# Machine Learning

CS 539

Worcester Polytechnic Institute

Department of Computer Science

Instructor: Prof. Kyumin Lee

# HW1 grading

# HW3

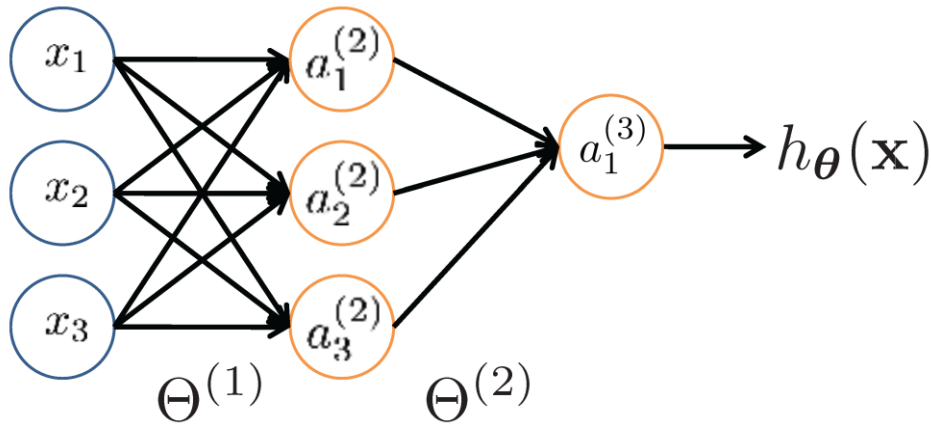- https://canvas.wpi.edu/courses/57384/assignments/340697?module_item_id=1073181

- Due date is Feb 20.

# Vectorization

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right) = g\left(z_1^{(2)}\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right) = g\left(z_2^{(2)}\right)$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right) = g\left(z_3^{(2)}\right)$$

$$h_\Theta(\mathbf{x}) = g\left(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)}\right) = g\left(z_1^{(3)}\right)$$
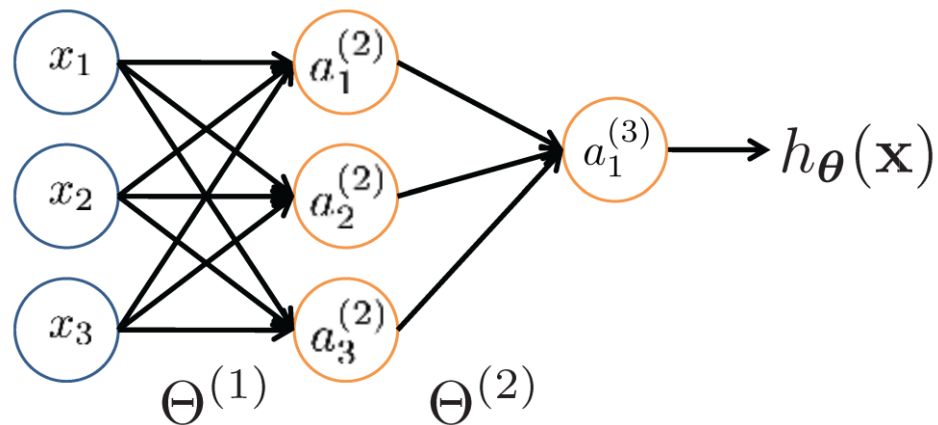


Feed-Forward Steps:
$$\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{x}$$
$$\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$
$$\text{Add } a_0^{(2)} = 1$$
$$\mathbf{z}^{(3)} = \Theta^{(2)}\mathbf{a}^{(2)}$$
$$h_\Theta(\mathbf{x}) = \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$

Feed-Forward Steps:
$$\mathbf{z}^{(2)} = \Theta^{(1)}\mathbf{x}$$
$$\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$$
Add $a_0^{(2)} = 1$
$$\mathbf{z}^{(3)} = \Theta^{(2)}\mathbf{a}^{(2)}$$
$$h_\Theta(\mathbf{x}) = \mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$$

where $\mathrm{cost}(\mathbf{x}_i) = -y_i \log h_\Theta(\mathbf{x}_i) - (1-y_i)\log(1 - h_\Theta(\mathbf{x}_i))$

$$\frac{\partial\, cost(x_i)}{\partial\, \theta^{(2)}} = \frac{\partial\, cost(x_i)}{\partial\, a^{(3)}} \cdot \frac{\partial\, a^{(3)}}{\partial\, z^{(3)}} \cdot \frac{\partial\, z^{(3)}}{\partial\, \theta^{(2)}} = \frac{a^{(3)}-y}{a^{(3)}(1-a^{(3)})} \cdot a^{(3)}(1-a^{(3)}) \cdot a^{(2)} = a^{(2)} \cdot (a^{(3)}-y)$$

$$\frac{\partial\, cost(x_i)}{\partial\, \theta^{(1)}} = \frac{\partial\, cost(x_i)}{\partial\, a^{(3)}} \cdot \frac{\partial\, a^{(3)}}{\partial\, z^{(3)}} \cdot \frac{\partial\, z^{(3)}}{\partial\, a^{(2)}} \cdot \frac{\partial\, a^{(2)}}{\partial\, z^{(2)}} \cdot \frac{\partial\, z^{(2)}}{\partial\, \theta^{(1)}} = (a^{(3)}-y)\cdot \theta^{(2)} \cdot a^{(2)}(1-a^{(2)}) \cdot x_i$$
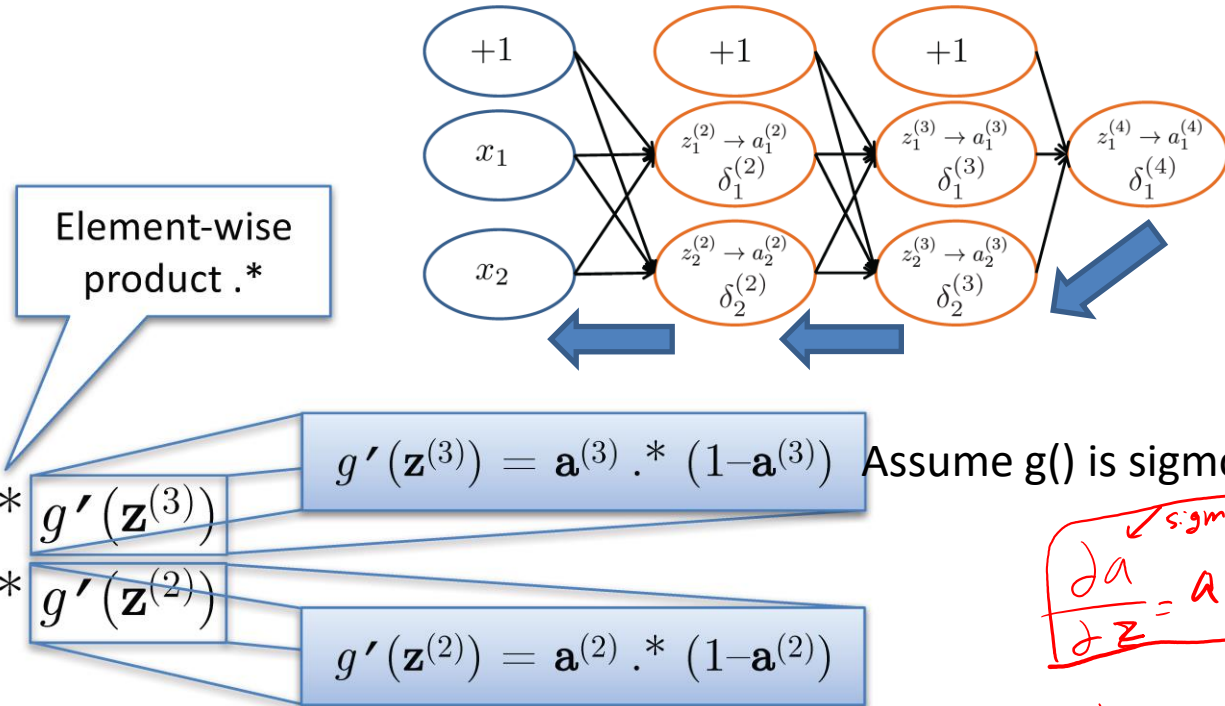
# Backpropagation: Gradient Computation (binary classification)

Let $\delta_j^{(l)}$ = "error" of node $j$ in layer $l$

(#layers $L$ = 4)



## Backpropagation

- $\boldsymbol{\delta}^{(4)} = \boldsymbol{a}^{(4)} - \mathbf{y}$
- $\boldsymbol{\delta}^{(3)} = (\Theta^{(3)})^{\mathsf{T}}\boldsymbol{\delta}^{(4)} .* g'(\mathbf{z}^{(3)})$
- $\boldsymbol{\delta}^{(2)} = (\Theta^{(2)})^{\mathsf{T}}\boldsymbol{\delta}^{(3)} .* g'(\mathbf{z}^{(2)})$
- (No $\boldsymbol{\delta}^{(1)}$)

Element-wise product .*

$g'(\mathbf{z}^{(3)}) = \mathbf{a}^{(3)} .* (1-\mathbf{a}^{(3)})$  Assume g() is sigmoid

$g'(\mathbf{z}^{(2)}) = \mathbf{a}^{(2)} .* (1-\mathbf{a}^{(2)})$

*(handwritten annotations in red:)*

$\dfrac{\partial \text{cost}(\lambda_i)}{\partial a^{(4)}} \times \dfrac{\partial a^{(4)}}{\partial z^{(4)}} =$

$\dfrac{\partial \text{cost}(\lambda_i)}{\partial a^{(4)}} \cdot \dfrac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \dfrac{\partial z^{(4)}}{\partial a^{(3)}} \cdot \dfrac{\partial a^{(3)}}{\partial z^{(3)}}$

same $\times \dfrac{\partial z^{(3)}}{\partial a^{(2)}} \times \dfrac{\partial a^{(2)}}{\partial z^{(2)}}$

$\checkmark$ sigmoid $\dfrac{\partial a}{\partial z} = a(1-a)$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

$= \dfrac{\partial \text{cost}(\lambda_i)}{\partial a_i^{(l+1)}} \cdot \dfrac{\partial a_i^{(l+1)}}{\partial z_i^{(l+1)}} \cdot \dfrac{\partial z_i^{(l+1)}}{\partial \Theta_{ij}^{(l)}}$

(ignoring $\lambda$; if $\lambda = 0$)

Mathematical Proof: https://web.cs.wpi.edu/~kmlee/cs539/cs229-notes-deep_learning.pdf

# Backpropagation

Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$ $\qquad$ (Used to accumulate gradient)

For each training instance $(\mathbf{x}_i, y_i)$:

$\qquad$ Set $\mathbf{a}^{(1)} = \mathbf{x}_i$

$\qquad$ Compute $\{\mathbf{a}^{(2)}, \ldots, \mathbf{a}^{(L)}\}$ via forward propagation

$\qquad$ Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y_i$

$\qquad$ Compute errors $\{\boldsymbol{\delta}^{(L-1)}, \ldots, \boldsymbol{\delta}^{(2)}\}$

$\qquad$ Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n} \Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

$\boldsymbol{D}^{(l)}$ is the matrix of partial derivatives of $J(\Theta)$

Note: Can vectorize $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ as $\boldsymbol{\Delta}^{(l)} = \boldsymbol{\Delta}^{(l)} + \boldsymbol{\delta}^{(l+1)} \mathbf{a}^{(l)\mathsf{T}}$

# Training a Neural Network via Gradient Descent with Backprop

Given: training set $\{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)\}$
Initialize all $\Theta^{(l)}$ randomly (NOT to 0!)
Loop // each iteration is called an epoch
    Set $\Delta_{ij}^{(l)} = 0 \quad \forall l, i, j$             (Used to accumulate gradient)
    For each training instance $(\mathbf{x}_i, y_i)$:
        Set $\mathbf{a}^{(1)} = \mathbf{x}_i$
        Compute $\{\mathbf{a}^{(2)}, \ldots, \mathbf{a}^{(L)}\}$ via forward propagation
        Compute $\boldsymbol{\delta}^{(L)} = \mathbf{a}^{(L)} - y_i$
        Compute errors $\{\boldsymbol{\delta}^{(L-1)}, \ldots, \boldsymbol{\delta}^{(2)}\}$
        Compute gradients $\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

    Compute avg regularized gradient $D_{ij}^{(l)} = \begin{cases} \frac{1}{n}\Delta_{ij}^{(l)} + \lambda\Theta_{ij}^{(l)} & \text{if } j \neq 0 \\ \frac{1}{n}\Delta_{ij}^{(l)} & \text{otherwise} \end{cases}$

    Update weights via gradient step $\Theta_{ij}^{(l)} = \Theta_{ij}^{(l)} - \alpha D_{ij}^{(l)}$
Until weights converge or max #epochs is reached

Backpropagation

# Multi-Class Classification
## (Softmax regression and Fully Connected Neural Network)

# Softmax Regression

:a generalization of logistic regression to the case where we want to handle multiple classes

# Logistic Regression

$x_1$

$w_1$

$\otimes$

...

