# Learning for vision II Neural networks

### Karel Zimmermann
http://cmp.felk.cvut.cz/~zimmerk/

Vision for Robotics and Autonomous Systems
https://cyber.felk.cvut.cz/vras/

Center for Machine Perception
https://cmp.felk.cvut.cz

Department for Cybernetics
Faculty of Electrical Engineering
Czech Technical University in Prague

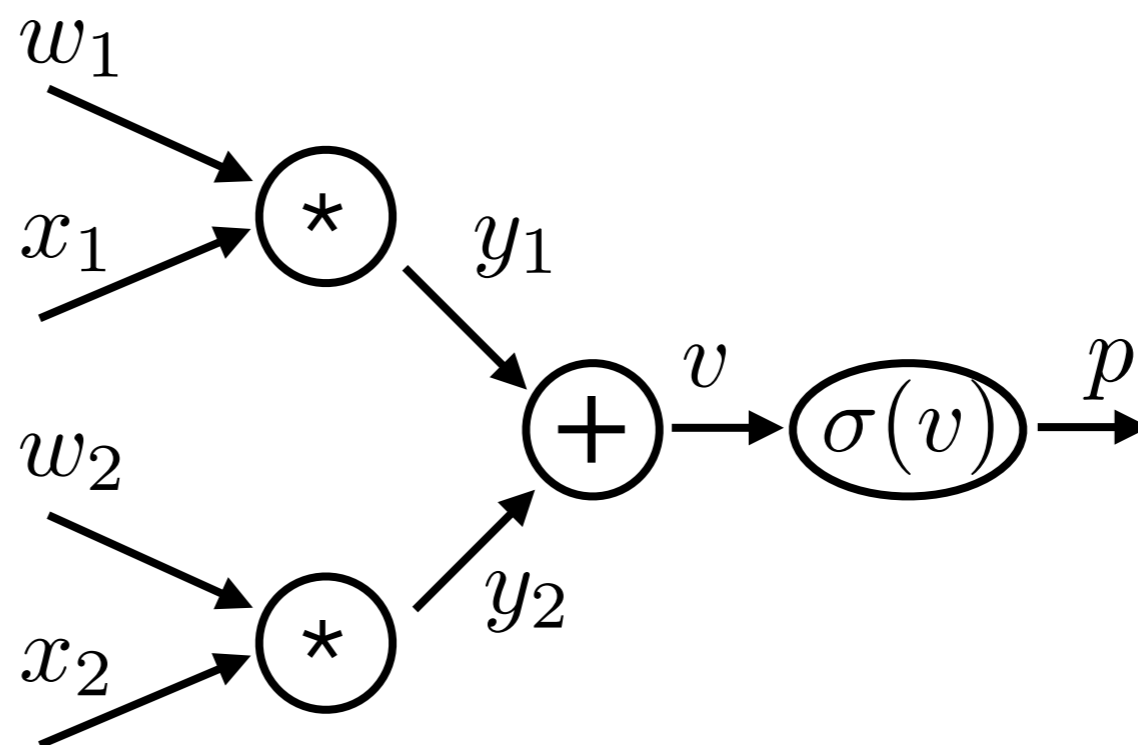# Linear classifier and neuron

Labels
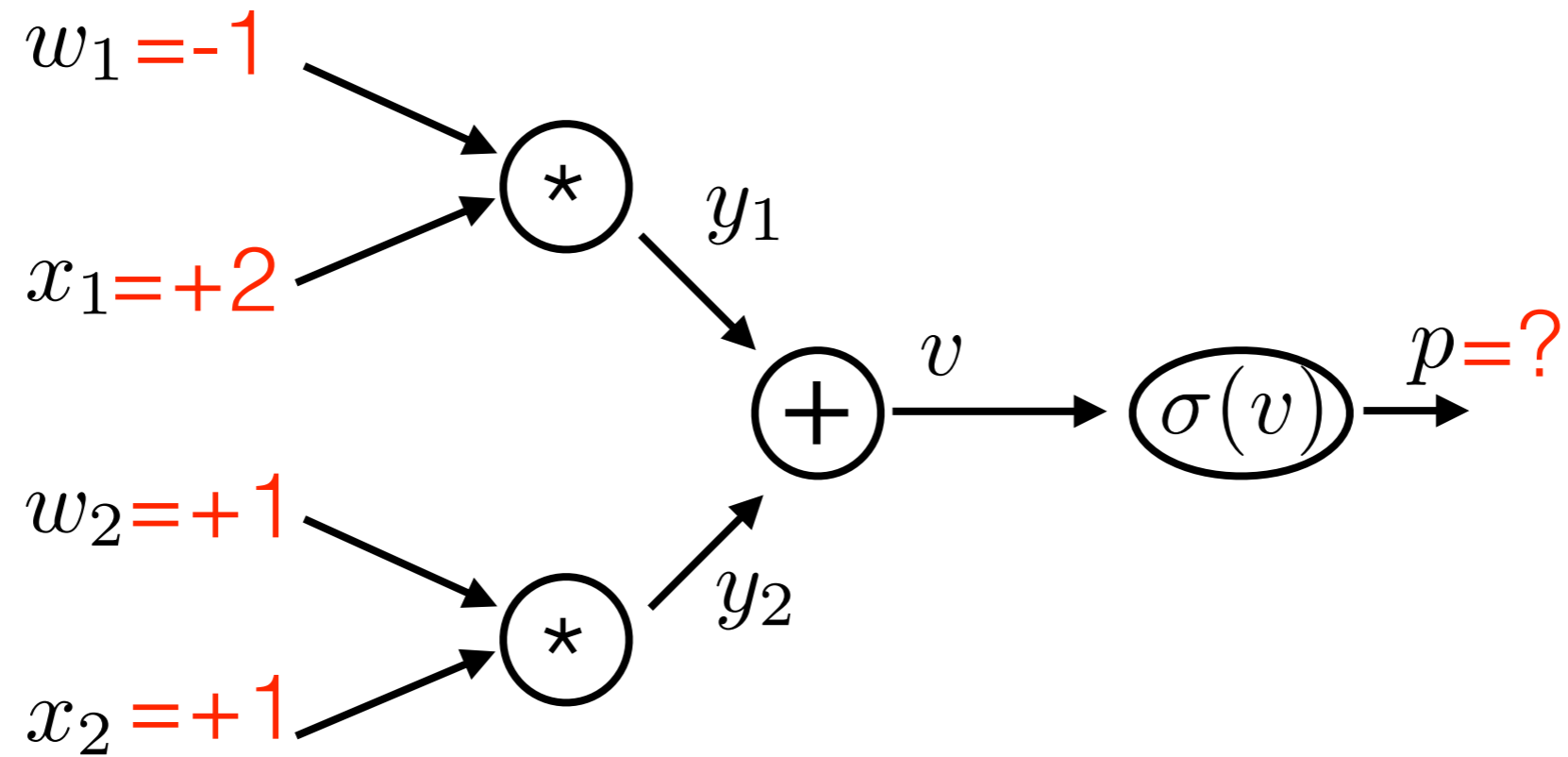
RGB images

+1



−1



```
def classify(    ):
    # Linear classifier
    x = vec(    )
    p = σ(wᵀx)
    return p
```

$$\mathbf{x} = \mathrm{vec}(\quad)$$

$$p = \sigma\left(\mathbf{w}^\top \mathbf{x}\right)$$

Computational graph of linear classifier



$w_1$

$x_1$

$*$ $y_1$

$w_2$

$x_2$

$*$ $y_2$

$+$ $v$ $\sigma(v)$ $p$

# Example I: given trained neuron, and input, what is output?



$w_1 = -1$

$x_1 = +2$

$y_1$

$v$

$p = ?$

$w_2 = +1$

$x_2 = +1$

$y_2$

$\sigma(v)$

# Example I: given trained classifier, and input, what is output?



$w_1 = -1$
$x_1 = +2$
$y_1 = -2$
$w_2 = +1$
$x_2 = +1$
$y_2 = 1$
$v = -1$
$\sigma(v)$
$p = 0.27$

# Relation to biological neuron

$w_1 = -1$

$x_1 = +2$

$*$   $y_1 = -2$

$+$   $v = -1$   $\sigma(v)$   $p = 0.27$

$w_2 = +1$

$x_2 = +1$

$*$   $y_2 = 1$

Dendrites

Soma

Axon

# Relation to biological neuron



$w_1 = -1$

$x_1 = +2$

$y_1 = -2$

$w_2 = +1$

$x_2 = +1$

$y_2 = 1$

$v = -1$

$\sigma(v)$

$p = 0.27$

Dendrites

Soma

Axon

# Relation to biological neuron



$w_1 = -1$

$x_1 = +2$

$y_1 = -2$

$w_2 = +1$

$x_2 = +1$

$y_2 = 1$

$v = -1$

$\sigma(v)$

$p = 0.27$

Dendrites

Soma

Axon

# Modeling dynamic neuron behaviour



http://jackterwilliger.com/biological-neural-networks-part-i-spiking-neurons/

# Activation functions $\sigma(v)$

**Sigmoid**

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

**tanh**

$\tanh(x)$

**ReLU**

$\max(0, x)$

**Leaky ReLU**

$\max(0.1x, x)$

**Maxout**

$\max(w_1^T x + b_1, w_2^T x + b_2)$

**ELU**

$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

# Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1$=-1+ ??

