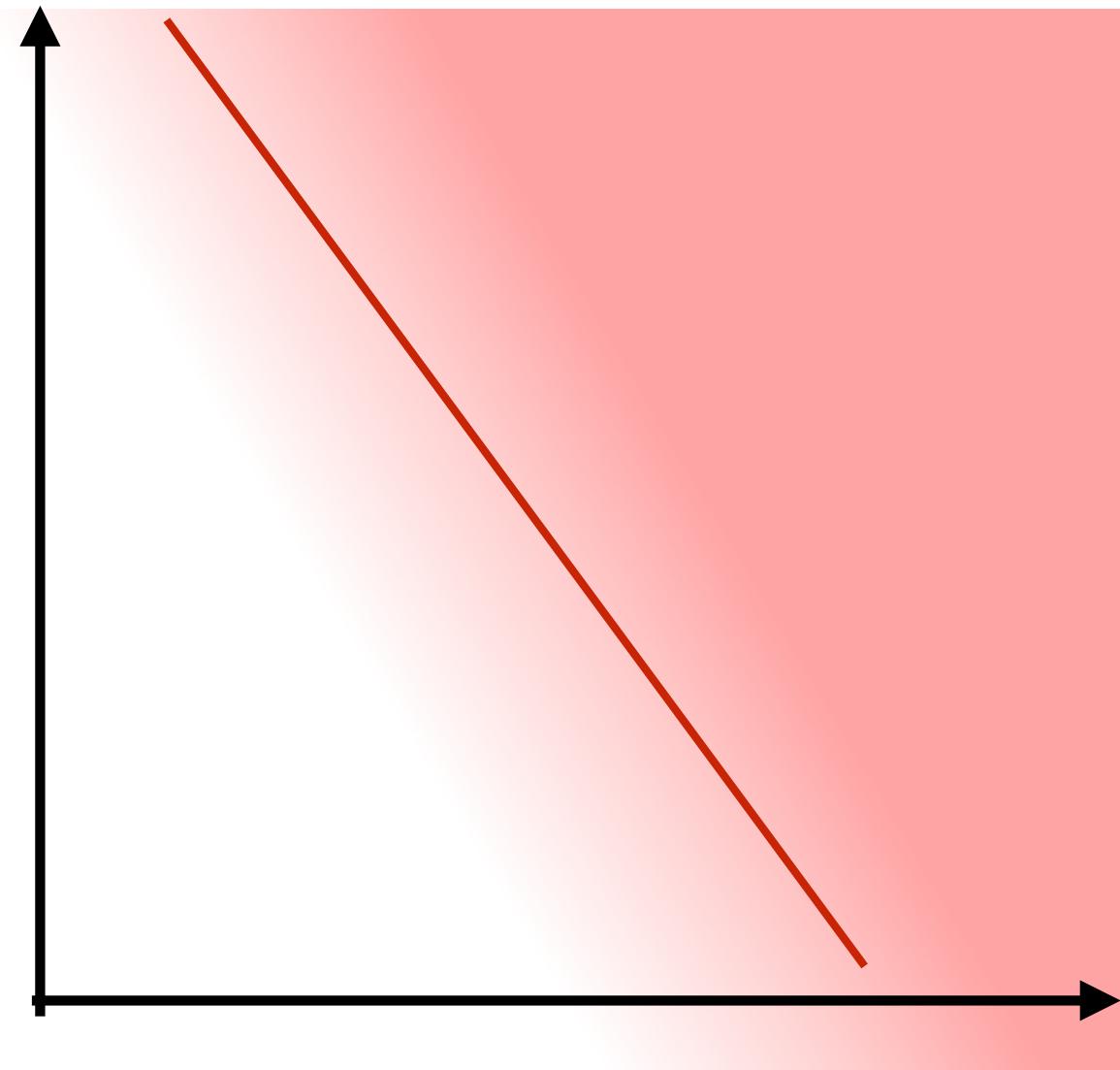
One training instance x and y



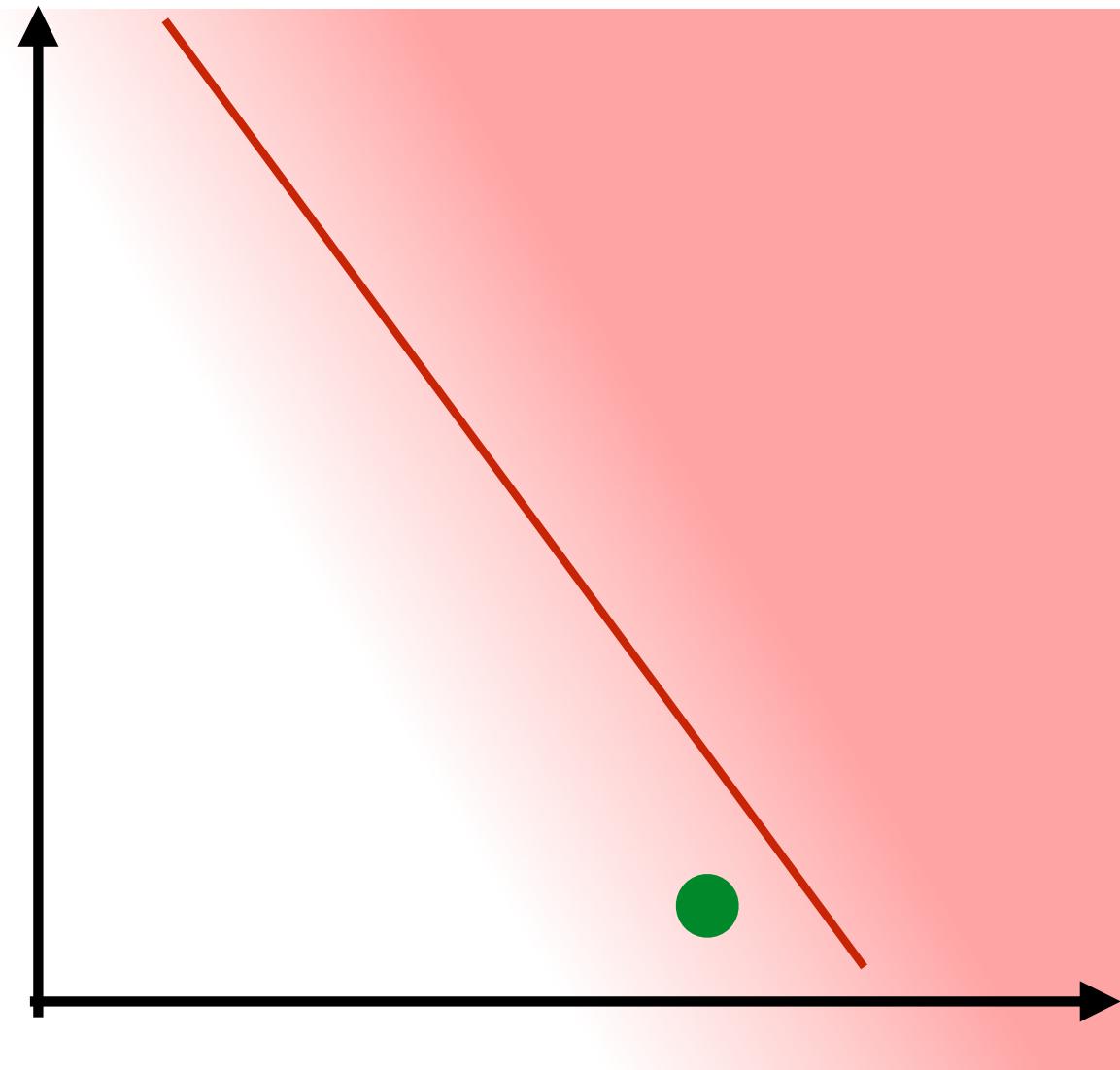
gradient of the loss w.r.t. parameters (w, b)



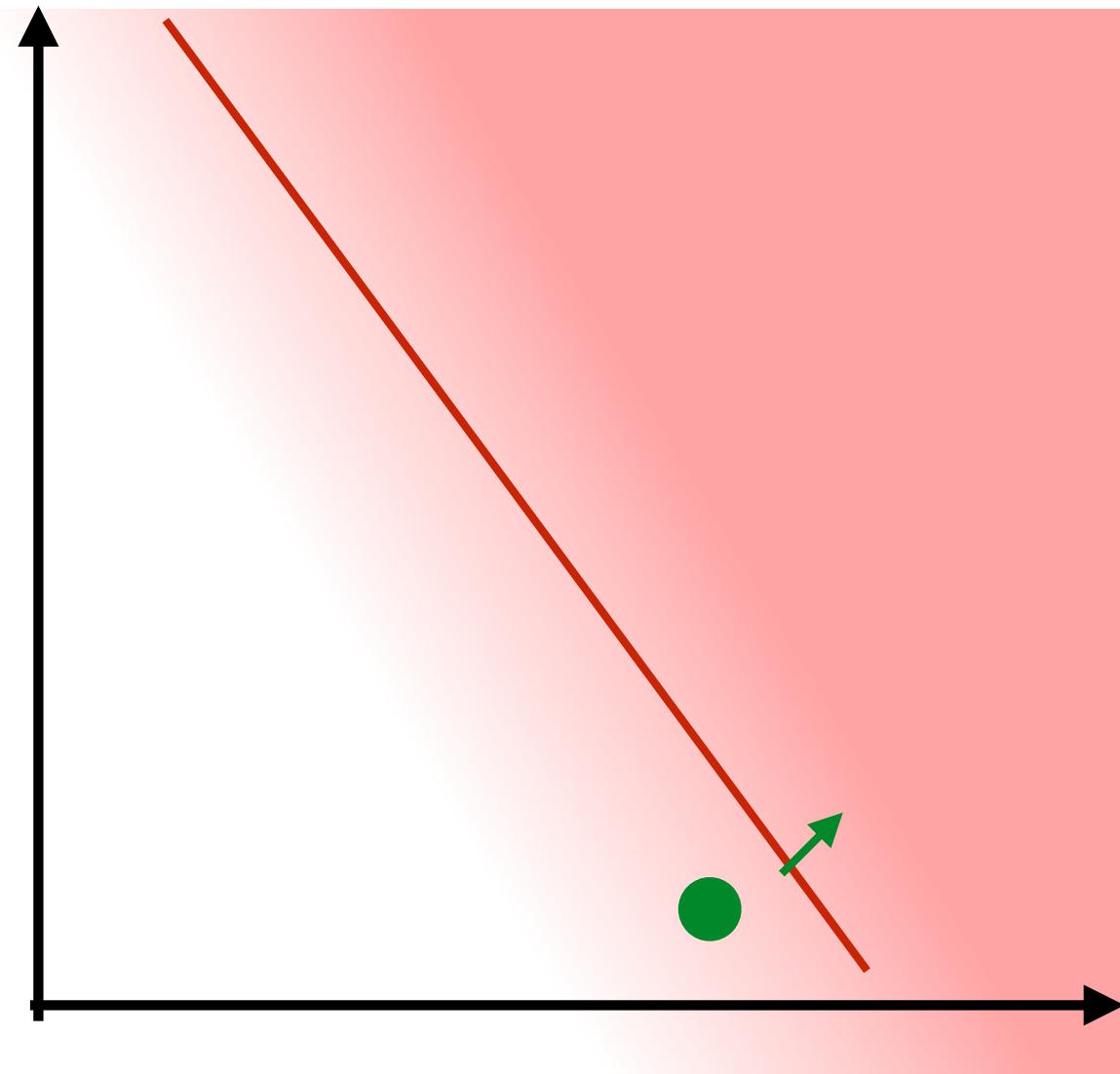
update model



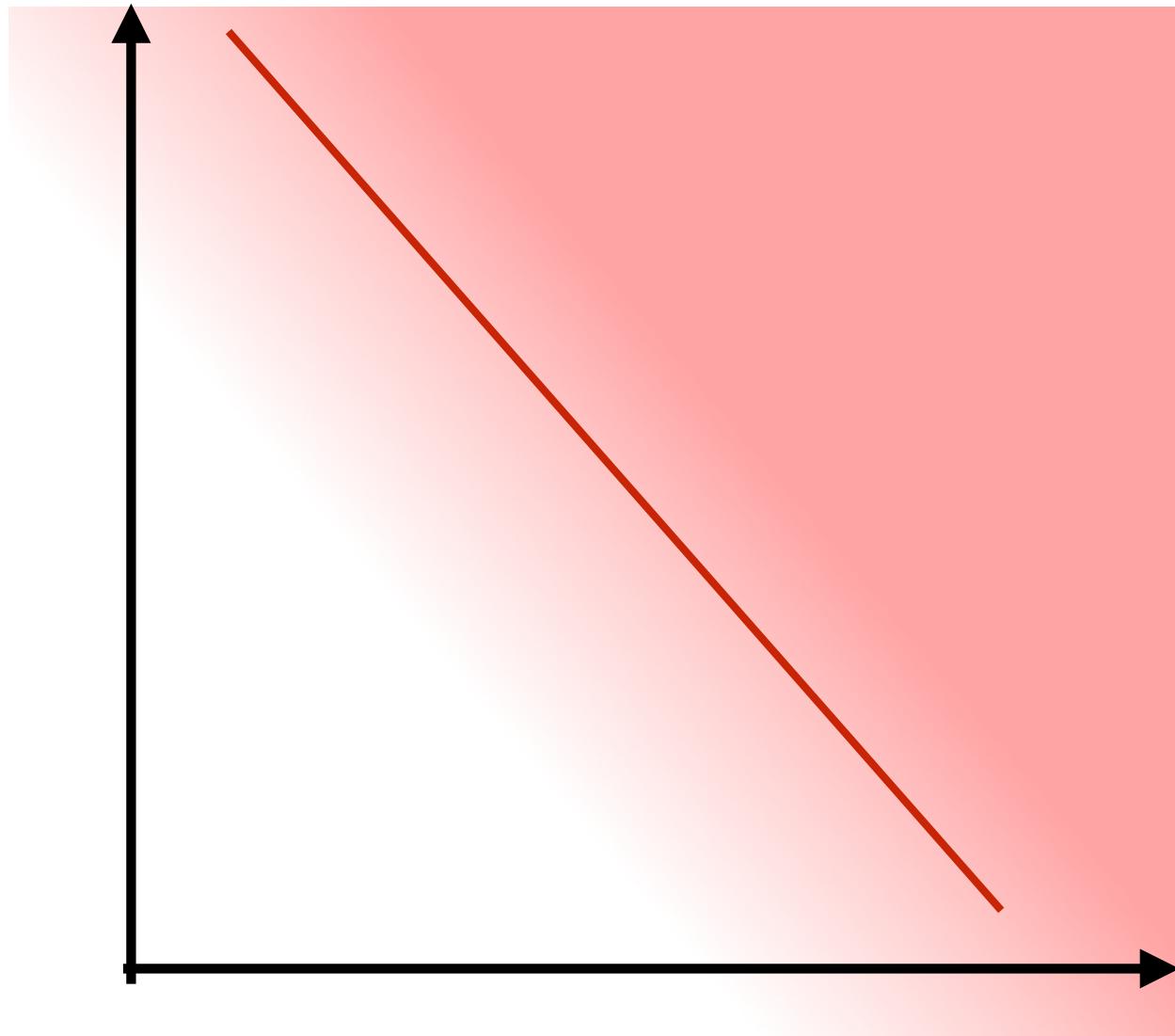
One training instance x and y



gradient of the loss w.r.t. parameters (w, b)

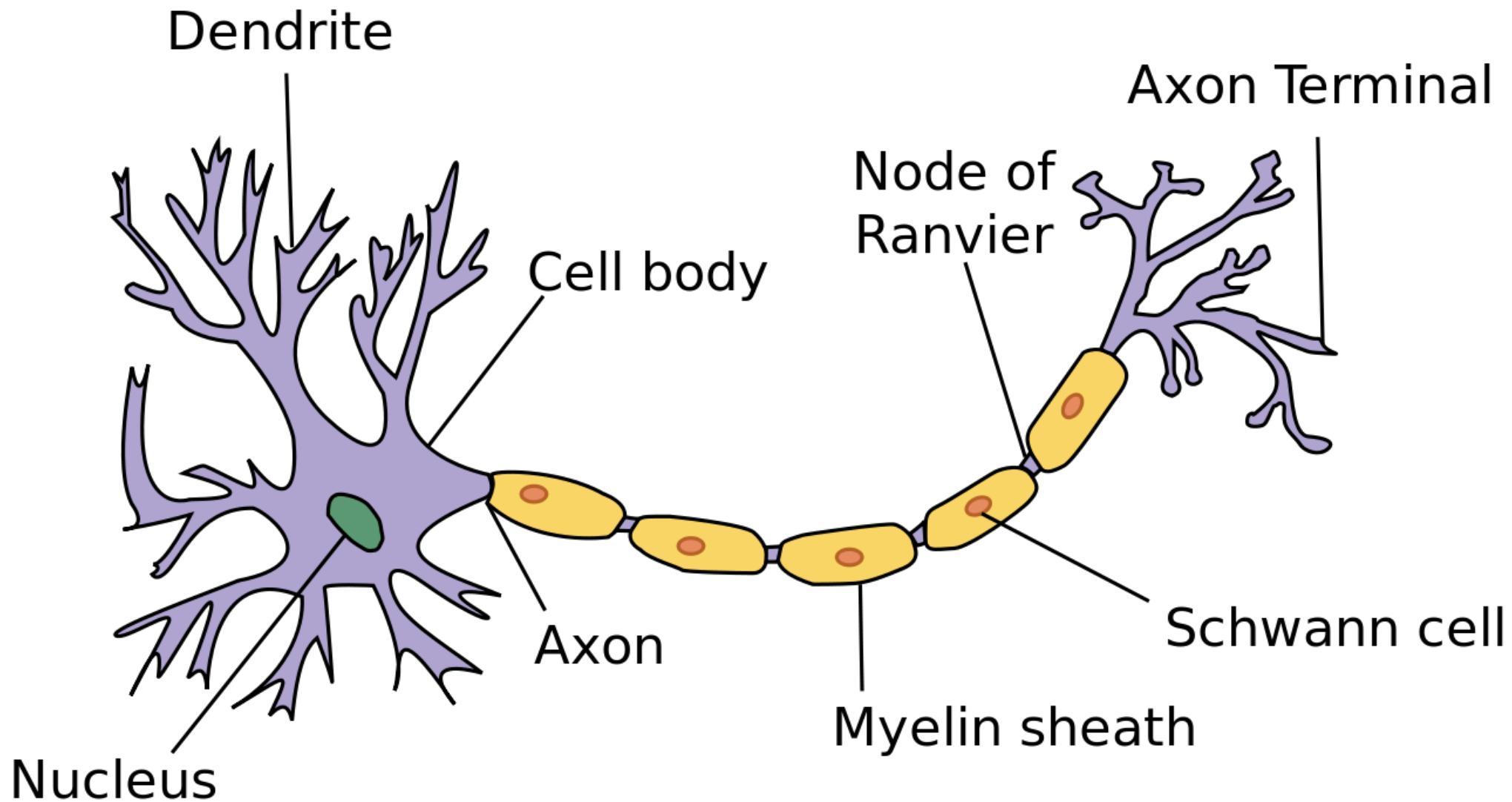


update model parameters (w , b)



Softmax Regression

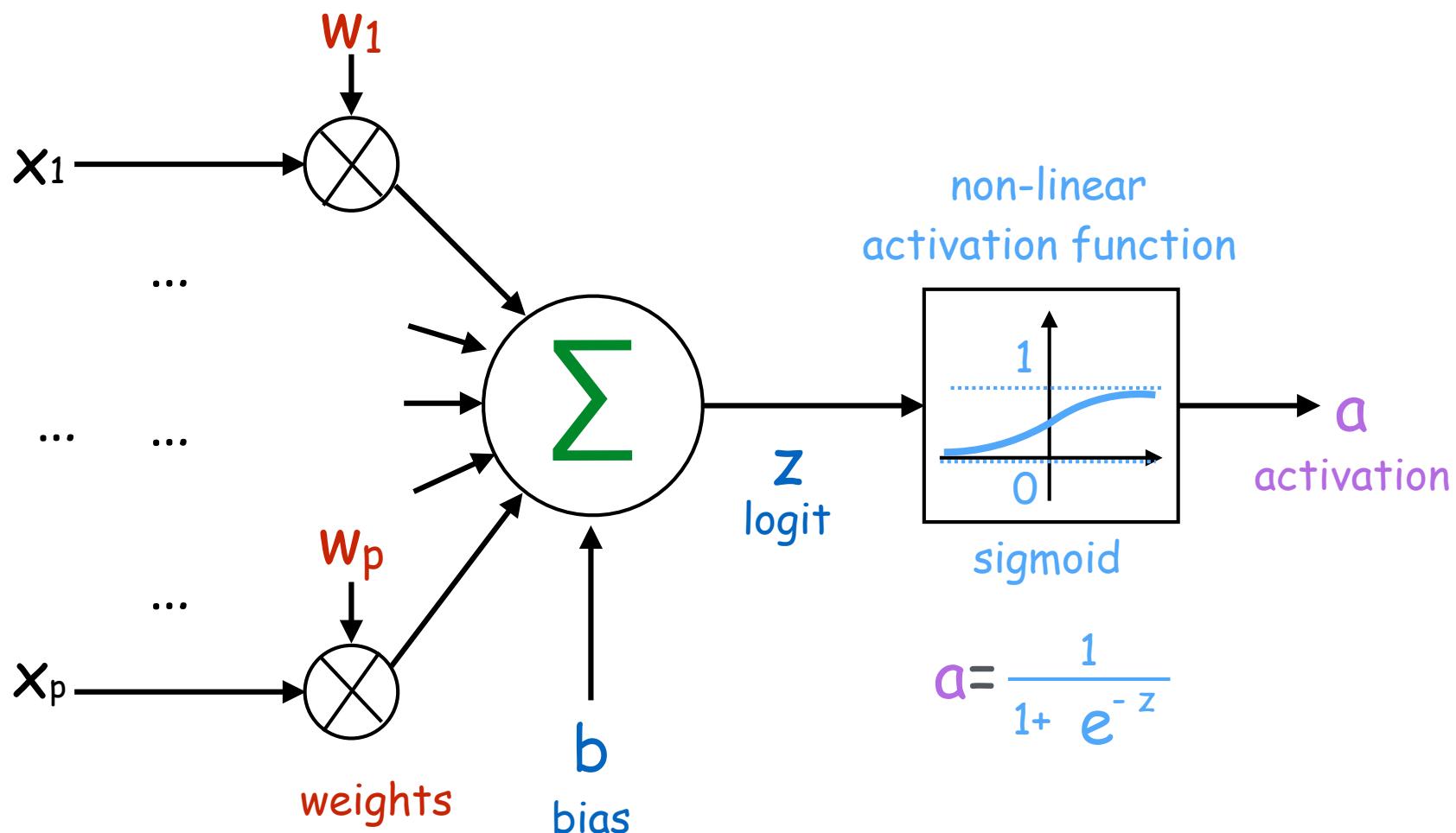
Neuron



<https://www.youtube.com/watch?v=fHRC8SILcH0>

Logistic Regression

Demo: <http://playground.tensorflow.org/>



Multi-class Classification Problem

data set

$2 \rightarrow 2, 5 \rightarrow 5, 4 \rightarrow 8, 0 \rightarrow 0, 2 \rightarrow 2, 7 \rightarrow 7, 5 \rightarrow 5, 1 \rightarrow 1,$
 $3 \rightarrow 3, 0 \rightarrow 0, 3 \rightarrow 3, 9 \rightarrow 9, 6 \rightarrow 6, 2 \rightarrow 2, 8 \rightarrow 8, 2 \rightarrow 2,$
 $0 \rightarrow 0, 4 \rightarrow 6, 6 \rightarrow 6, 1 \rightarrow 1, 1 \rightarrow 1, 7 \rightarrow 7, 8 \rightarrow 8, 5 \rightarrow 5,$
 $0 \rightarrow 0, 4 \rightarrow 4, 7 \rightarrow 7, 6 \rightarrow 6, 0 \rightarrow 0, 2 \rightarrow 2, 5 \rightarrow 5,$
 $3 \rightarrow 3, 1 \rightarrow 1, 5 \rightarrow 5, 6 \rightarrow 6, 7 \rightarrow 7, 5 \rightarrow 5, 4 \rightarrow 4, 1 \rightarrow 1,$
 $9 \rightarrow 9, 3 \rightarrow 3, 6 \rightarrow 6, 8 \rightarrow 8, 0 \rightarrow 0, 9 \rightarrow 9, 3 \rightarrow 3,$
 $0 \rightarrow 0, 3 \rightarrow 3, 7 \rightarrow 7, 4 \rightarrow 4, 4 \rightarrow 4, 3 \rightarrow 3, 8 \rightarrow 8, 0 \rightarrow 0,$
 $4 \rightarrow 4, 1 \rightarrow 1, 3 \rightarrow 3, 7 \rightarrow 7, 6 \rightarrow 6, 4 \rightarrow 4, 7 \rightarrow 7, 2 \rightarrow 2,$
 $7 \rightarrow 7, 2 \rightarrow 2, 5 \rightarrow 5, 2 \rightarrow 2, 0 \rightarrow 0, 9 \rightarrow 9, 8 \rightarrow 8, 9 \rightarrow 9,$
 ~~$8 \rightarrow 8, 1 \rightarrow 1, 6 \rightarrow 6, 4 \rightarrow 4, 8 \rightarrow 8, 5 \rightarrow 5, 8 \rightarrow 8,$~~
 $0 \rightarrow 0, 6 \rightarrow 6, 7 \rightarrow 7, 4 \rightarrow 4, 5 \rightarrow 5, 8 \rightarrow 8, 4 \rightarrow 4,$
 $3 \rightarrow 3, 1 \rightarrow 1, 5 \rightarrow 5, 1 \rightarrow 1, 9 \rightarrow 9, 9 \rightarrow 9, 9 \rightarrow 9, 2 \rightarrow 2,$
 $4 \rightarrow 4, 7 \rightarrow 7, 3 \rightarrow 3, 1 \rightarrow 1, 9 \rightarrow 9, 2 \rightarrow 2, 9 \rightarrow 9, 6 \rightarrow 6]$

input
(instance)



candidate labels
(classes)

0
1
2
⋮
9

output
(label)

0 or 1 or ... or 9

Feature Matrix X (n by p)

| Instance | Feature (pixel) | | | | | | | | | | | |
|----------|-----------------|---|-----|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | ... | p | | | | | | | | |
| X_1 2 | 1 | 3 | 4 | 3 | 8 | 3 | 5 | 7 | 9 | 7 | 3 | 4 |
| X_2 5 | 3 | 3 | 5 | 7 | 7 | 0 | 4 | 1 | 2 | 1 | 9 | 7 |
| 0 | 7 | 0 | 4 | 1 | 1 | 4 | 3 | 7 | 8 | 6 | 2 | 7 |
| 3 | 1 | 4 | 3 | 7 | 9 | 7 | 3 | 2 | 7 | 0 | 4 | 1 |
| 8 | 7 | 7 | 3 | 2 | 2 | 1 | 9 | 8 | 1 | 4 | 3 | 7 |
| 0 | 2 | 1 | 9 | 8 | 8 | 6 | 2 | 0 | 7 | 7 | 3 | 2 |
| 9 | 8 | 6 | 2 | 0 | 0 | 4 | 1 | 1 | 4 | 1 | 9 | 8 |
| 6 | 0 | 2 | 1 | 4 | 1 | 3 | 7 | 9 | 7 | 6 | 2 | 0 |
| 7 | 3 | 5 | 3 | 3 | 7 | 3 | 2 | 2 | 1 | 2 | 1 | 3 |
| X_n 6 | 1 | 7 | 2 | 3 | 2 | 2 | 1 | 2 | 3 | 5 | 3 | 1 |

Label Vector
 y

(length n)

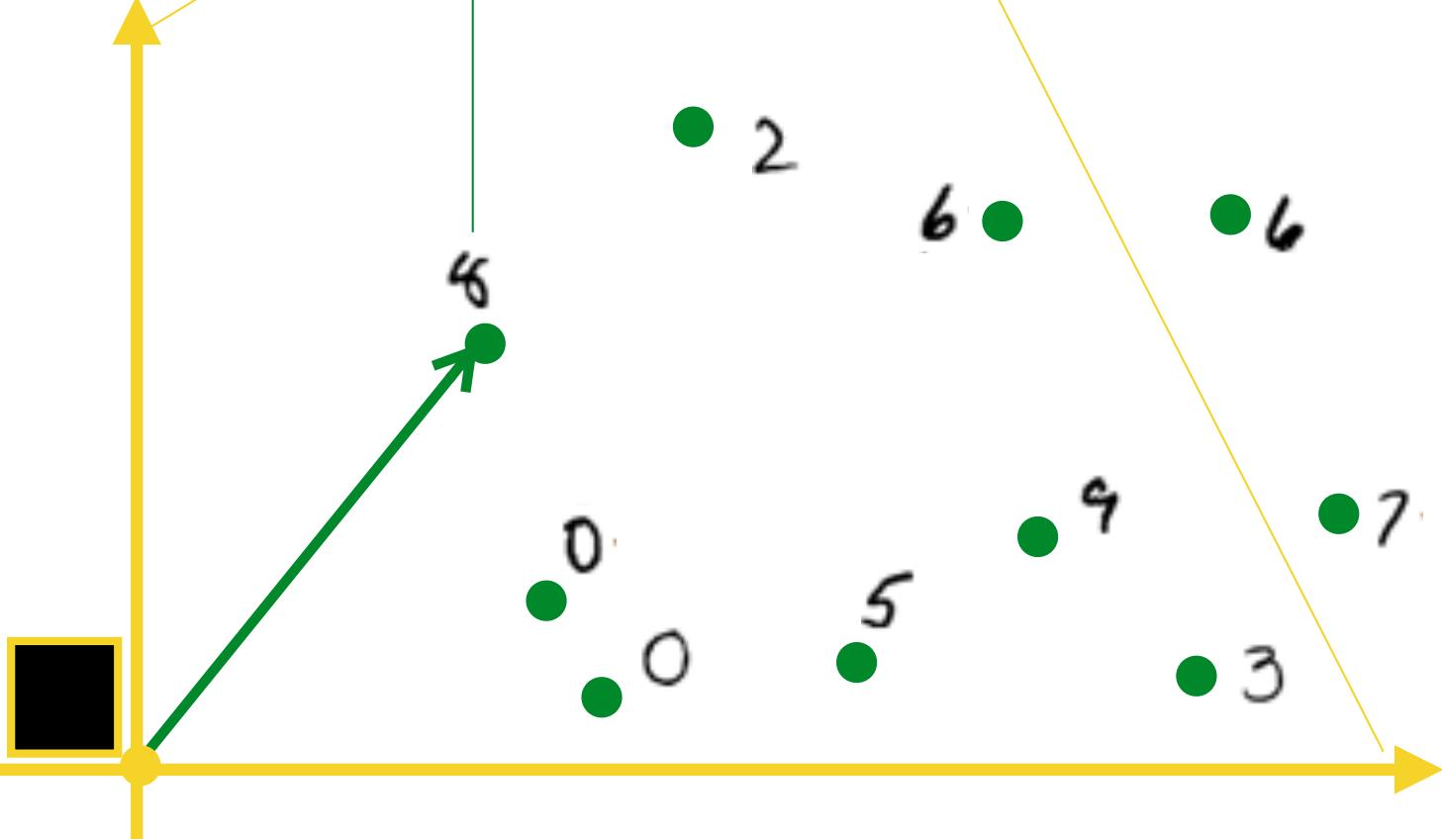
| | |
|-------|---|
| y_1 | 2 |
| y_2 | 5 |
| | 0 |
| | 3 |
| | 8 |
| | 0 |
| | 9 |
| | 6 |
| | 7 |
| y_n | 6 |

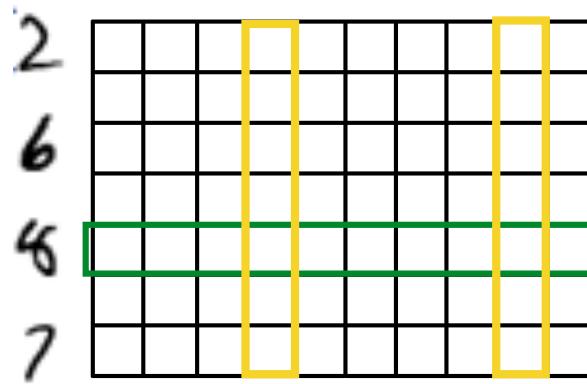
n - # instances p - # features

| | | | | | | | | |
|---|--|--|--|--|--|--|--|--|
| 2 | | | | | | | | |
| 6 | | | | | | | | |
| 8 | | | | | | | | |
| 7 | | | | | | | | |

X
(n by p)

y
(length n)