... ...

$w_p$

...

$x_p$

$\otimes$

weights

$\Sigma$

$b$

bias

$z$ logit

non-linear
activation function

1

0

sigmoid

$a$

activation

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

# Multi-class Classification Problem



data set

2→2, 5→5, 6→8, 0→0, 2→2, 7→7, 5→5, 1→1,
3→3, 0→0, 3→3, 9→9, 6→6, 2→2, 8→8, 2→2,
0→0, 6→6, 6→6, 1→1, 1→1, 7→7, 8→8, 5→5,
0→0, 4→4, 7→7, 6→6, 0→0, 2→2, 5→5,
3→3, 1→1, 5→5, 6→6, 7→7, 5→5, 4→4, 1→1,
9→9, 3→3, 6→6, 8→8, 0→0, 9→9, 3→3,
0→0, 3→3, 7→7, 4→4, 4→4, 3→3, 8→8, 0→0,
4→4, 1→1, 3→3, 7→7, 6→6, 4→4, 7→7, 2→2,
7→7, 2→2, 5→5, 2→2, 0→0, 9→9, 8→8, 9→9,
8→8, 1→1, 6→6, 4→4, 8→8, 5→5, 8→8,
0→0, 6→6, 7→7, 4→4, 5→5, 8→8, 4→4,
3→3, 1→1, 5→5, 1→1, 9→9, 9→9, 9→9, 2→2,
4→4, 7→7, 3→3, 1→1, 9→9, 2→2, 9→9, 6→6 }]

candidate labels
(classes)

0
1
2
⋮
9

input
(instance)

output
(label)

0 or 1 or … or 9

# Feature Matrix X
## (n by p)

## Label Vector y
## (length n)

**Feature (pixel)**

|  | 1 | 2 | ... | | | | | | | | | p | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $x_1$ | 2 | 1 | 3 | 4 | 3 | 8 | 3 | 5 | 7 | 9 | 7 | 3 | 4 | $y_1$ | 2 |
| $x_2$ | 5 | 3 | 3 | 5 | 7 | 7 | 0 | 4 | 1 | 2 | 1 | 9 | 7 | $y_2$ | 5 |
|  | 0 | 7 | 0 | 4 | 1 | 1 | 4 | 3 | 7 | 8 | 6 | 2 | 7 | | 0 |
|  | 3 | 1 | 4 | 3 | 7 | 9 | 7 | 3 | 2 | 7 | 0 | 4 | 1 | | 3 |
|  | 8 | 7 | 7 | 3 | 2 | 2 | 1 | 9 | 8 | 1 | 4 | 3 | 7 | | 8 |
|  | 0 | 2 | 1 | 9 | 8 | 8 | 6 | 2 | 0 | 7 | 7 | 3 | 2 | | 0 |
| $\vdots$ | 4 | 8 | 6 | 2 | 0 | 0 | 4 | 1 | 1 | 4 | 1 | 9 | 8 | $\vdots$ | 9 |
|  | 6 | 0 | 2 | 1 | 4 | 1 | 3 | 7 | 9 | 7 | 6 | 2 | 0 | | 6 |
|  | 7 | 3 | 5 | 3 | 3 | 7 | 3 | 2 | 2 | 1 | 2 | 1 | 3 | | 7 |
| $x_n$ | 6 | 1 | 7 | 2 | 3 | 2 | 2 | 1 | 2 | 3 | 5 | 3 | 1 | $y_n$ | 6 |

**Instance**

n – # instances   p – # features

$x_1$

$w_{11}$

$\Sigma$

$z_1$ logit

1

0

$a_1$ activation

$w_{1p}$

$b_1$ bias

$w_{c1}$

bias $b_c$

$w_{cp}$ weights

$x_p$

$\Sigma$

$z_c$ logit

1

0

$a_c$ activation

$W$

$x_1$

$w_{11}$

$w_{1p}$

$\sum$

$b_1$
bias

$z_1$
logit

1

0

$a_1$
activation

$w_{c1}$

$b_c$

$w_{cp}$
weights

$\sum$

$z_c$
logit

1

0

$a_c$
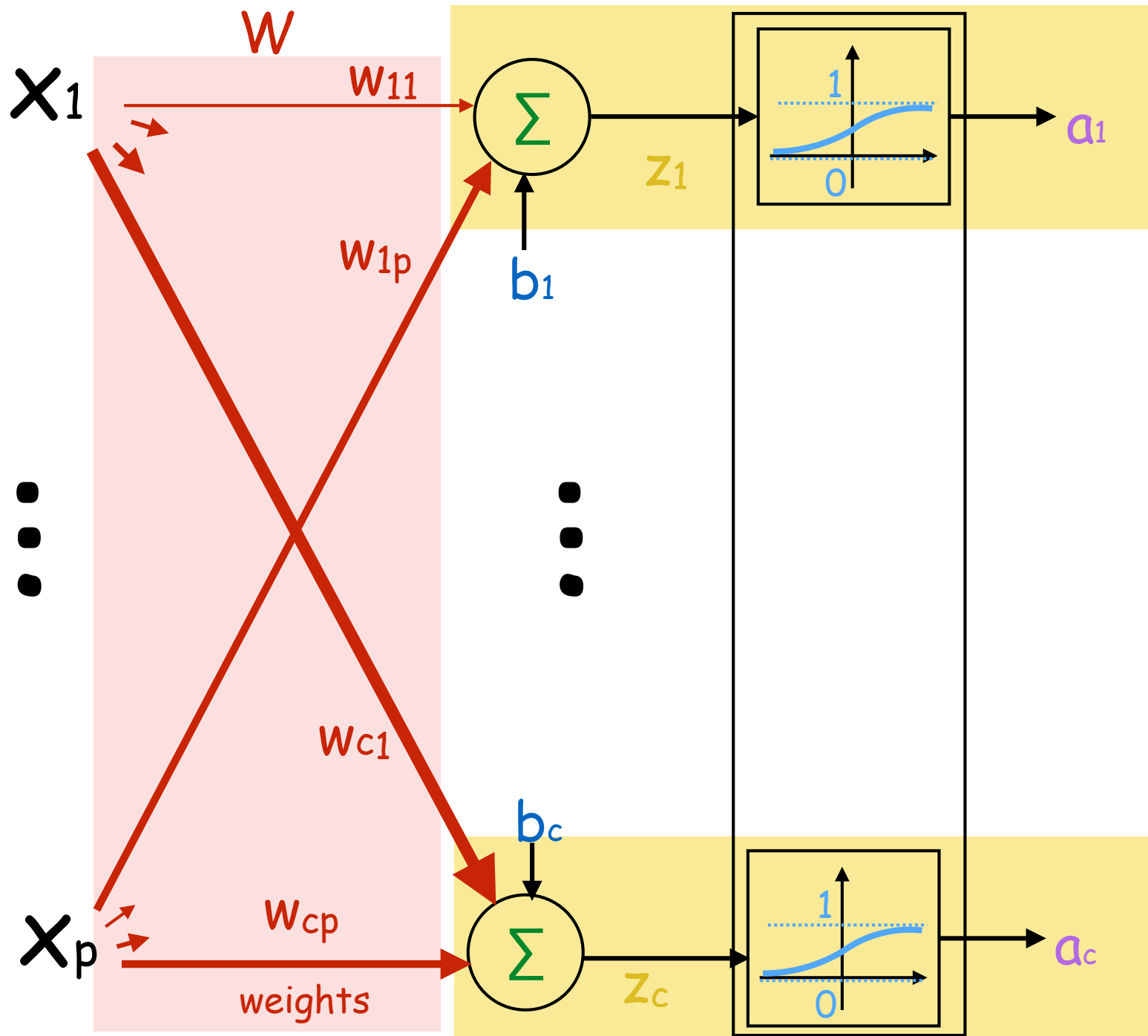activation

$x_p$

Weight Matrix

# Independent Outputs

# Independent Outputs

# The outputs we need



W

$x_1$

$w_{11}$

$w_{1p}$

$\Sigma$   $z_1$

$b_1$

$a_1$   0.6   ✔

0.1   ✘

0.1   ✘

$w_{c1}$

$b_c$

$w_{cp}$

weights

$\Sigma$   $z_c$

0.1   ✘

0.1   ✘

$a_c$   ✘

$x_p$

# The outputs we need

# Coordinated Outputs

# Softmax Activation

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

$W$

$X_1$

$w_{11}$

$w_{1p}$

$\Sigma$  $b_1$  $z_1$

$\frac{e^{z_1}}{\sum\limits_i e^{z_i}}$  $a_1$

$\frac{e^{z_2}}{\sum\limits_i e^{z_i}}$  $a_2$

softmax

$w_{c1}$

$b_c$

$w_{cp}$

$X_p$

weights

$\Sigma$  $z_c$

$\frac{e^{z_{c-1}}}{\sum\limits_i e^{z_i}}$  $a_{c-1}$

$\frac{e^{z_c}}{\sum\limits_i e^{z_i}}$  $a_c$

# Softmax Activation

$W$

$x_1$

$w_{11}$

$w_{1p}$

$\Sigma$ $z_1$

$b_1$

$x_p$

$w_{c1}$

$w_{cp}$

weights

$\Sigma$ $z_c$

$b_c$

softmax

$$\frac{e^{z_1}}{\sum\limits_i e^{z_i}}$$

$$\frac{e^{z_2}}{\sum\limits_i e^{z_i}}$$

$$\frac{e^{z_{c-1}}}{\sum\limits_i e^{z_i}}$$

$$\frac{e^{z_c}}{\sum\limits_i e^{z_i}}$$

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

$a_1$

$a_2$

$a_{c-1}$

$a_c$

# Softmax Activation

$z_1$   3   $e$   $e^{z_1}$   20   $\div$   0.88   $a_1$

$z_2$   1   $e$   $e^{z_2}$   2.7   $\div$   0.12   $a_2$

$z_3$   -3   $e$   $e^{z_3}$   0.05   $\div$   0.00   $a_3$

$\sum$   $\sum_i e^{z_i}$

22.75

logits

## Softmax

sum =1

$0 \le a_i \le 1$

Probabilities

# Parameters **W**, b

# Training Model

## training set

## learn parameters



Feature Matrix X (n by p)

Label Vector Y

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

Weight Matrix **W** (shape c by p)

bias vector b    (length c)

# negative log Likelihood



Label: Y    3                    Y    1
observed data

Outputs: $a_1$    .4                    $a_1$    .4
probabilities
         $a_2$    .2                    $a_2$    .2

         $a_3$    .1                    $a_3$    .1

         $a_c$    .1                    $a_c$    .1

Weights **w**
biases  b

- log Likelihood = $-\log a_y$

Multi-class <u>cross entropy</u> loss

# Multi-Class Cross Entropy Loss



$x_1$

$\vdots$

$x_p$

$w_{11}$

$w_{1p}$

$w_{c1}$

$w_{cp}$

W

weights

$\Sigma$

$\Sigma$

$b_1$

$b_c$

$z_1$

$z_c$

softmax

$a_1$

$a_c$

cross entropy

$-\log a_y$

L

Loss

label y

(1,2,..., or, c)

# Multi-Class Cross Entropy Loss

# Gradient Descent

L= loss (w)