$x_1$=+2

$y_1$=-2

$v$=-1

$\sigma(v)$

$p$=0.27  => max

$w_2$=+1

$x_2$=+1

$y_2$=1

$p(w_1)$

0.27

-1

$w_1$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1 = -1 +$ ??



$x_1 = +2$

$y_1 = -2$

$v = -1$

$\sigma(v)$

$p = 0.27 \Rightarrow$ max

$w_2 = +1$

$y_2 = 1$

$x_2 = +1$

$$w_1 = w_1 + \frac{\partial p}{\partial w1}$$

$p(w_1)$

0.27

-1

$w_1$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1 = -1 + ??$



$y_1 = -2$

$x_1 = +2$

$v = -1$ $\sigma(v)$ $p = 0.27$ => max

$w_2 = +1$

$y_2 = 1$

$x_2 = +1$

$$w_1 = w_1 + \frac{\partial p}{\partial w1}$$

Chain-rule in computational graph $\quad \dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v} \dfrac{\partial v}{\partial y_1} \dfrac{\partial y_1}{\partial w_1}$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1 = -1 + ??$

$x_1 = +2$

$y_1 = -2$

$w_2 = +1$

$y_2 = 1$

$x_2 = +1$

$v = -1$

$\sigma(v)$

$p = 0.27 \Rightarrow \max$

$\dfrac{\partial p}{\partial v} = 0.2$

$w_1 = w_1 + \dfrac{\partial p}{\partial w1}$

Local gradient:

$$\frac{\partial p}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v)(1 - \sigma(v))$$

Chain-rule in computational graph $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v}\dfrac{\partial v}{\partial y_1}\dfrac{\partial y_1}{\partial w_1}$

# Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

## Sigmoid Function



## Derivative of Sigmoid



Local gradient:
$$\frac{\partial p}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v)(1 - \sigma(v))$$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1$=-1+ ??



$y_1$=-2

$x_1$=+2

$v$=-1   $\sigma(v)$   $p$=0.27  => max

$\dfrac{\partial p}{\partial v}$=0.2

$w_2$=+1

$y_2$=1

$x_2$=+1

$w_1 = w_1 + \dfrac{\partial p}{\partial w1}$

Local gradient:

$$\frac{\partial p}{\partial v} = \frac{\partial \sigma(v)}{\partial v} = \sigma(v)(1 - \sigma(v))$$

Chain-rule in computational graph   $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v} \dfrac{\partial v}{\partial y_1} \dfrac{\partial y_1}{\partial w_1}$

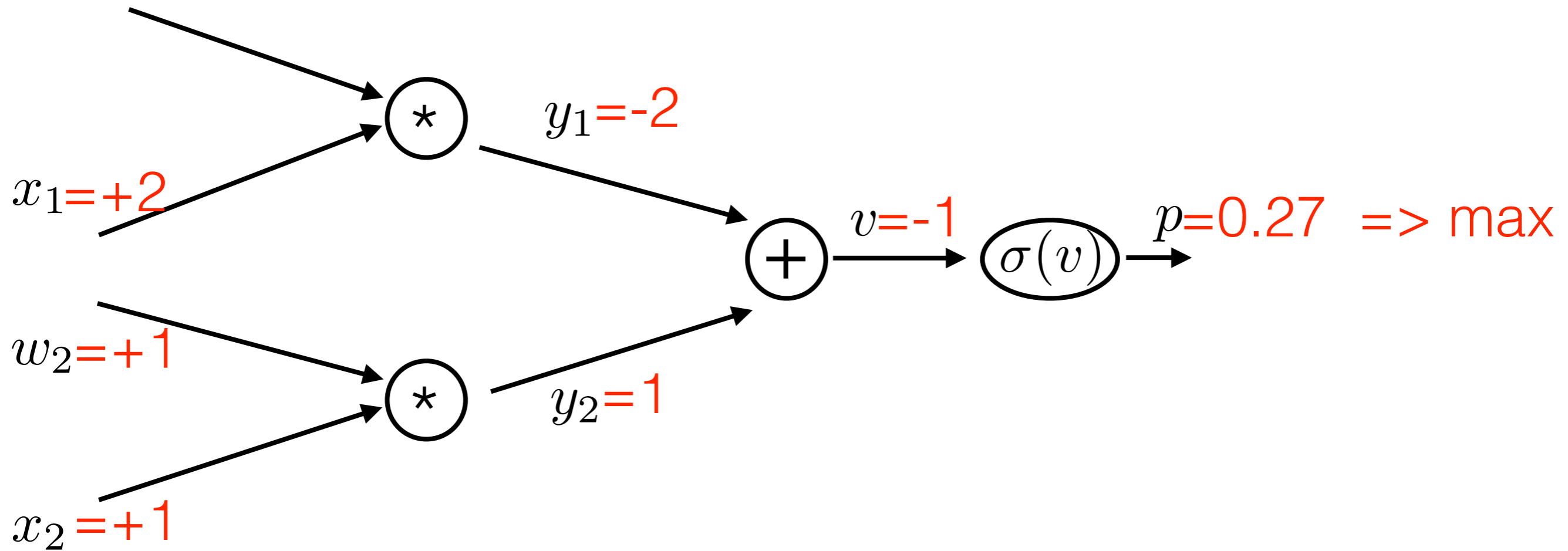Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1$ = -1+ ??

$x_1$ = +2

$y_1$ = -2

$\dfrac{\partial v}{\partial y_1}$ = 1

$w_2$ = +1

$y_2$ = 1

$x_2$ = +1

$v$ = -1

$\dfrac{\partial p}{\partial v}$ = 0.2

$\sigma(v)$

$p$ = 0.27  => max

$w_1 = w_1 + \dfrac{\partial p}{\partial w1}$
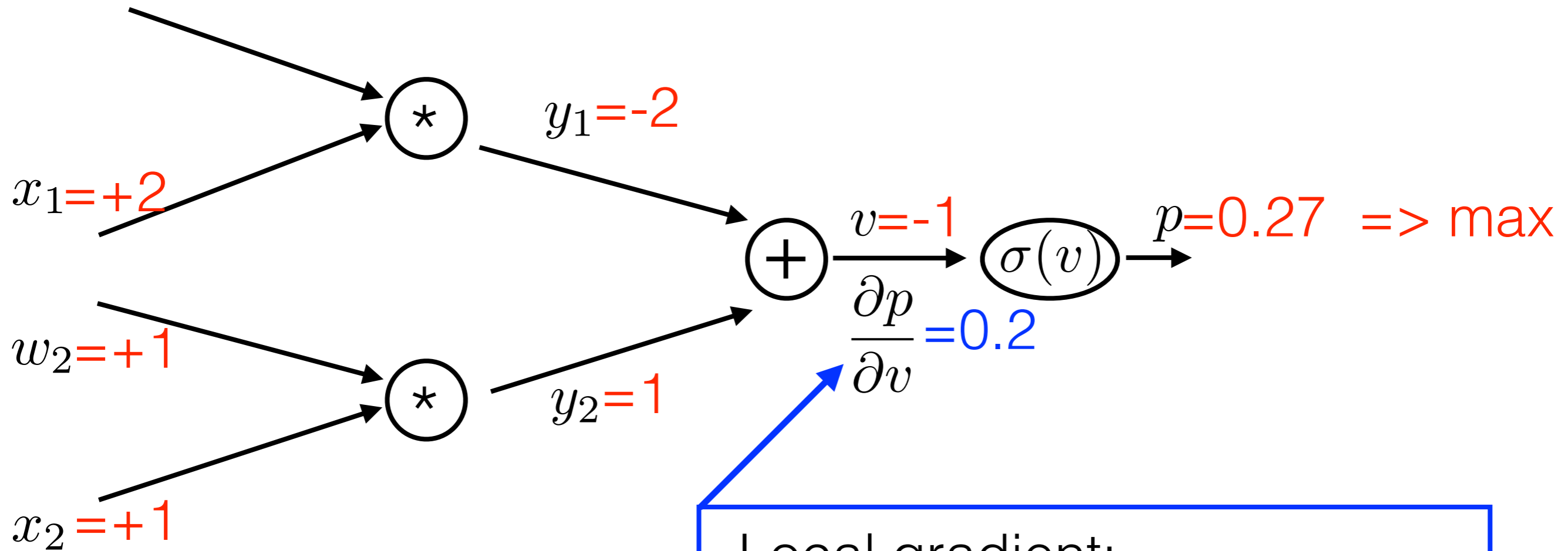
Local gradient:

$$\dfrac{\partial v}{\partial y_1} = \dfrac{\partial (y_1 + y_2)}{\partial y_1} = 1$$