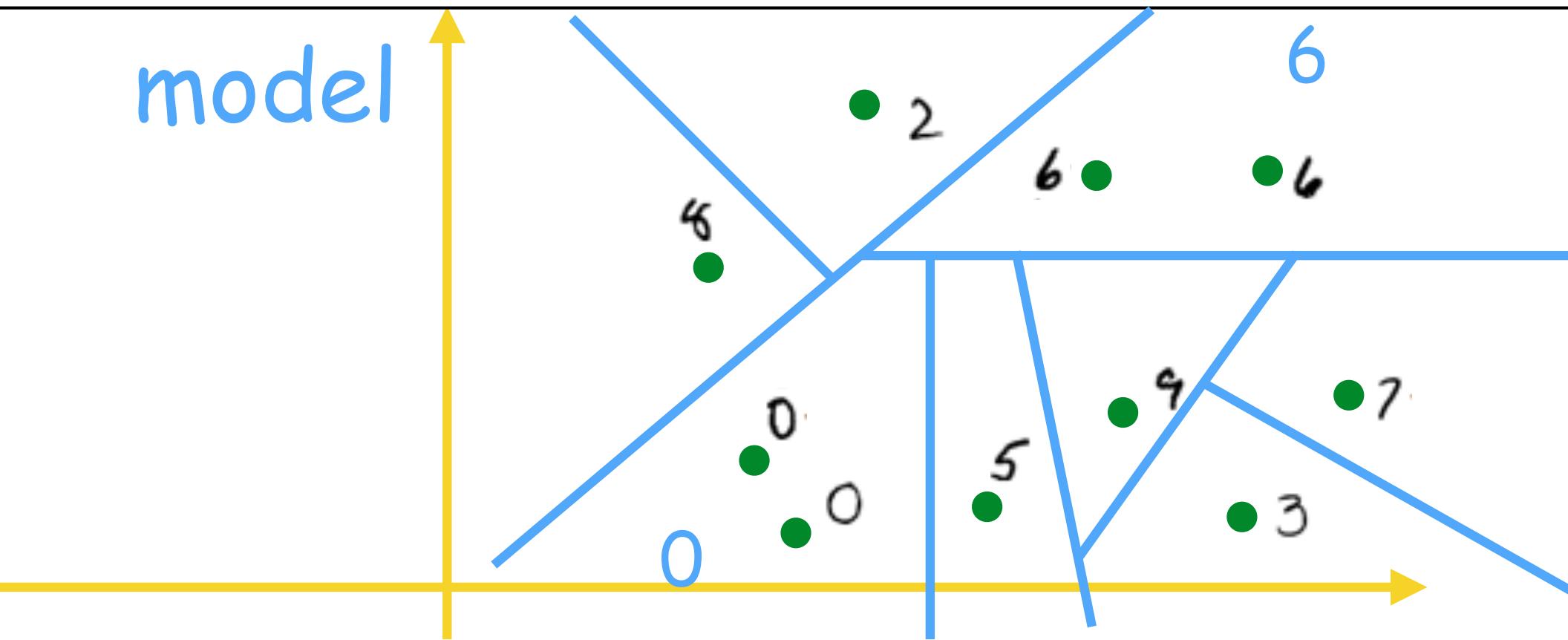


X
(n by p)

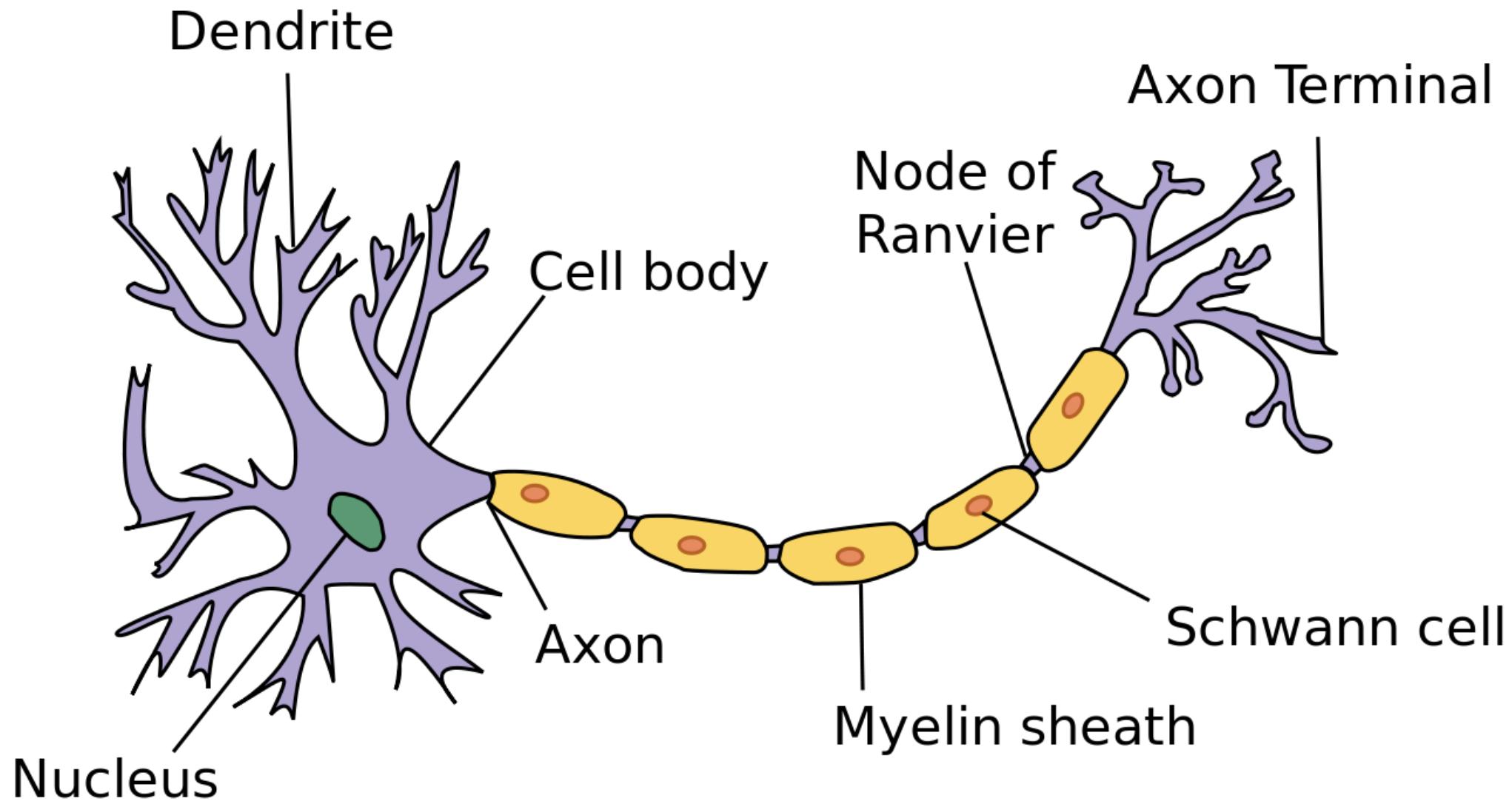
y

(length n)

model



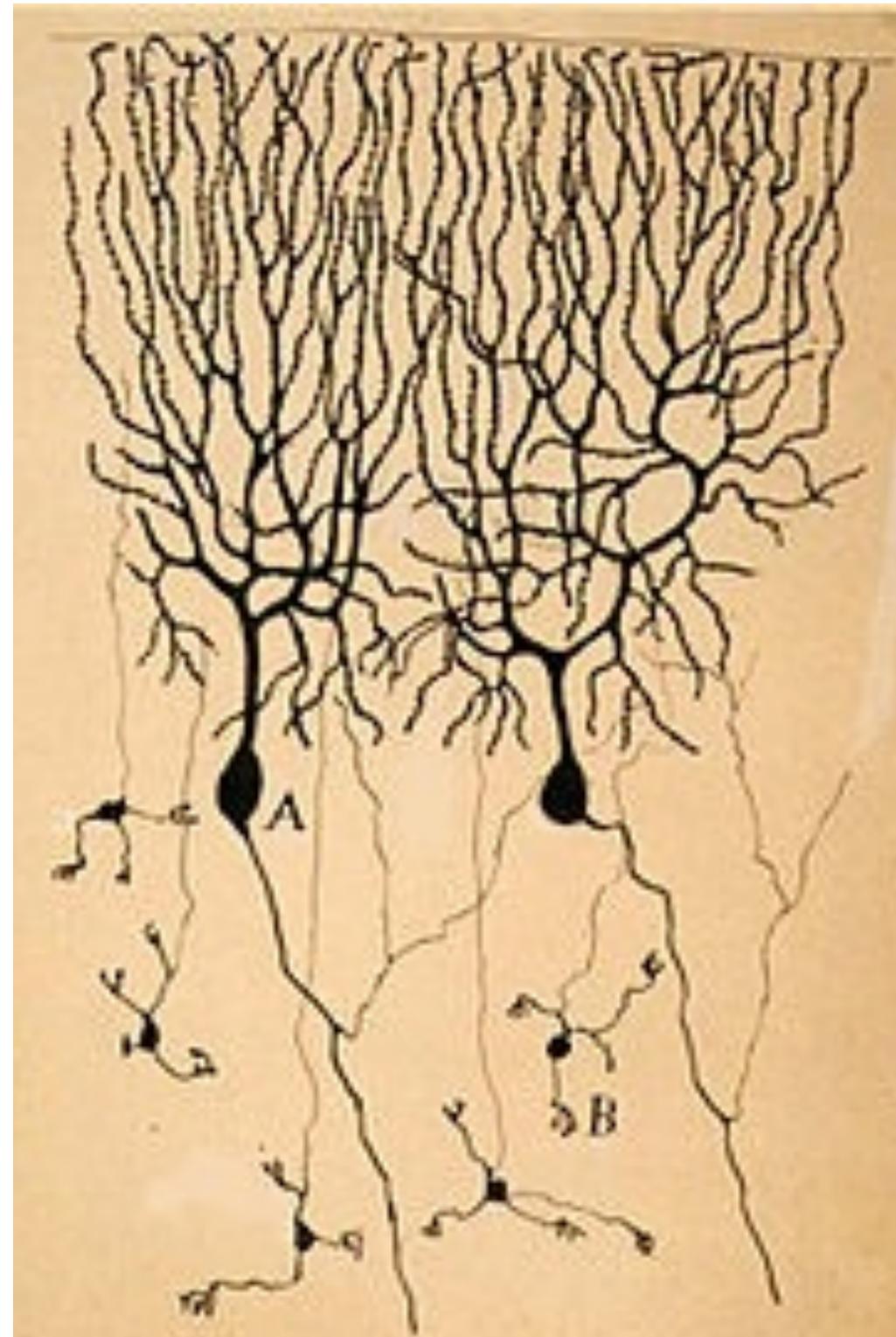
Neuron



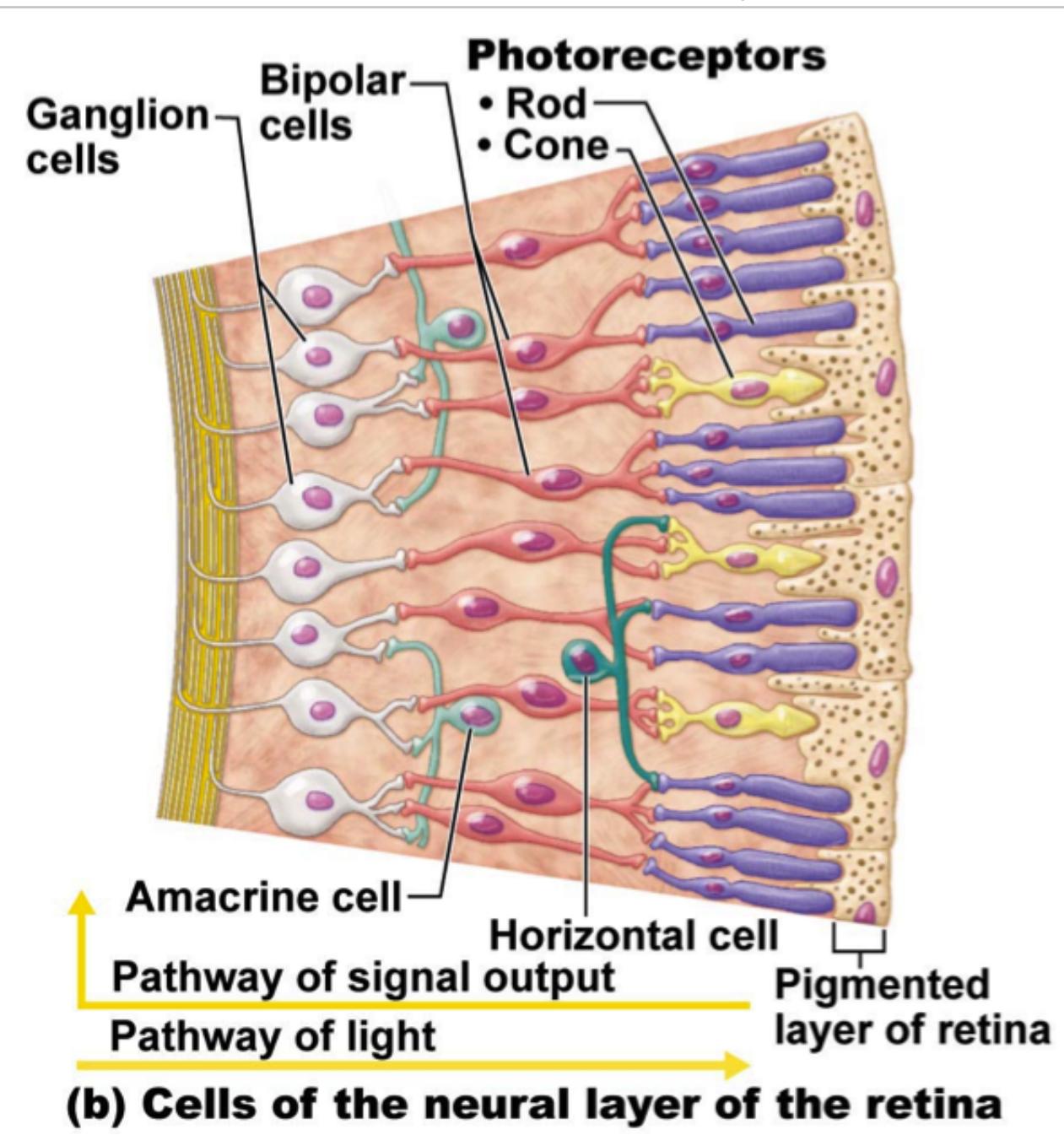
Neuron Layer



Santiago Ramón y Cajal
Nobel Prize in Medicine
1899

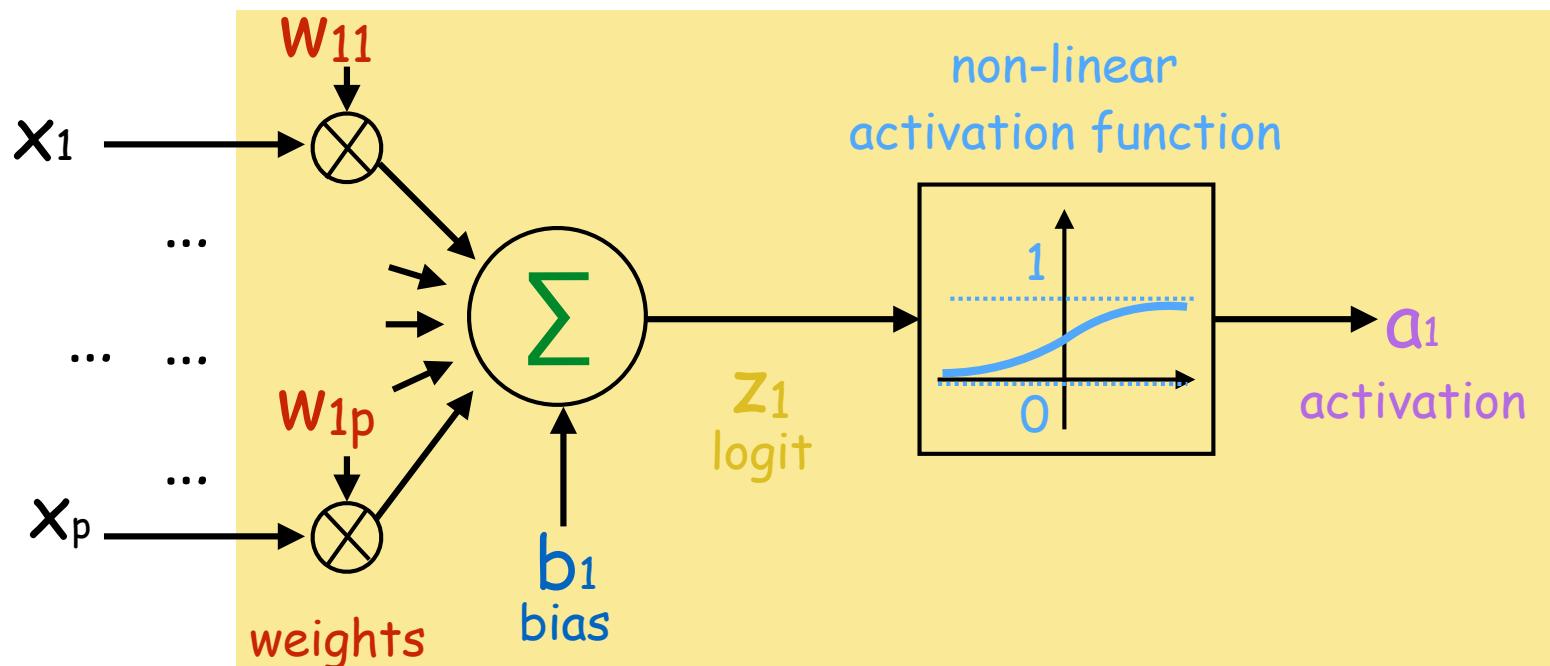


Neural Layer

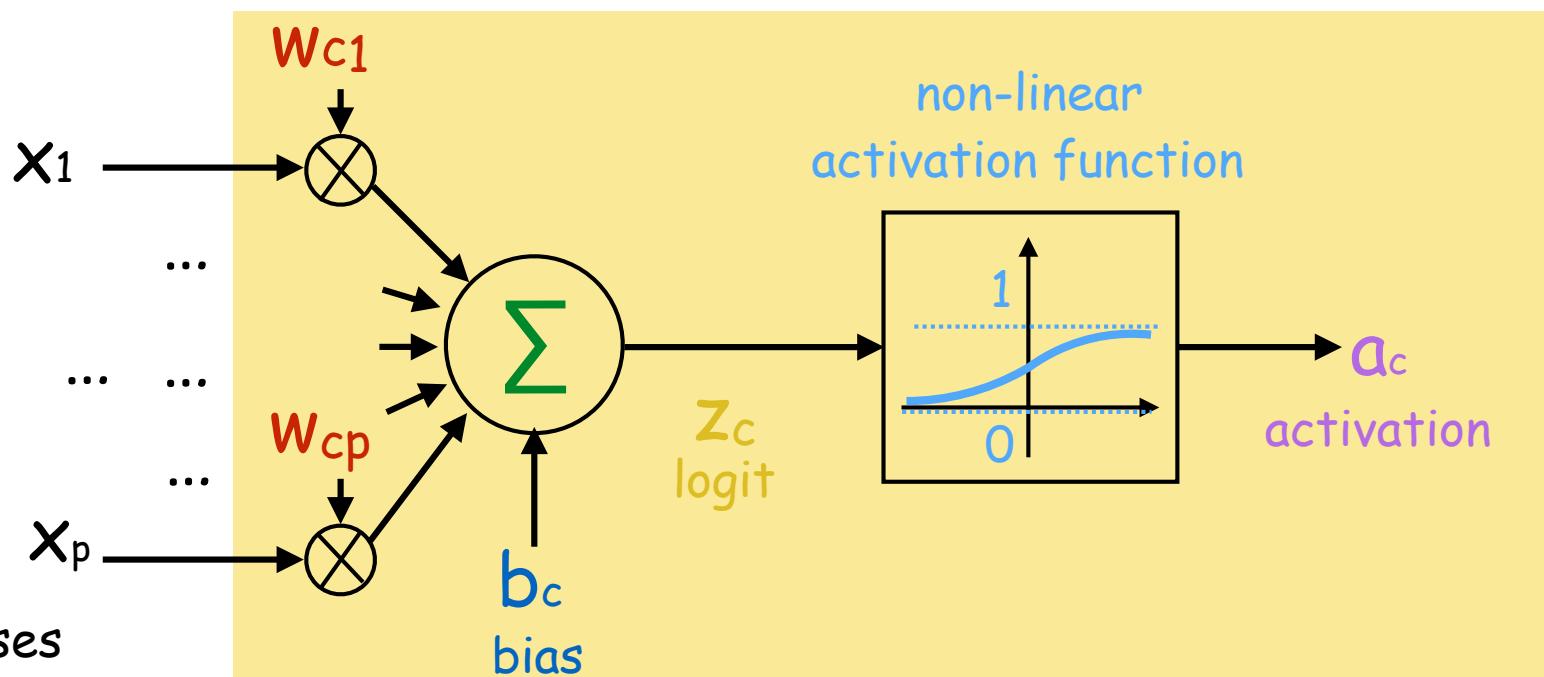


Demo

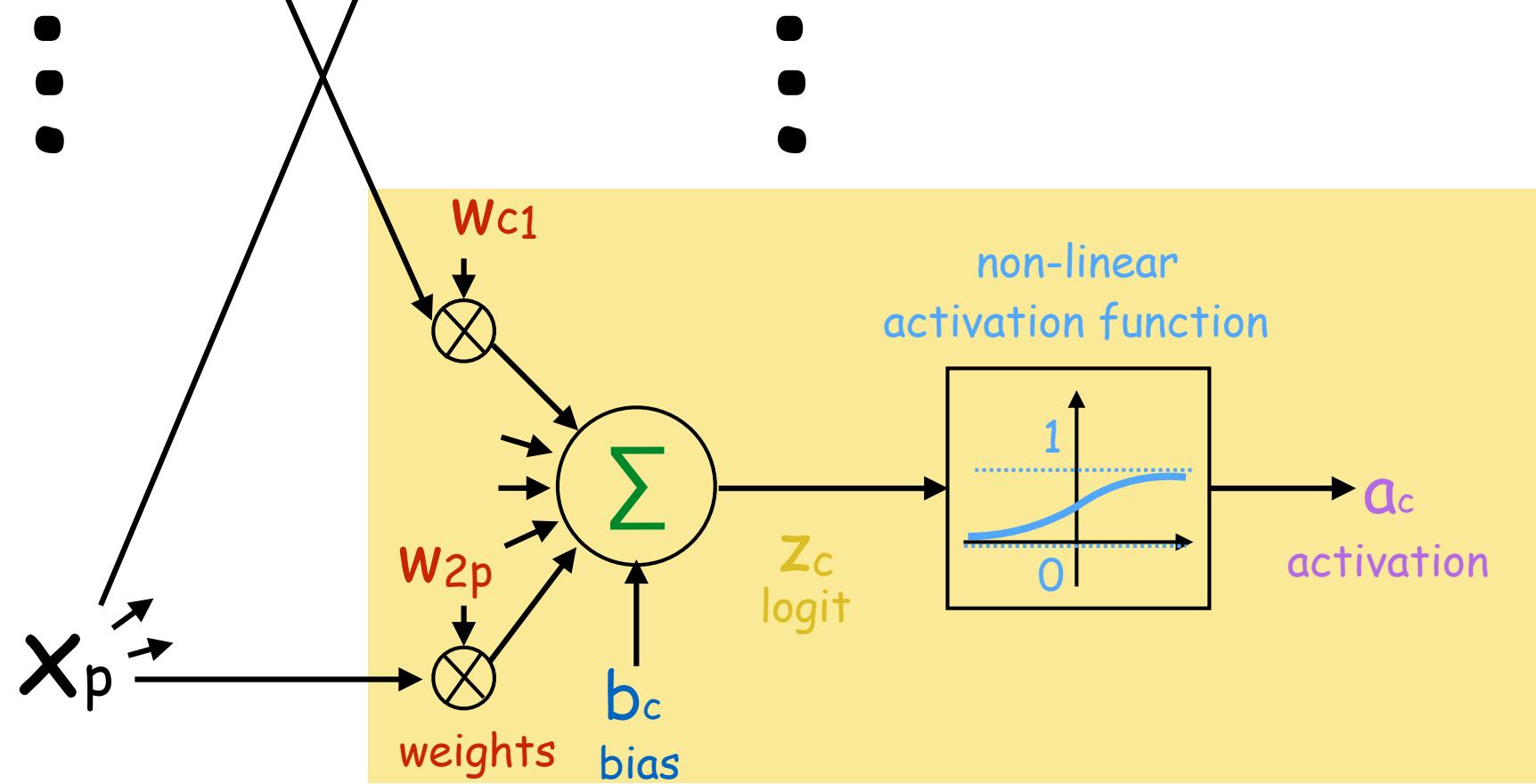
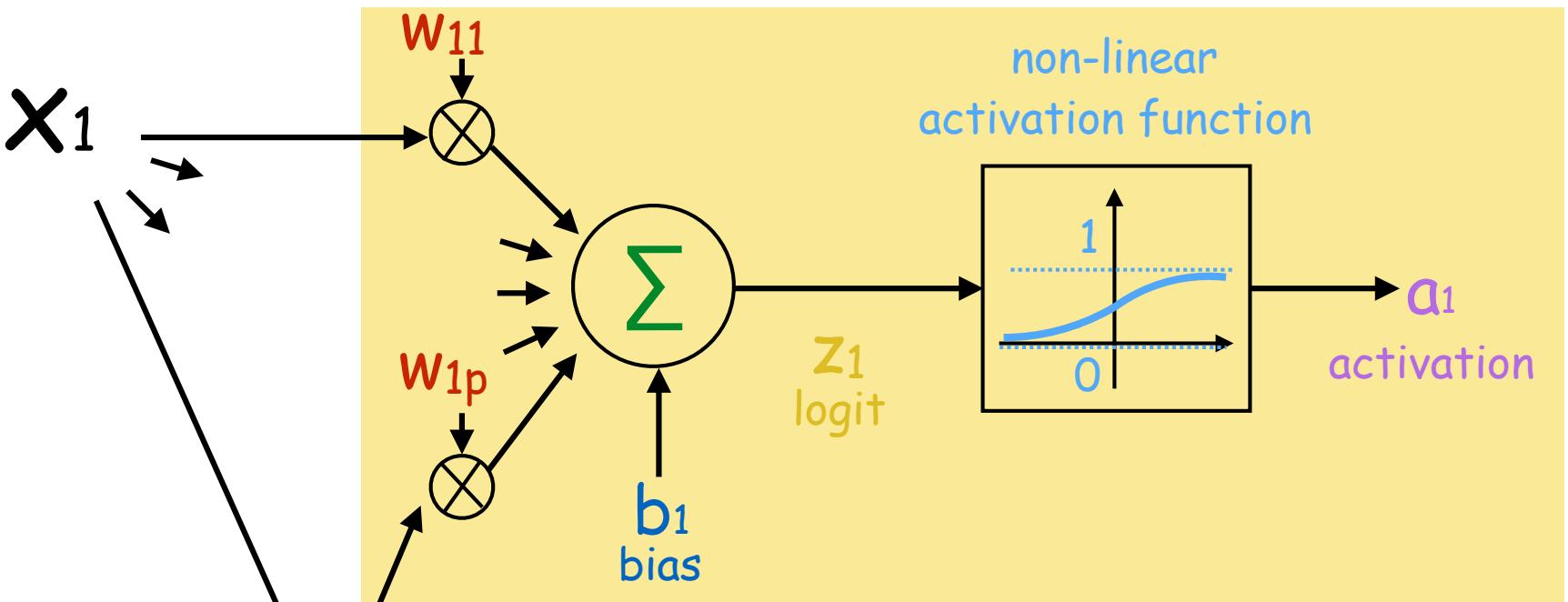
<https://youtu.be/3JQ3hYko51Y>

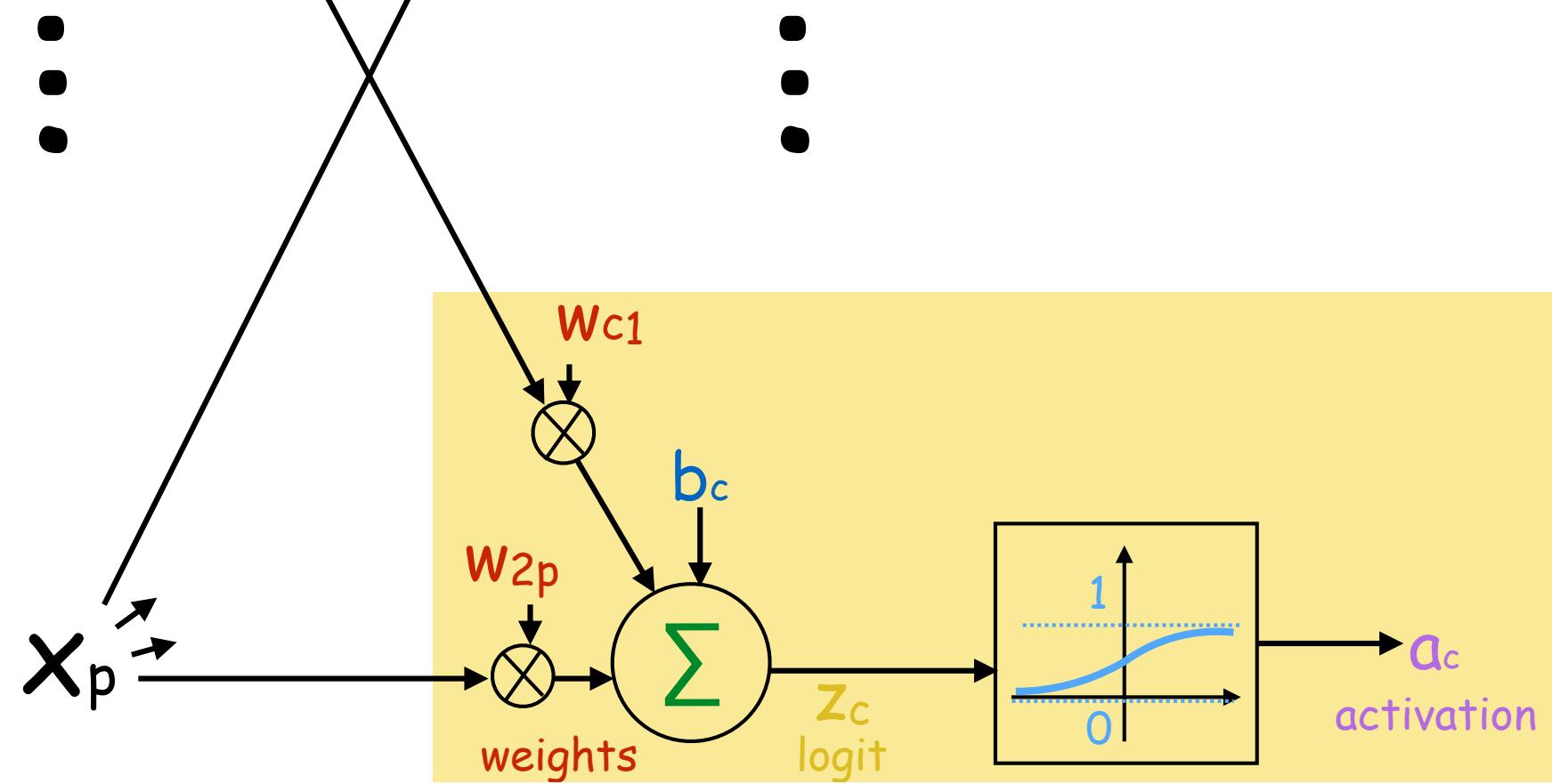
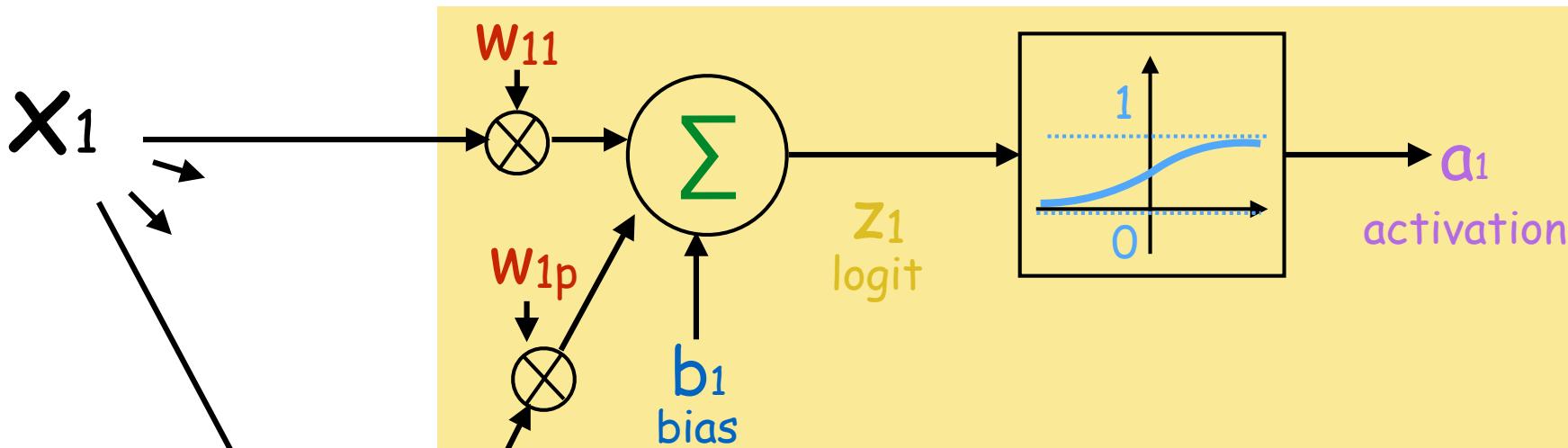