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

$$b \longleftarrow b - a\frac{\partial L}{\partial b}$$
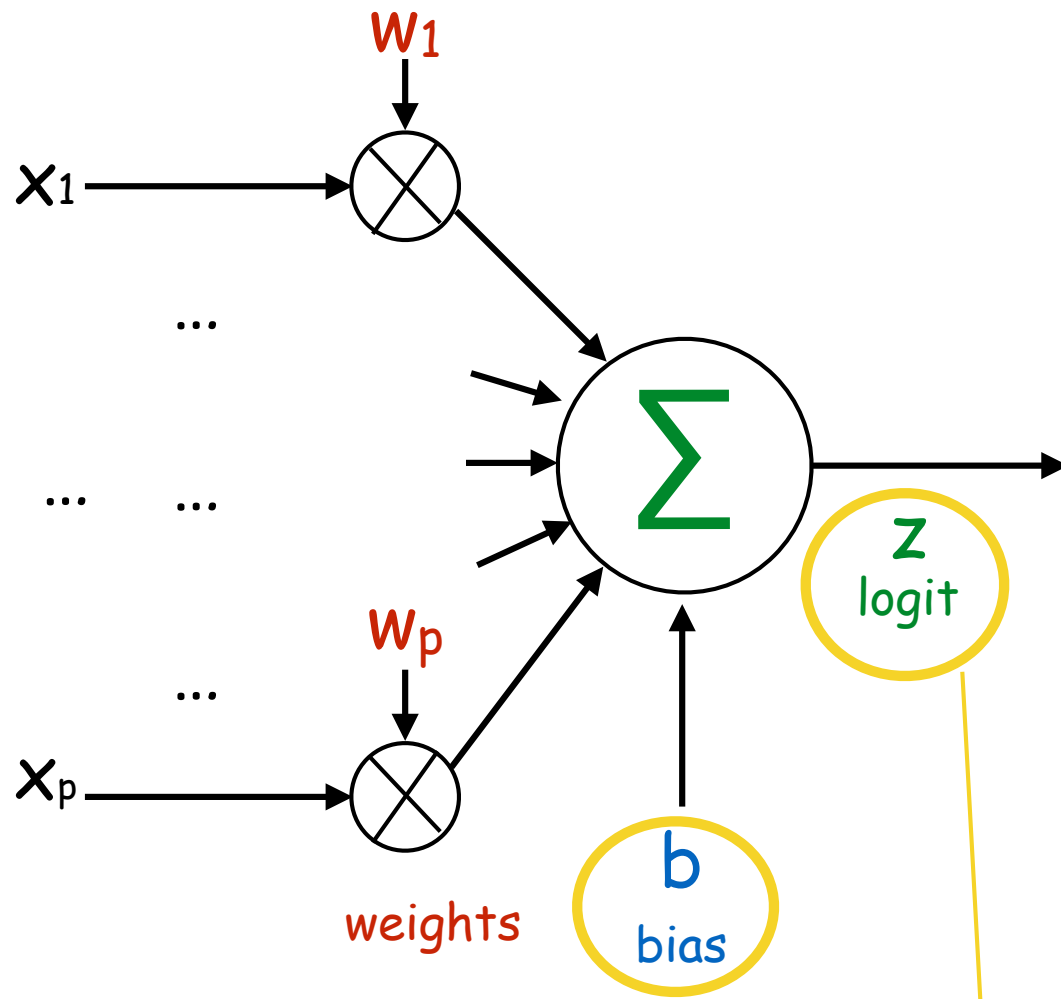
a   step size (a constant scalar)

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

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

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

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

Example

$x_1$

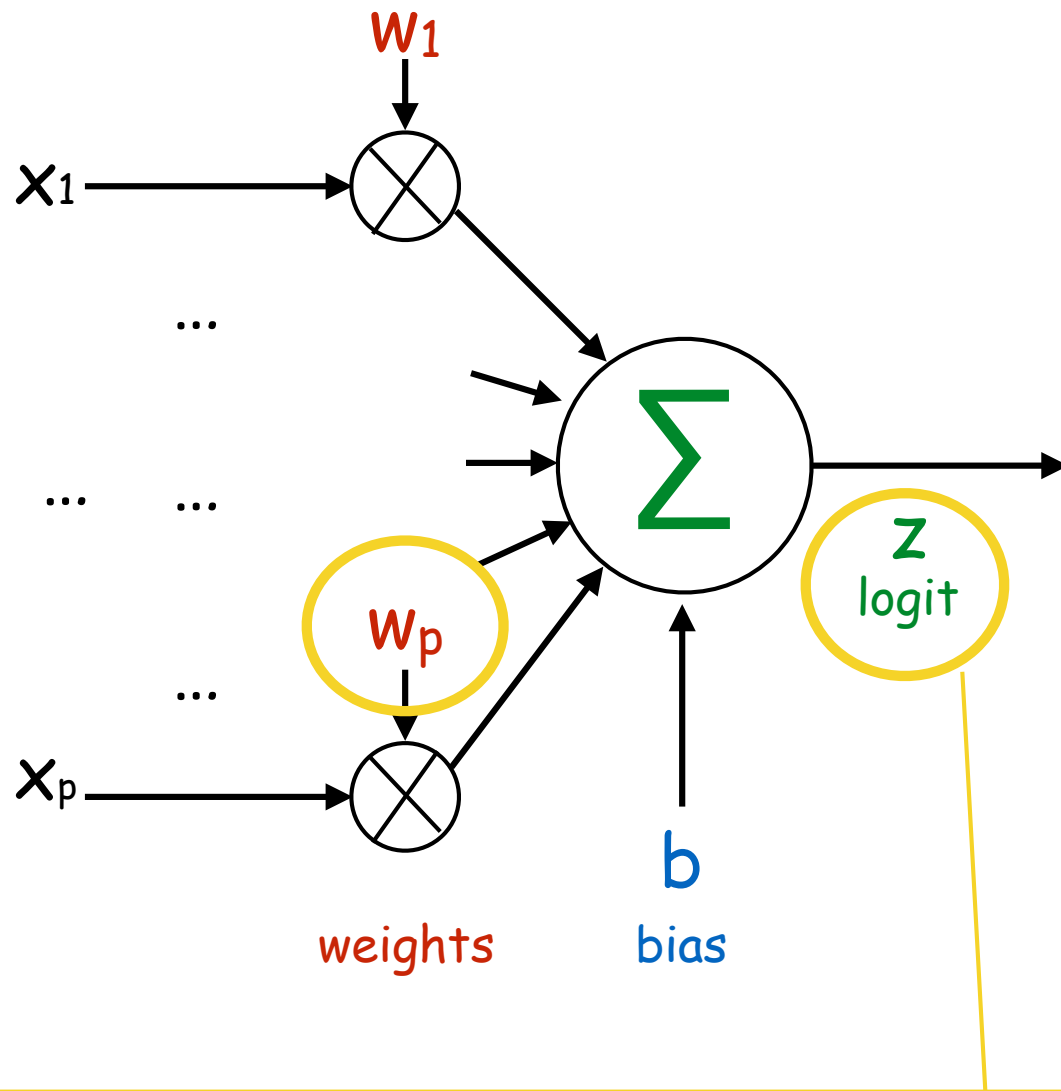$w_1$

...

... ...

$w_p$

...

$x_p$

weights

$\Sigma$

z
logit

b
bias

gradient

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

# Example

$x_1$

$w_1$

...

...     ...

$w_p$

...