Chain-rule in computational graph  $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v} \dfrac{\partial v}{\partial y_1} \dfrac{\partial y_1}{\partial w_1}$

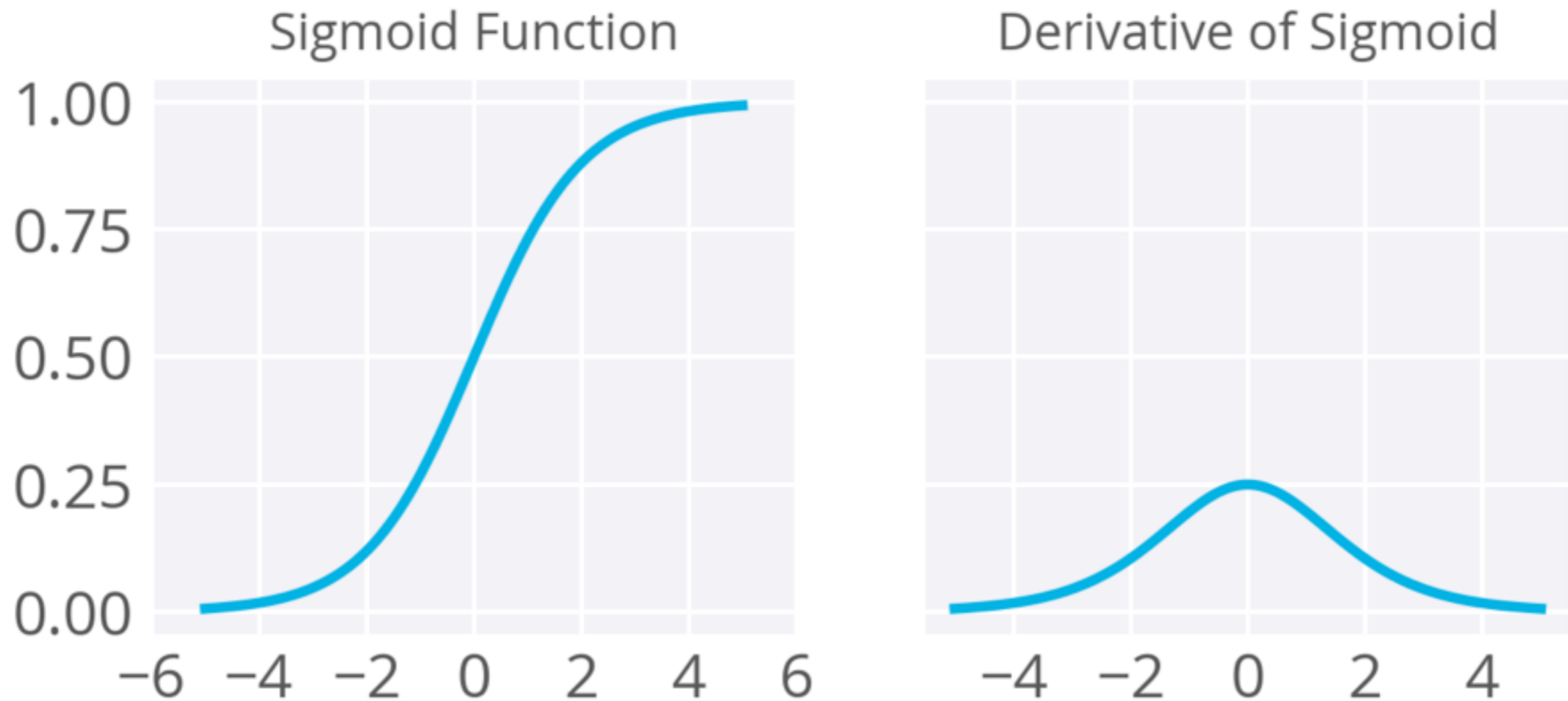Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1 = $ -1+ ??

$\frac{\partial y_1}{\partial w_1} = 2$

$x_1 = +2$

$y_1 = -2$

$\frac{\partial v}{\partial y_1} = 1$

$w_2 = +1$

$v = -1$

$\sigma(v)$

$p = 0.27$  => max

$\frac{\partial p}{\partial v} = 0.2$

$y_2 = 1$

$x_2 = +1$

Local gradient:
$$\frac{\partial y_1}{\partial w_1} = \frac{\partial (w_1 x_1)}{\partial w_1} = x_1$$

$$w_1 = w_1 + \frac{\partial p}{\partial w1}$$

Chain-rule in computational graph $\quad \frac{\partial p}{\partial w_1} = \frac{\partial p}{\partial v} \frac{\partial v}{\partial y_1} \frac{\partial y_1}{\partial w_1}$