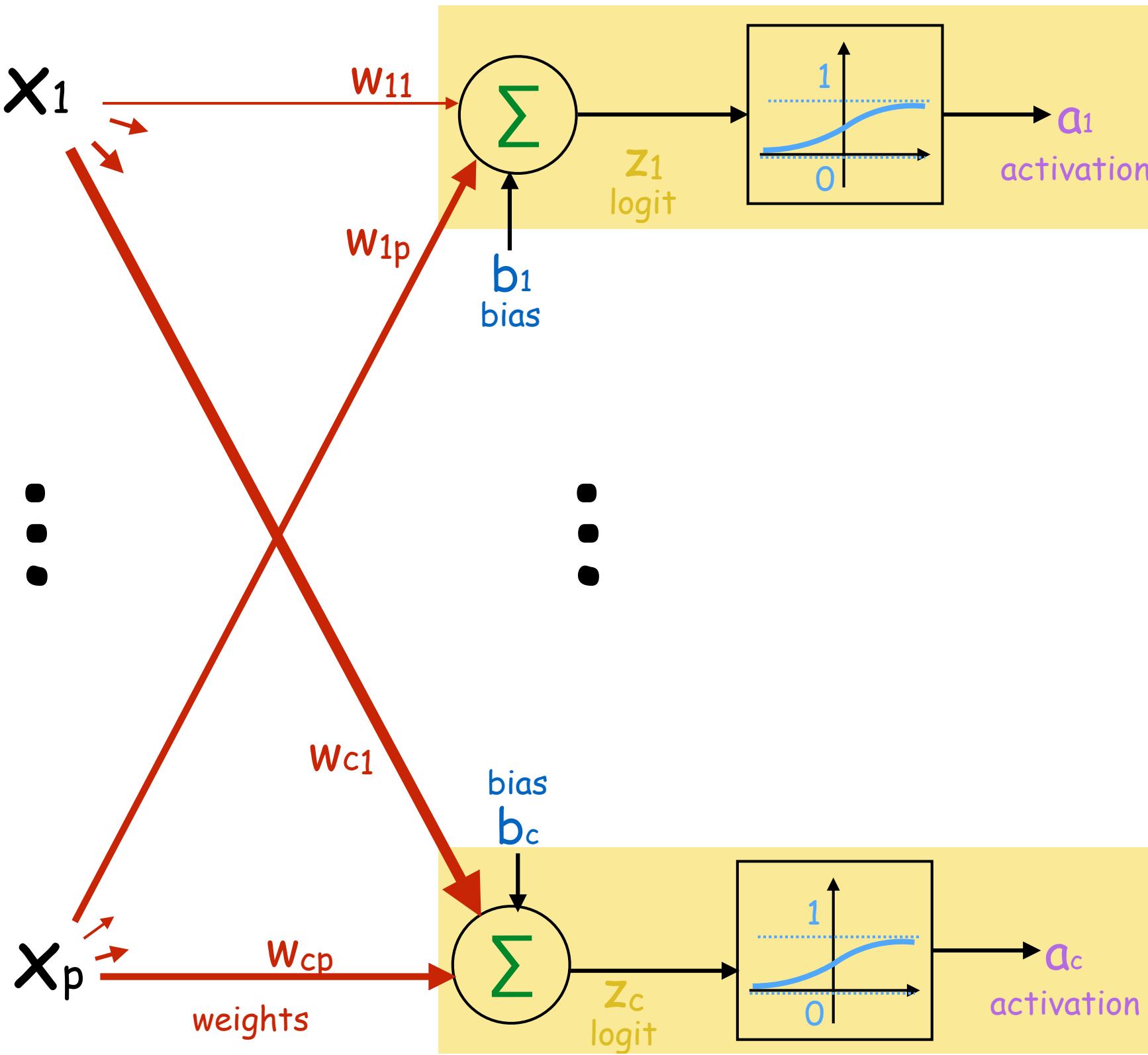
⋮

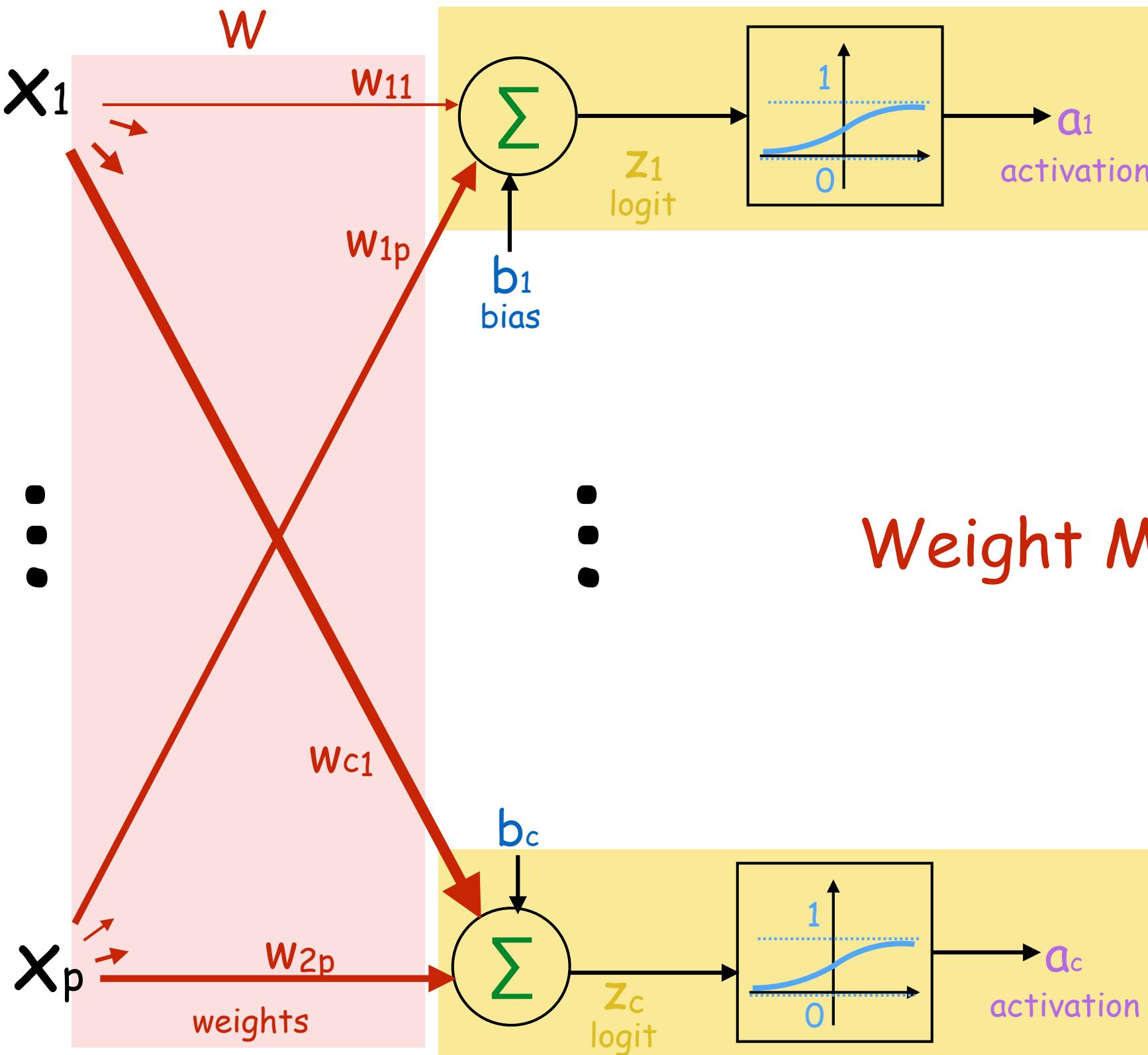


c : # classes





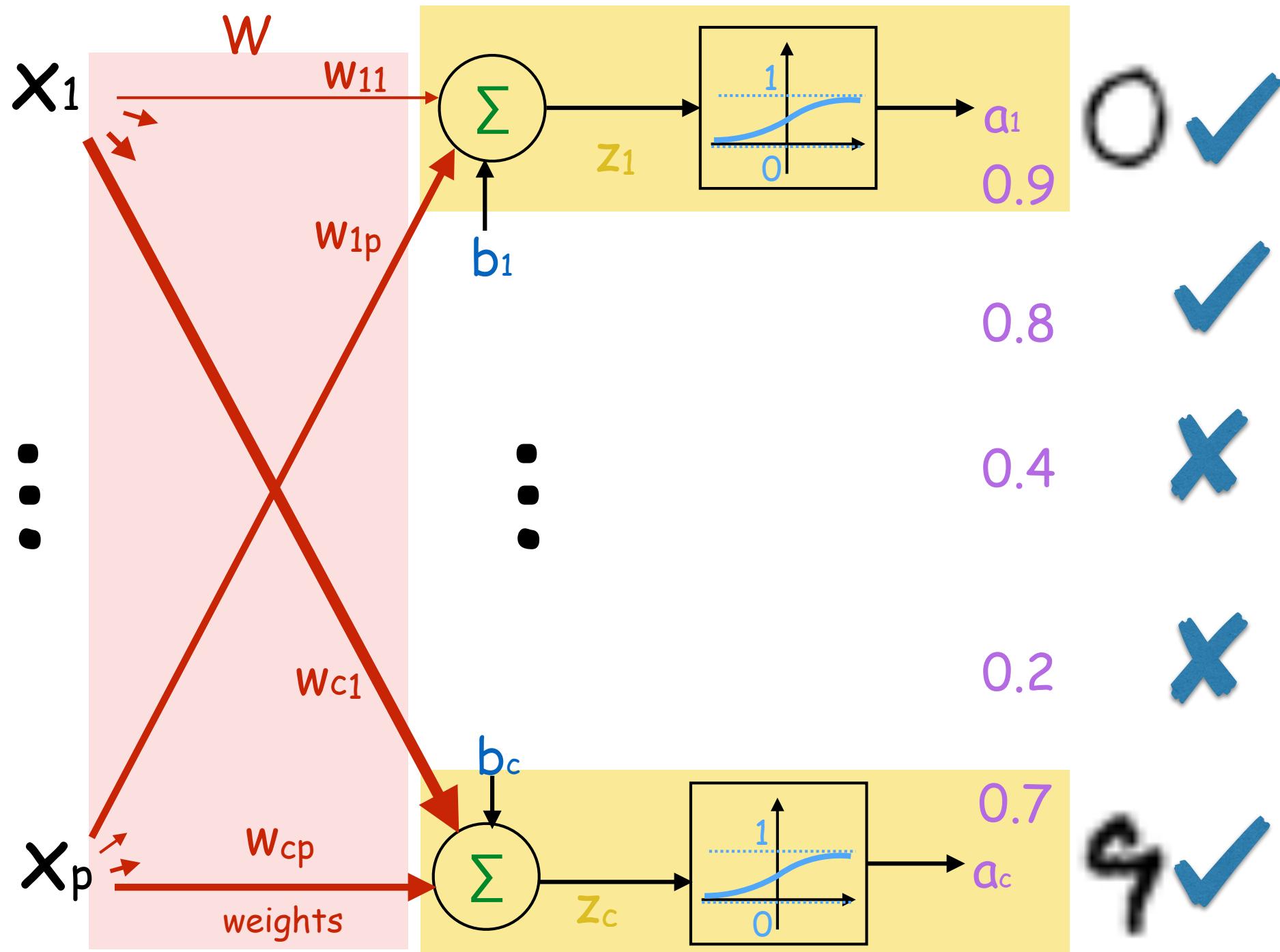




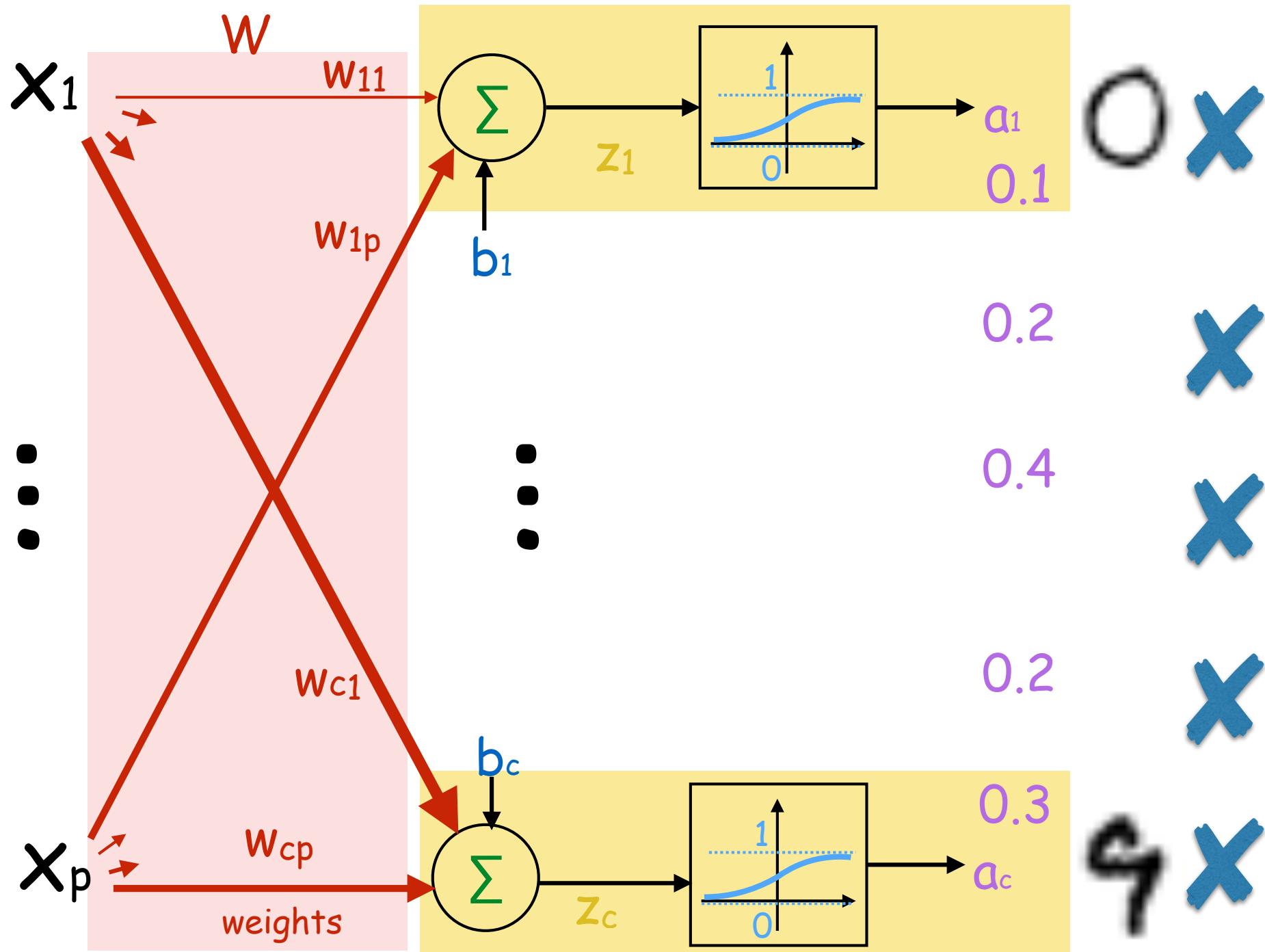
Weight Matrix



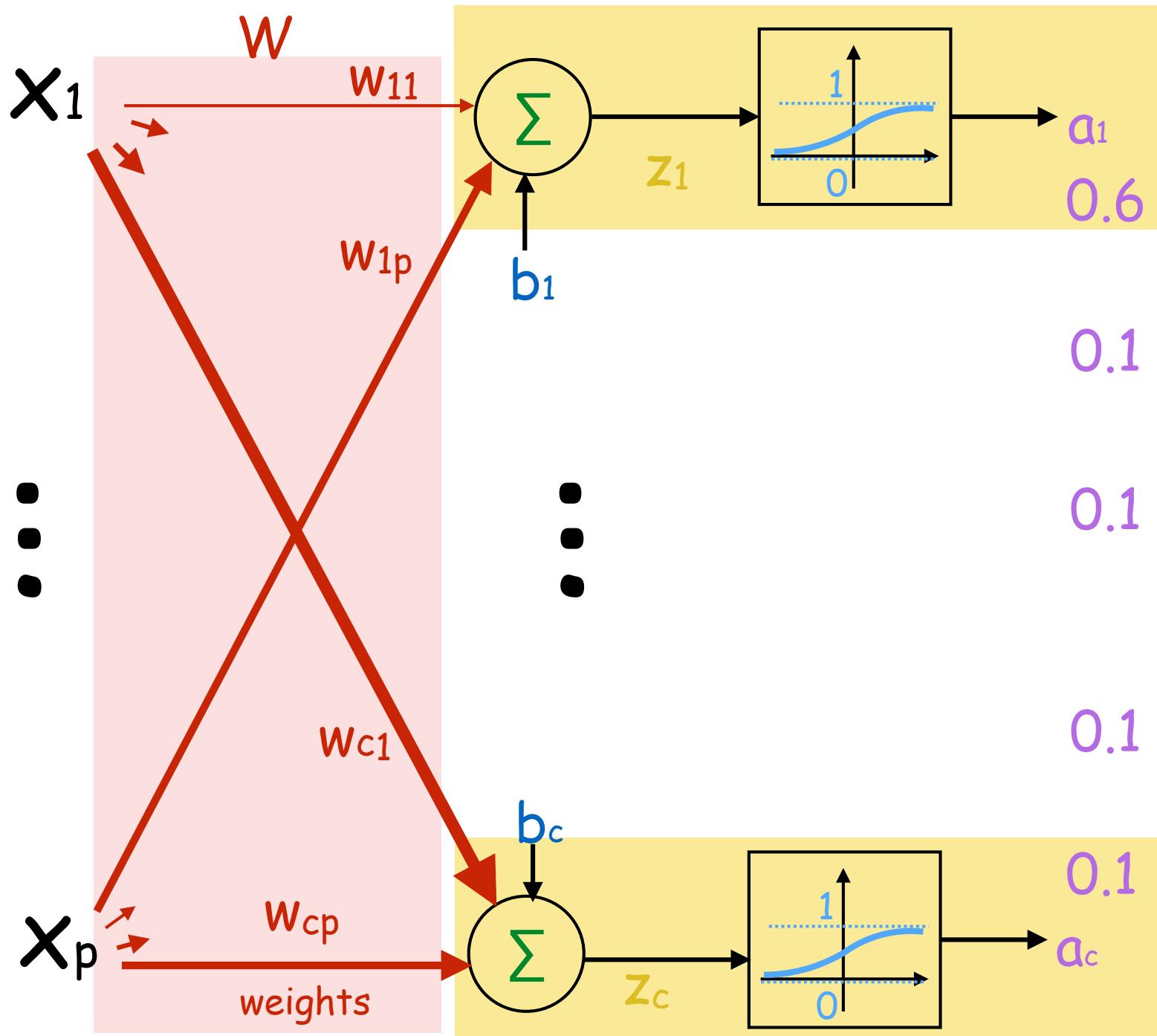
Independent Outputs



Independent Outputs



The outputs we need



0 ✓

0.1 ✗

0.1 ✗

0.1 ✗

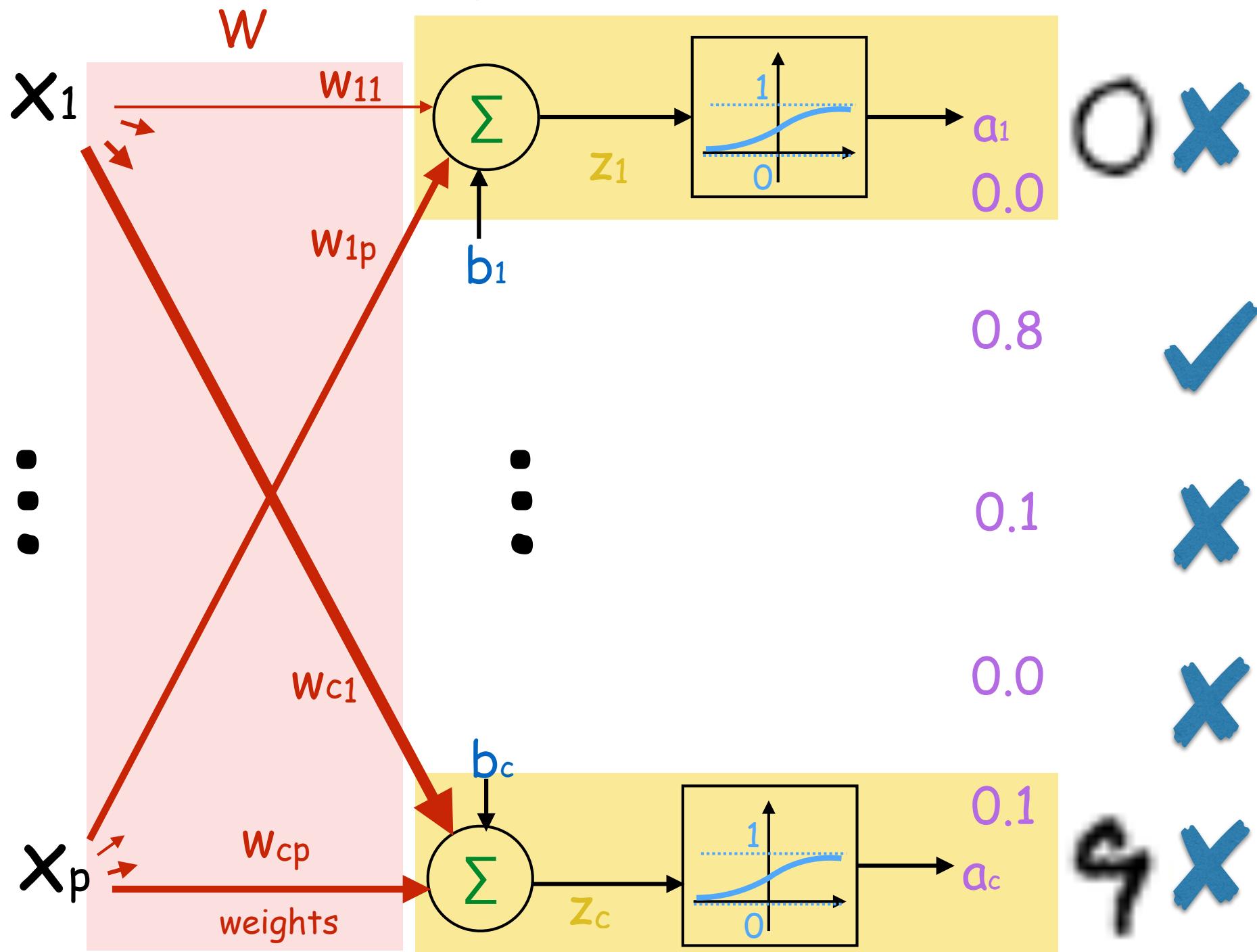
0.1 ✗

0.1 ✗

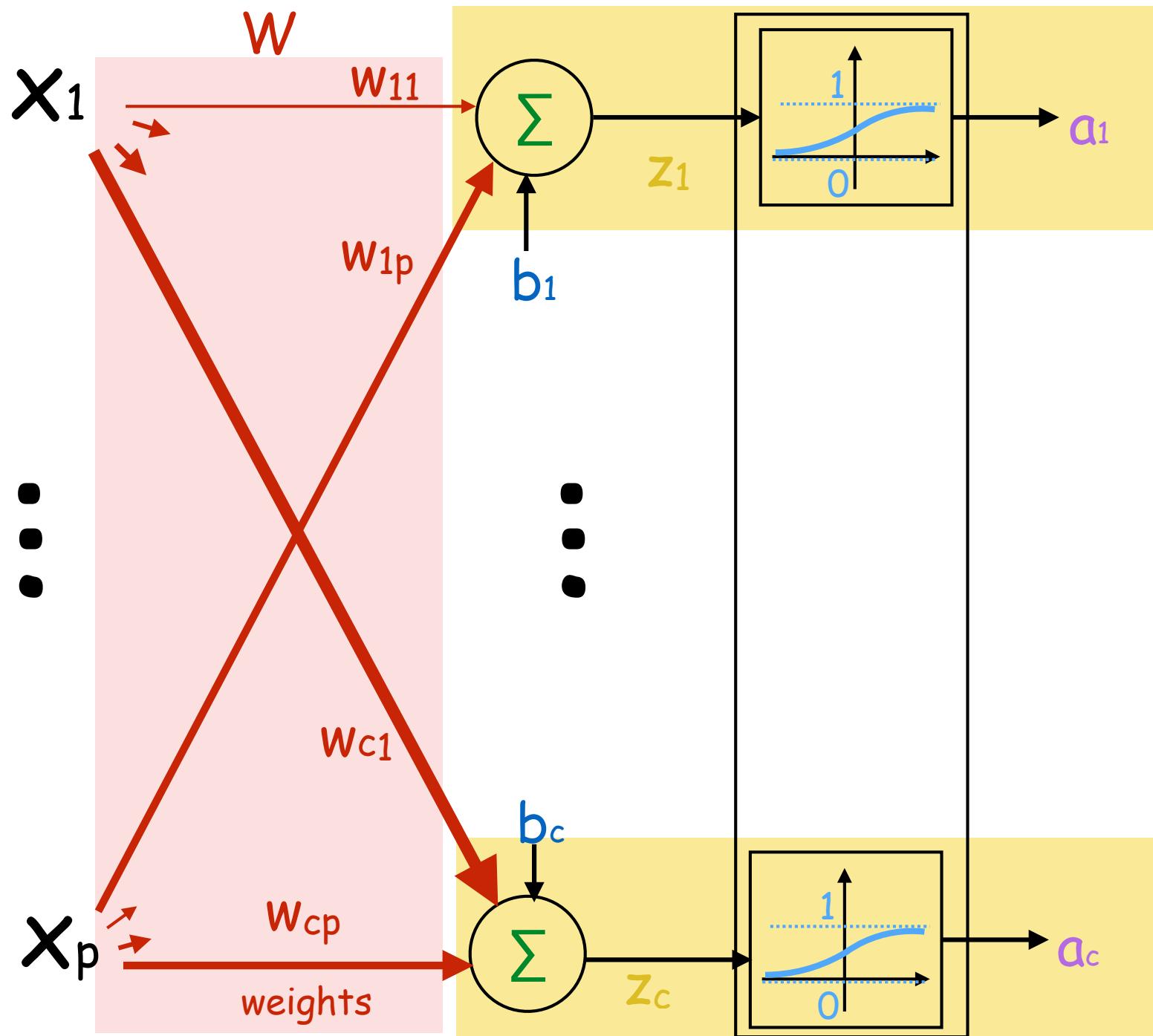
0.1 ✗

0.1 ✗ ✗

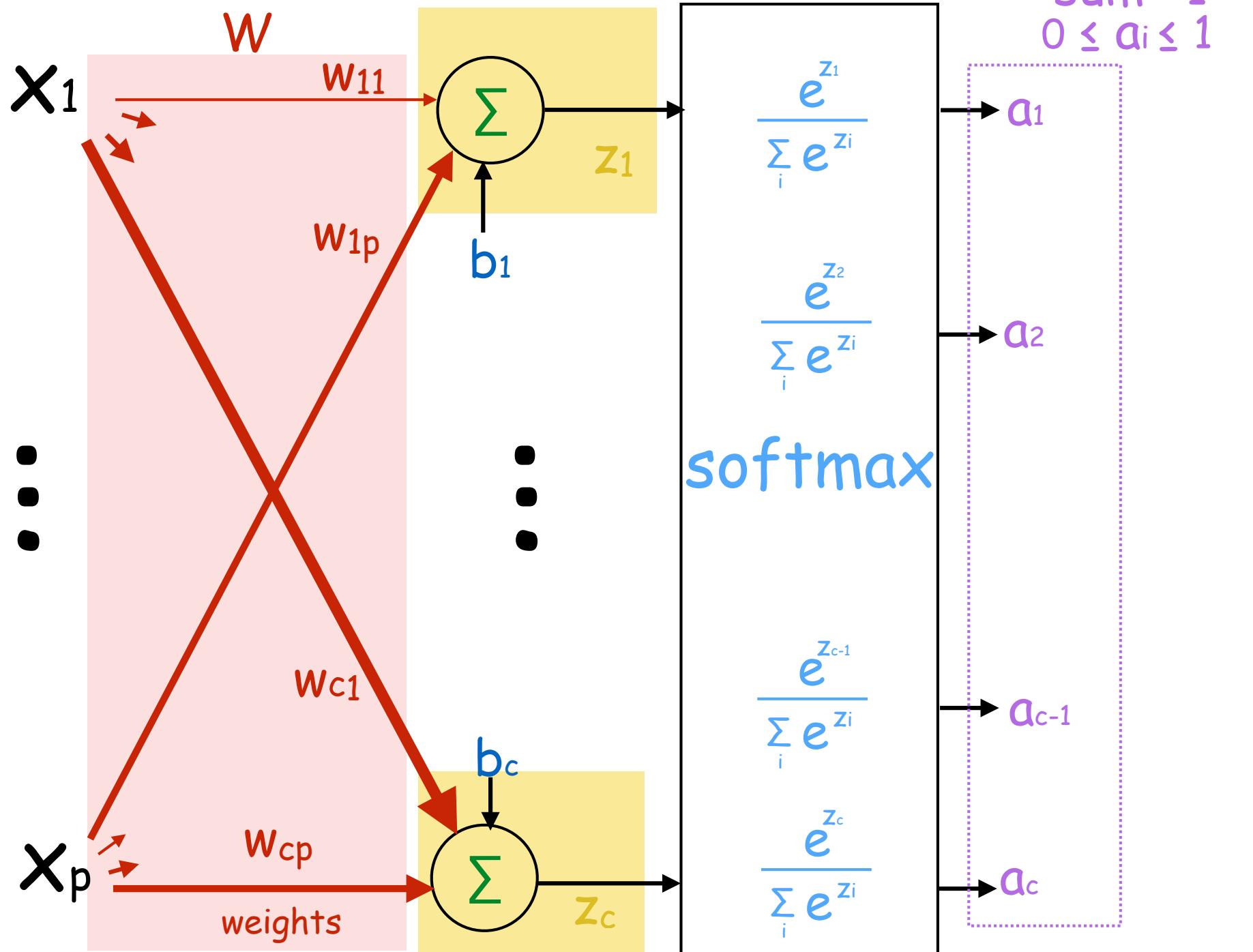
The outputs we need



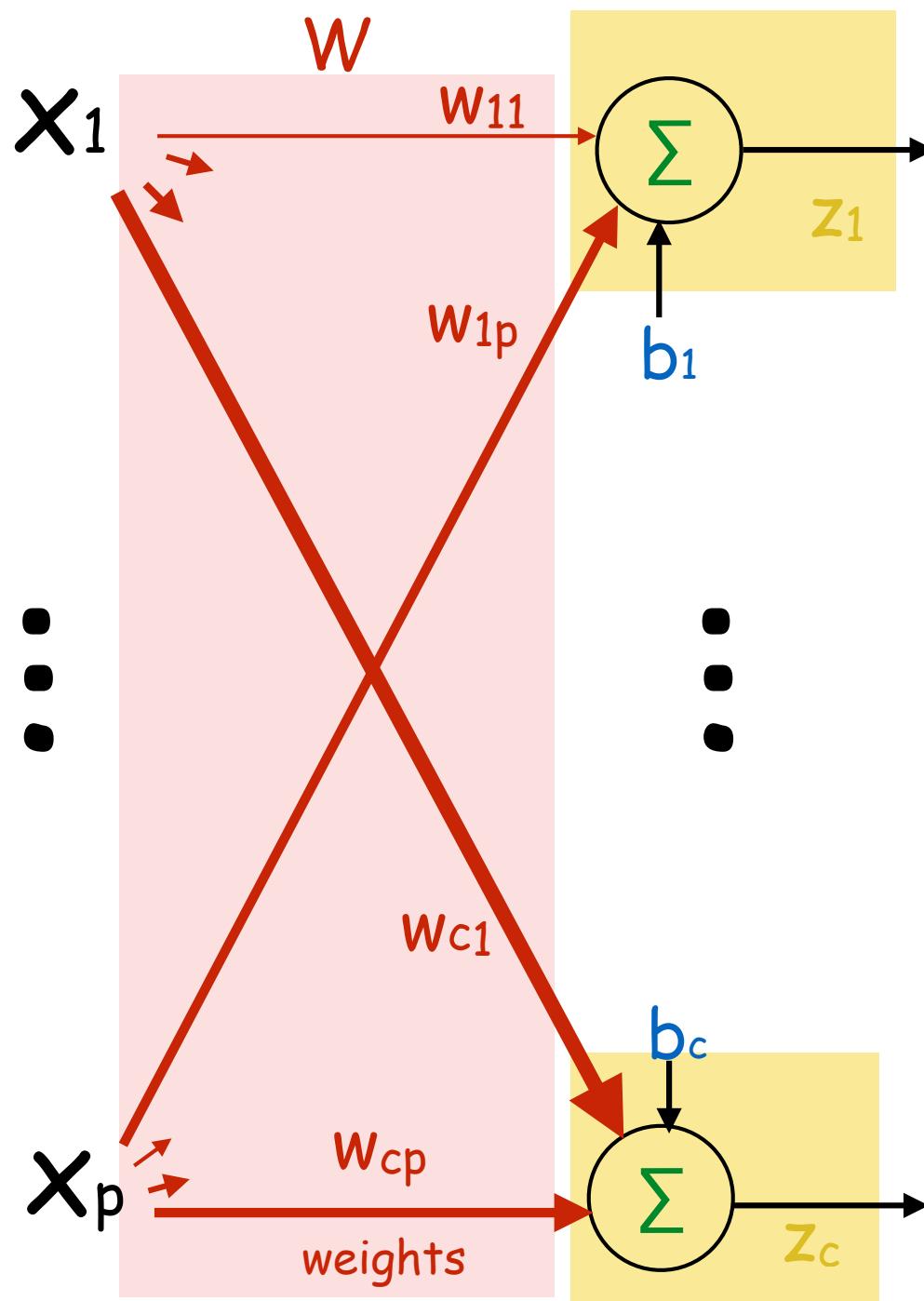
Coordinated Outputs



Softmax Activation



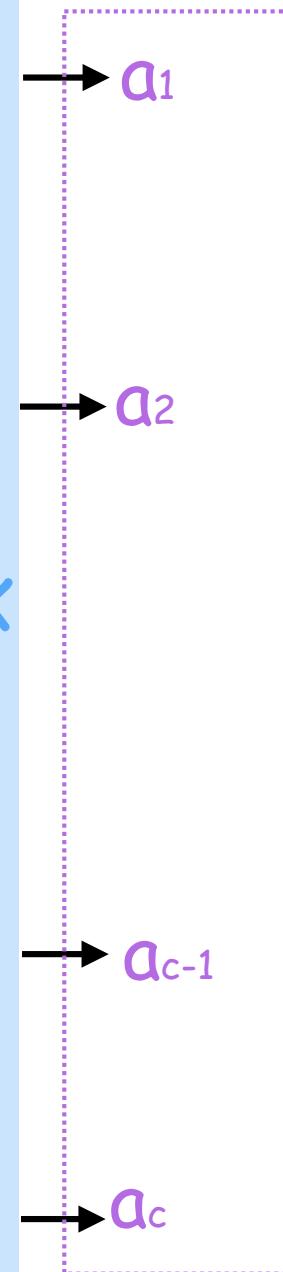
Softmax Activation



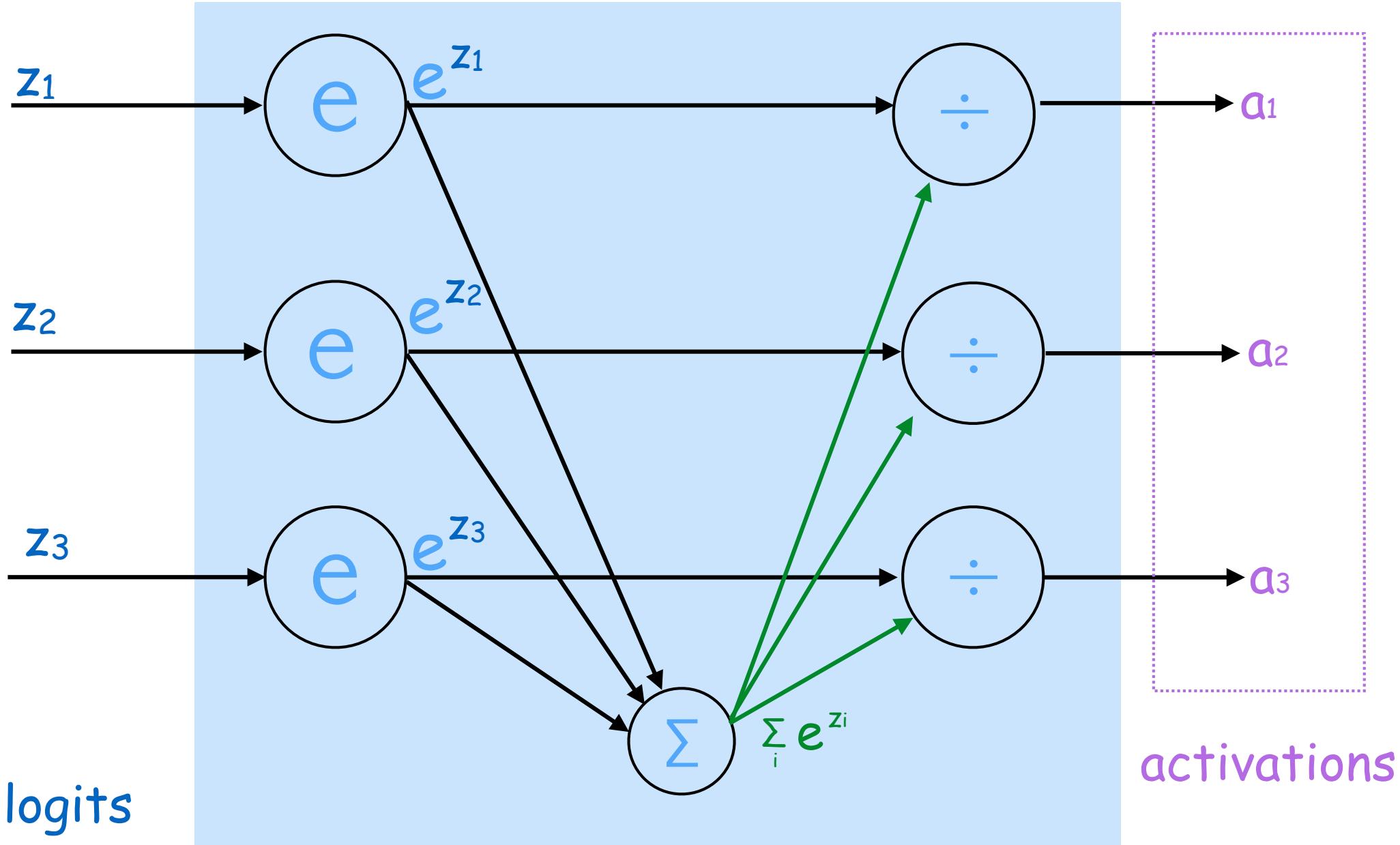
softmax

$$\frac{e^{z_1}}{\sum_i e^{z_i}}$$
$$\frac{e^{z_2}}{\sum_i e^{z_i}}$$
$$\frac{e^{z_{c-1}}}{\sum_i e^{z_i}}$$
$$\frac{e^{z_c}}{\sum_i e^{z_i}}$$