$x_p$

$\sum$

z
logit

b
bias

weights

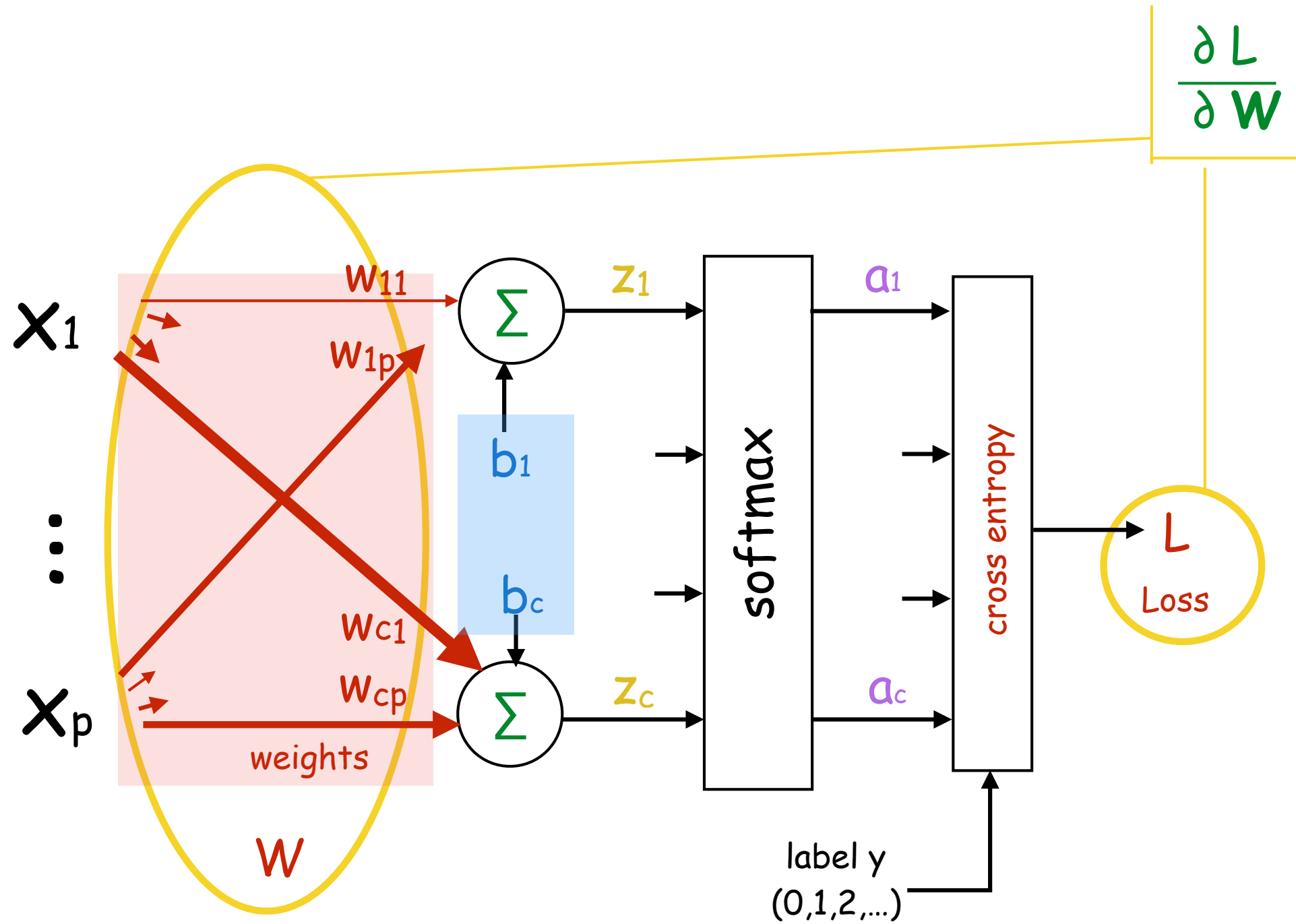gradient

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

is a function of $x_p$

Example

$$\frac{\partial z}{\partial \mathbf{w}} = \left( \frac{\partial z}{\partial w_1}, \frac{\partial z}{\partial w_2}, \ldots, \frac{\partial z}{\partial w_p} \right)$$
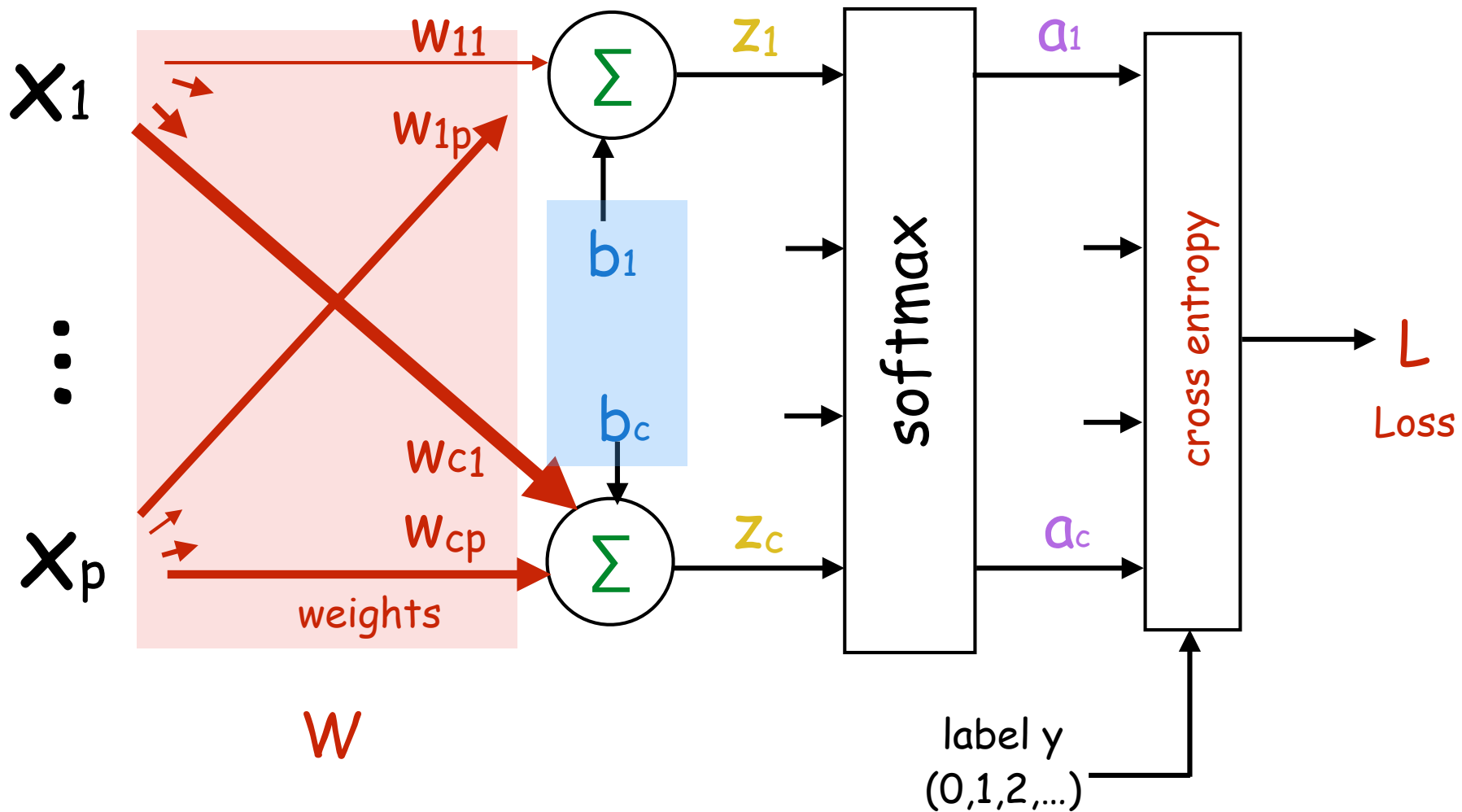
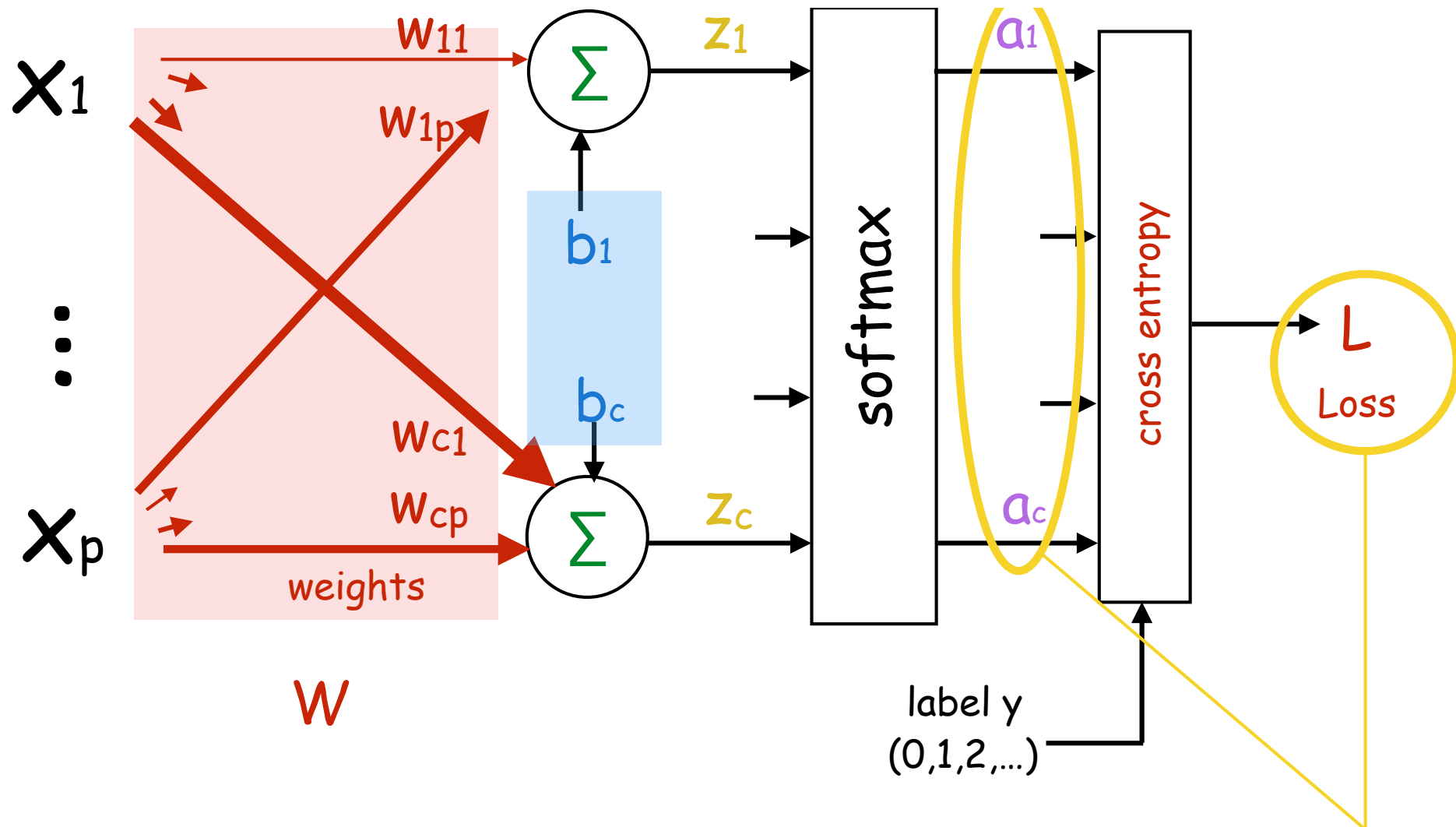$$= \left( x_1, x_2, \ldots, x_p \right) = \mathbf{x}$$

gradient

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

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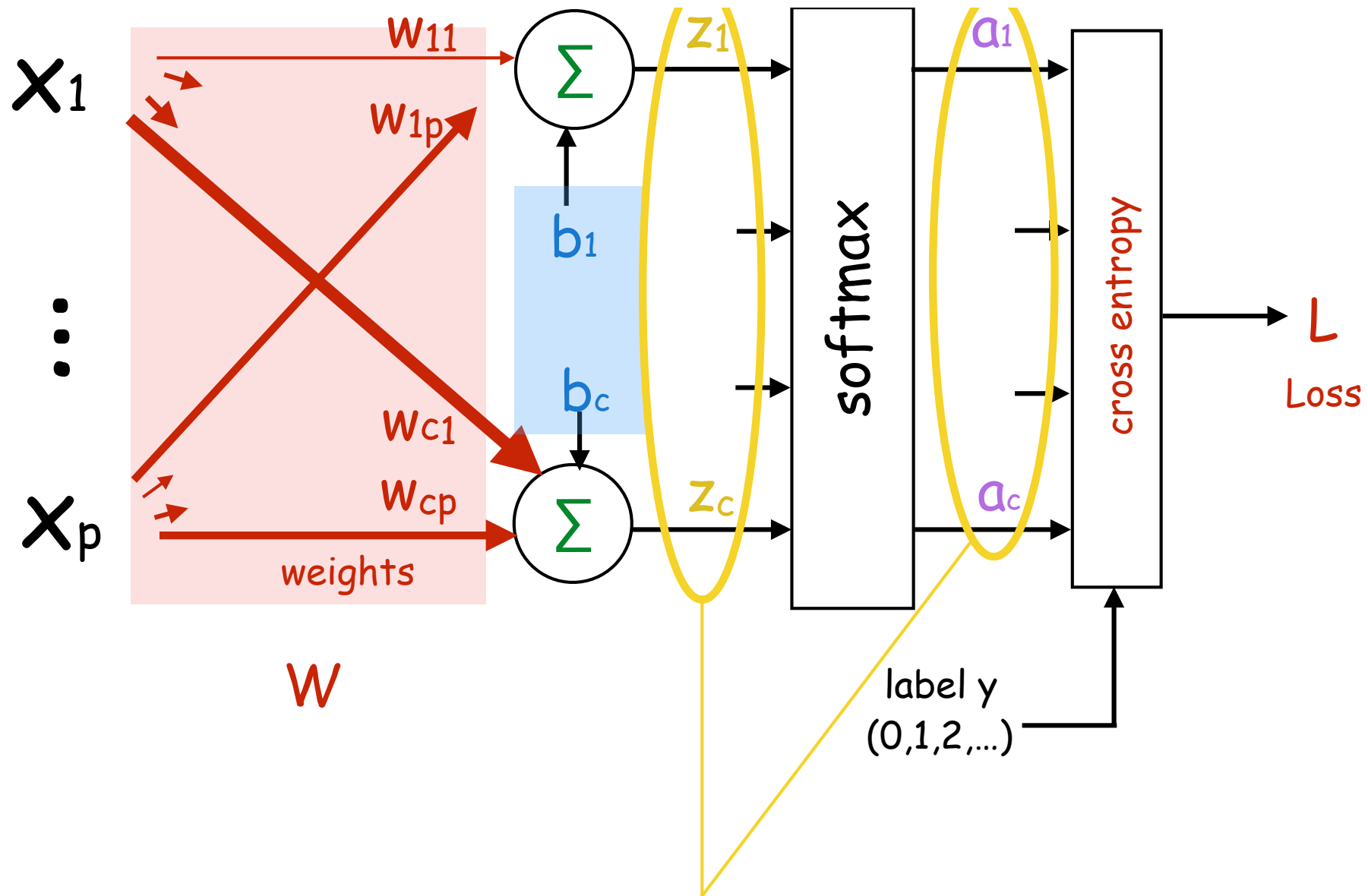
Chain Rule

$$\frac{\partial L}{\partial a} = \left( \frac{\partial L}{\partial a_1}, \dots, \frac{\partial L}{\partial a_y}, \dots, \frac{\partial L}{\partial a_c} \right)$$