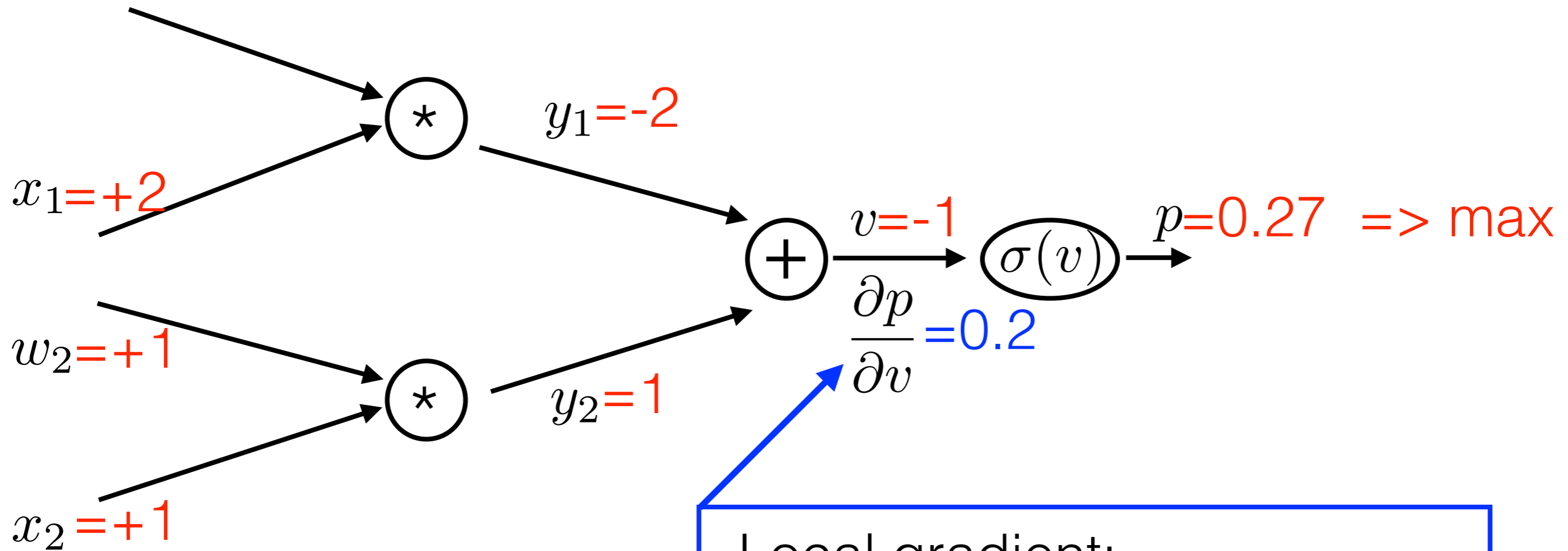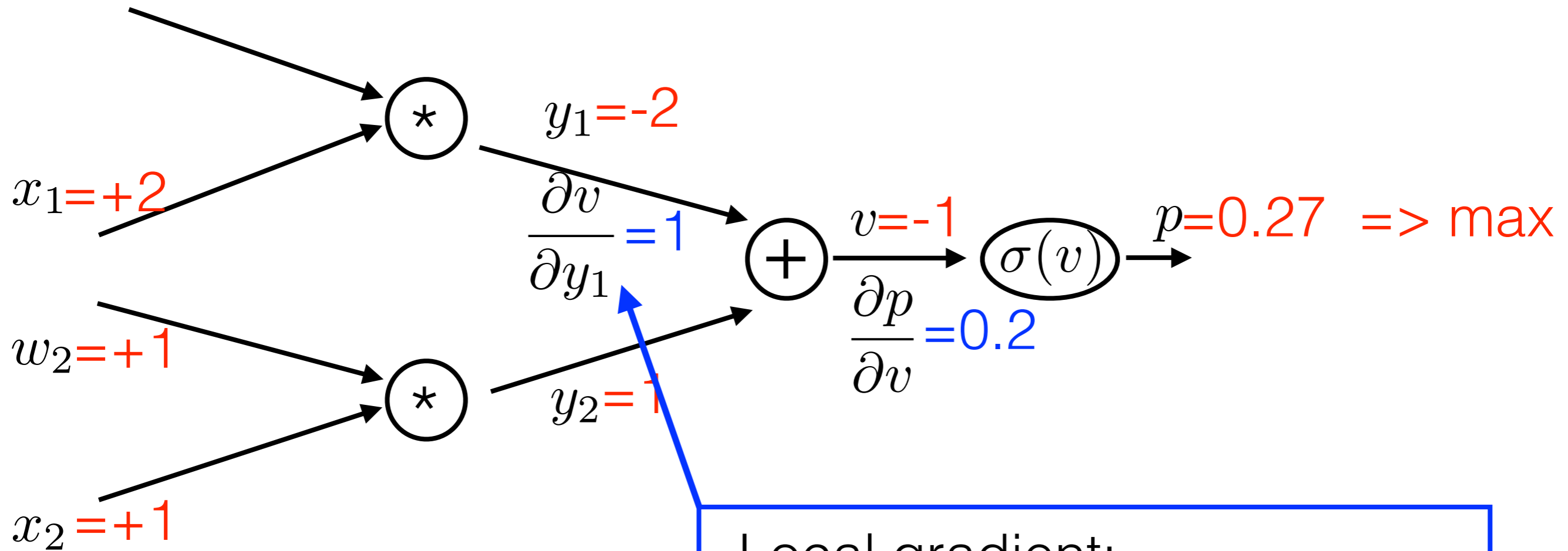Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1 = -1$

$\frac{\partial y_1}{\partial w_1} = 2$

$x_1 = +2$

$\frac{\partial p}{\partial y_1} = 1*0.2 = 0.2$

$y_1 = -2$

$\frac{\partial v}{\partial y_1} = 1$

$(*)$

$(+)$

$v = -1$

$\frac{\partial p}{\partial v} = 0.2$

$\sigma(v)$

$p = 0.27$   => max

$w_2 = +1$

$(*)$

$y_2 = 1$

$x_2 = +1$

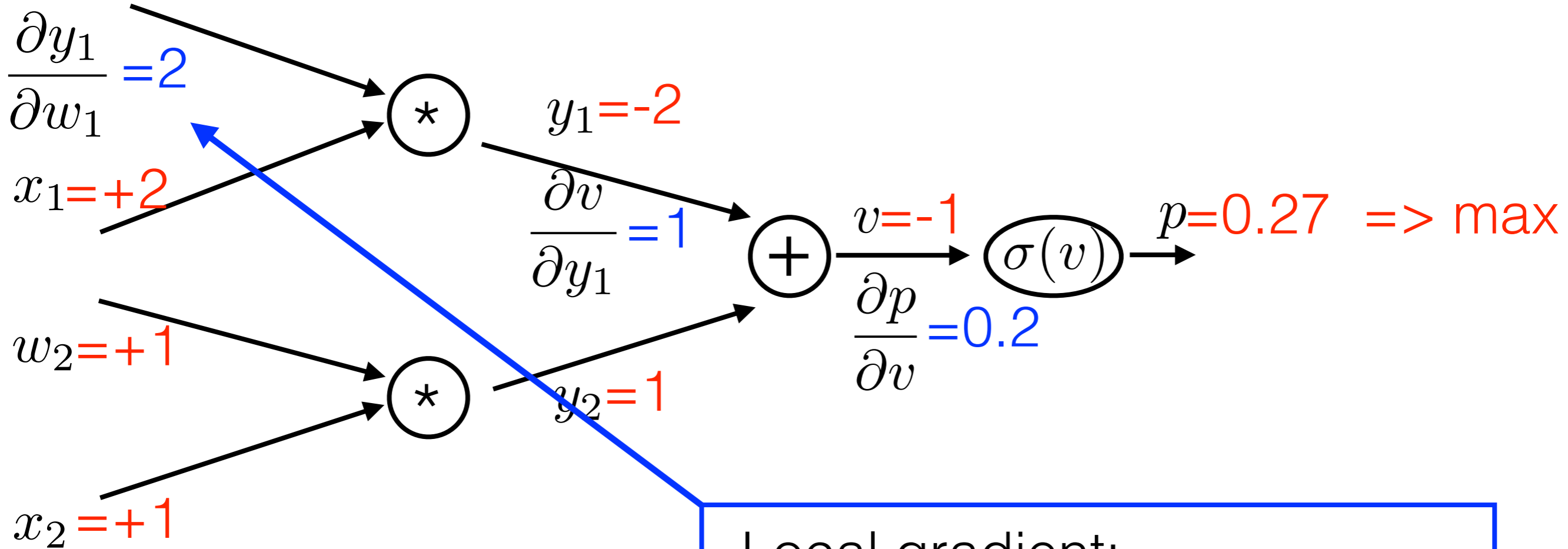$w_1 = w_1 + \frac{\partial p}{\partial w1}$

Edge gradient:
$$\frac{\partial p}{\partial y_1} = \frac{\partial p}{\partial v}\frac{\partial v}{\partial y_1}$$

Chain-rule in computational graph   $\frac{\partial p}{\partial w_1} = \frac{\partial p}{\partial v}\frac{\partial v}{\partial y_1}\frac{\partial y_1}{\partial w_1}$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1 = -1$

$\dfrac{\partial p}{\partial w_1} = 0.4$

$\dfrac{\partial y_1}{\partial w_1} = 2$

$\dfrac{\partial p}{\partial y_1} = 1*0.2 = 0.2$

$x_1 = +2$

$y_1 = -2$

$\dfrac{\partial v}{\partial y_1} = 1$

$v = -1$

$p = 0.27$ => max

$\sigma(v)$

$\dfrac{\partial p}{\partial v} = 0.2$

$w_2 = +1$

$y_2 = 1$

$x_2 = +1$

$w_1 = w_1 + \dfrac{\partial p}{\partial w1}$

Edge gradient:
$$\frac{\partial p}{\partial w_1} = \frac{\partial p}{\partial y_1}\frac{\partial y_1}{\partial w_1}$$

Chain-rule in computational graph $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v}\dfrac{\partial v}{\partial y_1}\dfrac{\partial y_1}{\partial w_1}$