Probabilities
sum = 1
 $0 \leq a_i \leq 1$

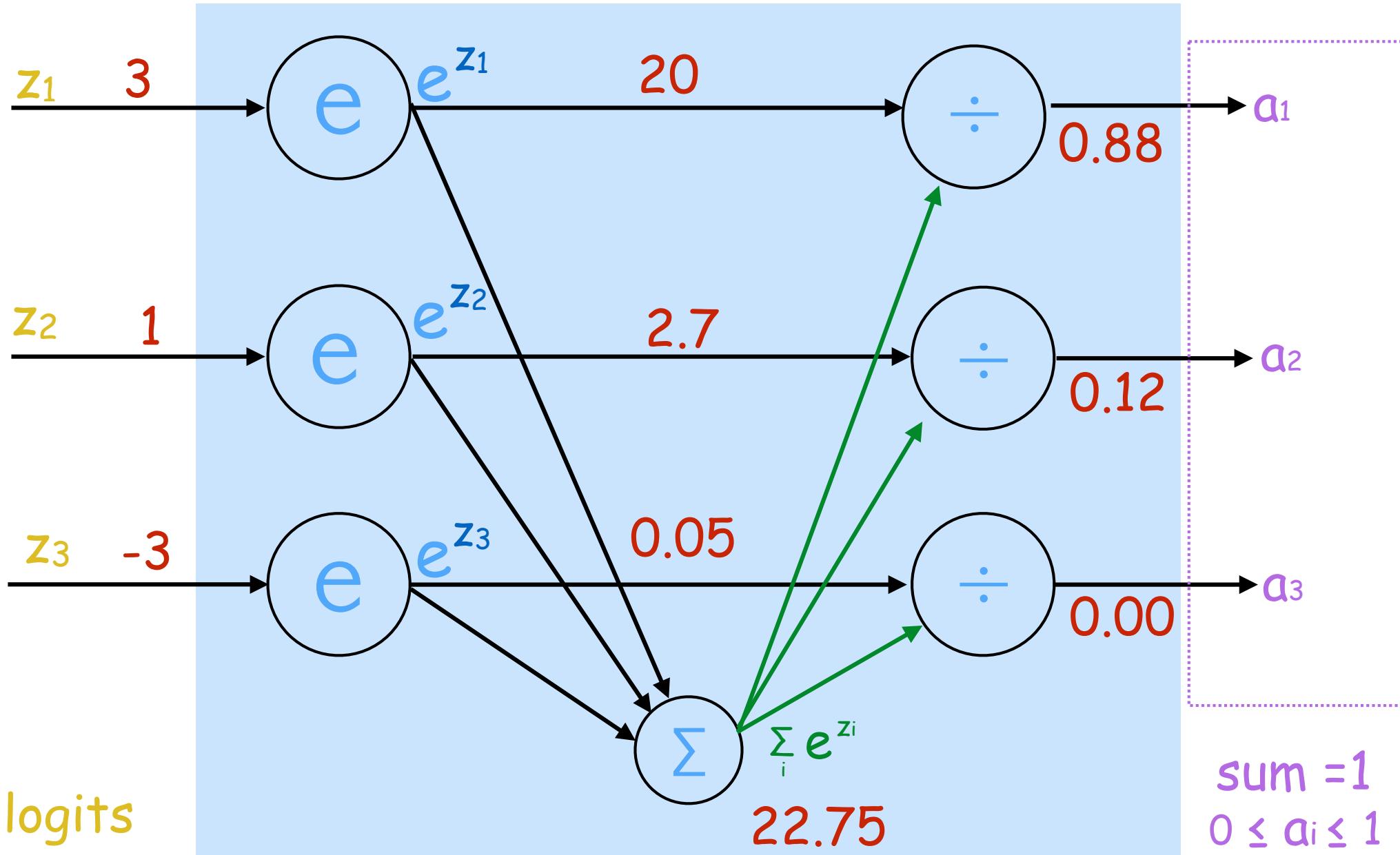


Softmax Activation



Softmax

Softmax Activation



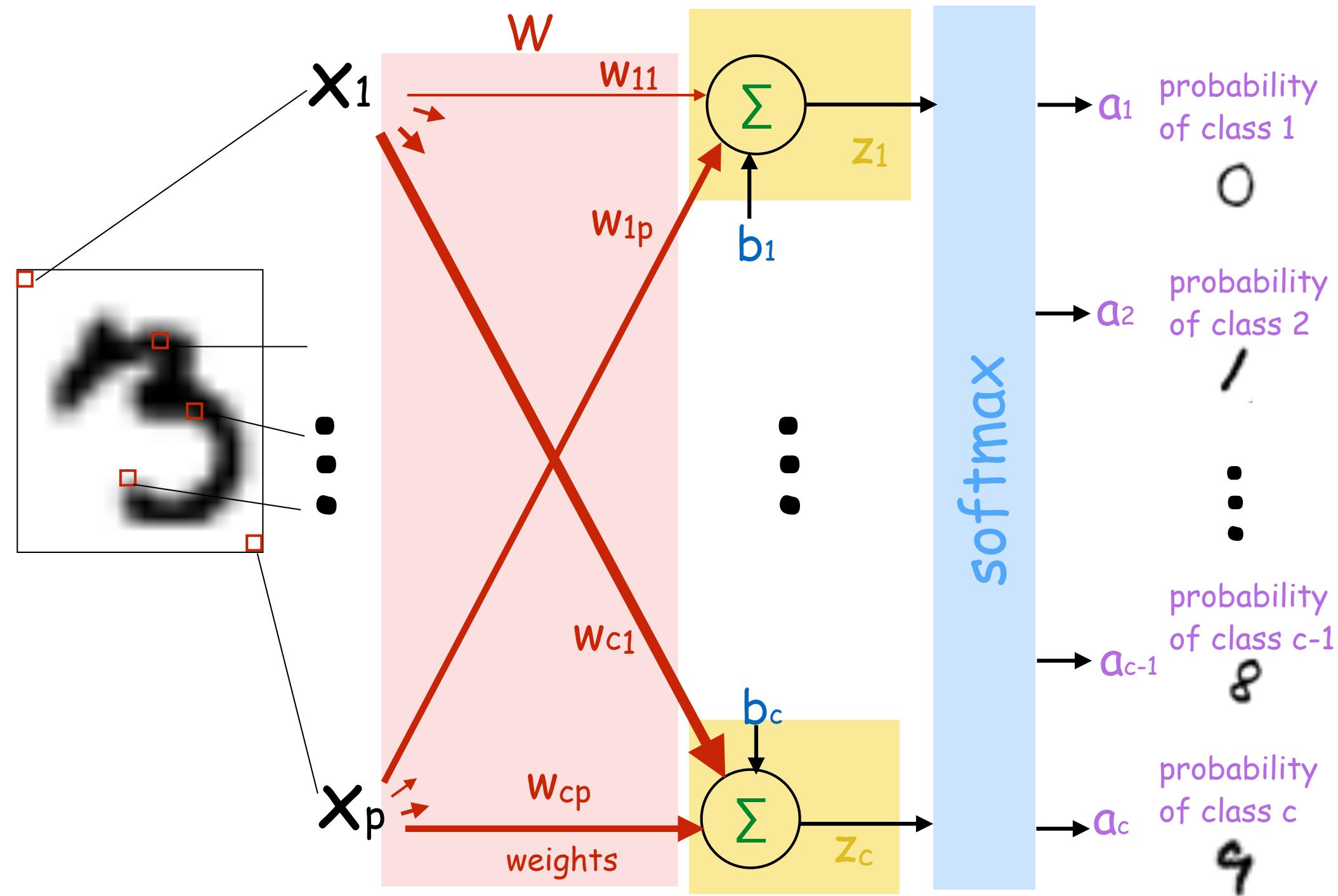
Softmax

Probabilities

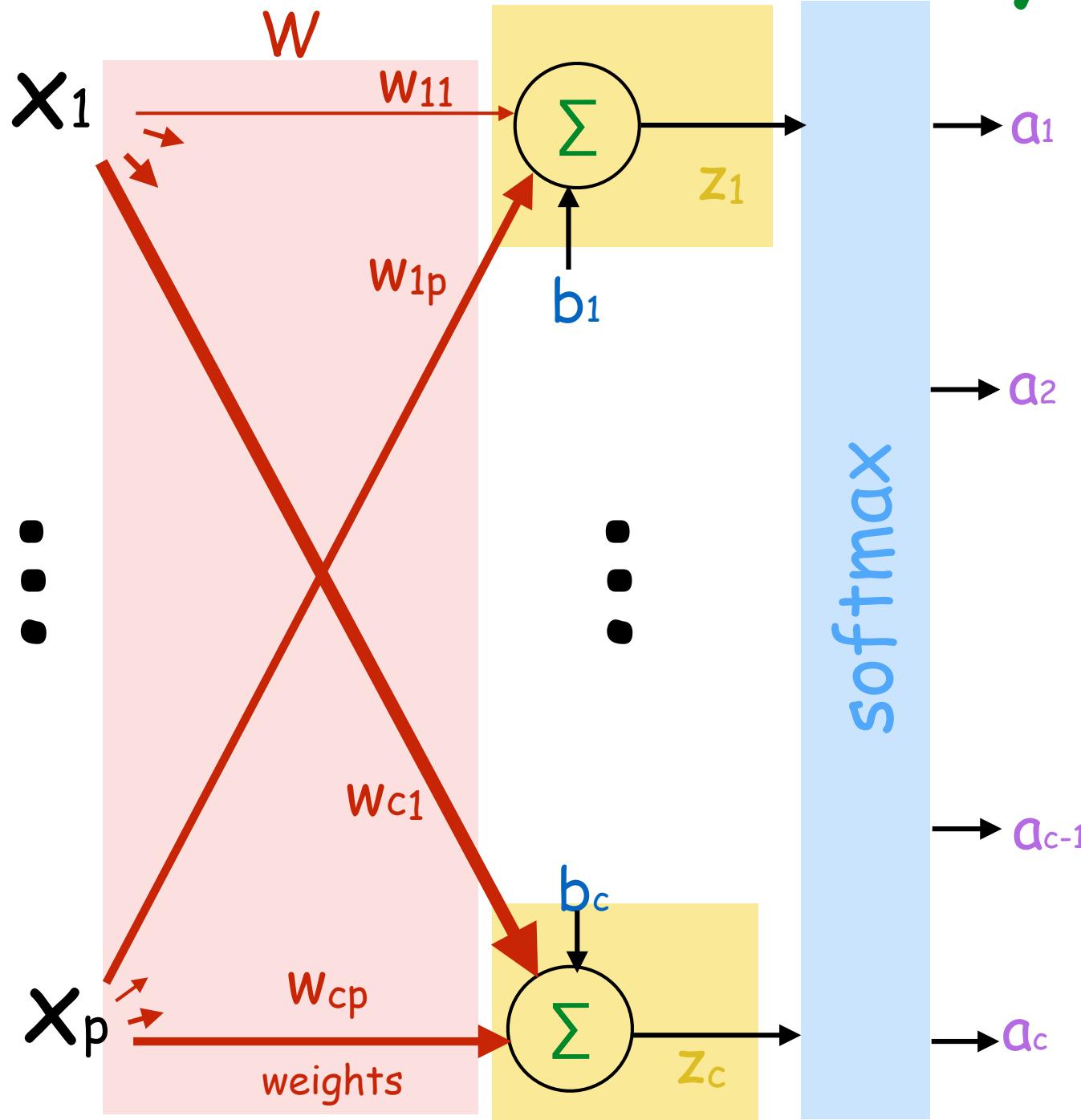
inputs

logits

activations

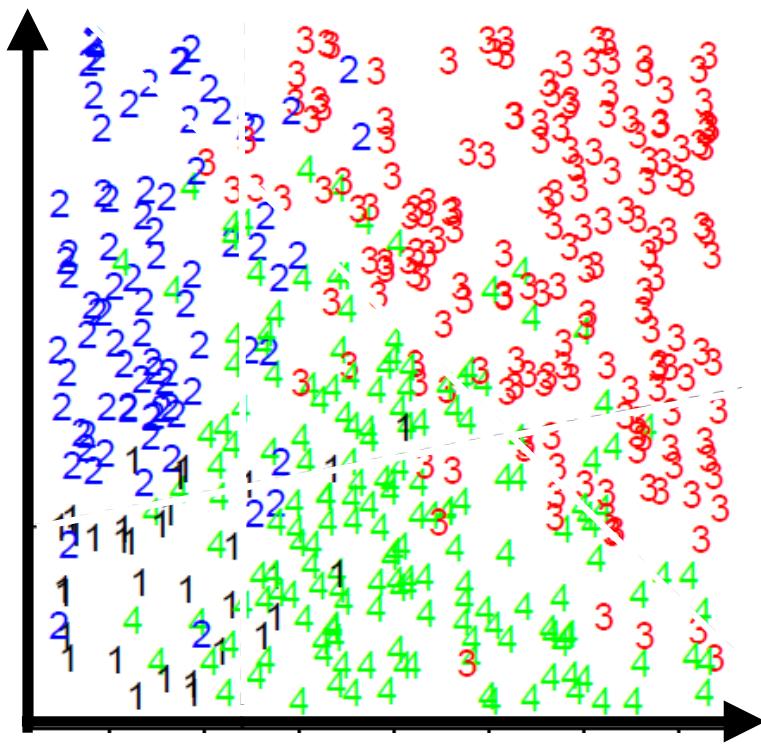


Parameters W, b



Training Model

training set

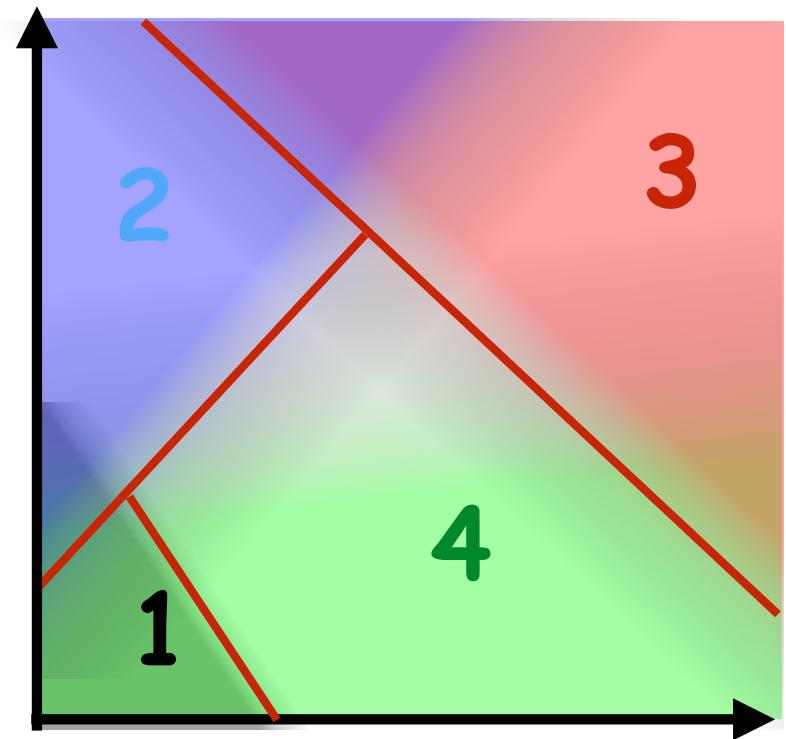
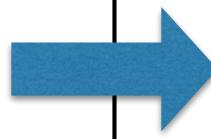


Feature Matrix X _(n by p)

Label Vector Y

(length n) (0, 1, 2, ... value)

learn parameters

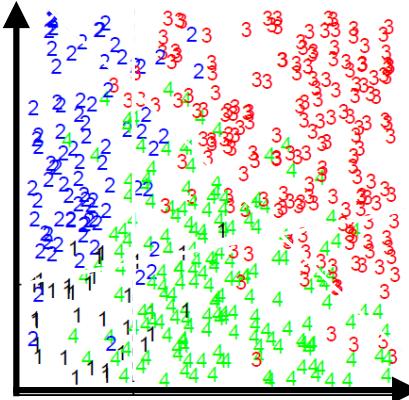


Weight Matrix W _(shape c by p)

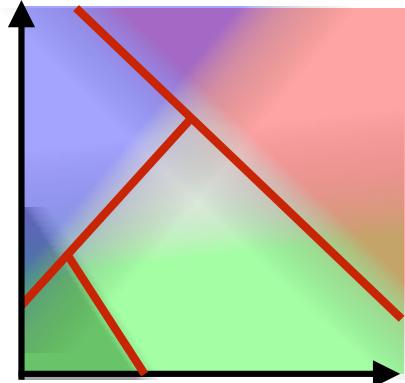
bias vector b

(length c)

negative log Likelihood



Label:
observed data



Weights w
biases b

Outputs:
probabilities

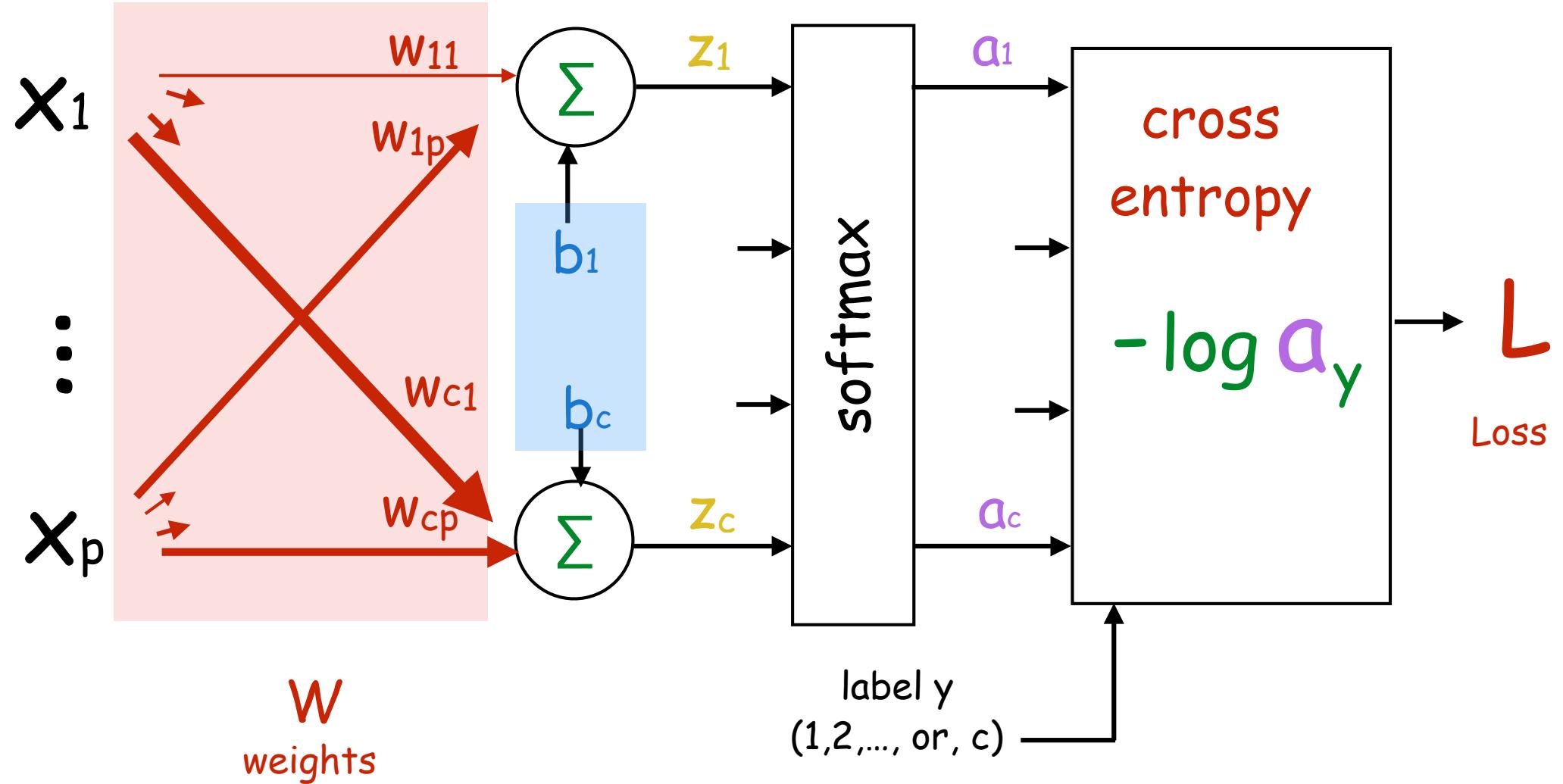
a_1 .4
 a_2 .2
 a_3 .1
 a_c .1

y 3
 y 1
 a_1 .4
 a_2 .2
 a_3 .1
 a_c .1

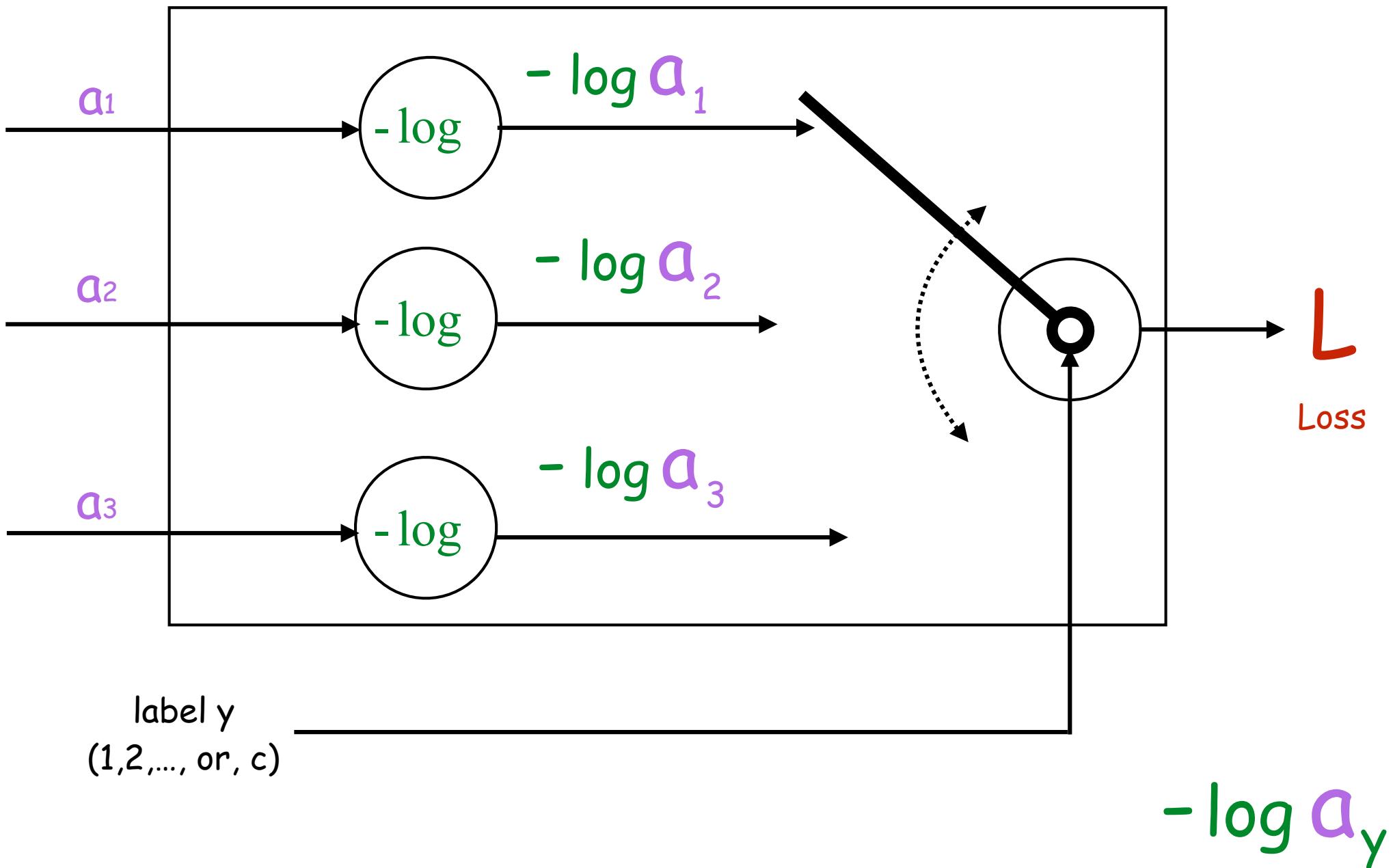
$$-\log \text{Likelihood} = -\log a_y$$

Multi-class cross entropy loss

Multi-Class Cross Entropy Loss

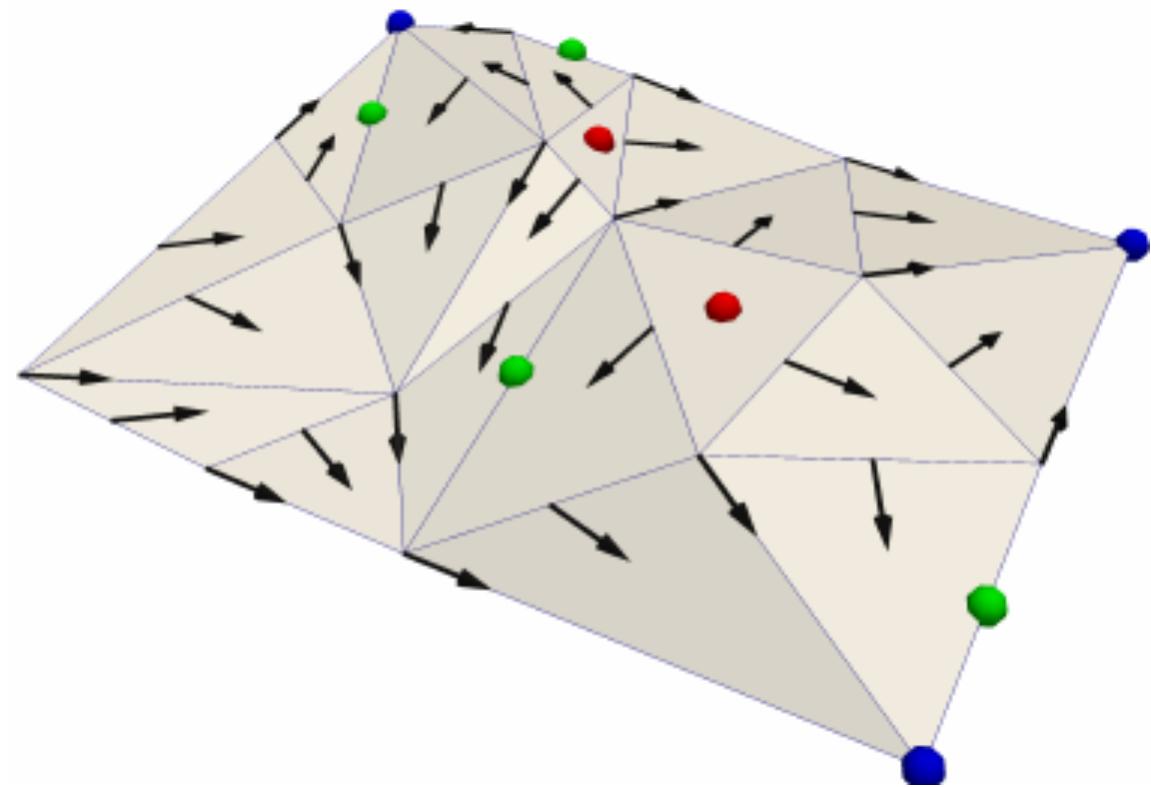
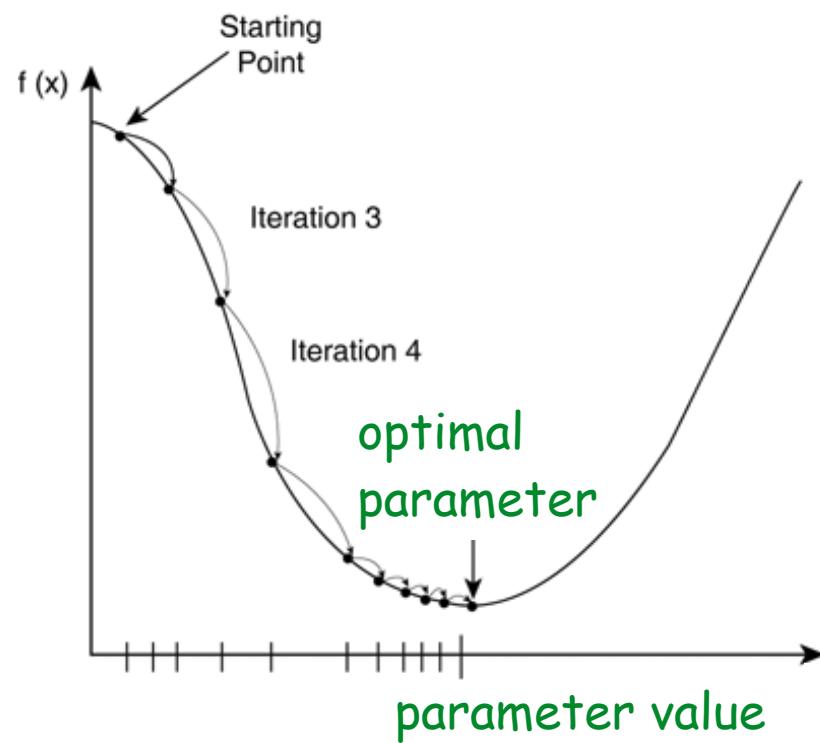