$$= \left( 0, \dots, -\frac{1}{a_y}, \dots, 0 \right)$$

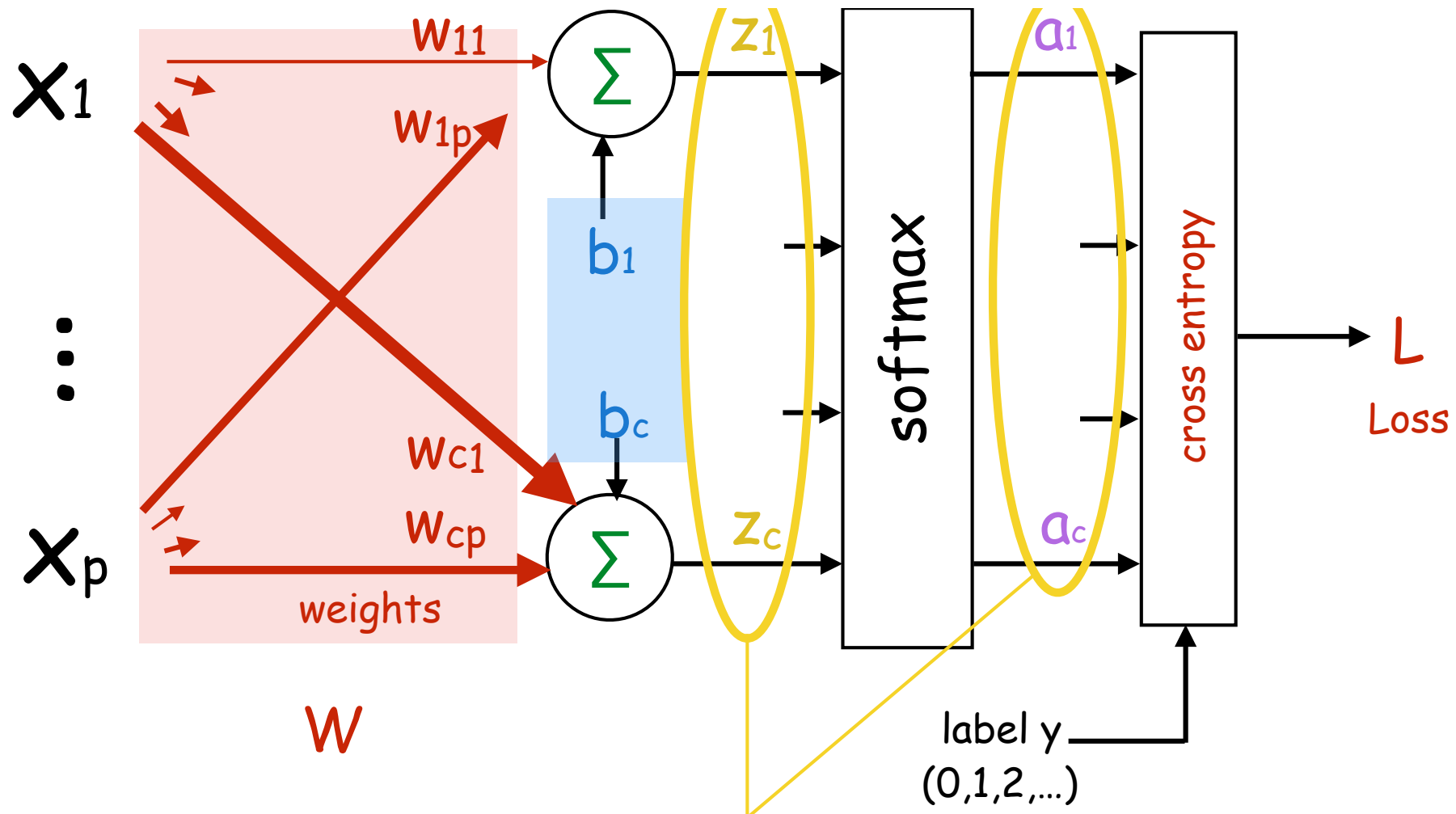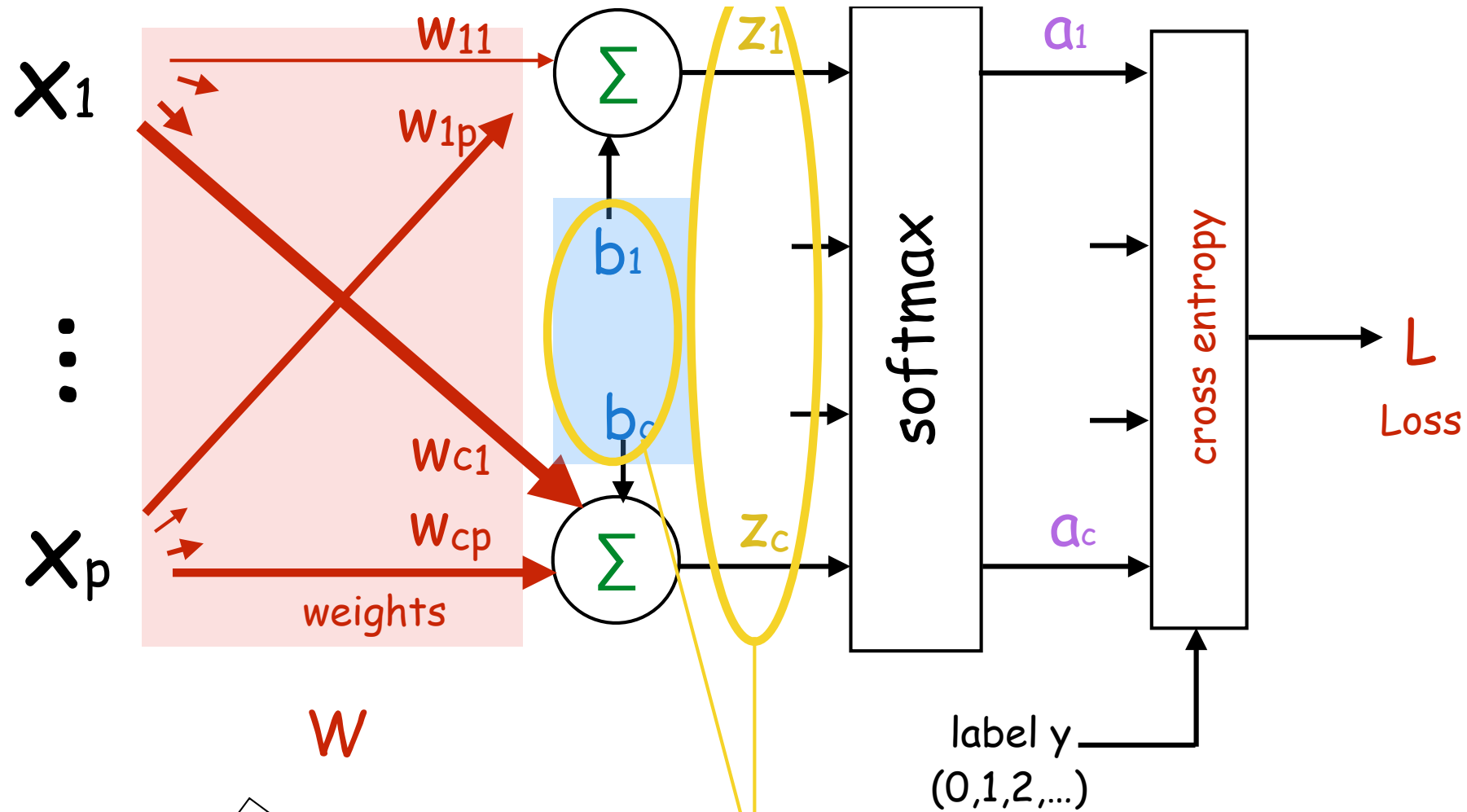$$\frac{\partial a}{\partial z} = ?$$

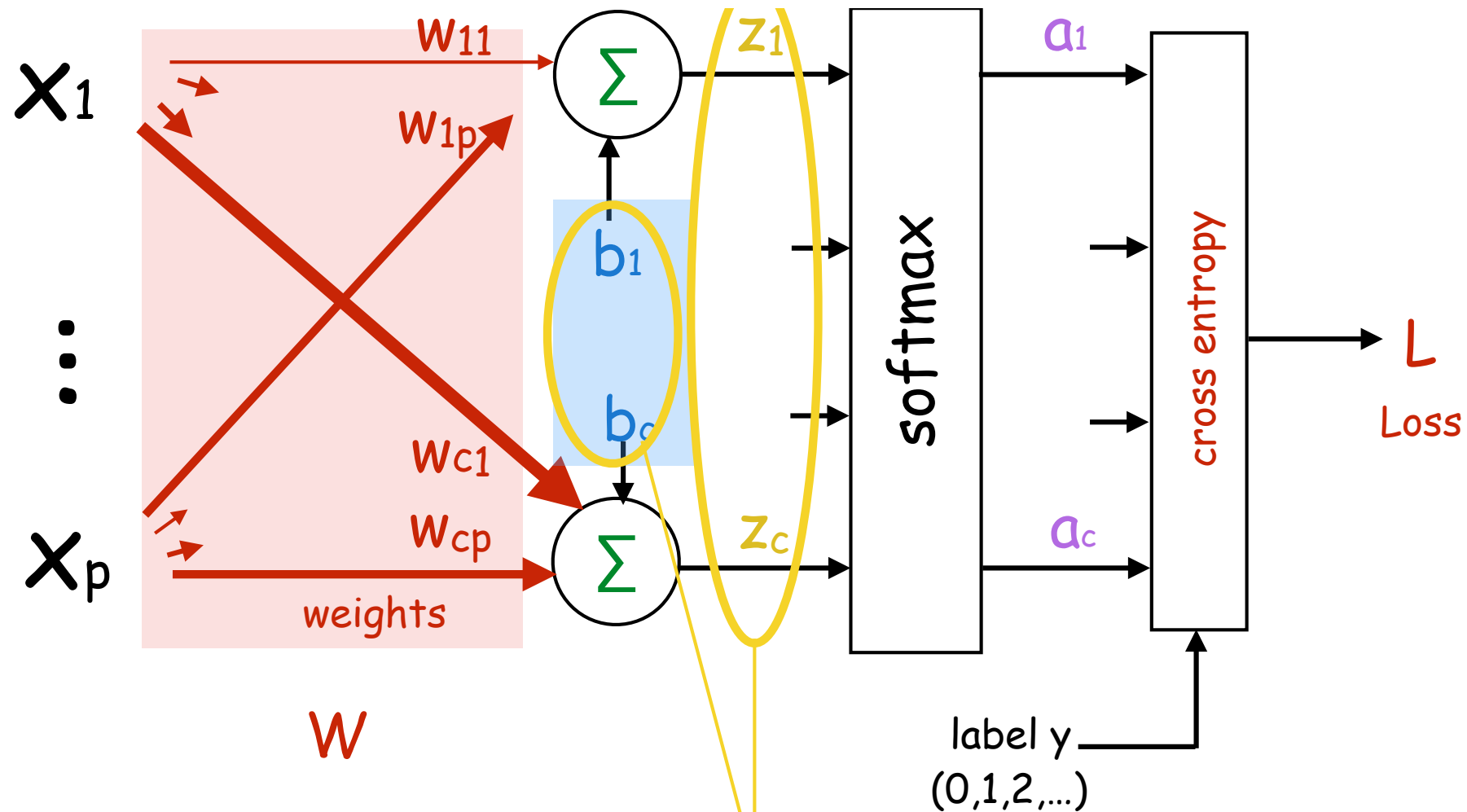Gradient of a vector w.r.t. a vector !!!

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

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

https://eli.thegreenplace.net/2016/the-softmax-function-and-its-derivative/

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

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

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

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

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

# 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)

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

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

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

vector (1 by c)    vector (1 by c)    matrix (c by c)

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

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

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

vector (1 by c)    vector (1 by c)    matrix (c by c)