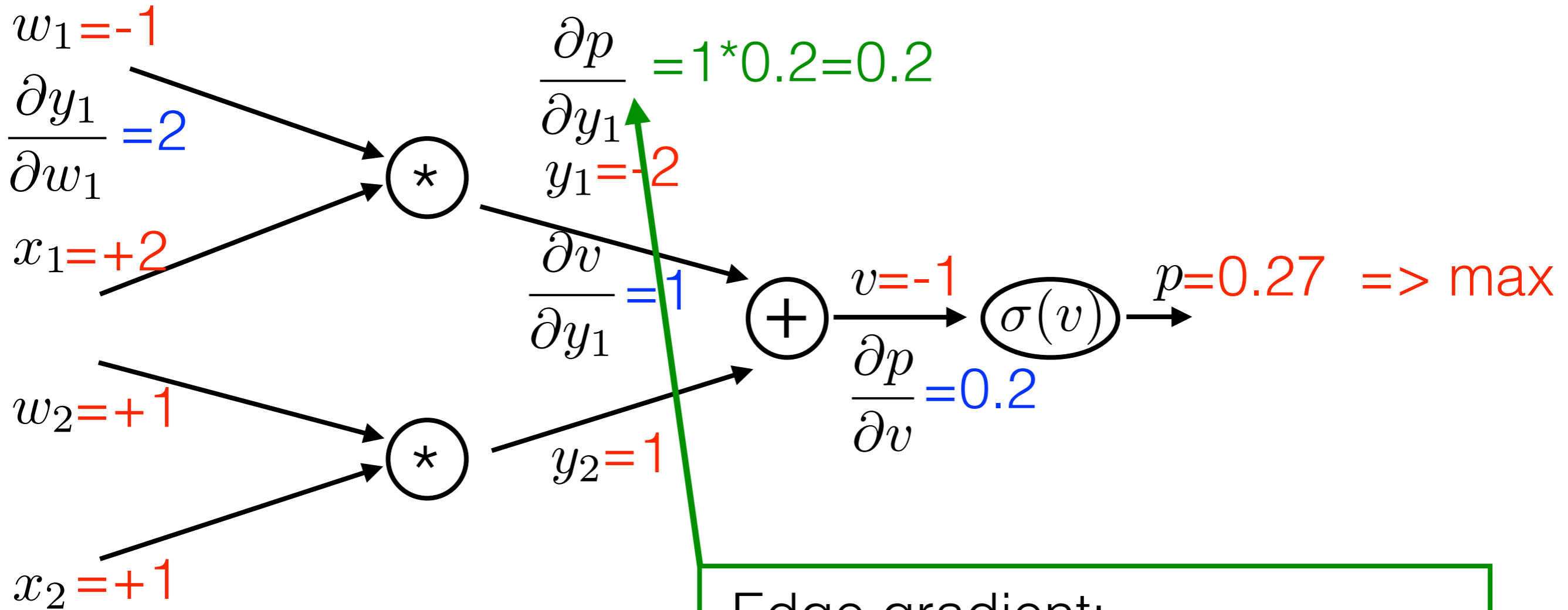Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1 = -1$

$\dfrac{\partial p}{\partial w_1} = 0.4$

$\dfrac{\partial y_1}{\partial w_1} = 2$

$\dfrac{\partial p}{\partial y_1} = 1*0.2 = 0.2$

$y_1 = -2$

$x_1 = +2$

$\dfrac{\partial v}{\partial y_1} = 1$

$v = -1$

$p = 0.27 \Rightarrow$ max

$\sigma(v)$

$w_2 = +1$

$\dfrac{\partial p}{\partial v} = 0.2$

$y_2 = 1$

$x_2 = +1$

$w_1 = w_1 + \dfrac{\partial p}{\partial w1}$

Chain-rule in computational graph $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v} \dfrac{\partial v}{\partial y_1} \dfrac{\partial y_1}{\partial w_1}$

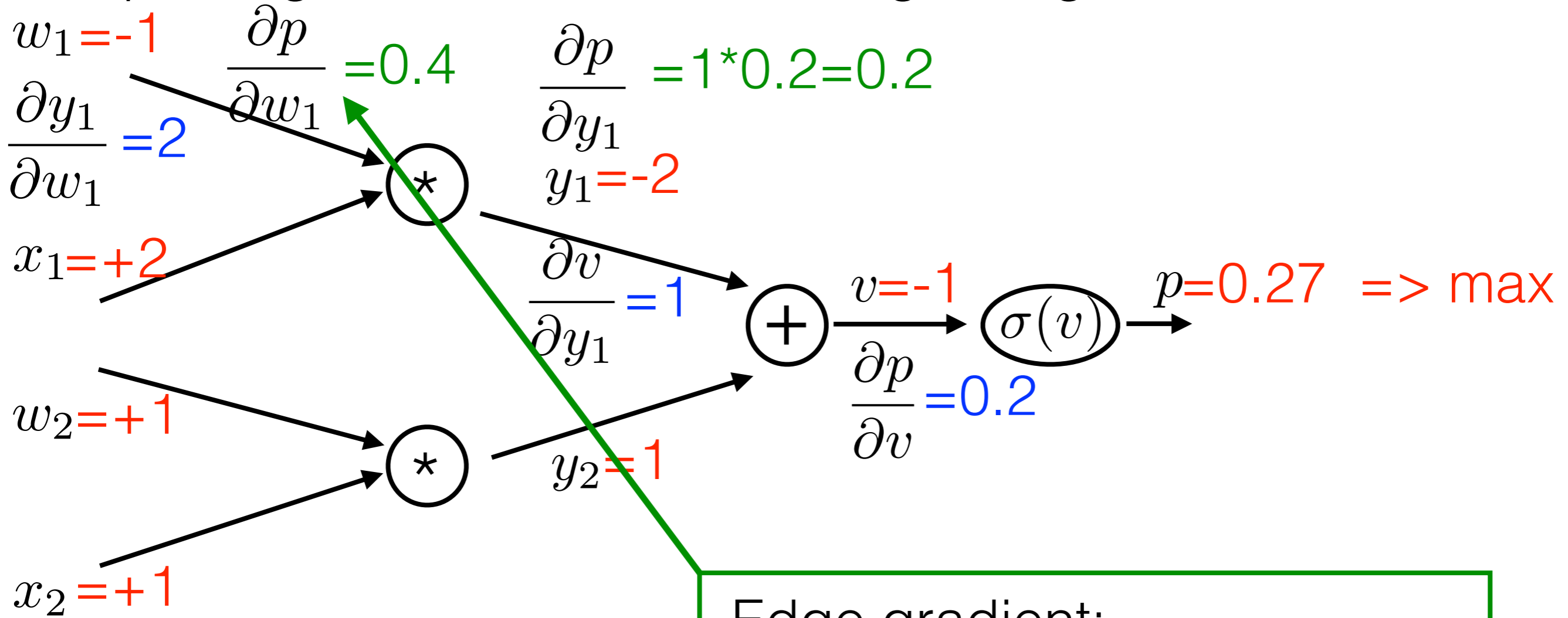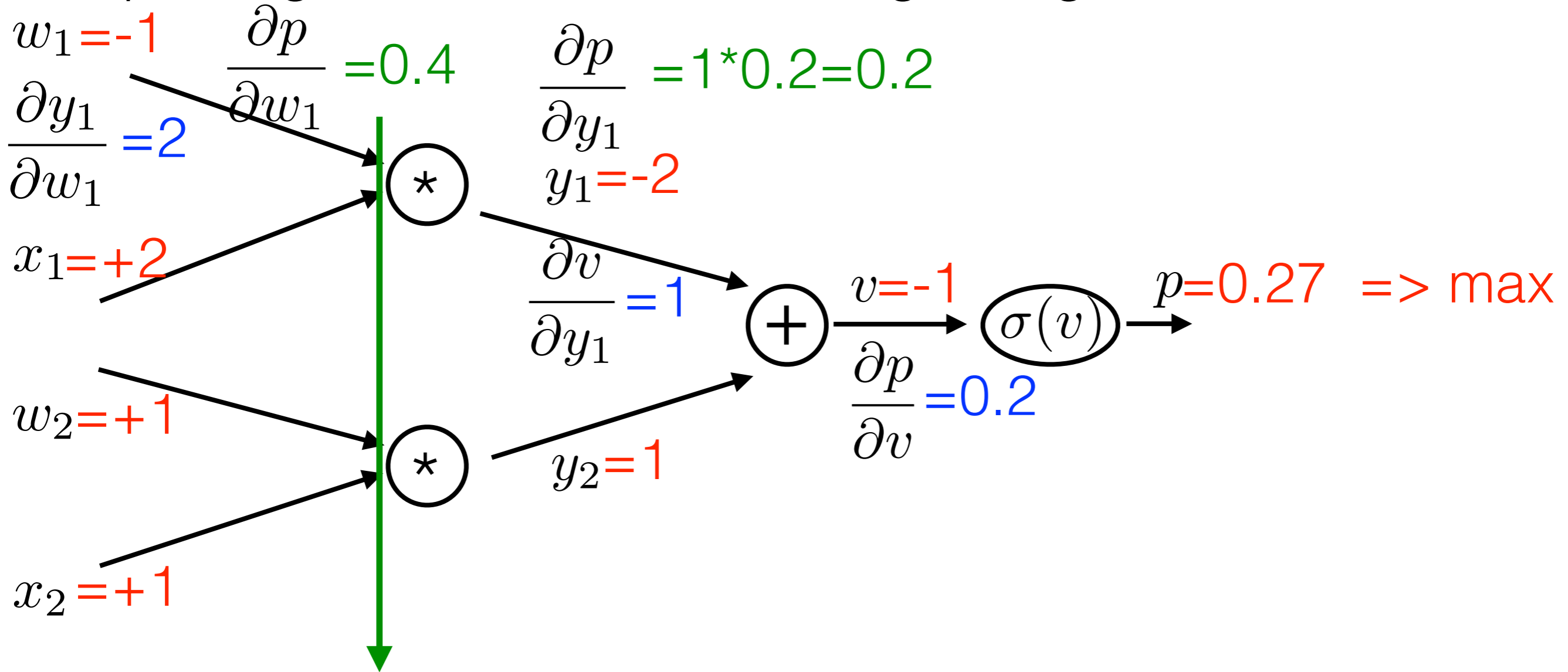Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$