Multi-Class Cross Entropy Loss



Gradient Descent

$L = \text{loss}(w)$



$$W \leftarrow W - a \frac{\partial L}{\partial W}$$

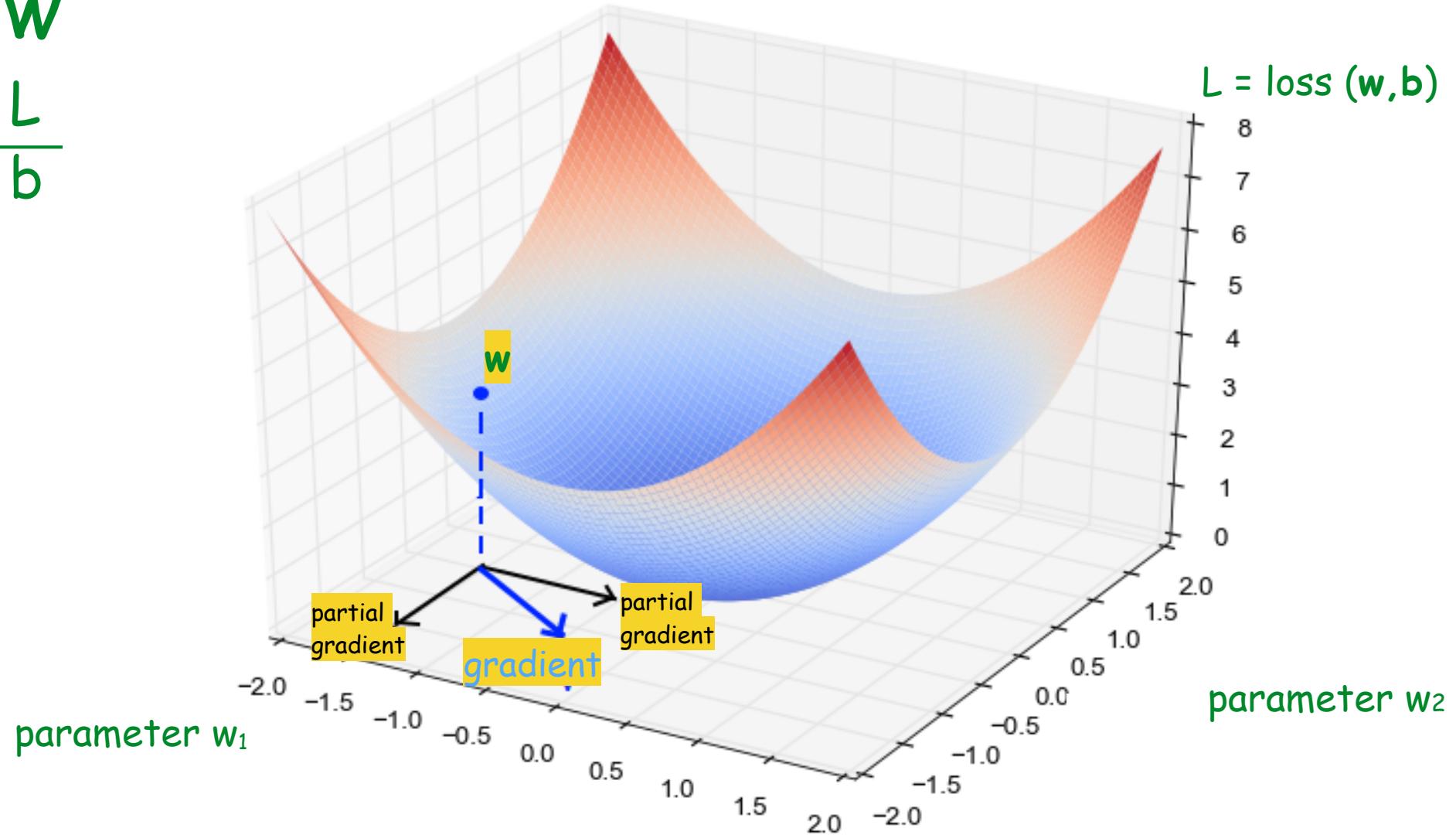
$$b \leftarrow b - a \frac{\partial L}{\partial b}$$

a step size (a constant scalar)

How to Compute Gradient

$$\frac{\partial L}{\partial W}$$

$$\frac{\partial L}{\partial b}$$



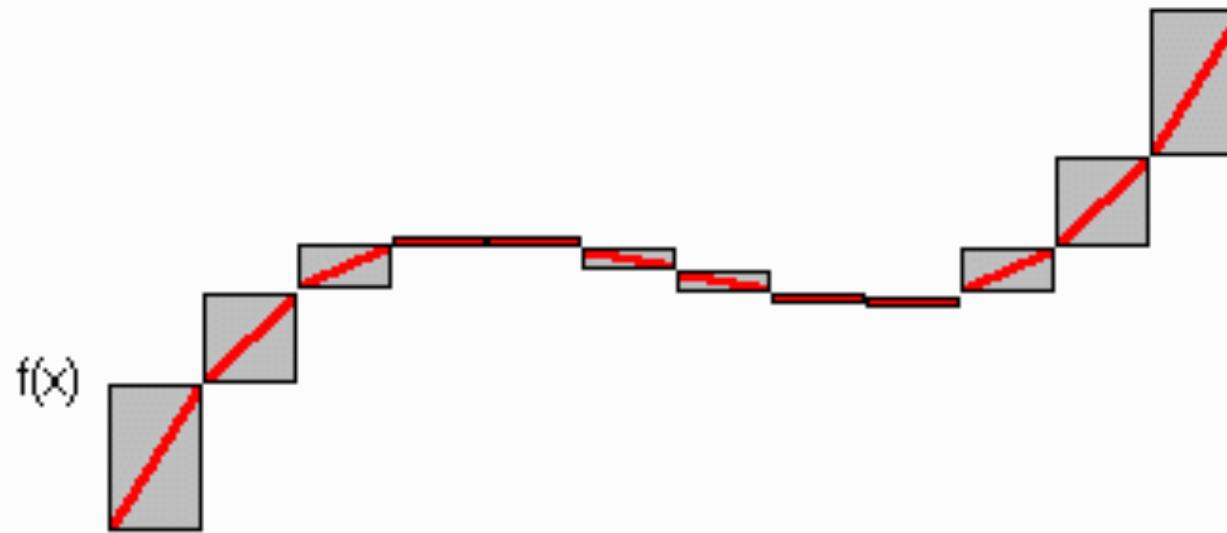
Gradient of L w.r.t. a **vector**?

$$\frac{\partial L}{\partial b}$$

Gradient of L w.r.t. a **matrix**?

$$\frac{\partial L}{\partial W}$$

Gradient of a function f with respect to (w.r.t.) a scalar x



differentiate



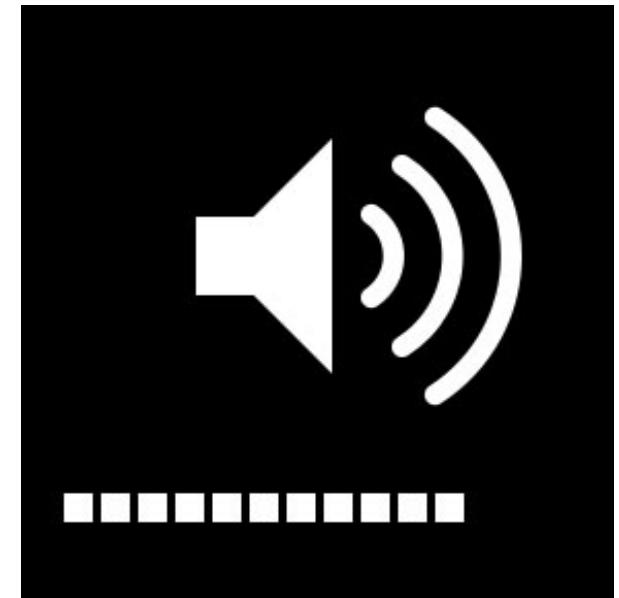
$$\frac{\partial f}{\partial x}$$

df/dx

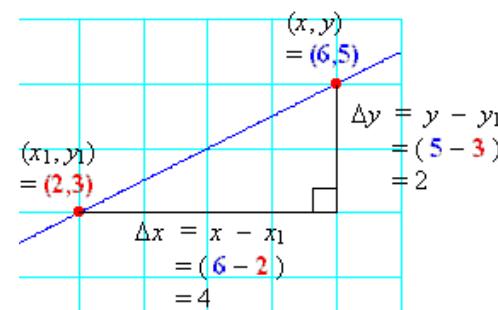


Gradient of f function w.r.t. input x

Gradient of L w.r.t. a scalar x



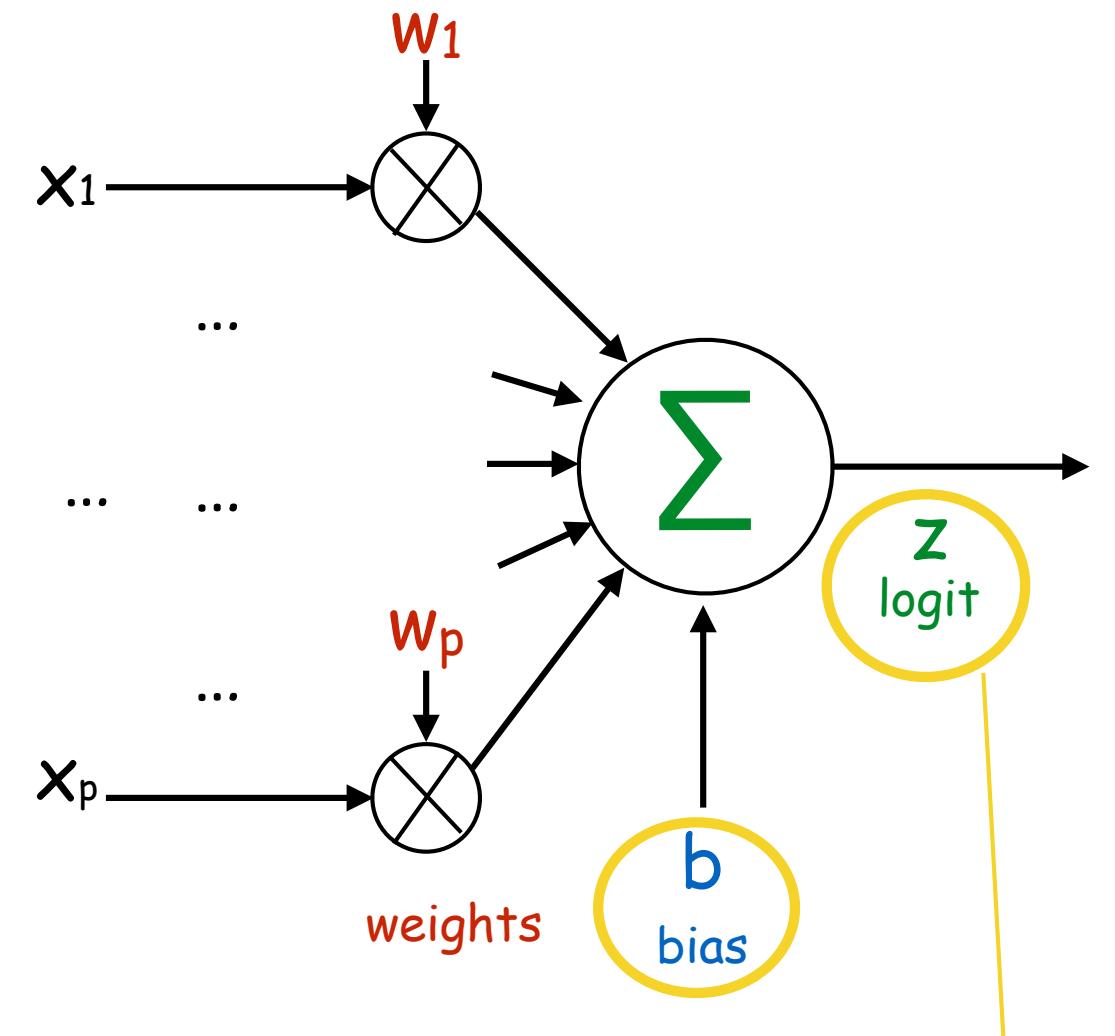
\times
control



L
response

$\frac{\partial L}{\partial x}$ a constance number

Example

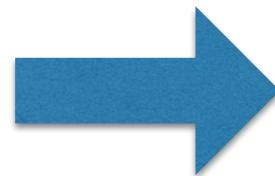


$$\frac{\partial z}{\partial b} = 1$$

Another Example



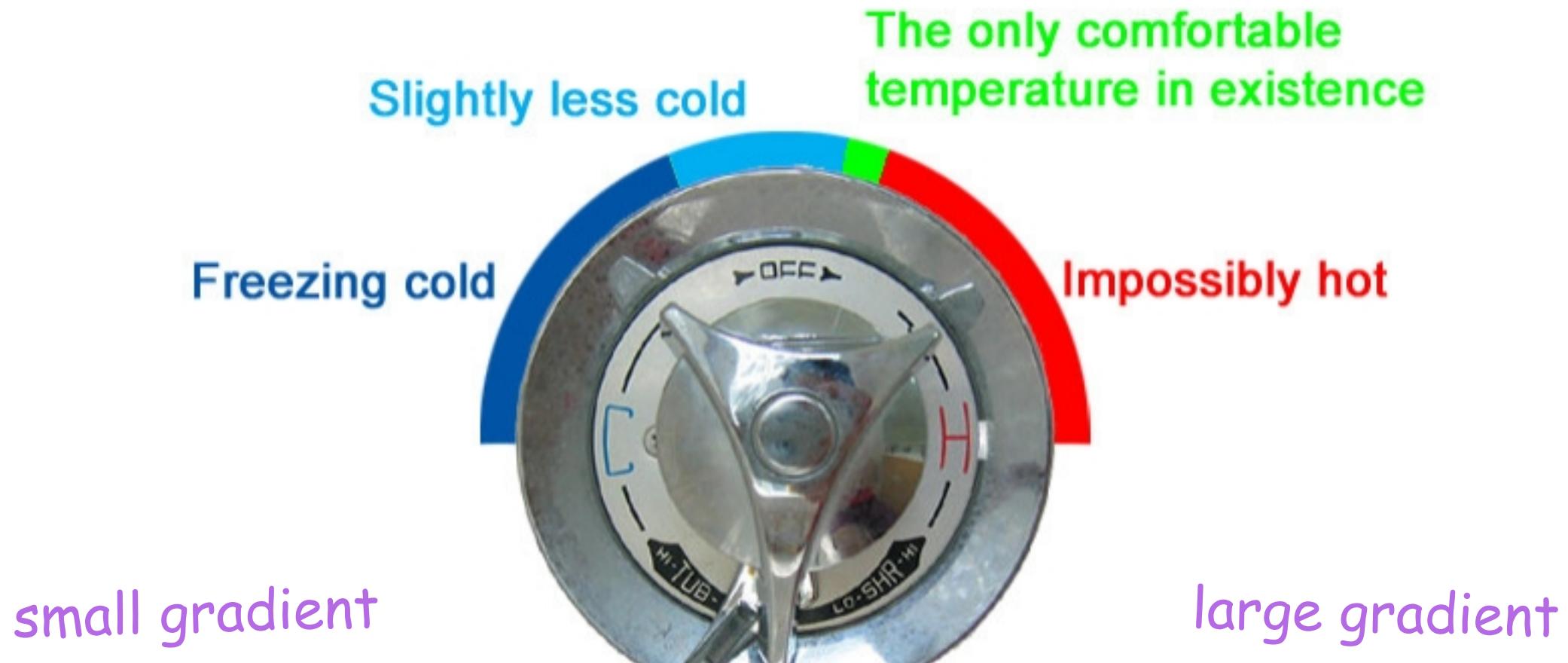
✗



L

water temperature

L is a function of x



Gradient $\frac{\partial L}{\partial x}$ is also a function of x

Partial Gradient w.r.t. a scalar $\frac{\partial L}{\partial x}$



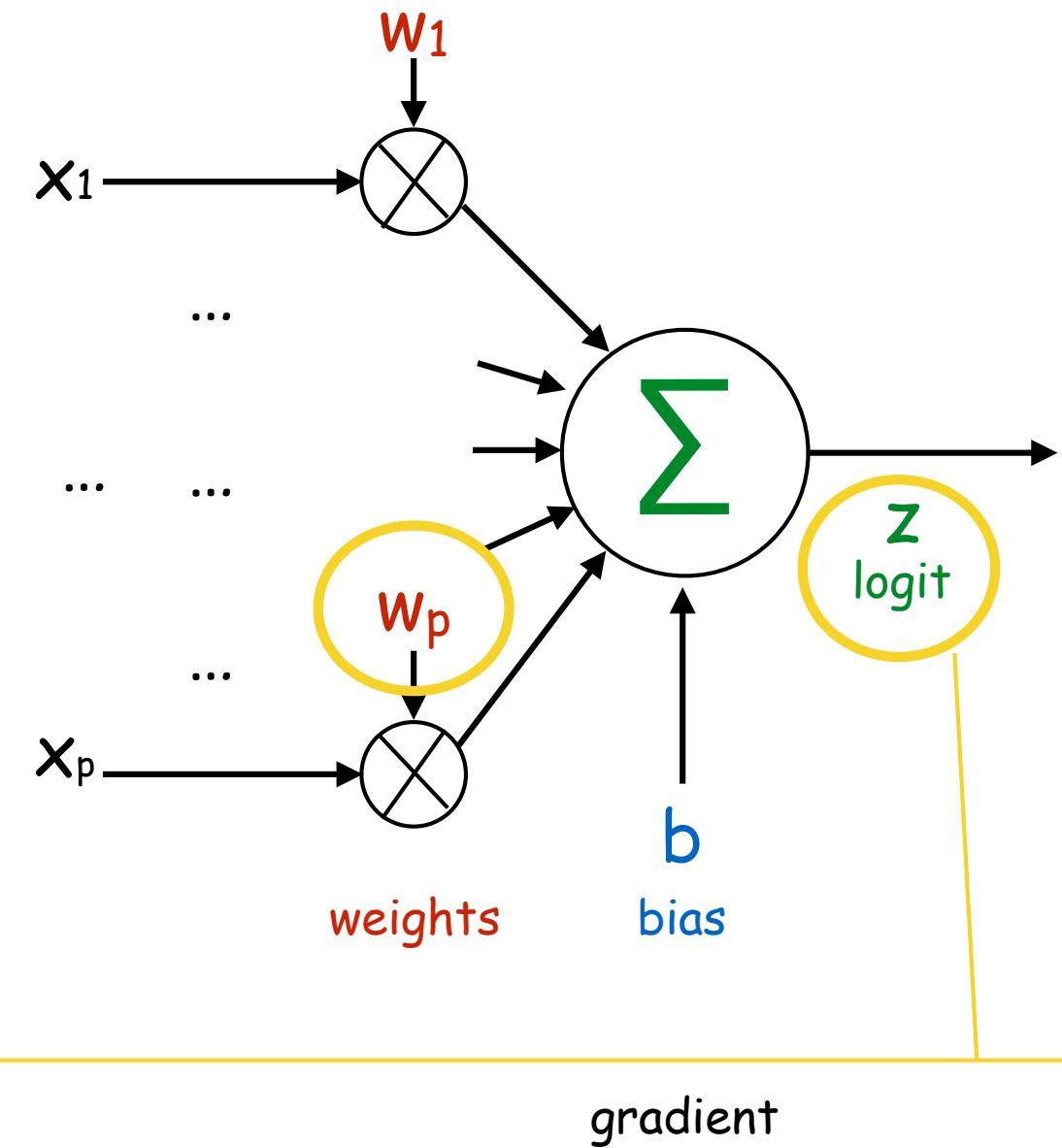
x
gas

L
speed

Gradient

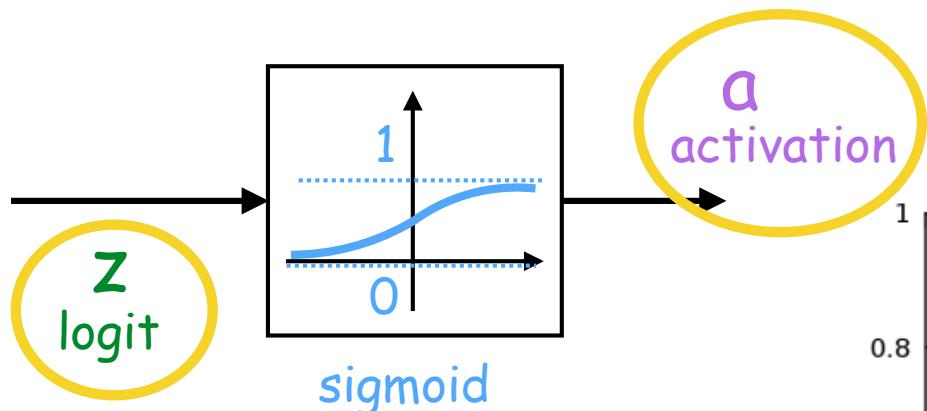
$\frac{\partial L}{\partial x}$ a function of x , speed, wind, gear ...

Example

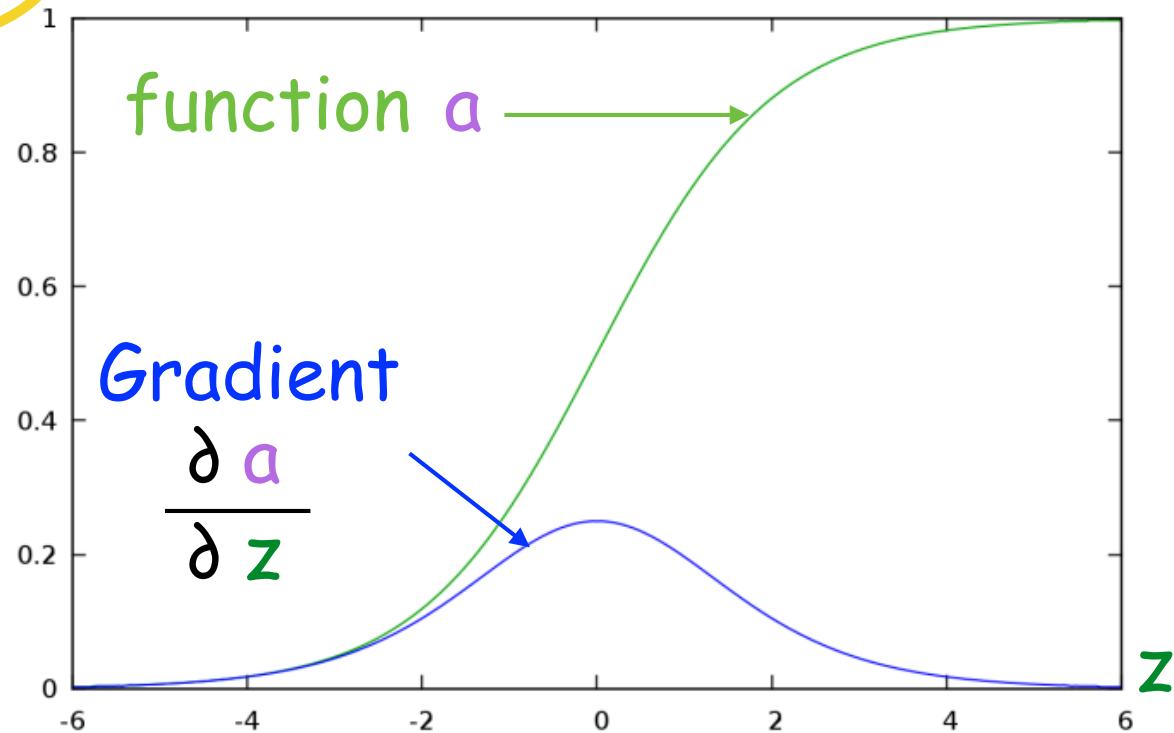


$$\frac{\partial z}{\partial w_p} = x_p$$

is a function of x_p



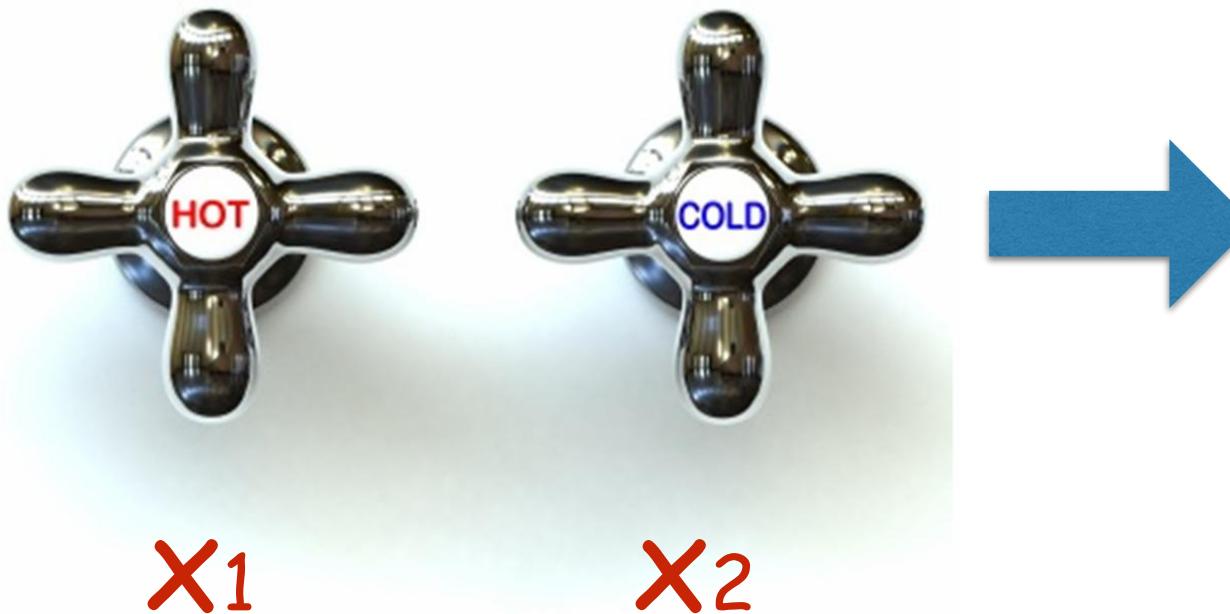
$$a = \frac{1}{1 + e^{-z}}$$



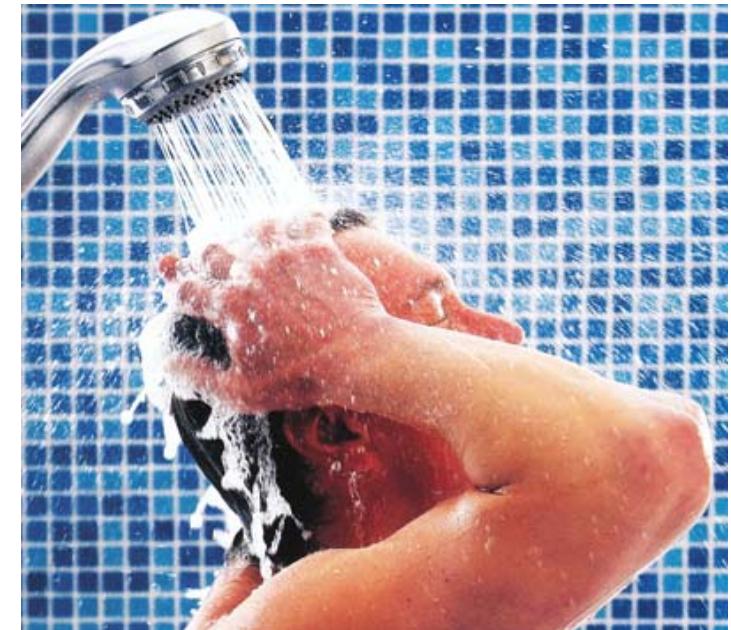
$$\begin{aligned} \frac{\partial a}{\partial z} &= a(1-a) \\ &= \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}}\right) \end{aligned}$$

is a function of z

Partial Gradients w.r.t. a vector $\frac{\partial L}{\partial x}$



$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2} \right)$$



$$L = f(x_1, x_2)$$

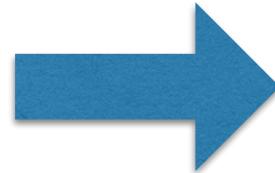
L
temperature

Partial Gradients w.r.t. a vector $\frac{\partial L}{\partial x}$

x_1

x_2

L



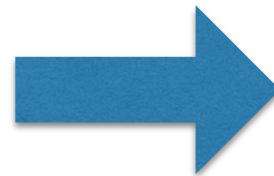
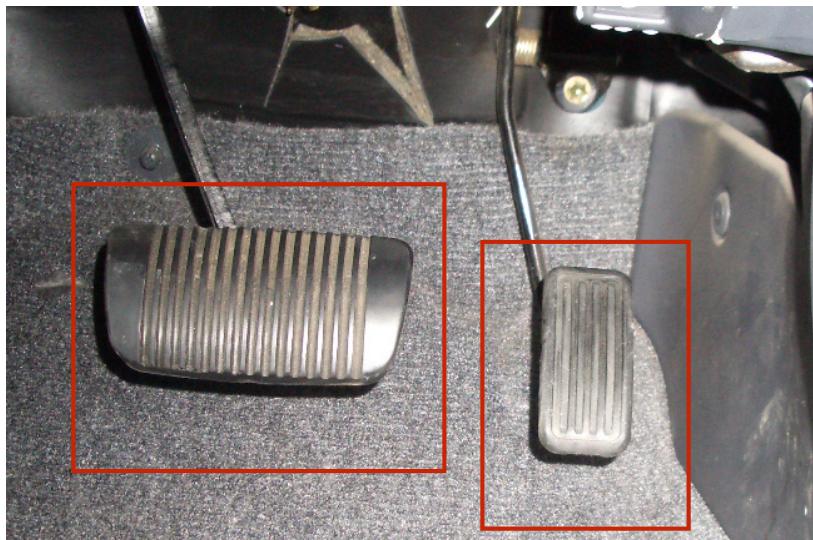
partial gradient – the gradient of f over x_1 while fixing x_2

$$\left(\frac{\partial L}{\partial x_1} \right)$$

$$\left. , \quad \frac{\partial L}{\partial x_2} \right)$$

$$L = f(x_1, x_2)$$

Partial Gradients w.r.t. a vector $\frac{\partial L}{\partial x}$



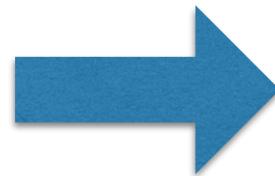
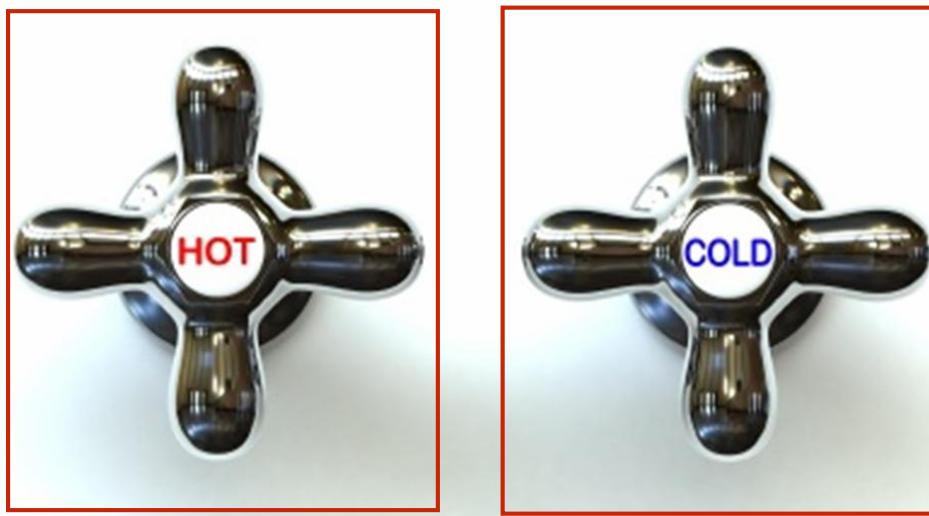
x_1
brake

x_2
gas

L
speed

$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2} \right)$$

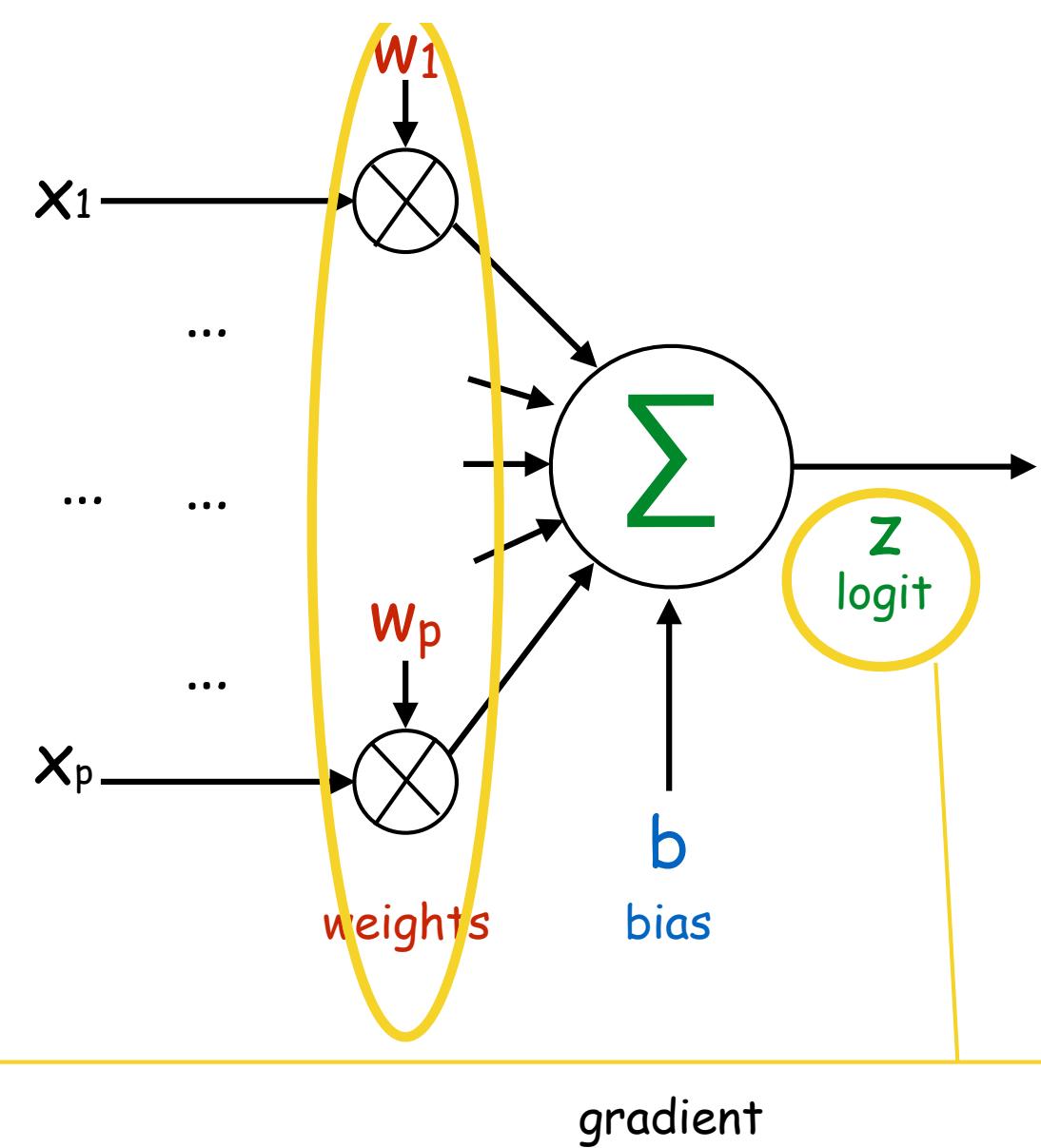
Partial Gradients w.r.t. a vector $\frac{\partial L}{\partial \mathbf{x}}$



$$\begin{array}{c} x_1 \quad x_2 \\ \left(\frac{\partial L}{\partial x_1}, \quad \frac{\partial L}{\partial x_2} \right) = \quad \frac{\partial L}{\partial \mathbf{x}} \\ 3 \quad -4 \quad \text{for example} \end{array}$$