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

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

$$\frac{\partial L}{\partial z_1}$$

$z_1$

$z_2$

$z_3$

softmax

$a_1$

$a_2$

$a_3$

$L$
Loss
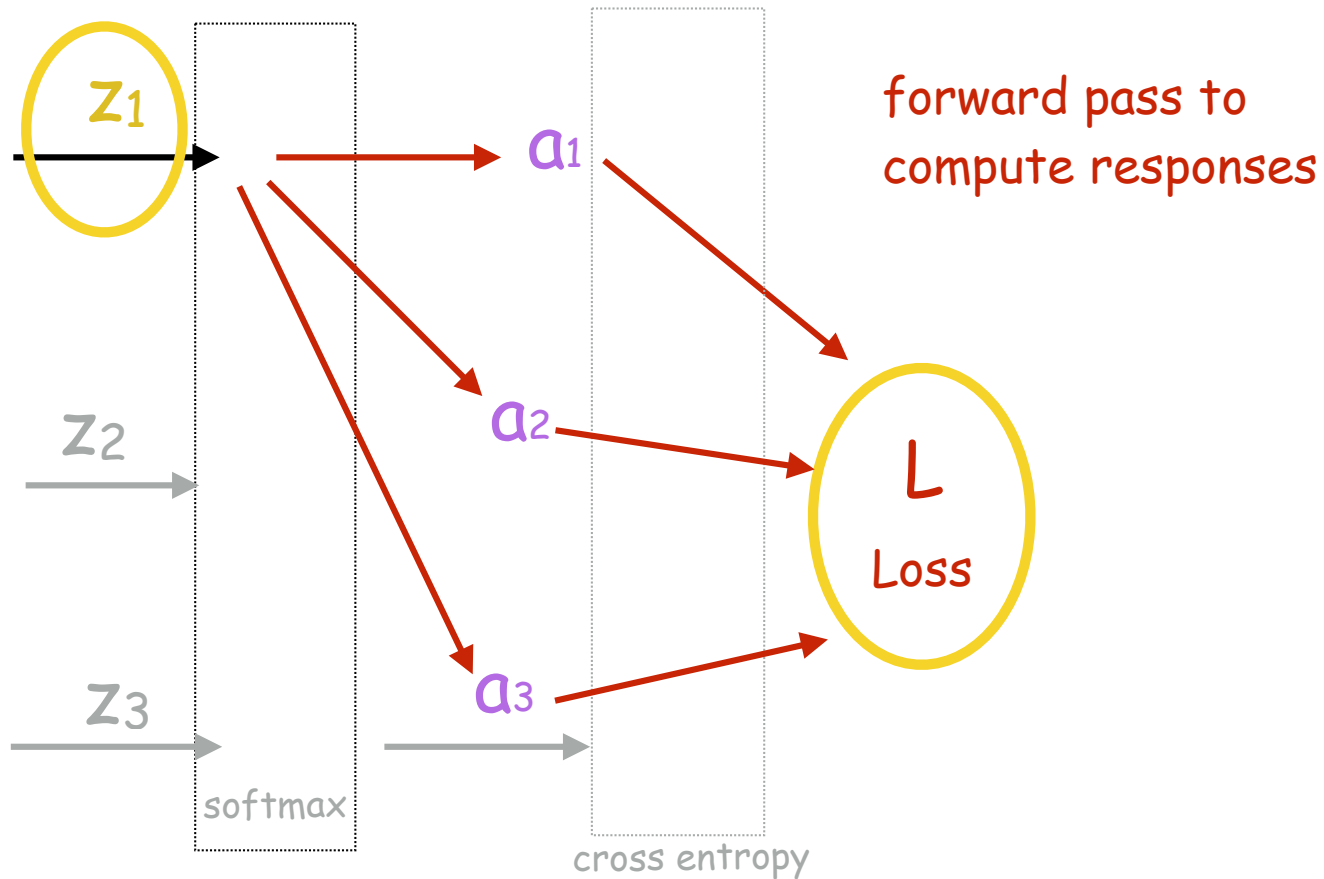
forward pass to
compute responses

cross entropy

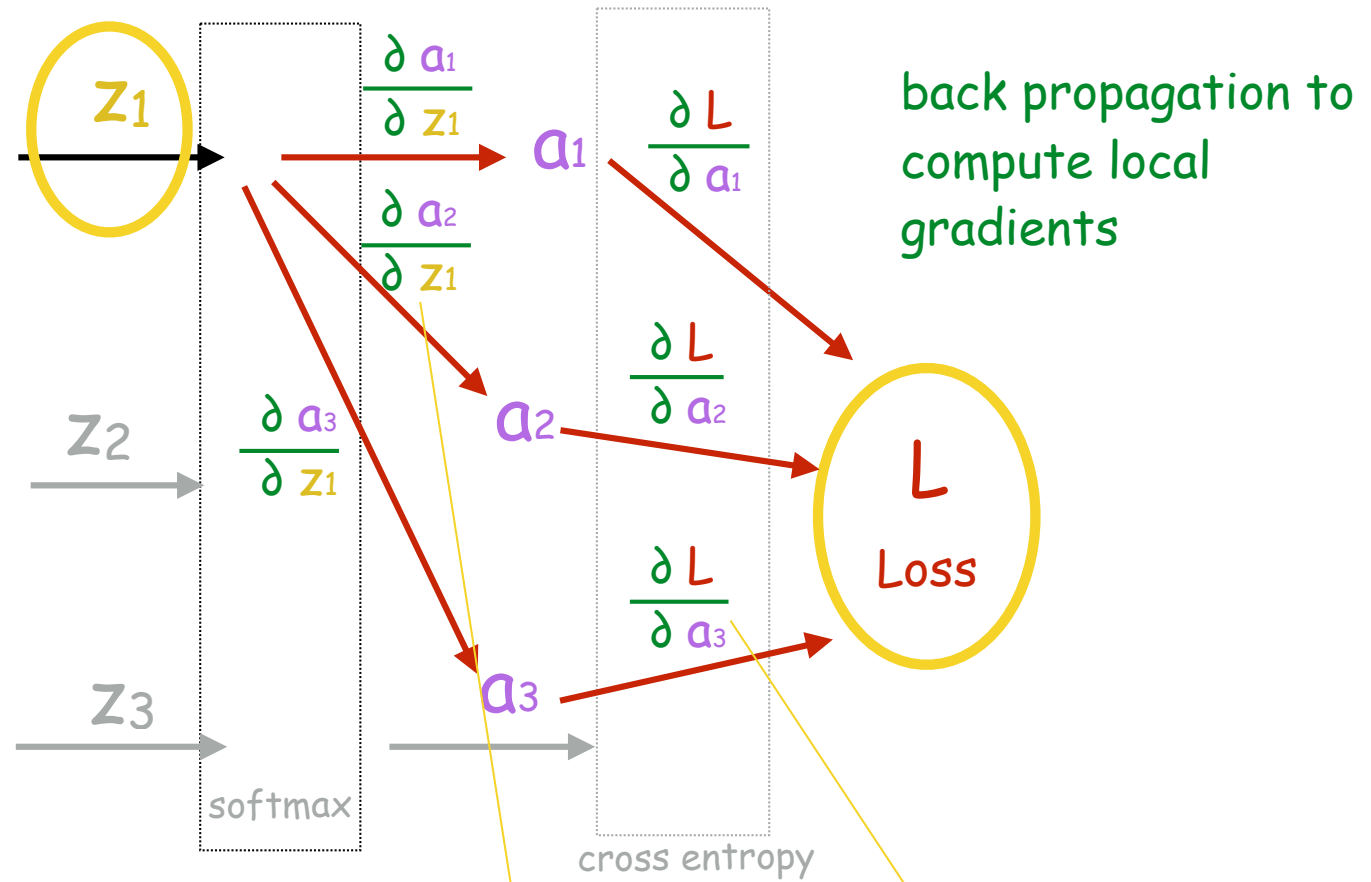$$\left( \frac{\partial L}{\partial a_1}, \ldots, \frac{\partial L}{\partial a_i}, \ldots, \frac{\partial L}{\partial a_c} \right) \times \begin{pmatrix} \frac{\partial a_1}{\partial z_1} & \cdots & \frac{\partial a_1}{\partial z_c} \\ \vdots & & \vdots \\ \frac{\partial a_c}{\partial z_1} & & \frac{\partial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

$$\frac{\partial L}{\partial z_1}$$

$$z_1$$

$$\frac{\partial a_1}{\partial z_1}$$

back propagation to compute local gradients

$$a_1 \qquad \frac{\partial L}{\partial a_1}$$

$$\frac{\partial a_2}{\partial z_1}$$

$$z_2 \qquad \frac{\partial a_3}{\partial z_1}$$

$$a_2 \qquad \frac{\partial L}{\partial a_2}$$

$$L$$

Loss

$$\frac{\partial L}{\partial a_3}$$

$$z_3$$

softmax

$$a_3$$

cross entropy

$$\frac{\partial L}{\partial a} = \left( \frac{\partial L}{\partial a_1} \quad , \quad \frac{\partial L}{\partial a_2} \quad , \quad \frac{\partial L}{\partial a_3} \right)$$

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

$$\frac{\partial L}{\partial z_1}$$

$z_1$

$$\frac{\partial a_1}{\partial z_1}$$

$a_1$

$$\frac{\partial L}{\partial a_1}$$

$$\frac{\partial a_2}{\partial z_1}$$

along each path:
compute product of
local gradients

$$\frac{\partial a_3}{\partial z_1}$$

$z_2$

$a_2$

$$\frac{\partial L}{\partial a_2}$$

$$L$$
Loss

$z_3$

$$\frac{\partial L}{\partial a_3}$$

$a_3$

softmax

cross entropy
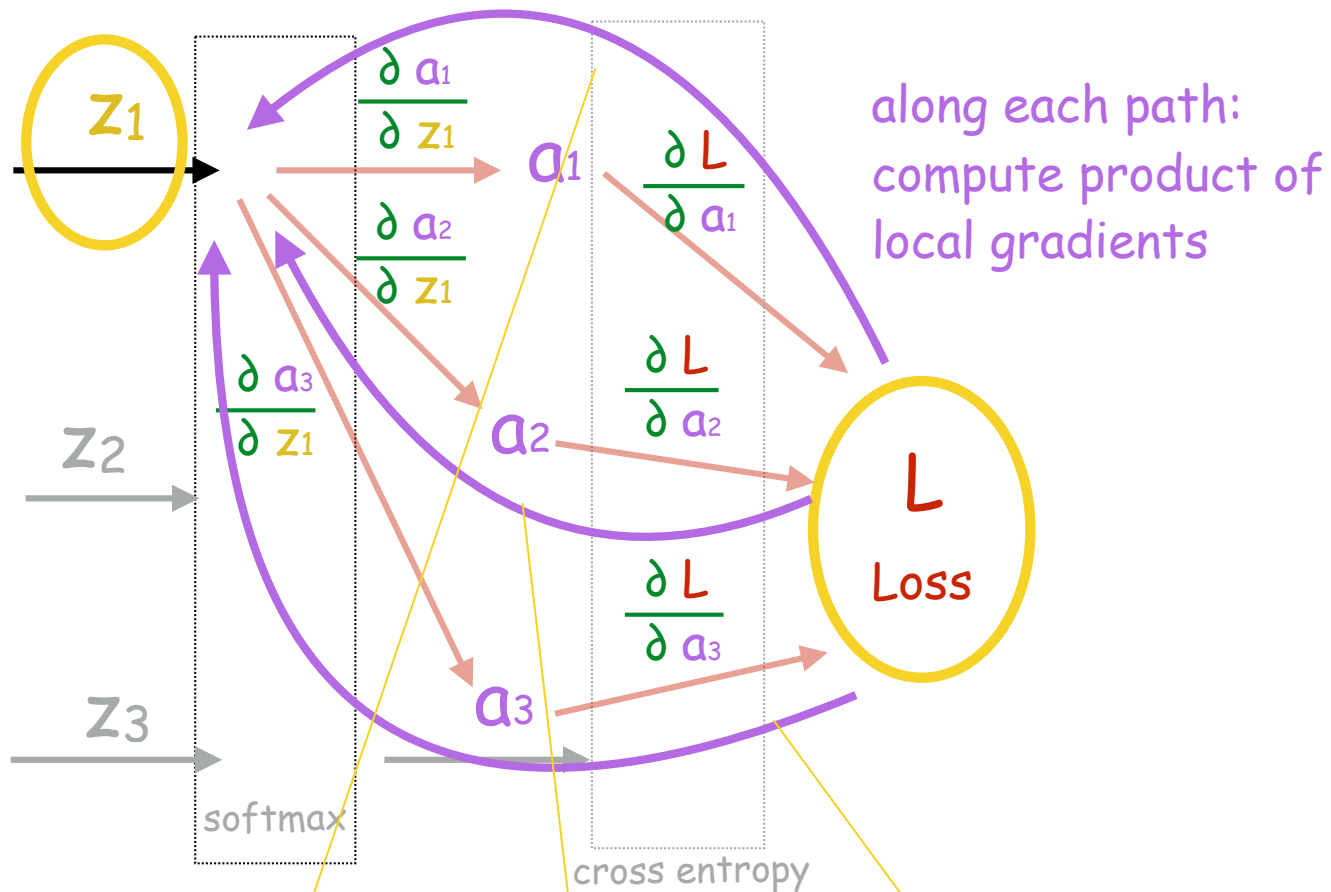
$$\frac{\partial L}{\partial a} = \left( \frac{\partial L}{\partial a_1} , \frac{\partial L}{\partial a_2} , \frac{\partial L}{\partial a_3} \right)$$
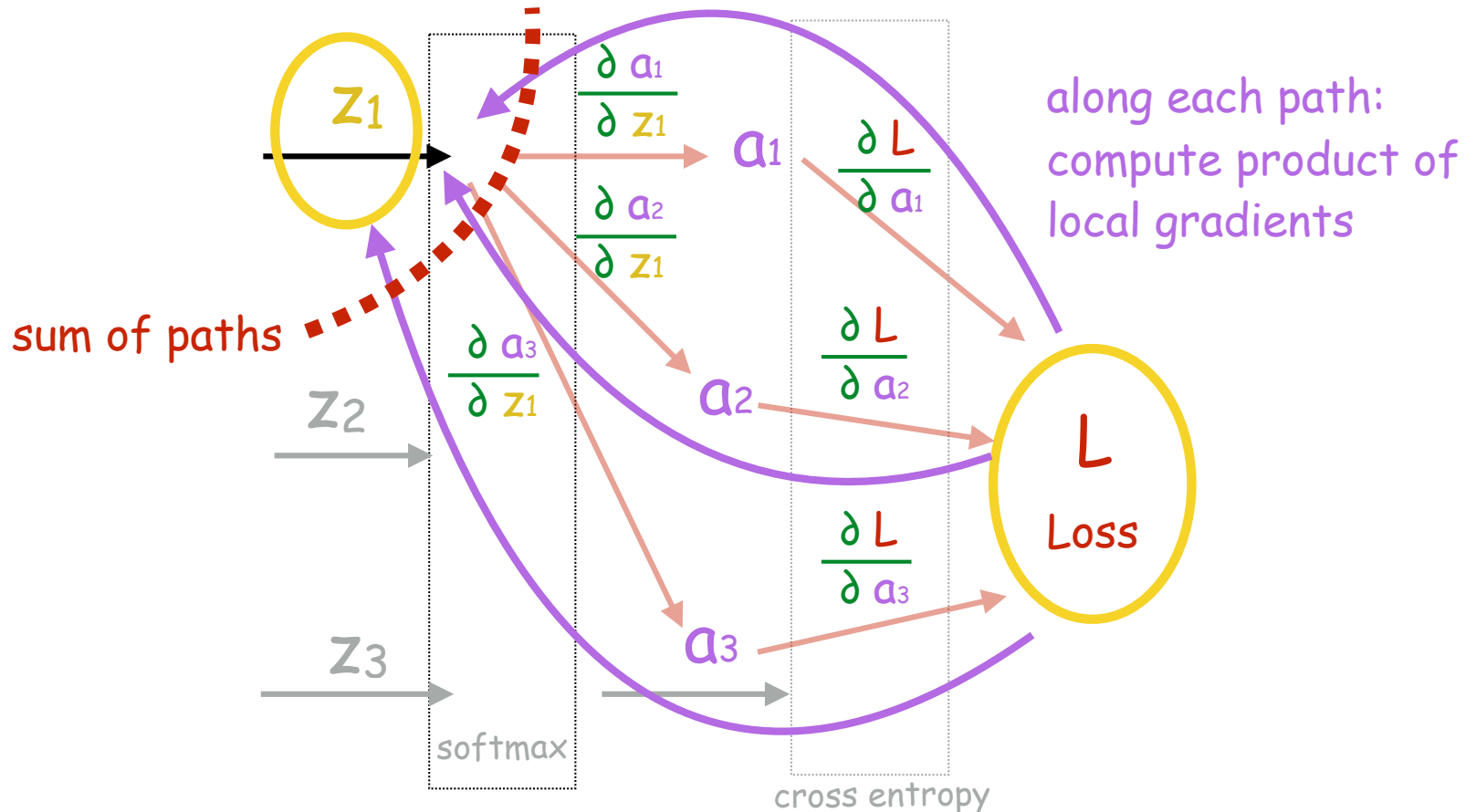
$$\times \qquad \times \qquad \times$$