$w_1$=-1

$\dfrac{\partial p}{\partial w_1}$ =0.4

$x_1$=+2

$y_1$=-2

$w_2$=+1

$x_2$=+1

$v$=-1

$\sigma(v)$

$p$=0.27  => max

$y_2$=1

$$w_1 = w_1 + \frac{\partial p}{\partial w1}$$

Chain-rule in computational graph  $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v}\dfrac{\partial v}{\partial y_1}\dfrac{\partial y_1}{\partial w_1}$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1$=-0.6

$x_1$=+2



$y_1$=-2

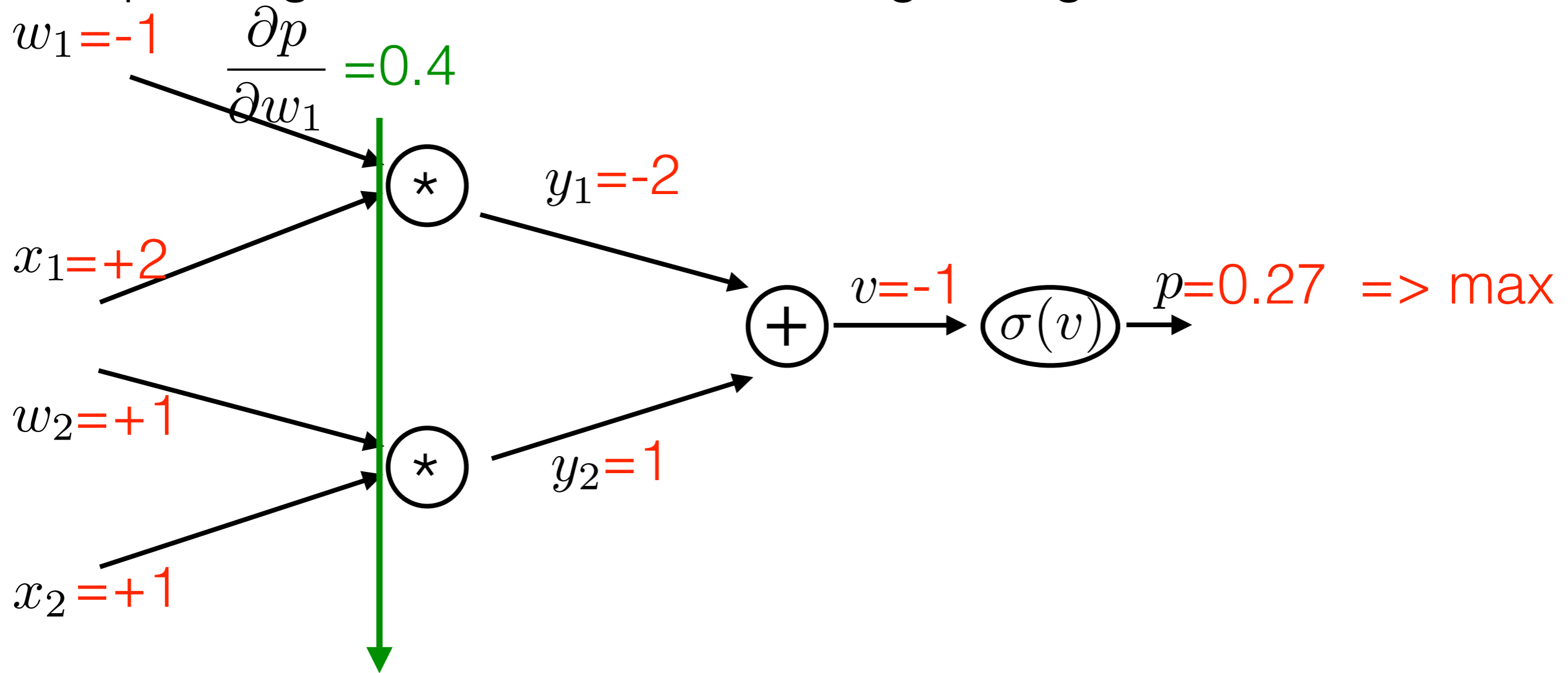$v$=-1

$p$=0.27  => max

$w_2$=+1

$x_2$=+1

$y_2$=1

$$w_1 = w_1 + \frac{\partial p}{\partial w1}$$

Chain-rule in computational graph  $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v}\dfrac{\partial v}{\partial y_1}\dfrac{\partial y_1}{\partial w_1}$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1$=-0.6