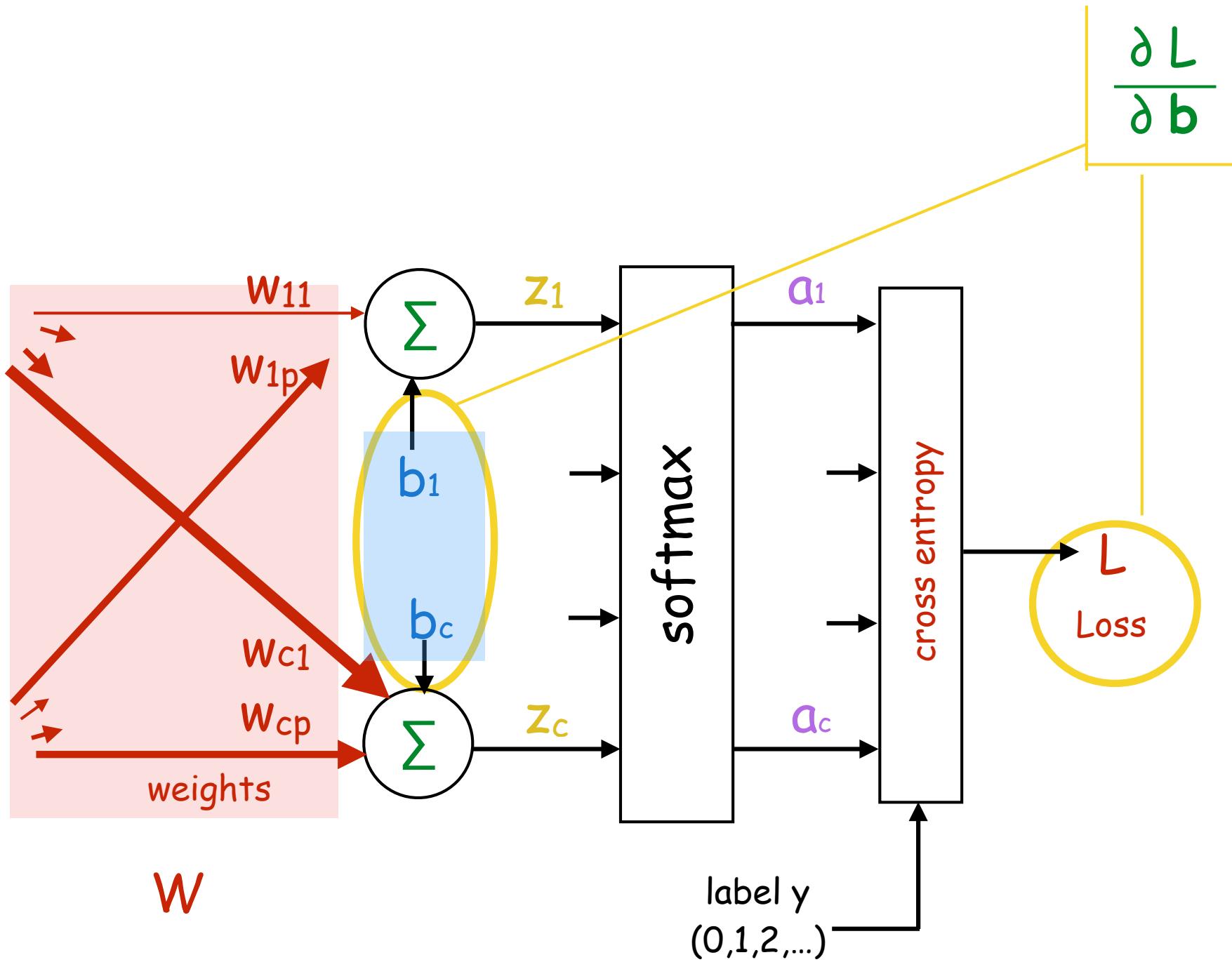
temperature

Example



$$\begin{aligned}\frac{\partial z}{\partial w} &= \left(\frac{\partial z}{\partial w_1}, \frac{\partial z}{\partial w_2}, \dots, \frac{\partial z}{\partial w_p} \right) \\ &= (x_1, x_2, \dots, x_p) = x\end{aligned}$$

Example



Gradient w.r.t. a Matrix?



height

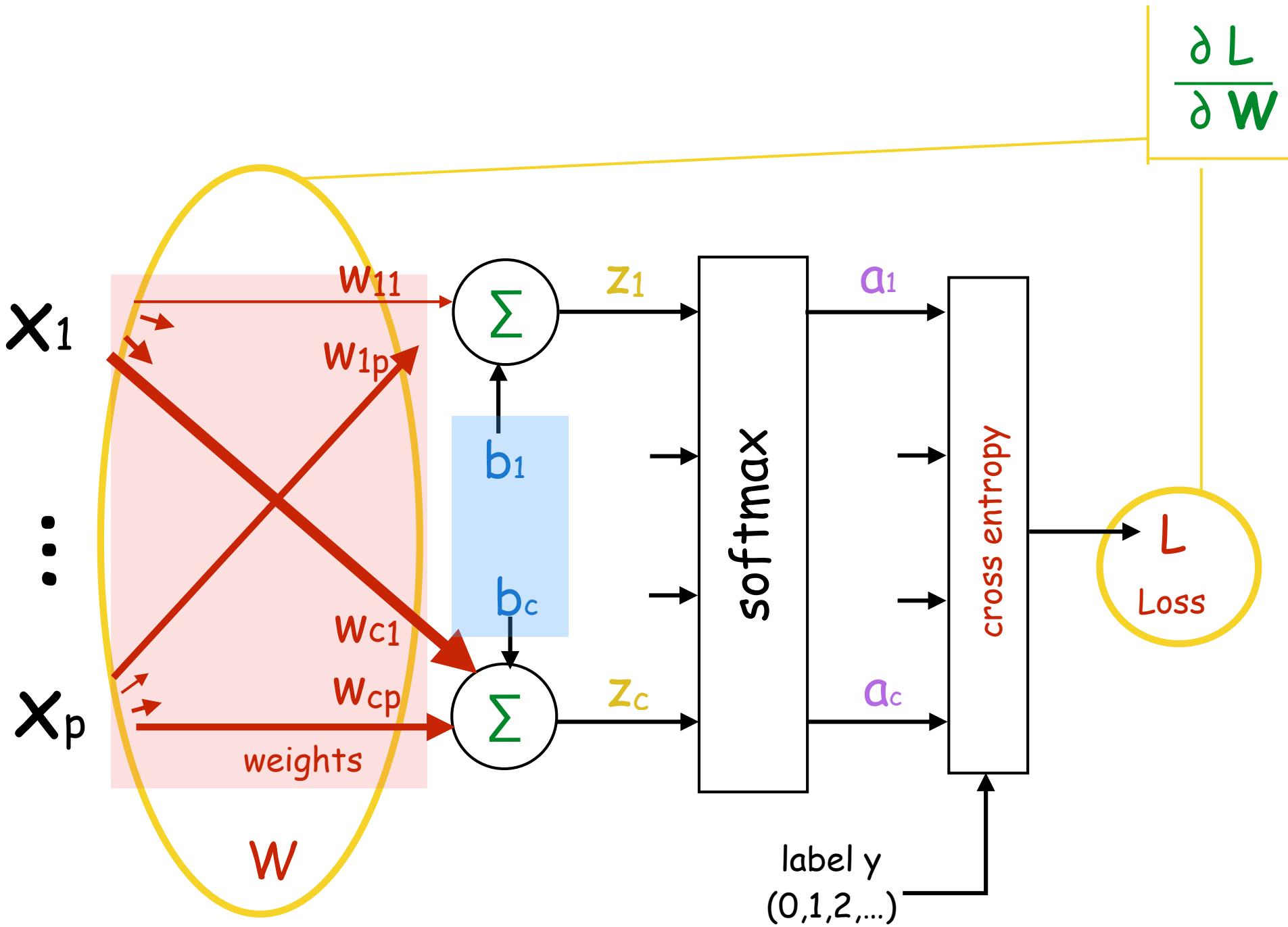
a matrix of controls

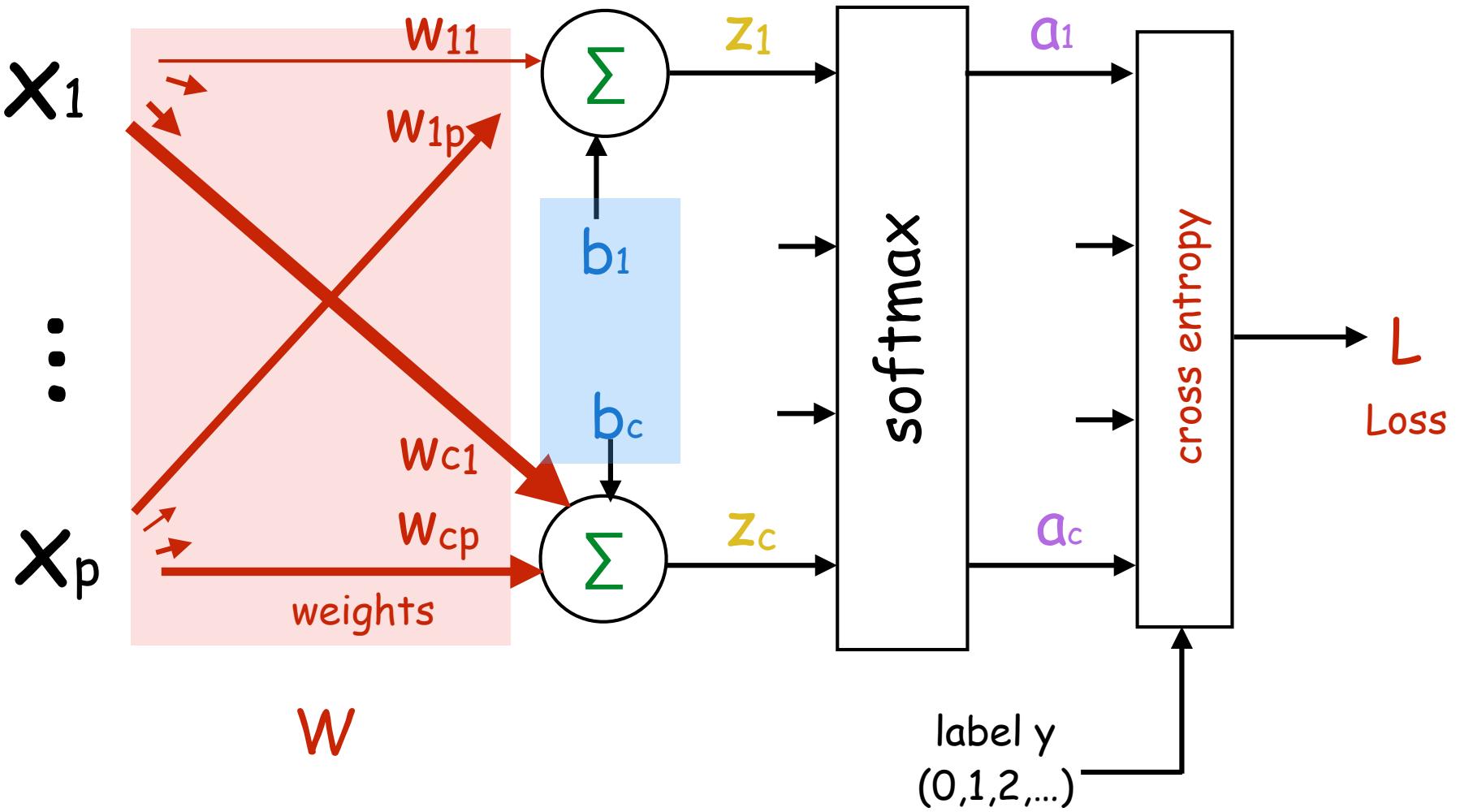
Gradient w.r.t. a Matrix?



$$\frac{\partial L}{\partial W} = \begin{matrix} \frac{\partial L}{\partial w_{11}} & \dots & \frac{\partial L}{\partial x_{1p}} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial w_{c1}} & & \frac{\partial L}{\partial w_{cp}} \end{matrix}$$

Example

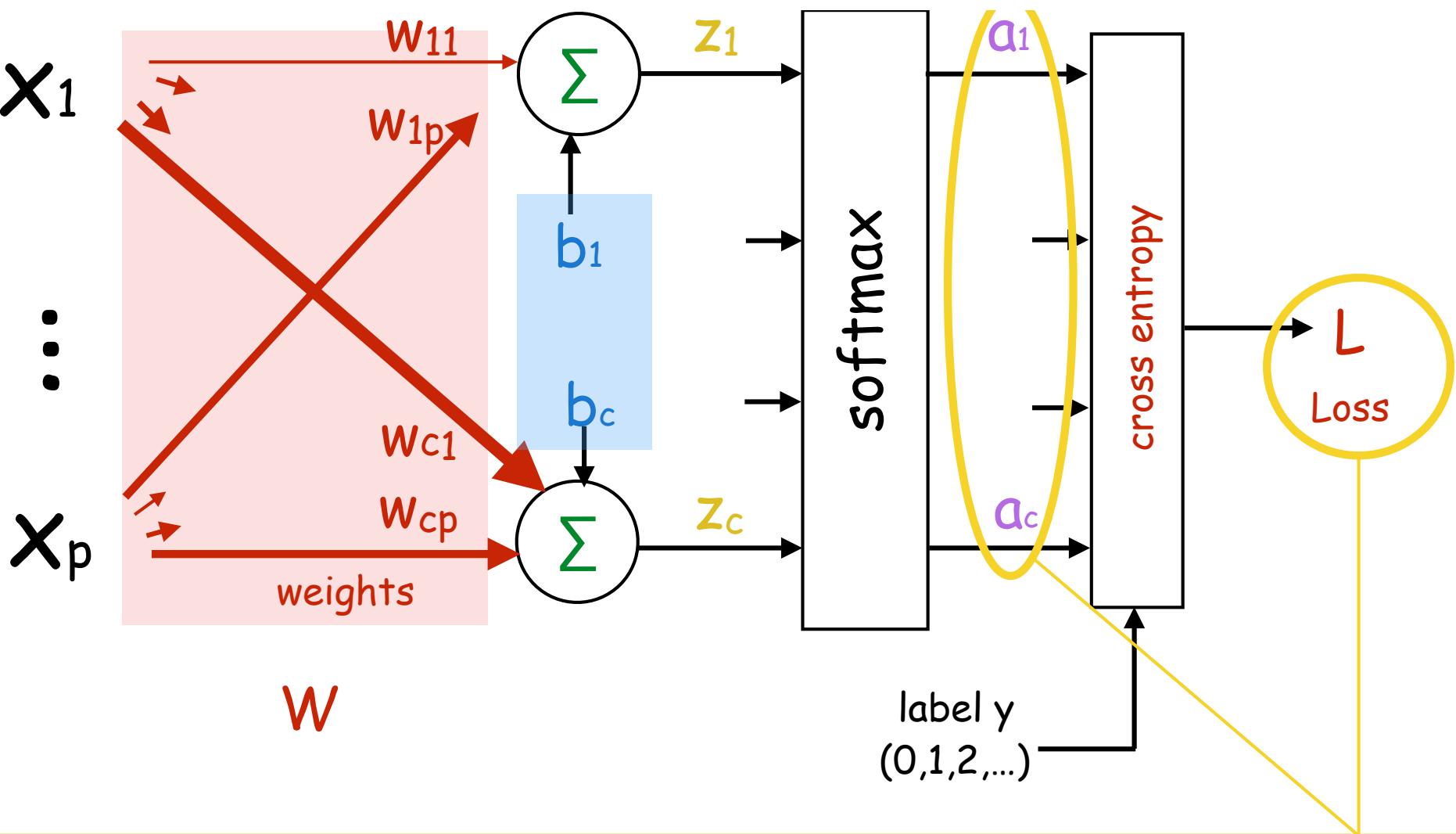




$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial W}$$

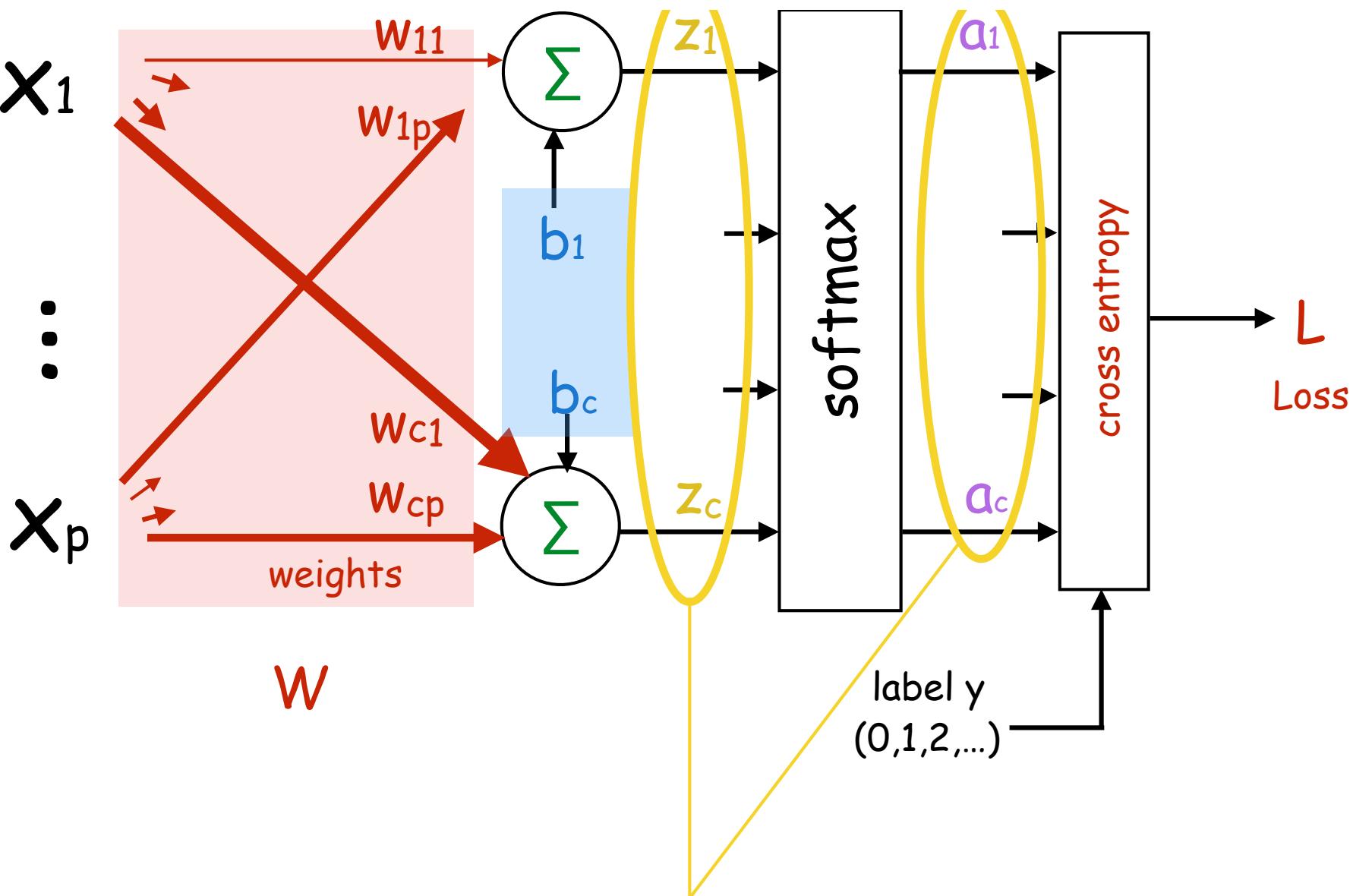
$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial b}$$

Chain Rule



$$\begin{aligned}
 \frac{\partial L}{\partial a} &= \left(\frac{\partial L}{\partial a_1}, \dots, \frac{\partial L}{\partial a_y}, \dots, \frac{\partial L}{\partial a_c} \right) \\
 &= \left(0, \dots, -\frac{1}{a_y}, \dots, 0 \right)
 \end{aligned}$$

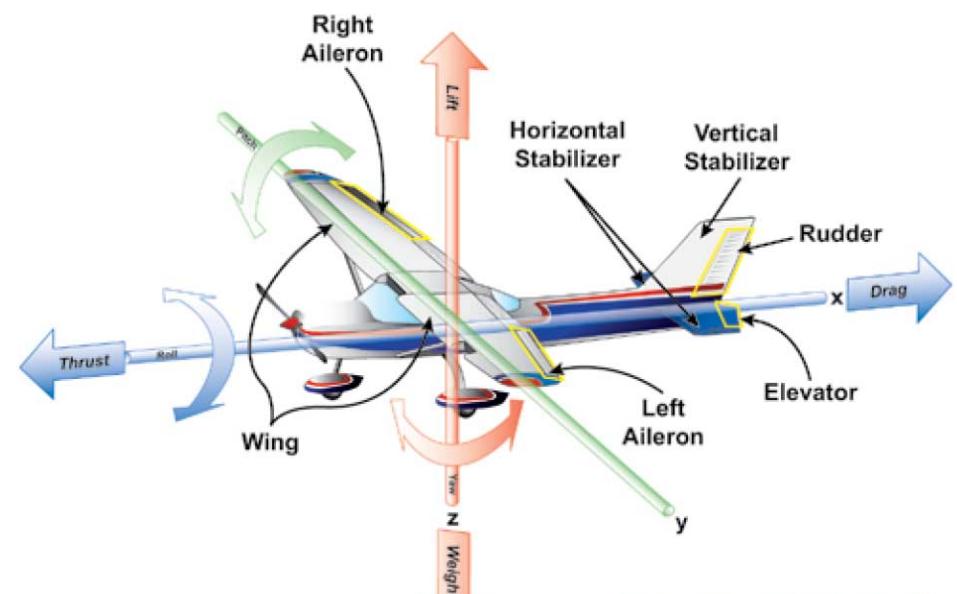
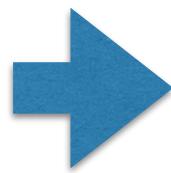
y



$$\frac{\partial a}{\partial z} = ?$$

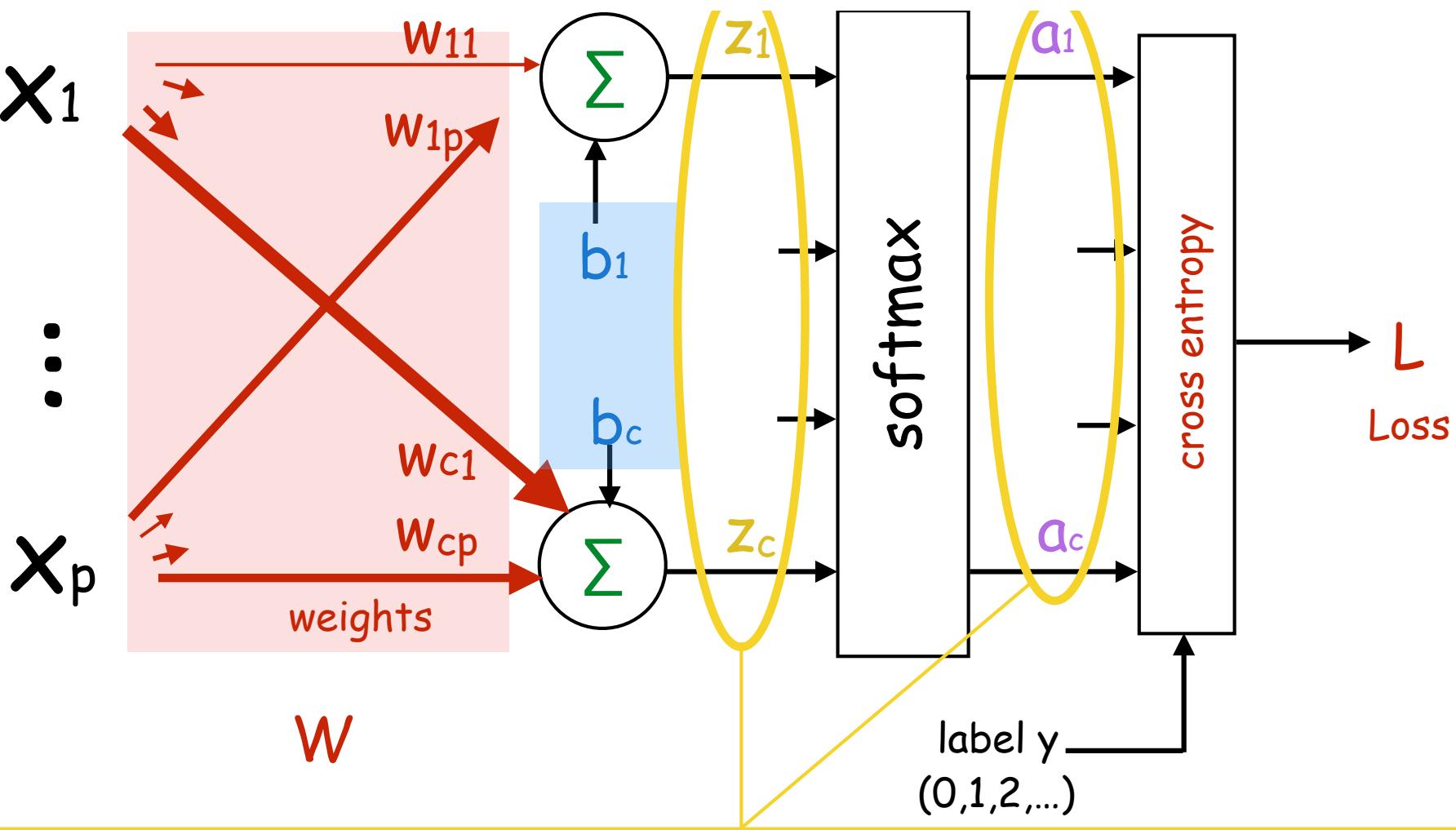
Gradient of a vector w.r.t. a vector !!!

Gradient of a vector w.r.t. a vector



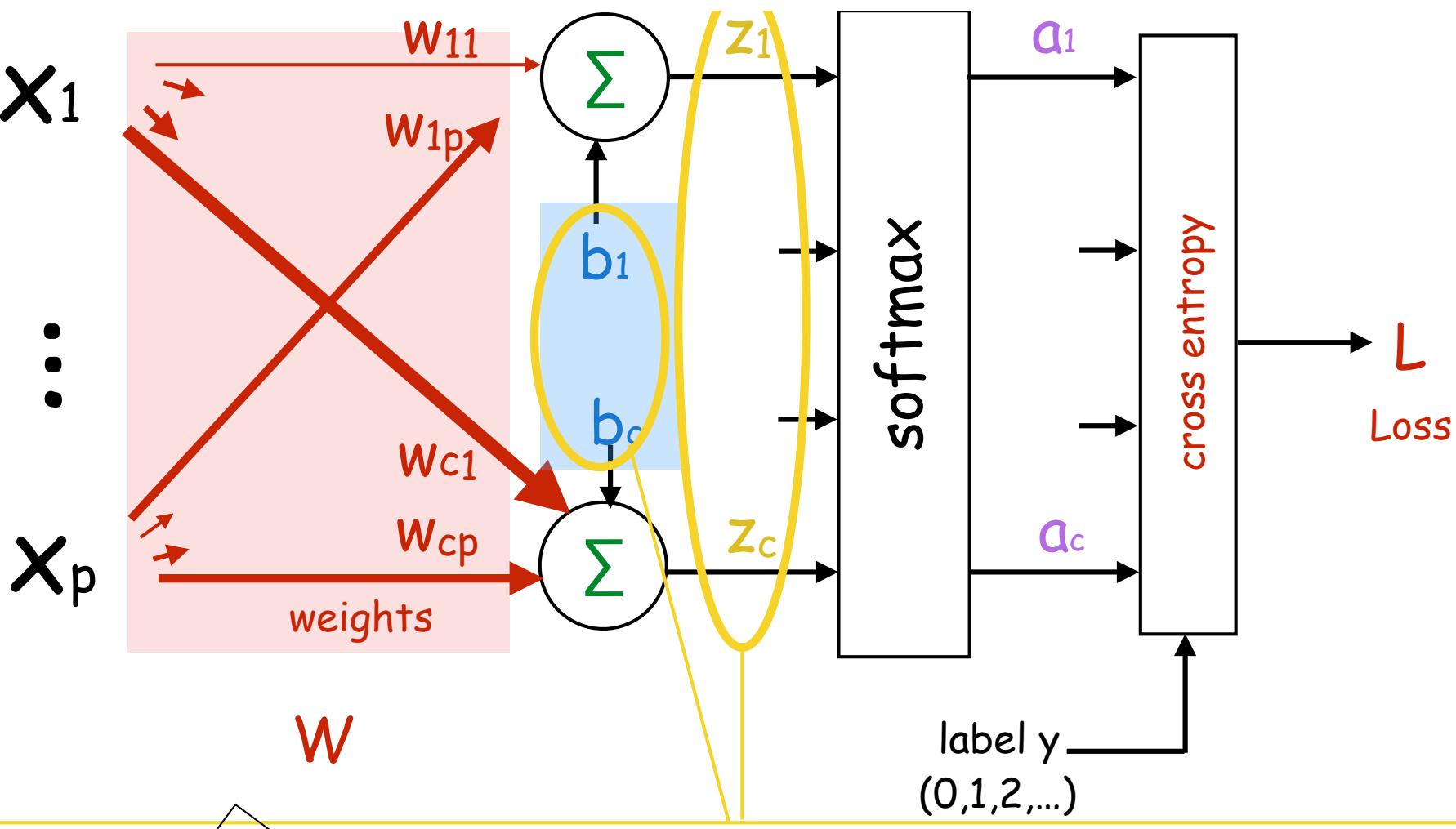
a vector of controls

a vector of response

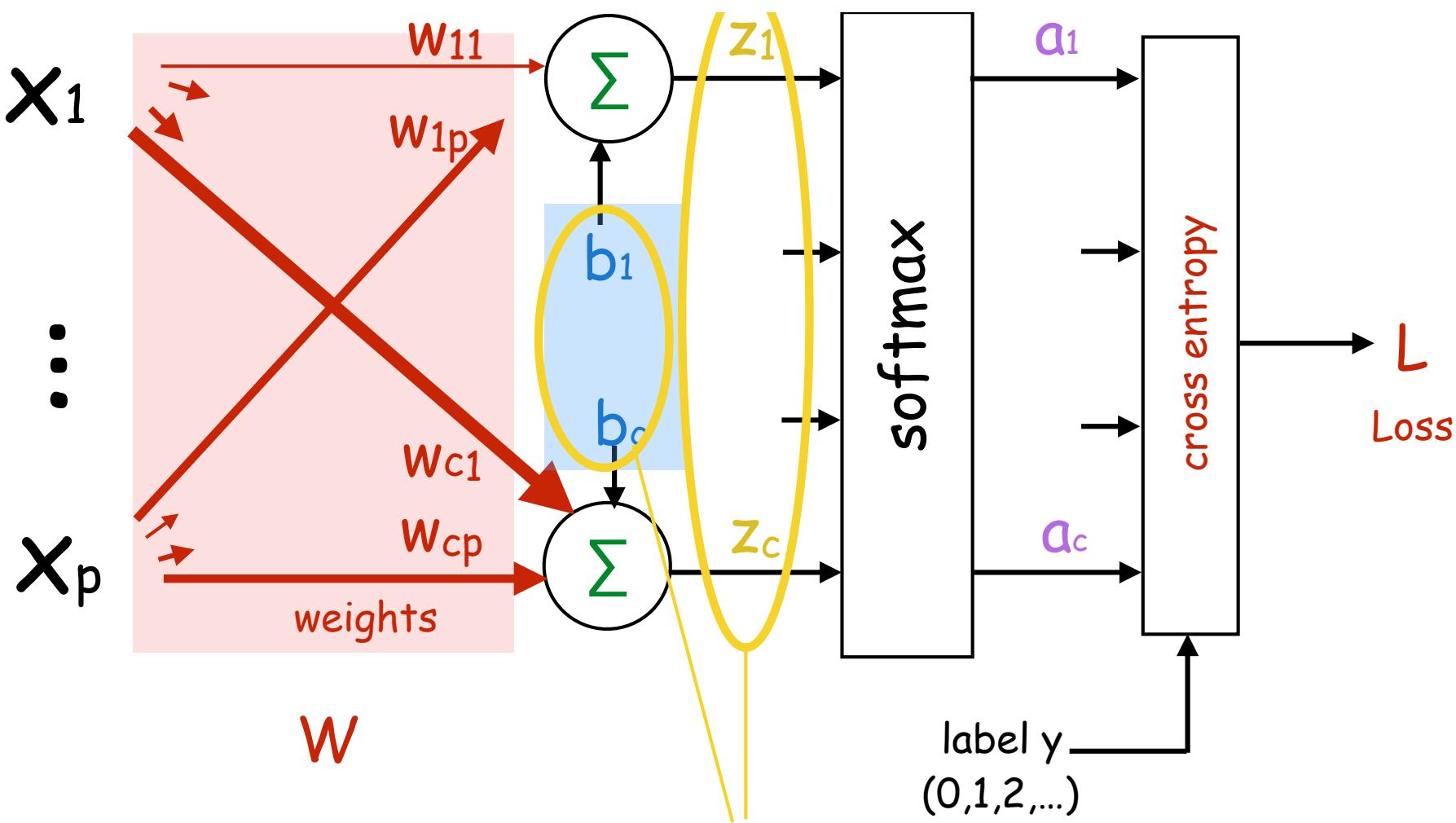


$$\frac{\partial a}{\partial z} = \begin{pmatrix} \frac{\partial a_1}{\partial z_1} & \cdots & \frac{\partial a_1}{\partial z_c} \\ \vdots & & \vdots \\ \frac{\partial a_c}{\partial z_1} & \cdots & \frac{\partial a_c}{\partial z_c} \end{pmatrix}$$

$$\frac{\partial a_i}{\partial z_j} = \begin{cases} a_i (1 - a_i) & \text{if } i=j \\ -a_i a_j & \text{if } i \neq j \end{cases}$$



$$\frac{\partial z}{\partial b} = \begin{pmatrix} 1 & 1 & \cdots & 0 \\ 1 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 1 \\ 0 & \cdots & 1 & 1 \end{pmatrix}$$