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

$$\frac{\partial L}{\partial z_1}$$

along each path: compute product of local gradients

sum of paths

$$\frac{\partial a_1}{\partial z_1}$$

$$\frac{\partial a_2}{\partial z_1}$$

$$\frac{\partial a_3}{\partial z_1}$$

$$\frac{\partial L}{\partial a_1}$$

$$\frac{\partial L}{\partial a_2}$$

$$\frac{\partial L}{\partial a_3}$$

$z_1$

$z_2$

$z_3$

$a_1$

$a_2$

$a_3$

$L$ Loss

softmax

cross entropy

product along each path

sum of paths

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

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

matrix (c by p)   vector (c by 1)   matrix (c by p)

$$\left( \frac{\partial L}{\partial z_1}, \ldots, \frac{\partial L}{\partial z_i}, \ldots, \frac{\partial L}{\partial z_c} \right)^{T} \times \begin{pmatrix} \frac{\partial z_1}{\partial W_{11}}, & \cdots & , \frac{\partial z_1}{\partial W_{1p}} \\ \frac{\partial z_i}{\partial W_{i1}}, & \cdots & , \frac{\partial z_i}{\partial W_{ip}} \\ \frac{\partial z_c}{\partial W_{c1}}, & \cdots & , \frac{\partial z_c}{\partial W_{cp}} \end{pmatrix}$$

transpose

element-wise product

$$= \begin{pmatrix} \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial W_{11}}, & \cdots & , \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial W_{1p}} \\ \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial W_{i1}}, & \cdots & , \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial W_{ip}} \\ \frac{\partial L}{\partial z_c} \frac{\partial z_c}{\partial W_{c1}}, & \cdots & , \frac{\partial L}{\partial z_c} \frac{\partial z_c}{\partial W_{cp}} \end{pmatrix}$$

# 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

$$\frac{\partial L}{\partial a} \quad \frac{\partial a}{\partial z} \quad \frac{\partial z}{\partial b} \quad \frac{\partial z}{\partial W}$$
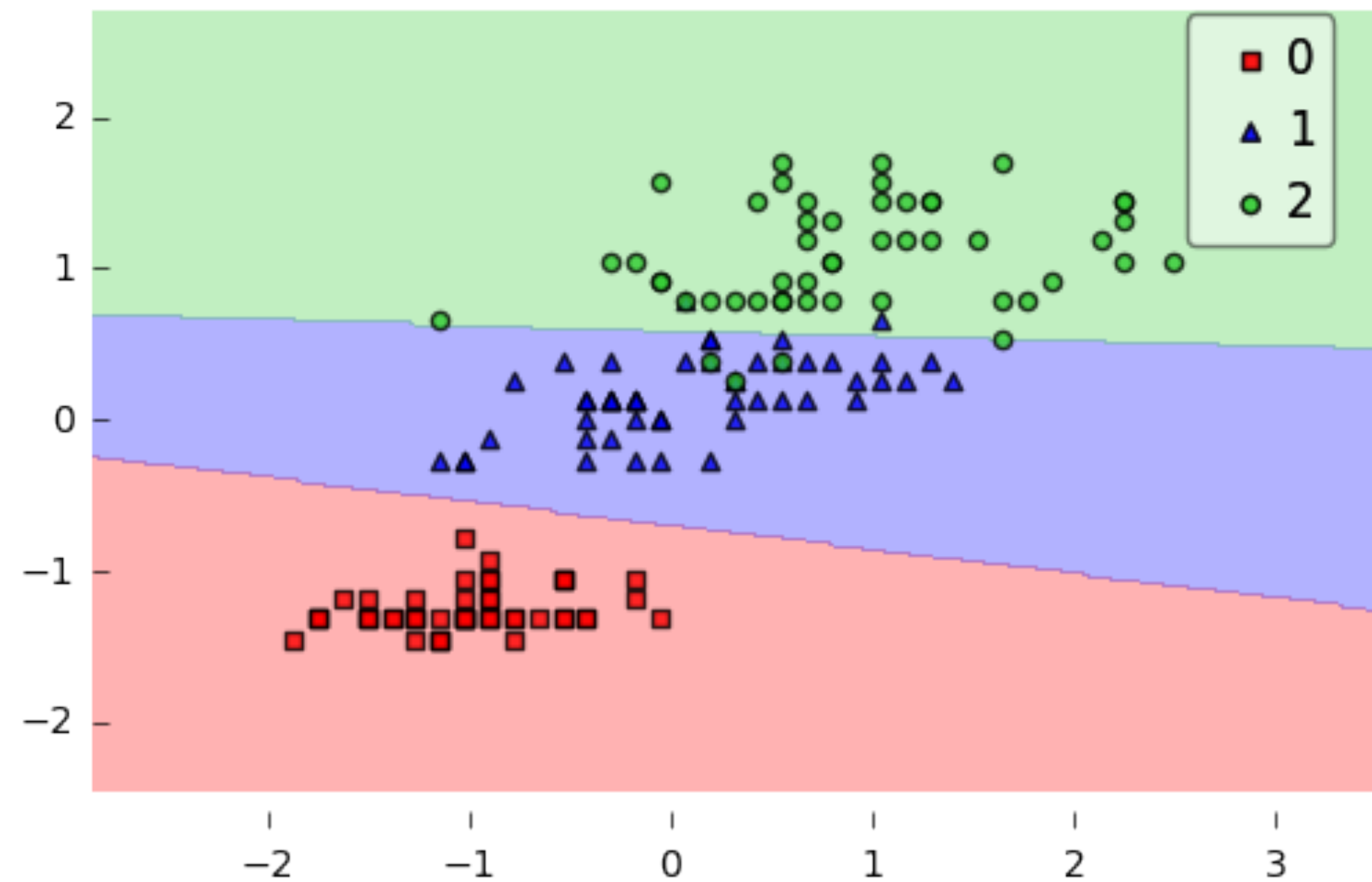
      compute global gradients using chain rule

$$\frac{\partial L}{\partial W} \quad \frac{\partial L}{\partial b}$$
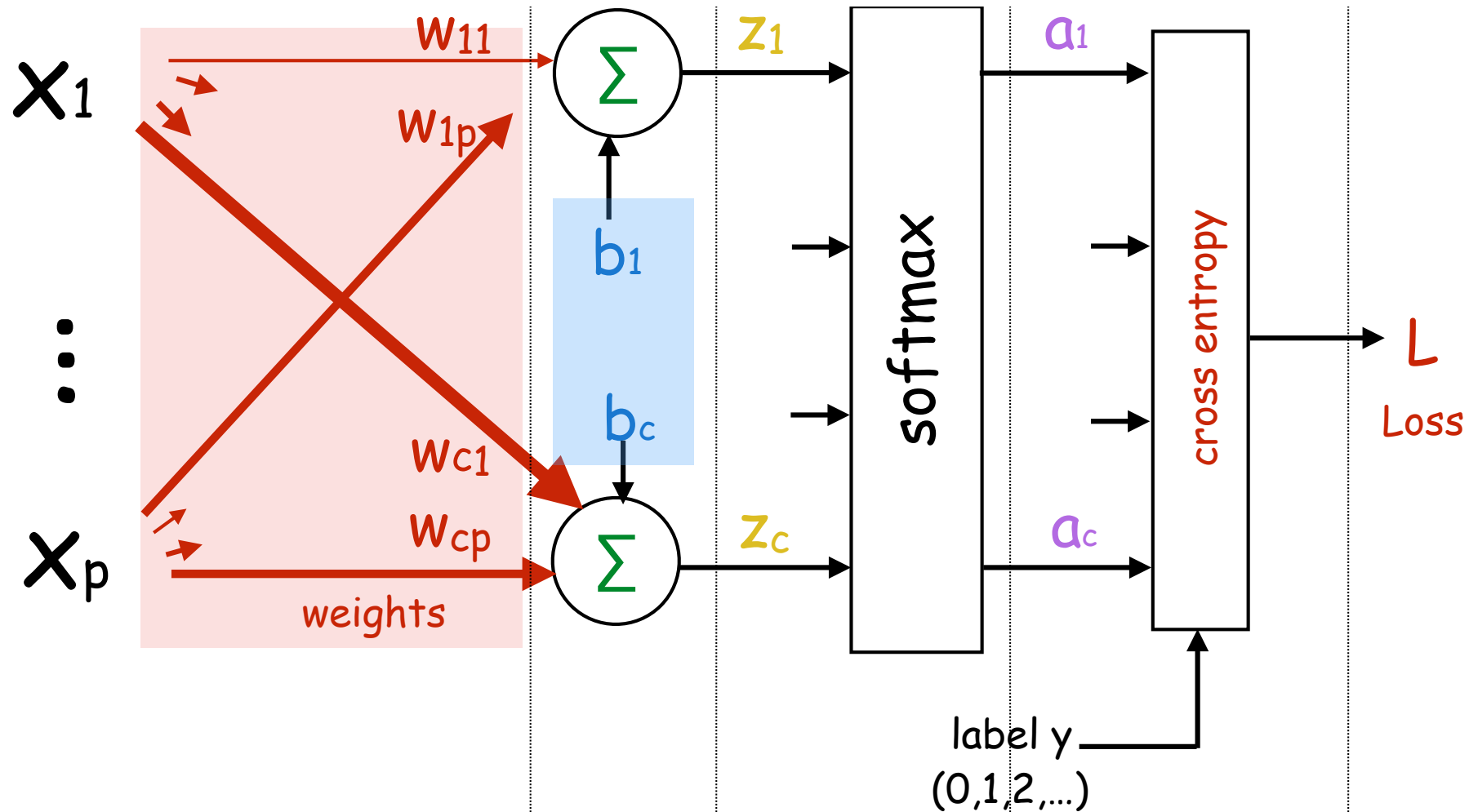
      update the parameters **W** and **b**