$x_1$=+2

$y_1$=-1.2
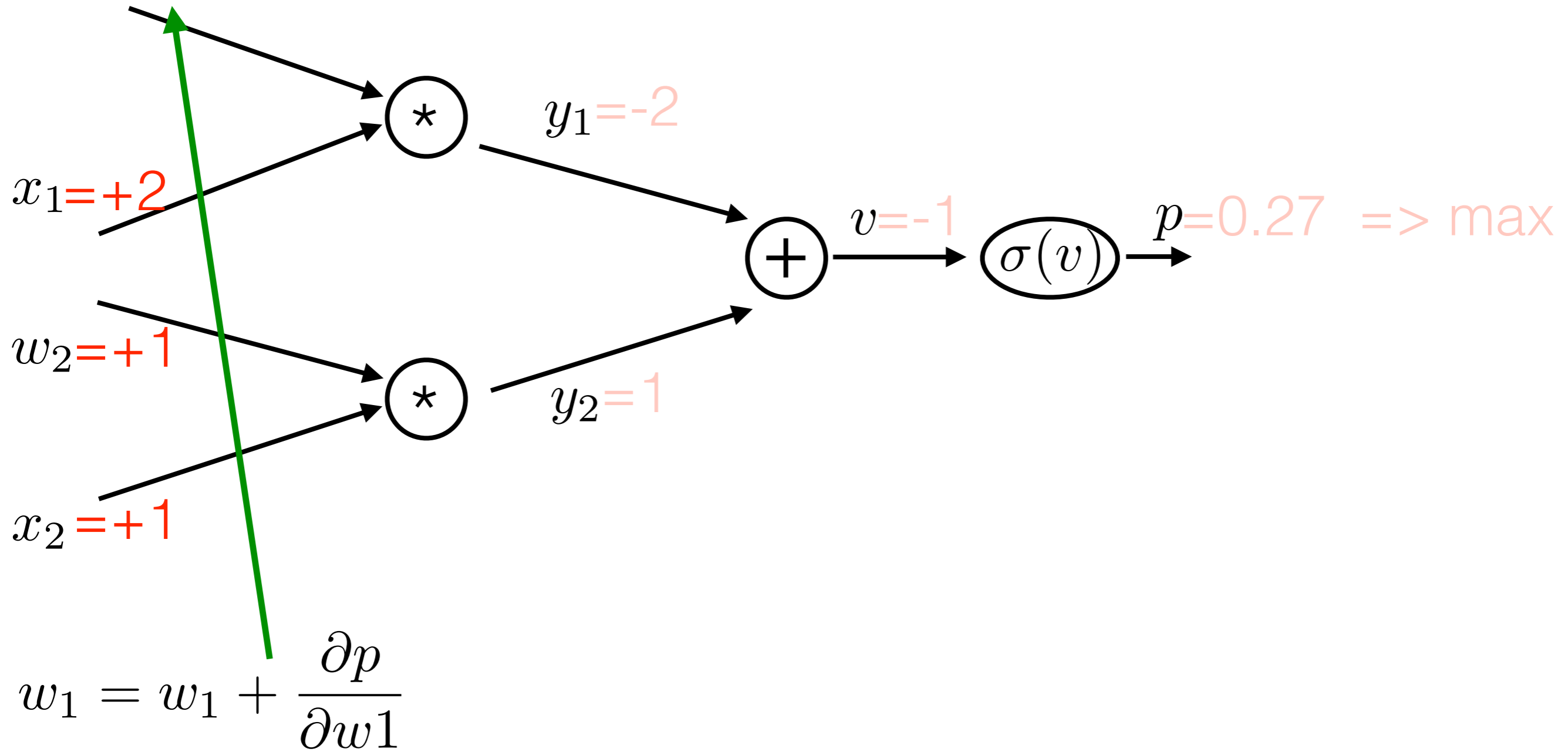
$v$=-1

$p$=0.27  => max

$w_2$=+1

$y_2$=1

$x_2$=+1

$$w_1 = w_1 + \frac{\partial p}{\partial w1}$$

Chain-rule in computational graph  $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v}\dfrac{\partial v}{\partial y_1}\dfrac{\partial y_1}{\partial w_1}$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1$=-0.6

$x_1$=+2

$y_1$=-1.2

$v$=-0.2

$p$=0.45

$\sigma(v)$

$w_2$=+1

$x_2$=+1

$y_2$=1

$p(w_1)$

0.45

0.27

-1   -0.6   $w_1$

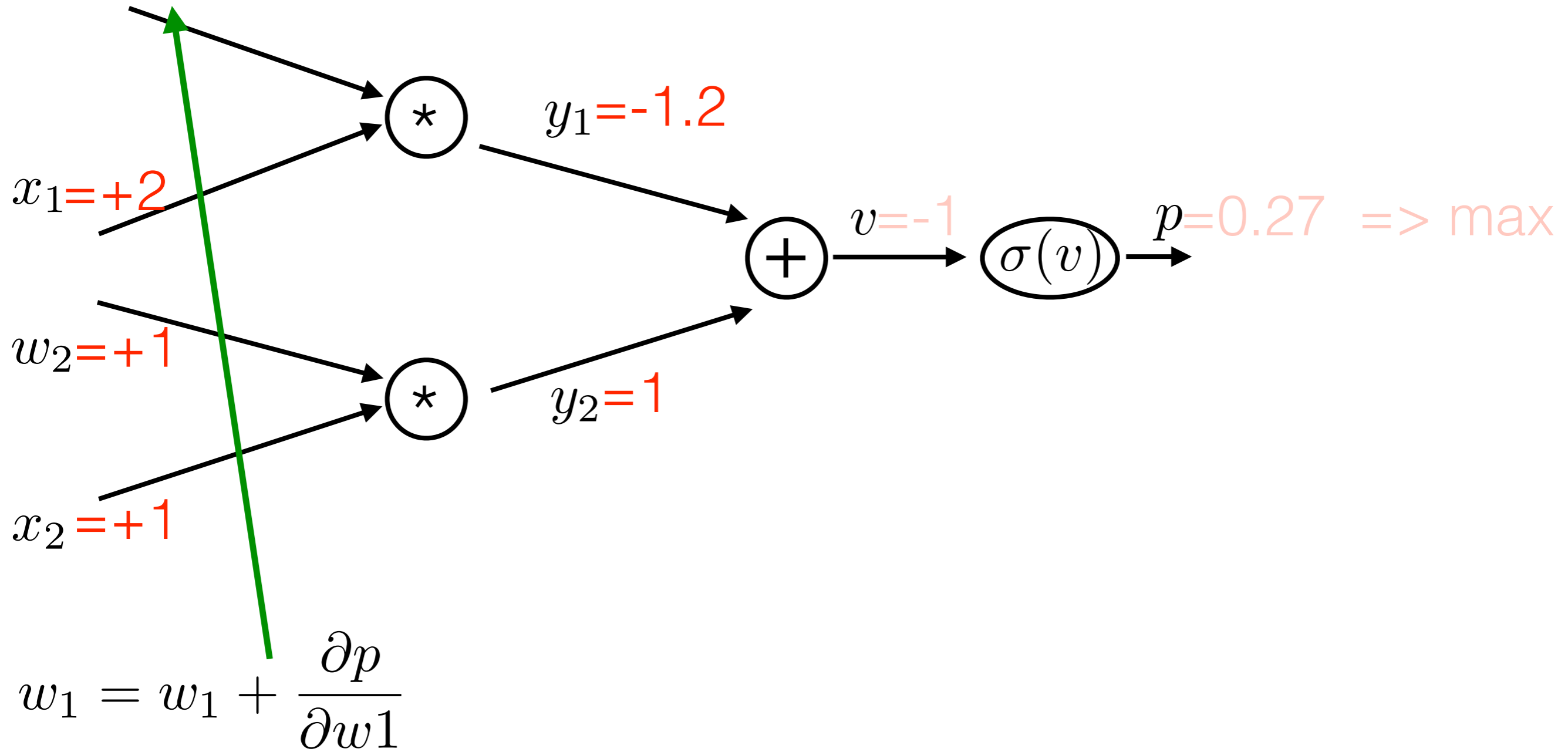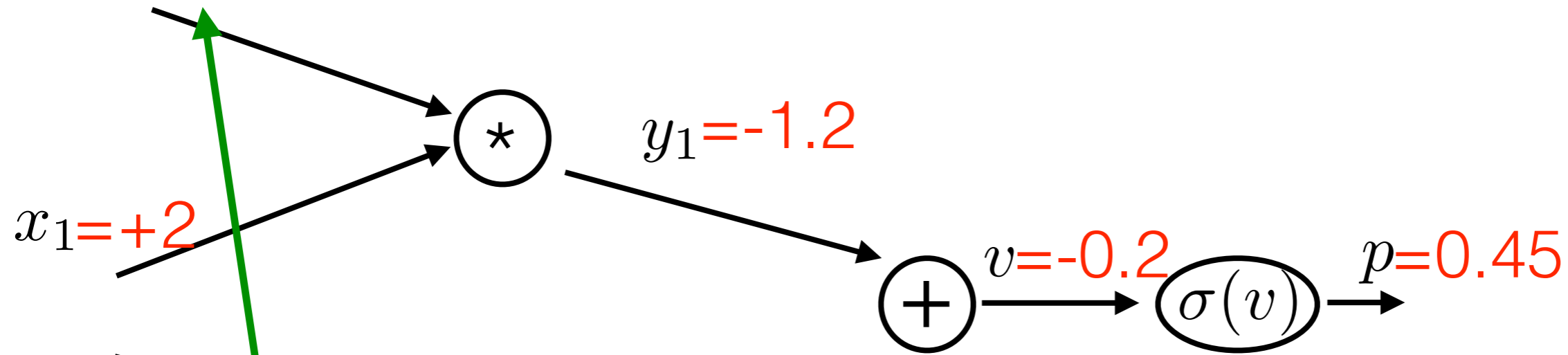$w_1 = w_1 + \dfrac{\partial p}{\partial w1}$

Chain-rule in computational graph   $\dfrac{\partial p}{\partial w_1} = \dfrac{\partial p}{\partial v}\dfrac{\partial v}{\partial y_1}\dfrac{\partial y_1}{\partial w_1}$