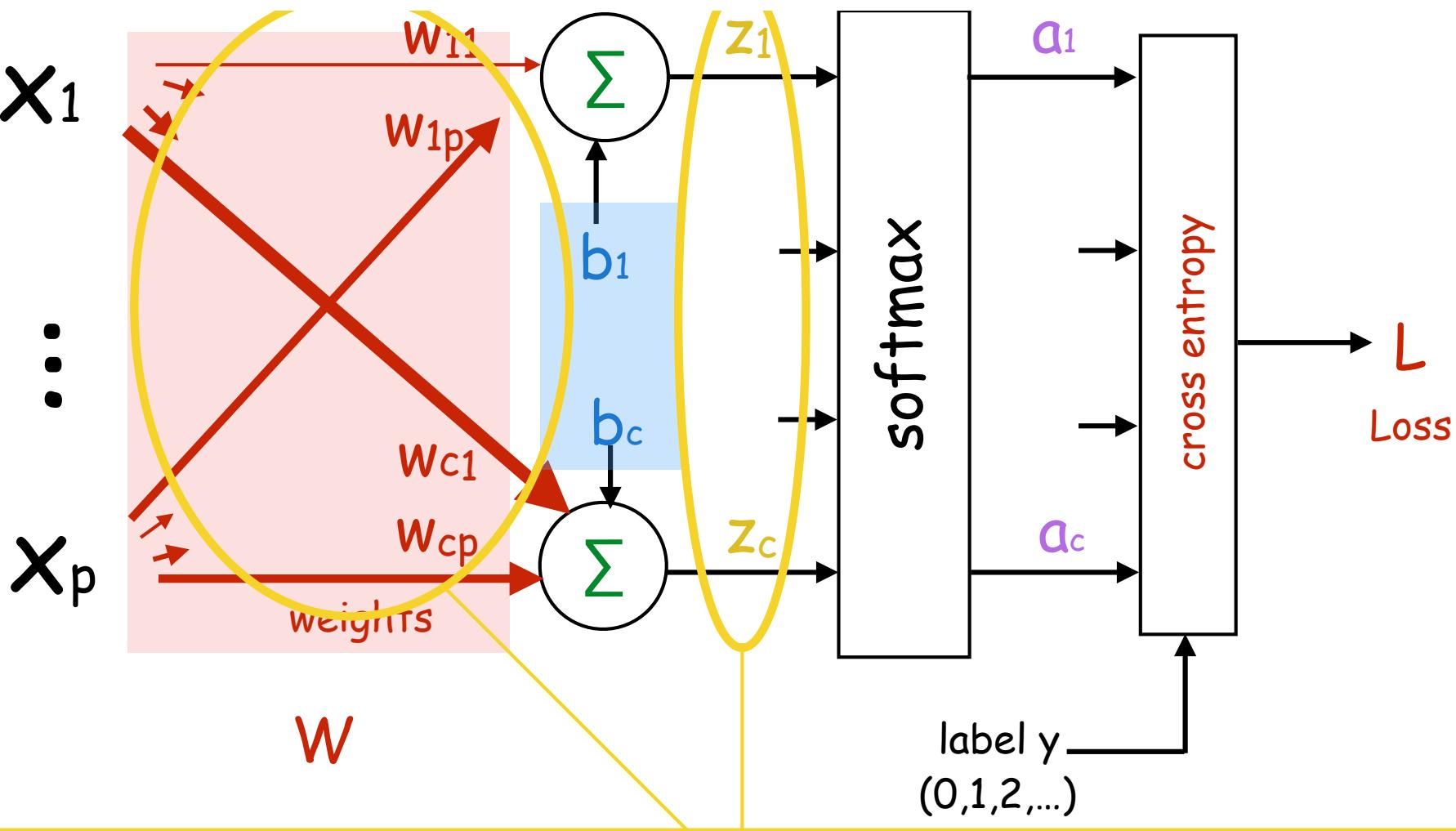


$$\frac{\partial z_1}{\partial b_1} = 1$$

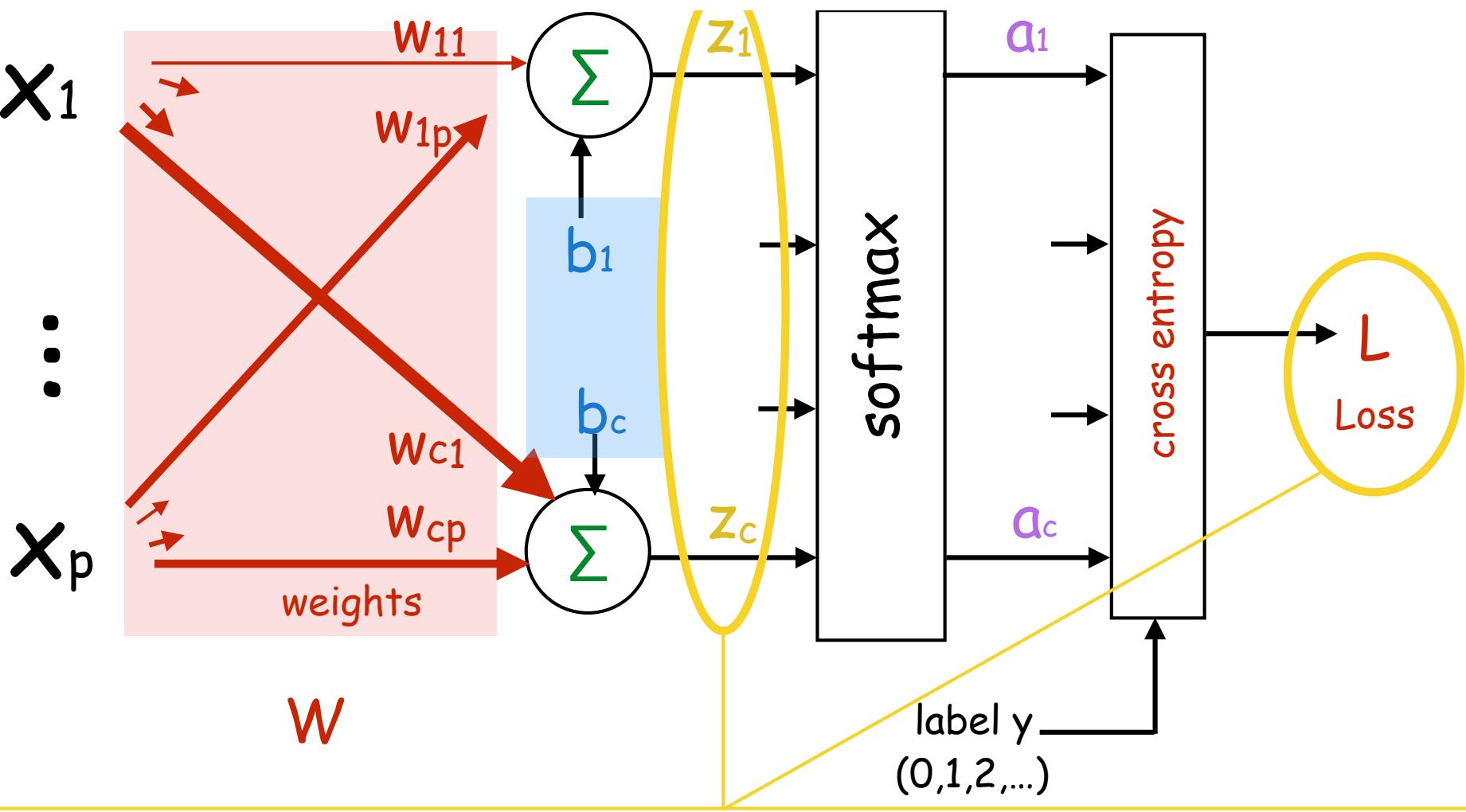
$$\frac{\partial z}{\partial b}$$

$$\frac{\partial z_i}{\partial b_i} = 1$$

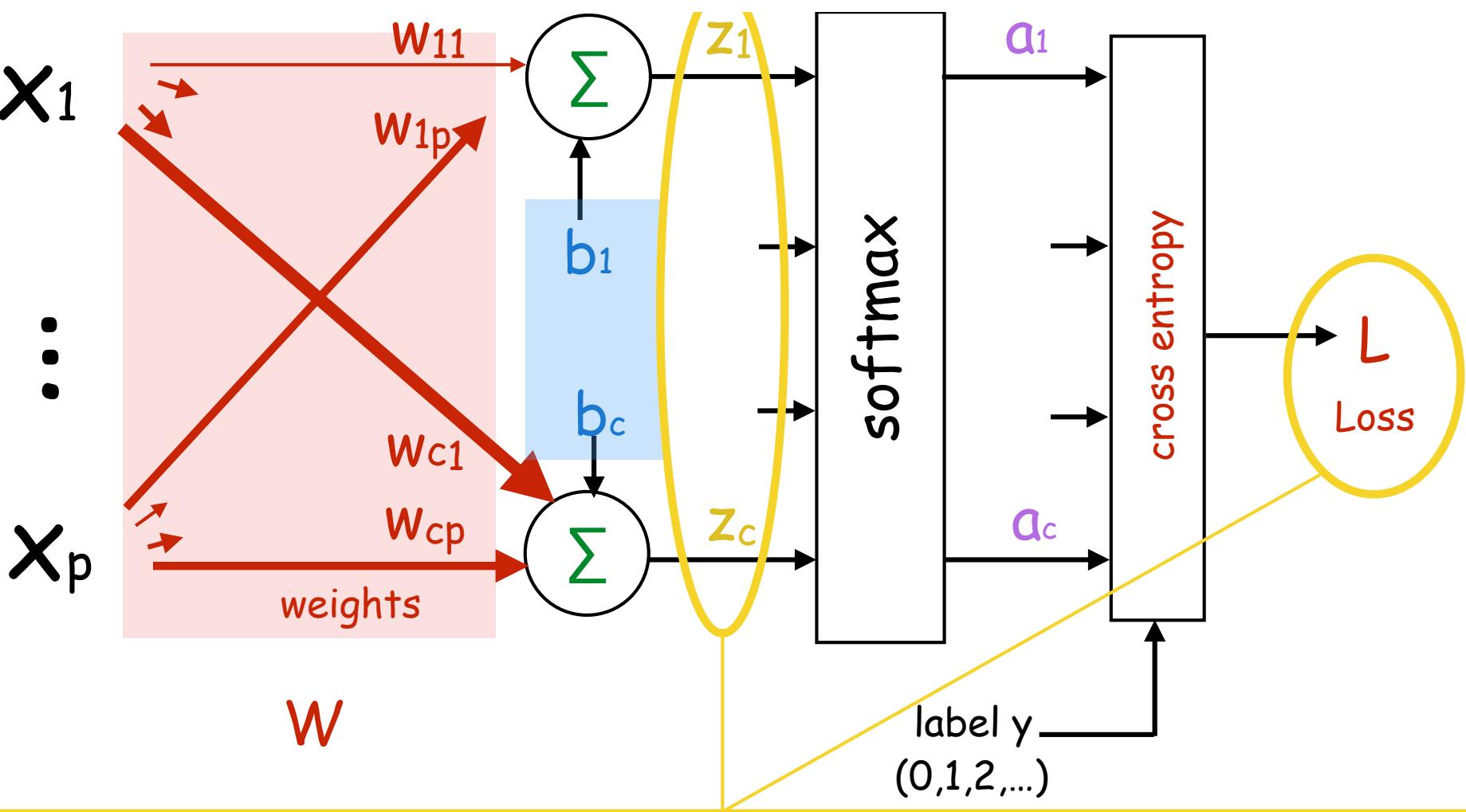
$$\frac{\partial z_c}{\partial b_c} = 1$$



$$\frac{\partial Z}{\partial W} = \begin{pmatrix} \frac{\partial Z_1}{\partial W_{1*}} \\ \frac{\partial Z_i}{\partial W_{i*}} \\ \frac{\partial Z_c}{\partial W_{c*}} \end{pmatrix} = \begin{pmatrix} \frac{\partial Z_1}{\partial W_{11}} & \dots & \frac{\partial Z_1}{\partial W_{1p}} \\ \frac{\partial Z_i}{\partial W_{i1}} & \dots & \frac{\partial Z_i}{\partial W_{ip}} \\ \frac{\partial Z_c}{\partial W_{c1}} & \dots & \frac{\partial Z_c}{\partial W_{cp}} \end{pmatrix} = \begin{pmatrix} x_1 & x_2 & \dots & x_p \\ x_1 & x_2 & \dots & x_p \\ x_1 & x_2 & \dots & x_p \end{pmatrix}$$



$$\frac{\partial L}{\partial z} = \left(\frac{\partial L}{\partial z_1}, \dots, \frac{\partial L}{\partial z_i}, \dots, \frac{\partial L}{\partial z_c} \right) = ?$$



Chain Rule

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z}$$

vector (1 by c)

vector (1 by c)

matrix (c by c)

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z}$$

vector (1 by c)

vector (1 by c)

matrix (c by c)

$$\left(\frac{\partial L}{\partial z_1}, \dots, \frac{\partial L}{\partial z_i}, \dots, \frac{\partial L}{\partial z_c} \right) =$$

$$\left(\frac{\partial L}{\partial a_1}, \dots, \frac{\partial L}{\partial a_i}, \dots, \frac{\partial L}{\partial a_c} \right) \times$$

| | | |
|-------------------------------------|-----|-------------------------------------|
| $\frac{\partial a_1}{\partial z_1}$ | ... | $\frac{\partial a_1}{\partial z_c}$ |
| \vdots | | \vdots |
| $\frac{\partial a_c}{\partial z_1}$ | | $\frac{\partial a_c}{\partial z_c}$ |

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z}$$

vector (1 by c)

vector (1 by c)

matrix (c by c)

$$\left(\frac{\partial L}{\partial z_1}, \dots, \frac{\partial L}{\partial z_i}, \dots, \frac{\partial L}{\partial z_c} \right) =$$

$$\left(\frac{\partial L}{\partial a_1}, \dots, \frac{\partial L}{\partial a_i}, \dots, \frac{\partial L}{\partial a_c} \right) \times$$

$$\begin{array}{c} \frac{\partial a_1}{\partial z_1} \\ \vdots \\ \frac{\partial a_c}{\partial z_1} \end{array} \cdots \begin{array}{c} \dots \\ \vdots \\ \dots \end{array} \begin{array}{c} \frac{\partial a_1}{\partial z_c} \\ \vdots \\ \frac{\partial a_c}{\partial z_c} \end{array}$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z}$$

vector (1 by c)

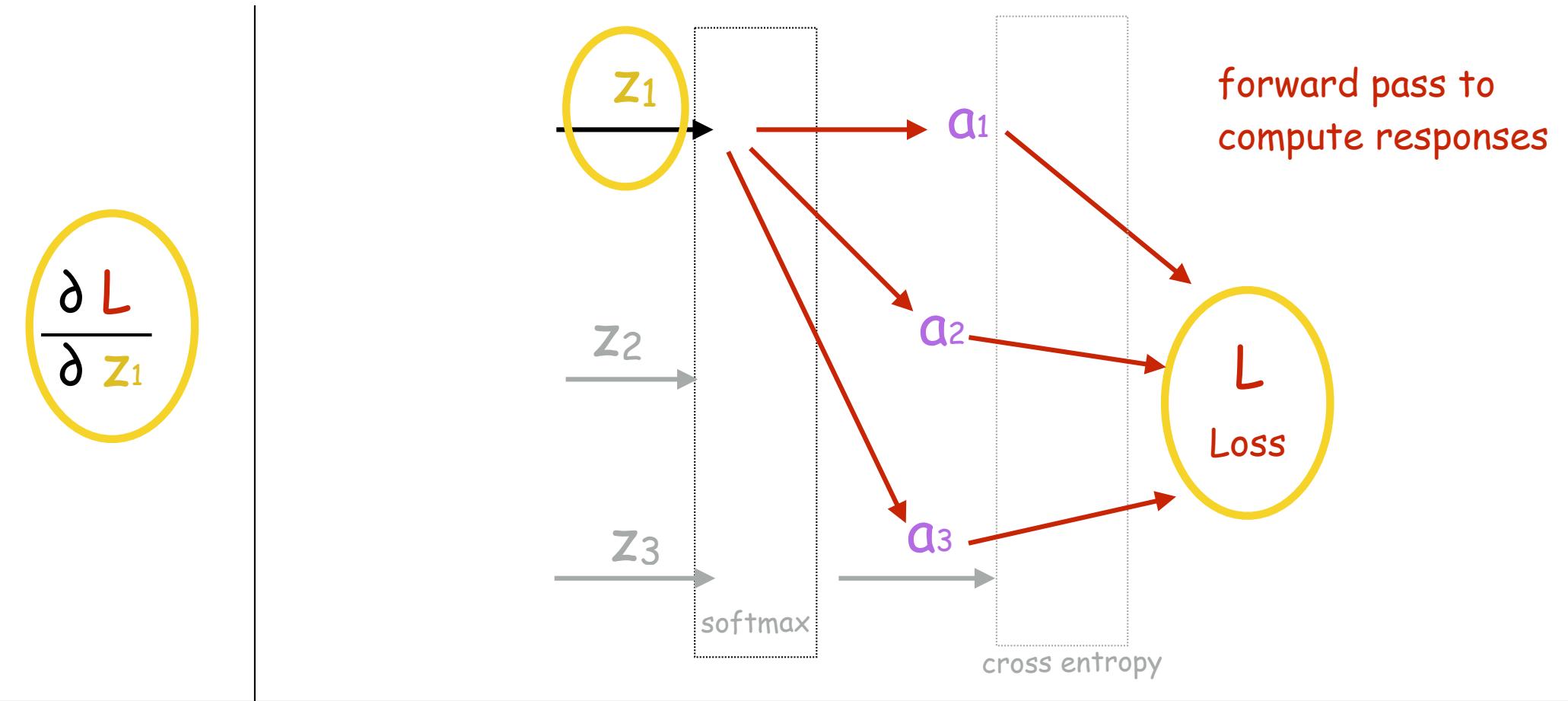
vector (1 by c)

matrix (c by c)

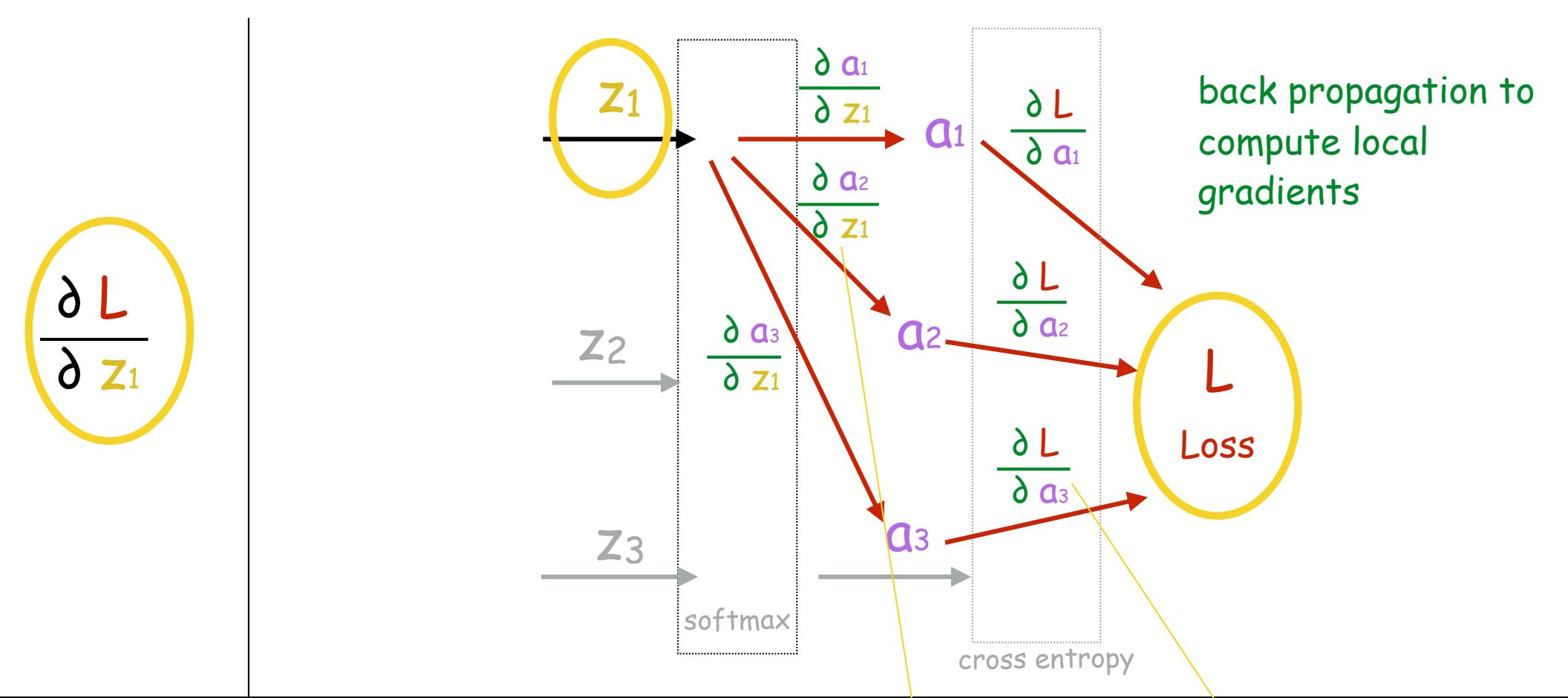
$$\left(\frac{\partial L}{\partial z_1}, \dots, \frac{\partial L}{\partial z_i}, \dots, \frac{\partial L}{\partial z_c} \right) =$$

$$\left(\frac{\partial L}{\partial a_1}, \dots, \frac{\partial L}{\partial a_i}, \dots, \frac{\partial L}{\partial a_c} \right) \times$$

$$\begin{array}{c} \frac{\partial a_1}{\partial z_1} \\ \vdots \\ \frac{\partial a_c}{\partial z_1} \end{array} \cdots \begin{array}{c} \frac{\partial a_1}{\partial z_c} \\ \vdots \\ \frac{\partial a_c}{\partial z_c} \end{array}$$

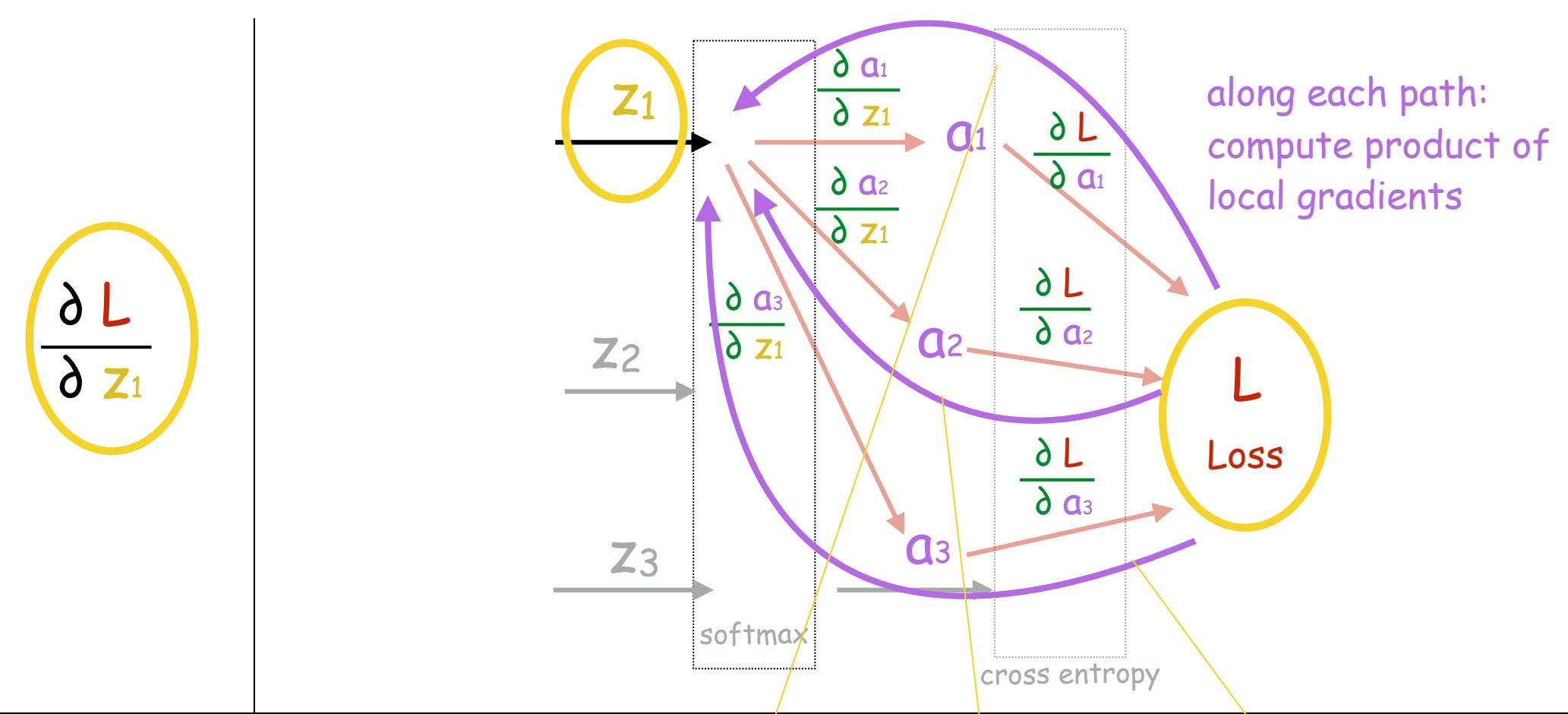


$$\left(\frac{\partial L}{\partial a_1}, \dots, \frac{\partial L}{\partial a_i}, \dots, \frac{\partial L}{\partial a_c} \right) \times \begin{pmatrix} \frac{\partial a_1}{\partial z_1} & \dots & \frac{\partial a_1}{\partial z_c} \\ \vdots & & \vdots \\ \frac{\partial a_c}{\partial z_1} & \dots & \frac{\partial a_c}{\partial z_c} \end{pmatrix}$$



$$\frac{\partial L}{\partial a} = \left(\frac{\partial L}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} + \frac{\partial L}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_1} + \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_1} \right)$$

$$\frac{\partial a}{\partial z_1} = \left(\frac{\partial a_1}{\partial z_1} \cdot \frac{\partial a_1}{\partial z_1} + \frac{\partial a_2}{\partial z_1} \cdot \frac{\partial a_2}{\partial z_1} + \frac{\partial a_3}{\partial z_1} \cdot \frac{\partial a_3}{\partial z_1} \right)$$



$$\frac{\partial L}{\partial a}$$

=

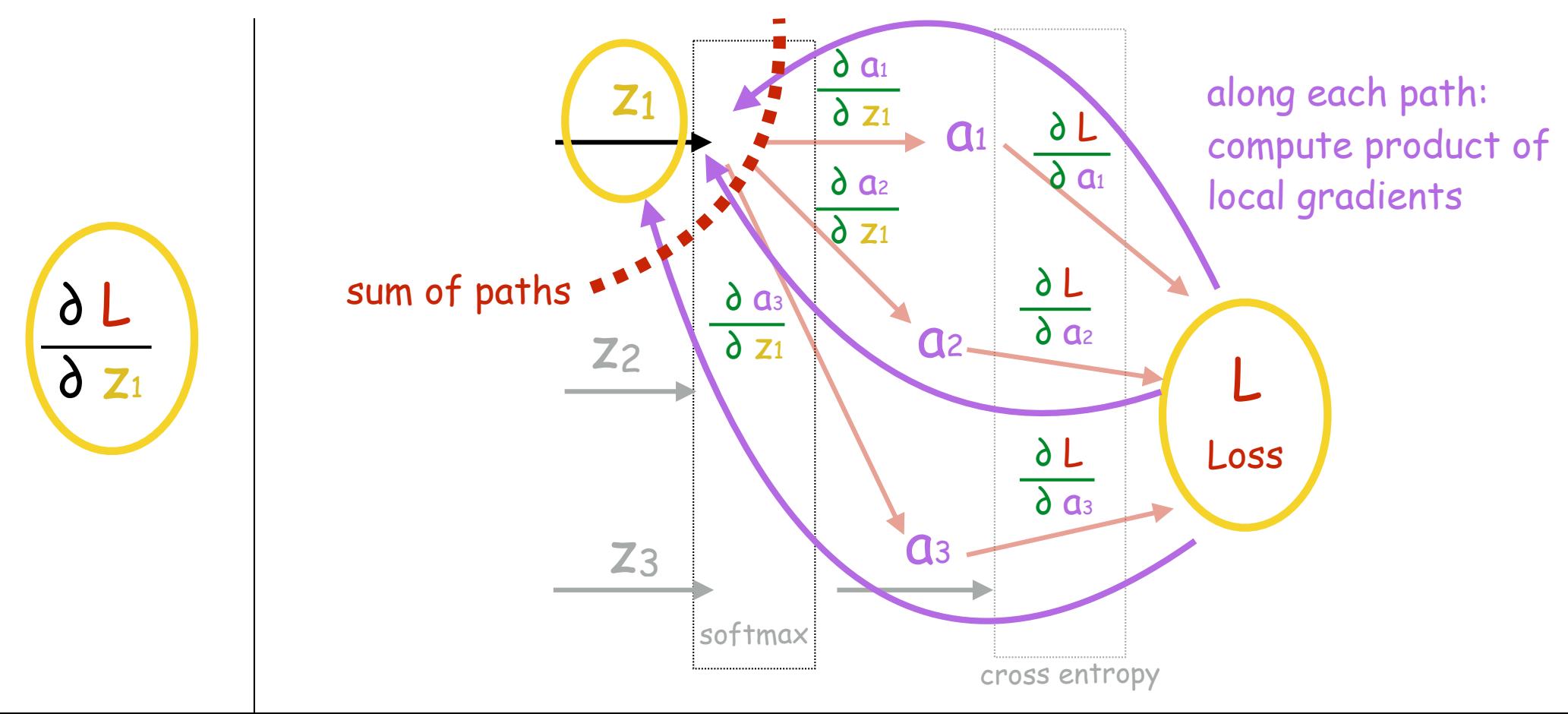
$$\left(\frac{\partial L}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \right)$$

$$- \left(\frac{\partial L}{\partial a_2} \times \frac{\partial a_2}{\partial z_1} \right)$$

$$- \left(\frac{\partial L}{\partial a_3} \times \frac{\partial a_3}{\partial z_1} \right)$$

$$\frac{\partial a}{\partial z_1}$$

=

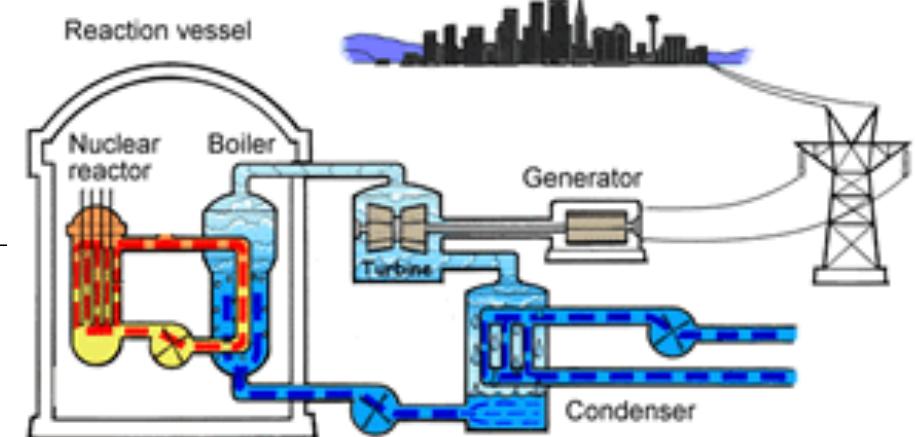
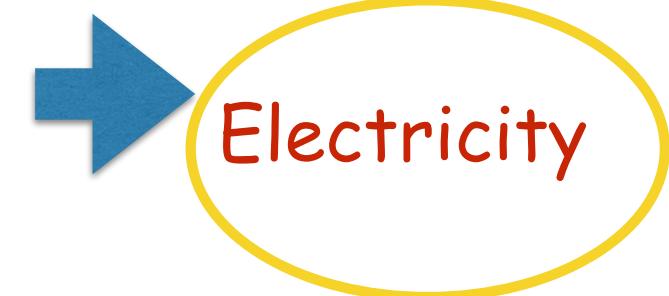
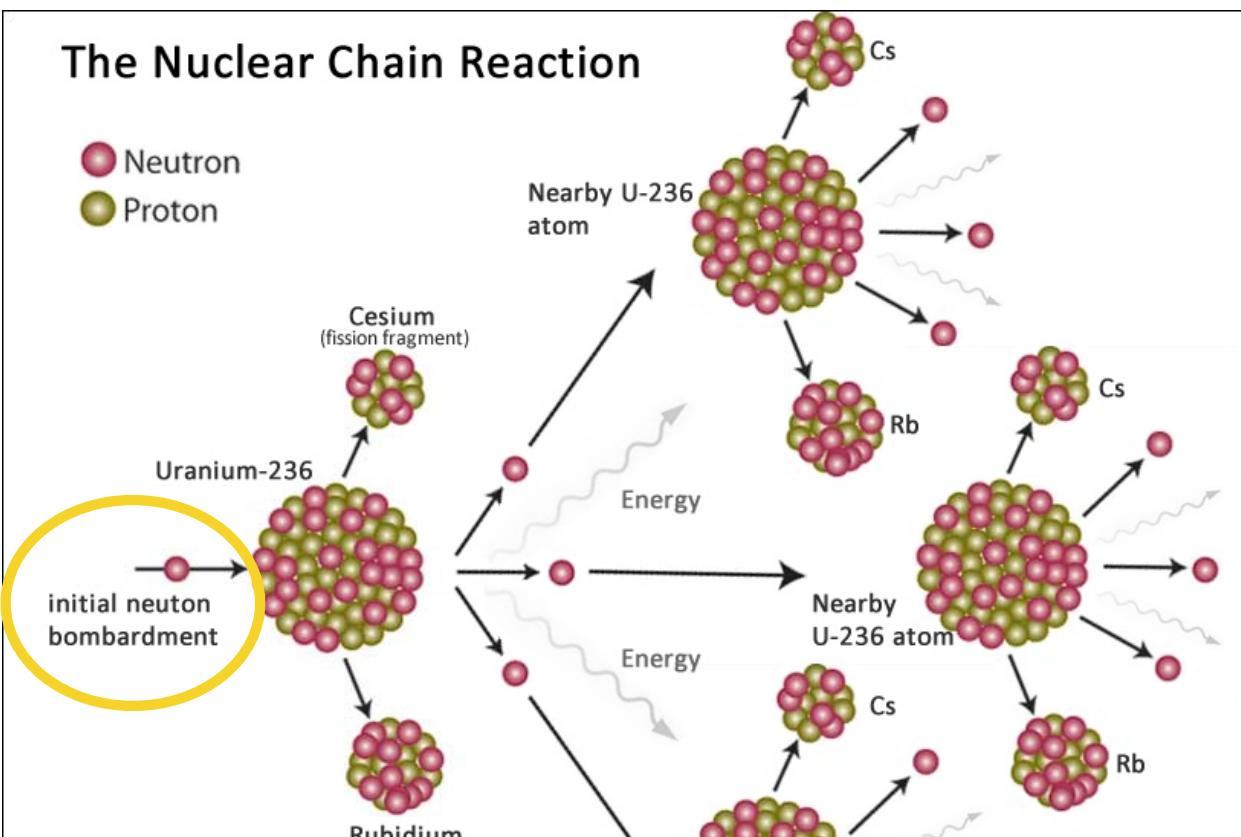


product along each path
sum of paths

$$\frac{\partial L}{\partial z_1} = \left(\frac{\partial L}{\partial a_1} \times \frac{\partial a_1}{\partial z_1} \right) + \left(\frac{\partial L}{\partial a_2} \times \frac{\partial a_2}{\partial z_1} \right) + \left(\frac{\partial L}{\partial a_3} \times \frac{\partial a_3}{\partial z_1} \right)$$