$$\mathbf{W} \longleftarrow \mathbf{W} - a\frac{\partial L}{\partial W}$$
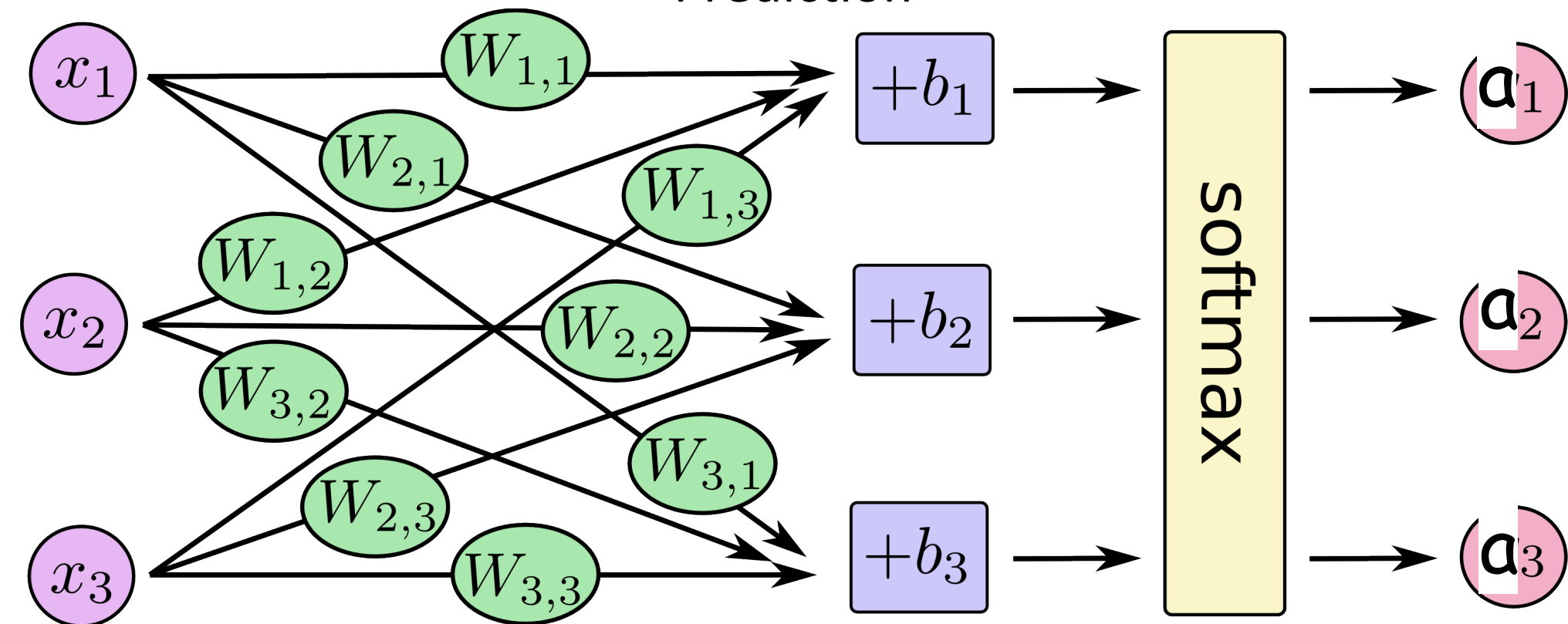
$$b \longleftarrow b - a\frac{\partial L}{\partial b}$$
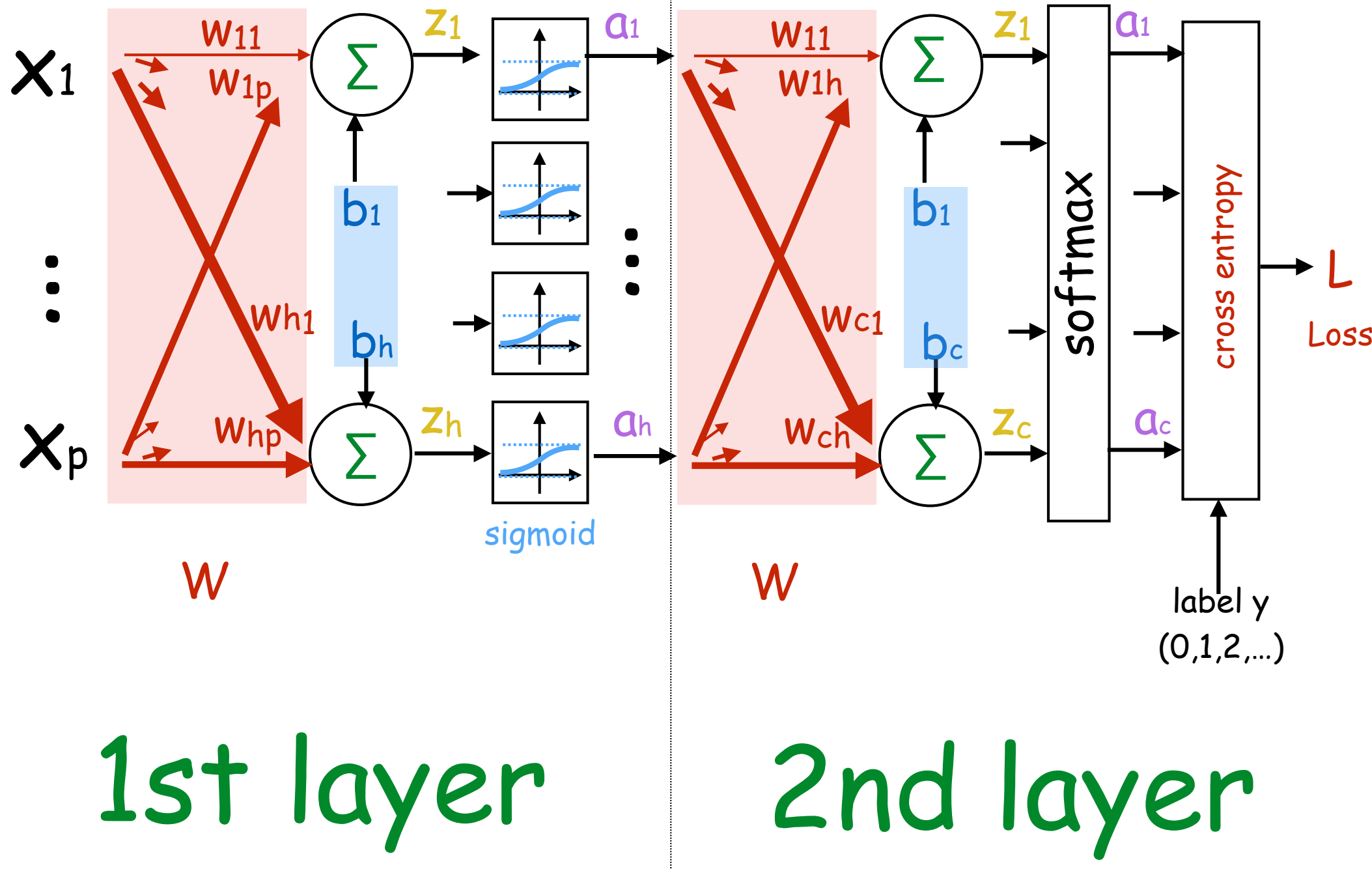
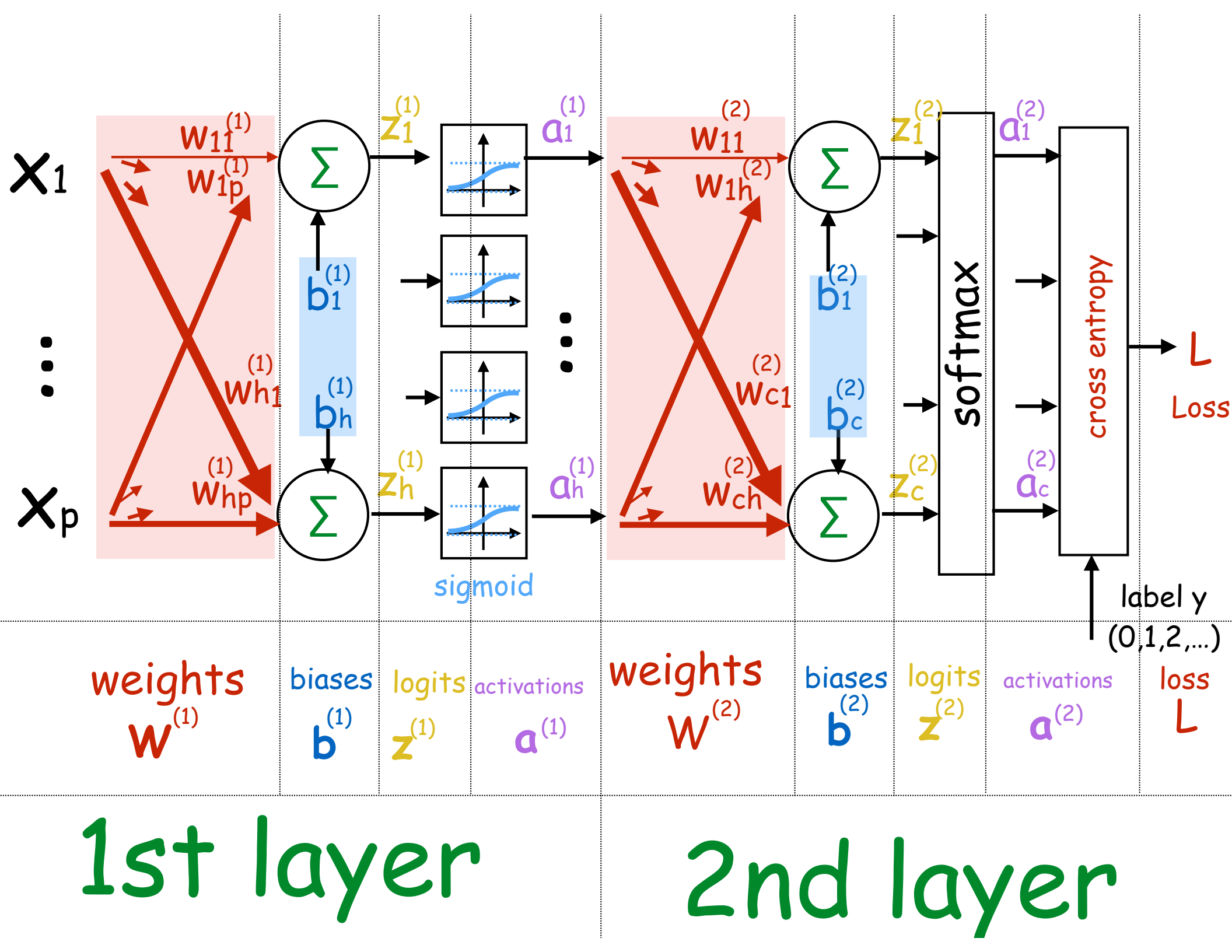Softmax Regression - Gradient Descent

Softmax Regression

Prediction



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

# Fully Connected Neural Network

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

vector length h          vector of length c          matrix (c by h)

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

$$= \left( \frac{\partial L}{\partial z_1^{(2)}}, \ldots, \frac{\partial L}{\partial z_c^{(2)}} \right) \times \begin{pmatrix} \frac{\partial z_1^{(2)}}{\partial a_1^{(1)}}, & \cdots & , \frac{\partial z_1^{(2)}}{\partial a_h^{(1)}} \\ & \cdots & \\ \frac{\partial z_c^{(2)}}{\partial a_1^{(1)}}, & \cdots & , \frac{\partial z_c^{(2)}}{\partial a_h^{(1)}} \end{pmatrix}$$