The Nuclear Chain Reaction

● Neutron
● Proton



gradient of Electricity w.r.t. Initial Neutron

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial z}^T \times \frac{\partial z}{\partial W}$$

matrix (c by p)

vector (c by 1)

matrix (c by p)

$$\left(\frac{\partial L}{\partial z_1}, \dots, \frac{\partial L}{\partial z_i}, \dots, \frac{\partial L}{\partial z_c} \right)^T \times$$

transpose
 element-wise product

| | | |
|--|---------|--|
| $\frac{\partial z_1}{\partial w_{11}}$ | \dots | $\frac{\partial z_1}{\partial w_{1p}}$ |
| $\frac{\partial z_i}{\partial w_{i1}}$ | \dots | $\frac{\partial z_i}{\partial w_{ip}}$ |
| $\frac{\partial z_c}{\partial w_{c1}}$ | \dots | $\frac{\partial z_c}{\partial w_{cp}}$ |

$$= \left(\frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial w_{11}}, \dots, \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial w_{1p}} \right)^T, \dots, \left(\frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial w_{i1}}, \dots, \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial w_{ip}} \right)^T, \dots, \left(\frac{\partial L}{\partial z_c} \frac{\partial z_c}{\partial w_{c1}}, \dots, \frac{\partial L}{\partial z_c} \frac{\partial z_c}{\partial w_{cp}} \right)^T$$

Softmax Regression (train)

initialize W and b

Loop for n_{epoch} iterations:

Loop for each training instance (x, y) in training set

forward pass to compute z , a and L for the instance

backward pass to compute local gradients

$$\frac{\partial L}{\partial a} \quad \frac{\partial a}{\partial z} \quad \frac{\partial z}{\partial b} \quad \frac{\partial z}{\partial W}$$

compute global gradients using chain rule

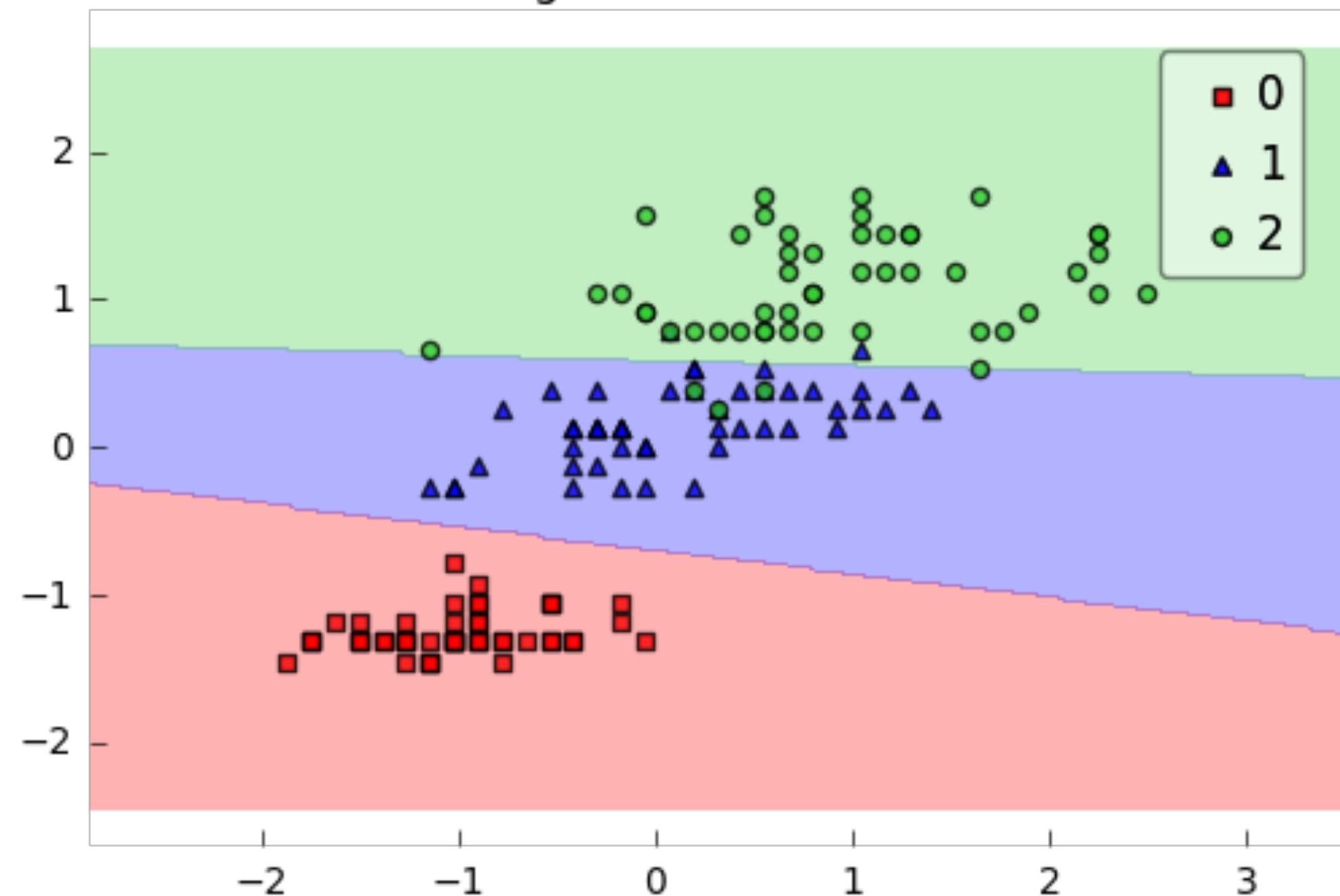
$$\frac{\partial L}{\partial W} \quad \frac{\partial L}{\partial b}$$

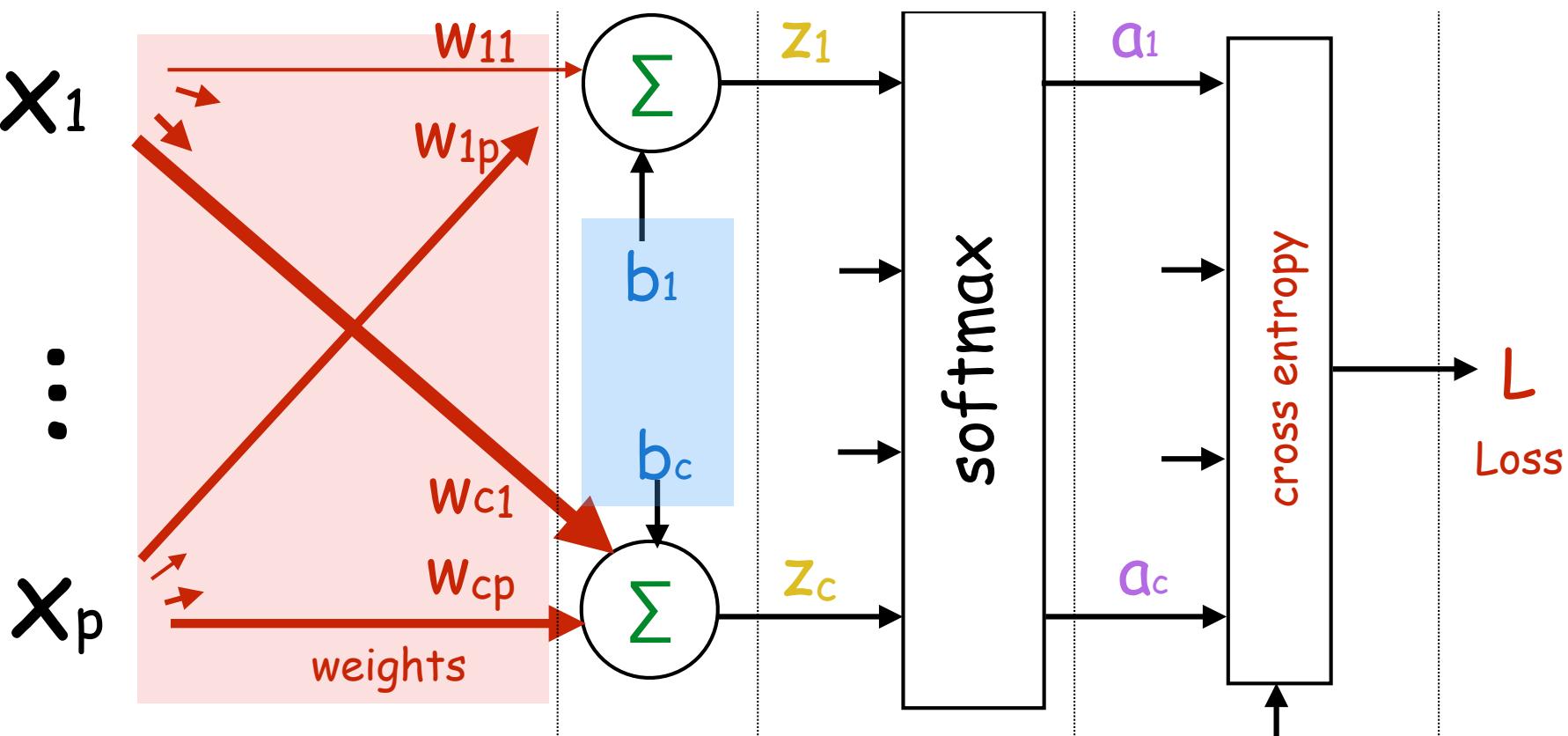
update the parameters W and b

$$W \leftarrow W - a \frac{\partial L}{\partial W}$$

$$b \leftarrow b - a \frac{\partial L}{\partial b}$$

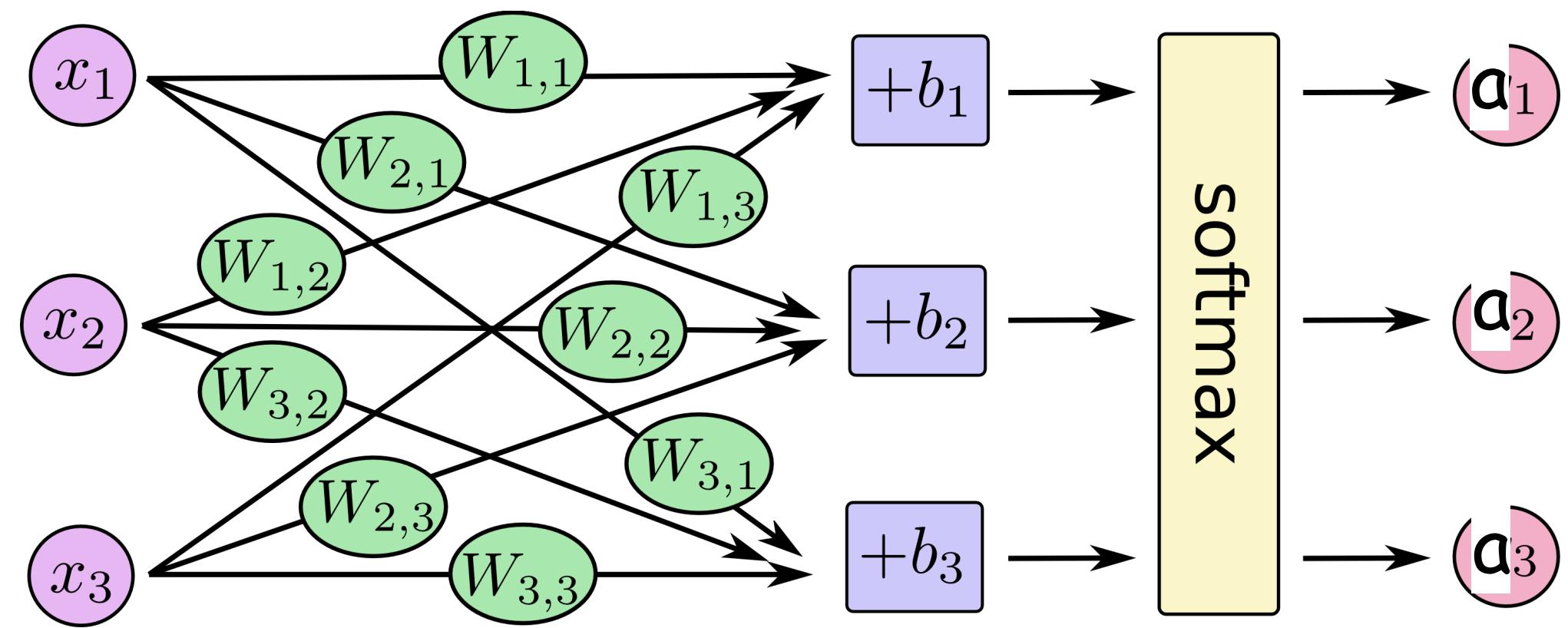
Softmax Regression - Gradient Descent





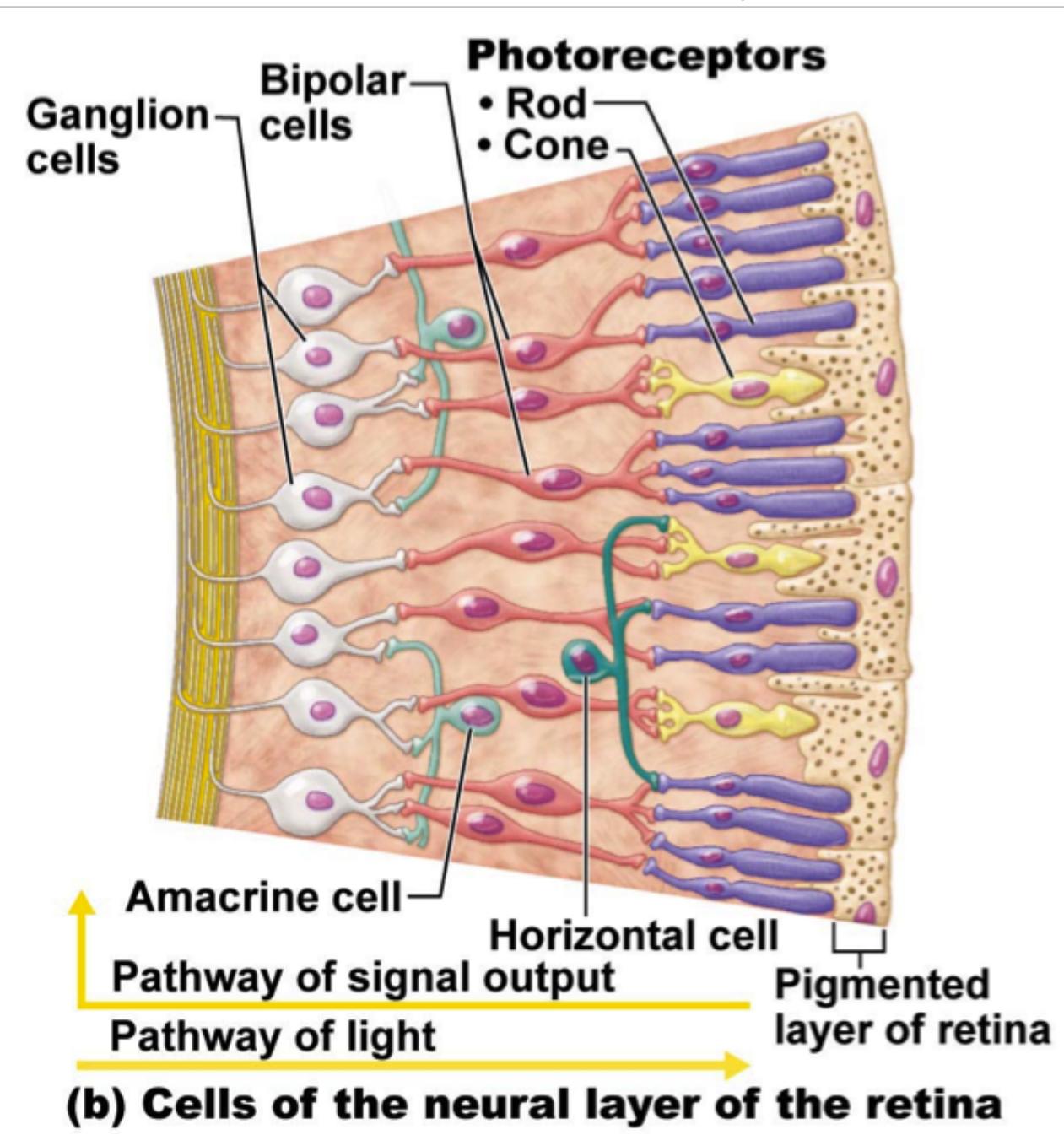
| weights | biases | logits | activations | loss |
|---------|--------|--------|-------------|------|
| w | b | z | a | L |

Softmax Regression

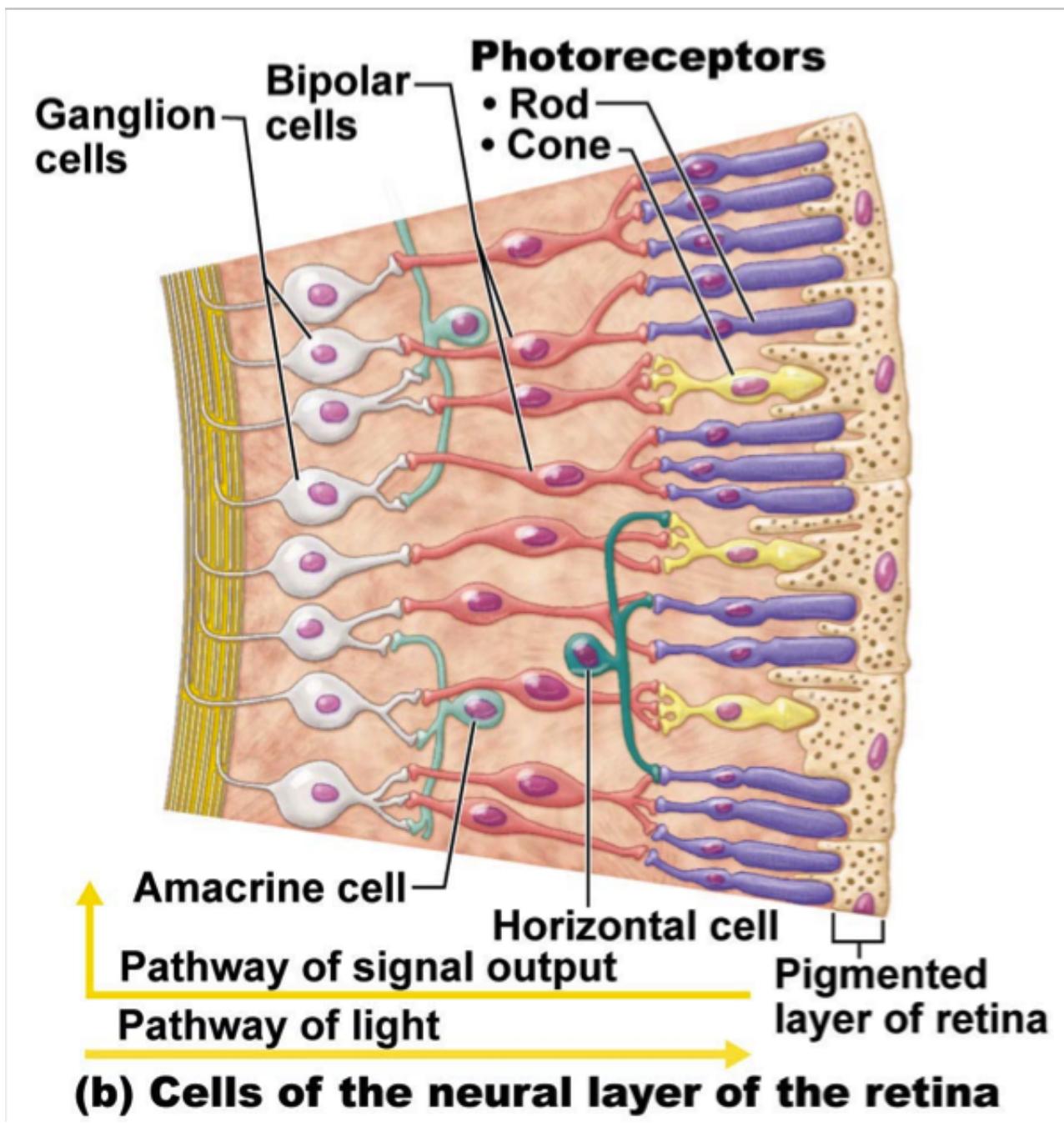


$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \text{softmax} \left[\begin{bmatrix} W_{1,1} & W_{1,2} & W_{1,3} \\ W_{2,1} & W_{2,2} & W_{2,3} \\ W_{3,1} & W_{3,2} & W_{3,3} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \right]$$

Neural Layer

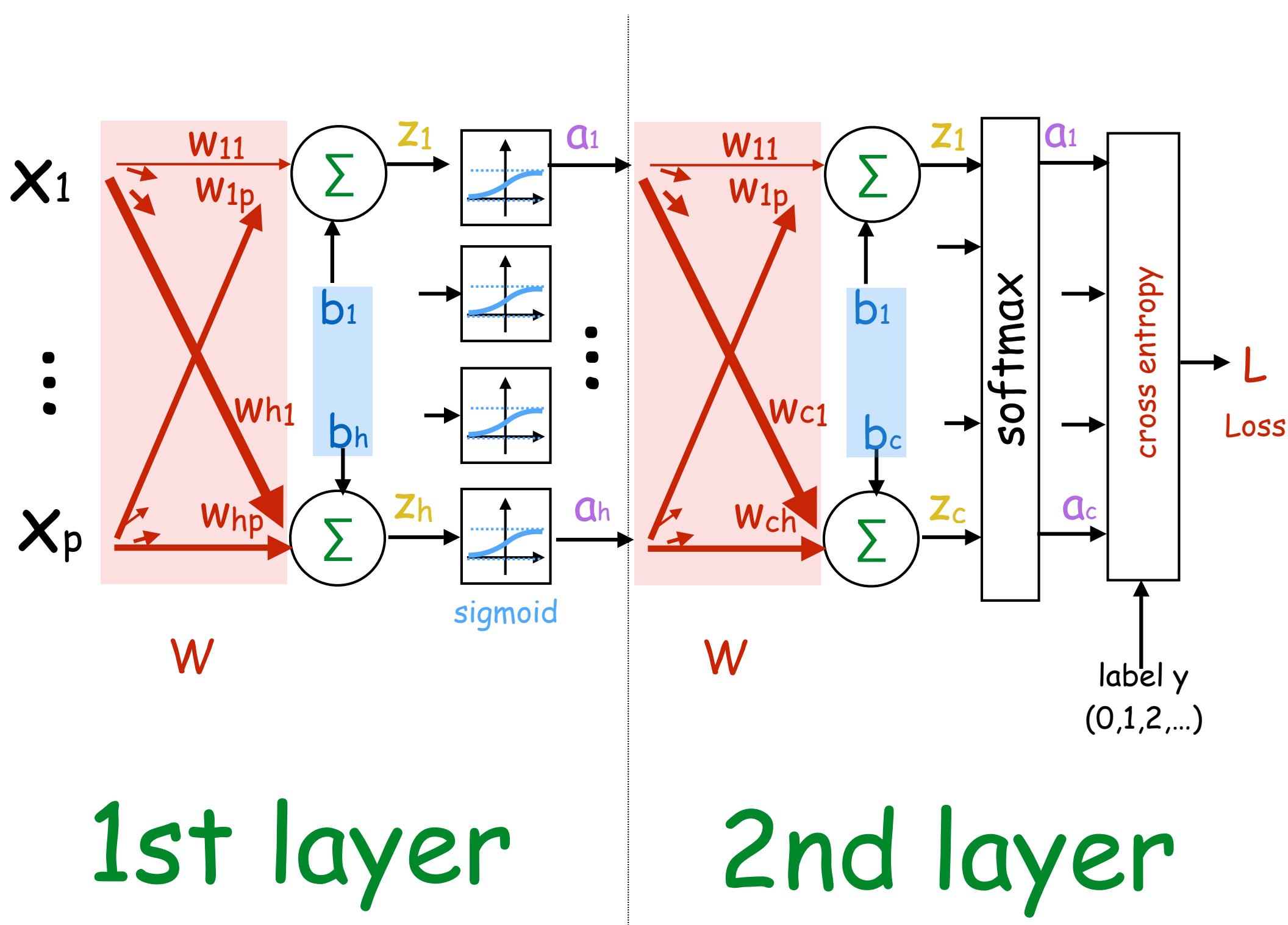


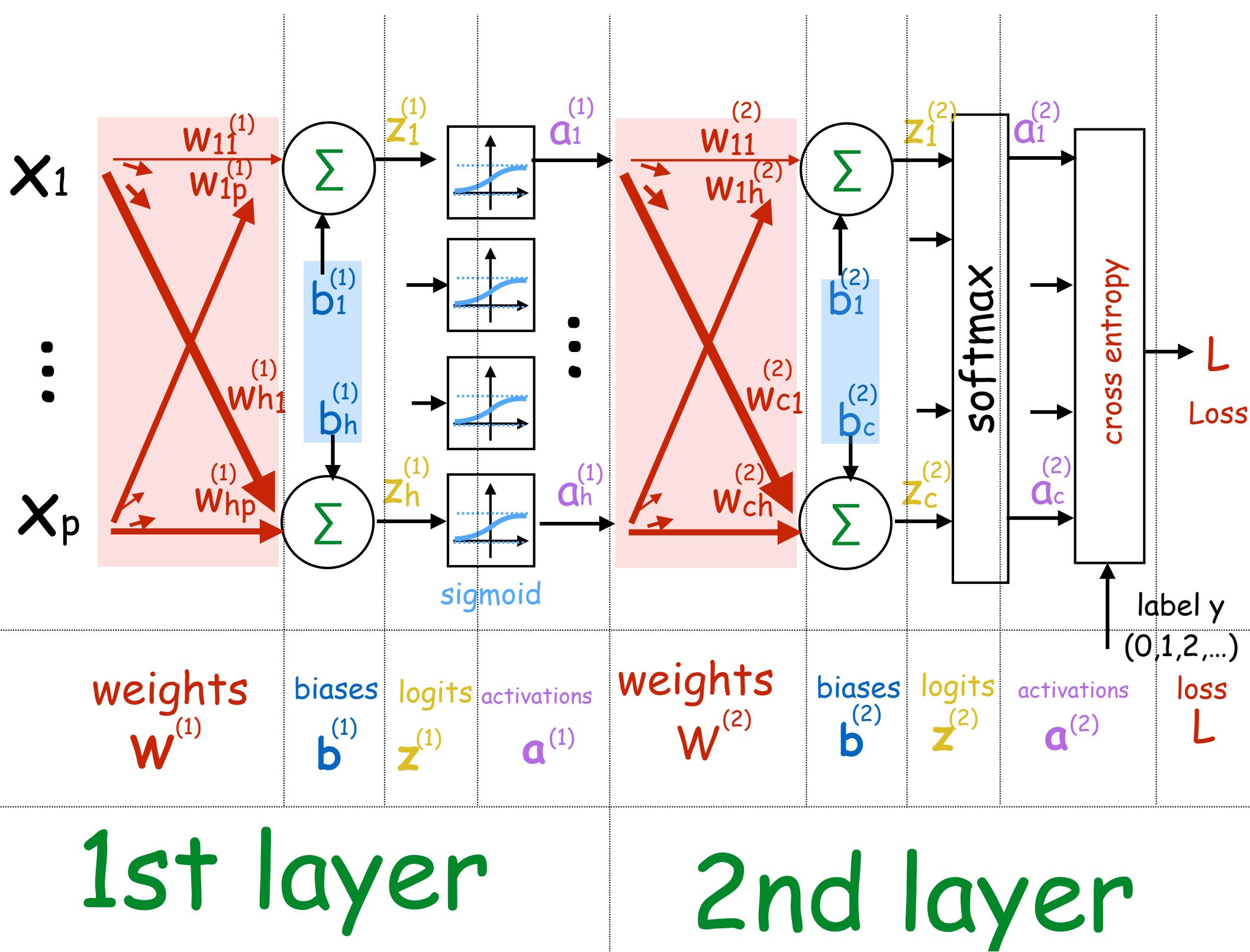
Multiple Neural Layers ?



Fully Connected Neural Network

Demo: <http://playground.tensorflow.org/>





$$\frac{\partial L}{\partial a^{(1)}} = \frac{\partial L}{\partial z^{(2)}} \times \frac{\partial z^{(2)}}{\partial a^{(1)}}$$

vector length h

vector of length c

matrix (c by h)

$$\left(\frac{\partial L}{\partial a_1^{(1)}}, \dots, \frac{\partial L}{\partial a_c^{(1)}} \right)$$

$$= \left(\frac{\partial L}{\partial z_1^{(2)}}, \dots, \frac{\partial L}{\partial z_c^{(2)}} \right) \times$$

| | | |
|---|---------|---|
| $\frac{\partial z_1^{(2)}}{\partial a_1^{(1)}}$ | \dots | $\frac{\partial z_1^{(2)}}{\partial a_h^{(1)}}$ |
| \dots | | \dots |
| $\frac{\partial z_c^{(2)}}{\partial a_1^{(1)}}$ | \dots | $\frac{\partial z_c^{(2)}}{\partial a_h^{(1)}}$ |