Example II: given $x_1, x_2, w_1, w_2$, change weight $w_1$ to increase $p$

$w_1$=-0.6

$y_1$=-1.2

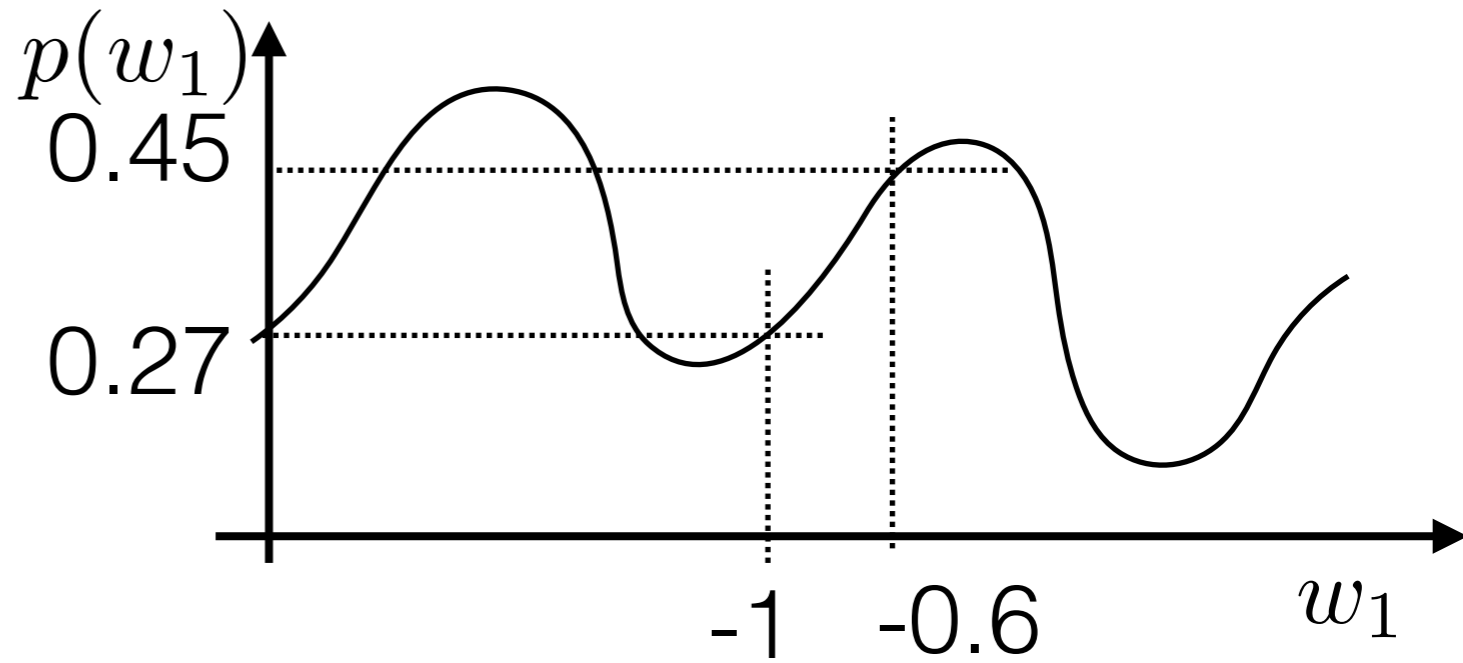$x_1$=+2

$v$=-0.2

$\sigma(v)$

$p$=0.45

$w_2$=+1

$y_2$=1

$x_2$=+1

Discussion:
- edge_gradient = upstream_gradient * local_gradient
- what is maximum $p$ (bounds)? can I also update $x_1$?
- relation to learning (max $p$ for positive samples)

# Computational graph of the learning- **adding loss layer**

$w_1$=-1

$x_1$=+2

$\ast$

$y_1$=-2

$w_2$=+1

$\ast$

$y_2$=1

$+$

$v$=-1

$\sigma(v)$

$p$=0.27

MAP estimate actually says:

Positive sample => p should be huge => minimize $-\log(p)$

# Computational graph of the learning from a **positive** sample

$w_1$ = -1

$y_1$ = -2

$x_1$ = +2

$v$ = -1    $p$ = 0.27    $\mathcal{L}$ = 0.57

$\ast$    $+$    $\sigma(v)$    $-\log(p)$

$w_2$ = +1

$y_2$ = 1

$x_2$ = +1

MAP estimate actually says:

Positive sample => p should be huge => minimize $-\log(p)$

# Computational graph of the learning from a **positive** sample

$w_1$=-1

$y_1$=-2

$x_1$=+2

$v$=-1    $p$=0.27    $\mathcal{L}$=0.57

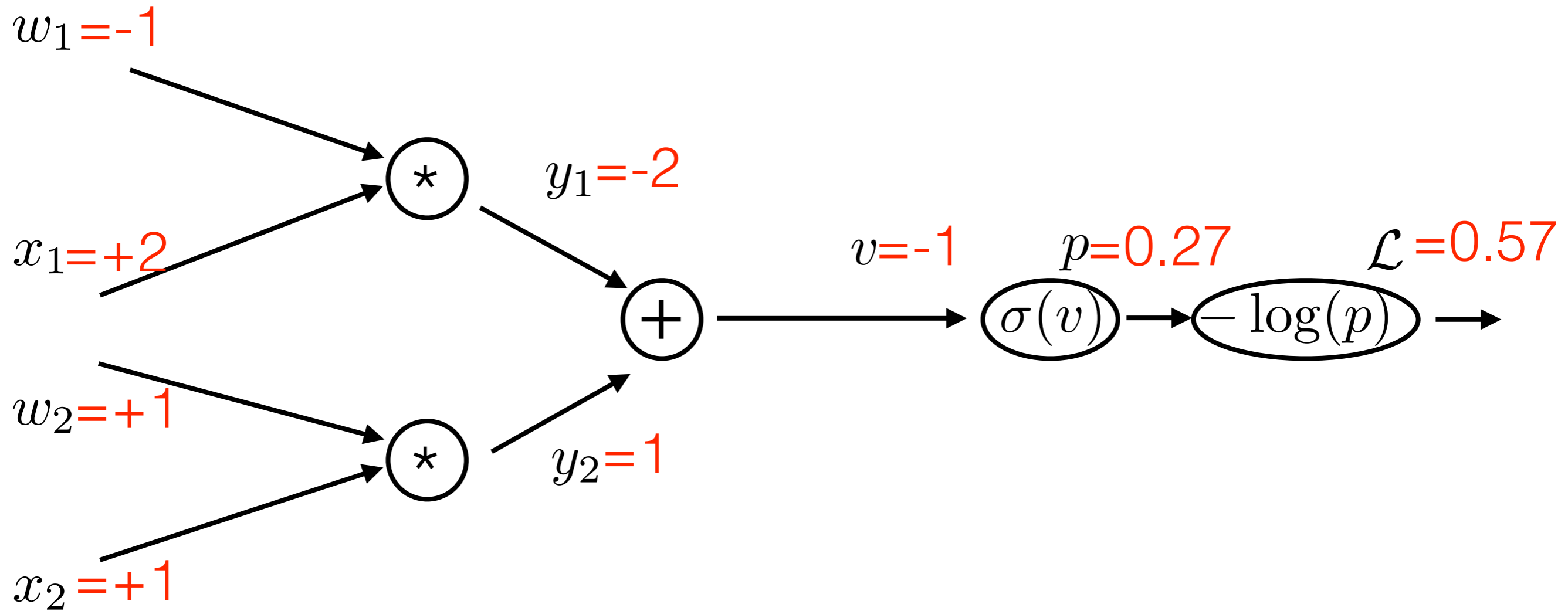$\sigma(v)$  →  $-\log(p)$

$w_2$=+1

$y_2$=1

$x_2$=+1

small p =>big loss

MAP estimate actually says:

Positive sample => p should be huge => minimize $-\log(p)$

# Computational graph of the learning from a **negative** sample

$w_1$=-1

$y_1$=-2

$x_1$=+2

$v$=-1    $p$=0.27    $\mathcal{L}$=0.13

$\sigma(v)$  →  $-\log(1-p)$

$w_2$=+1

$y_2$=1

$x_2$=+1

$*$  $+$

## MAP estimate actually says:

Positive sample => p should be huge => minimize $-\log(p)$

Negative sample => p should be small => minimize $-\log(1-p)$

# Computational graph of the learning from a **positive** sample

$w_1$=-1

$x_1$=+2

$y_1$=-2

$*$

$w_2$=+1

$x_2$=+1

$*$

$y_2$=1

$+$

$v$=-1

$\sigma(v)$

$p$=0.27

$-\log(1-p)$

$\mathcal{L}$=0.13

small p =>small loss

MAP estimate actually says:

Positive sample => p should be huge => minimize $-\log(p)$

Negative sample => p should be small => minimize $-\log(1-p)$

# Computational graph of the learning- adding loss layer

$w_1 = -1$

$x_1 = +2$

$*$

$y_1 = -2$

$v = -1$

$p = 0.27$

$\sigma(v)$

$w_2 = +1$

$*$

$y_2 = 1$

$+$

$x_2 = +1$

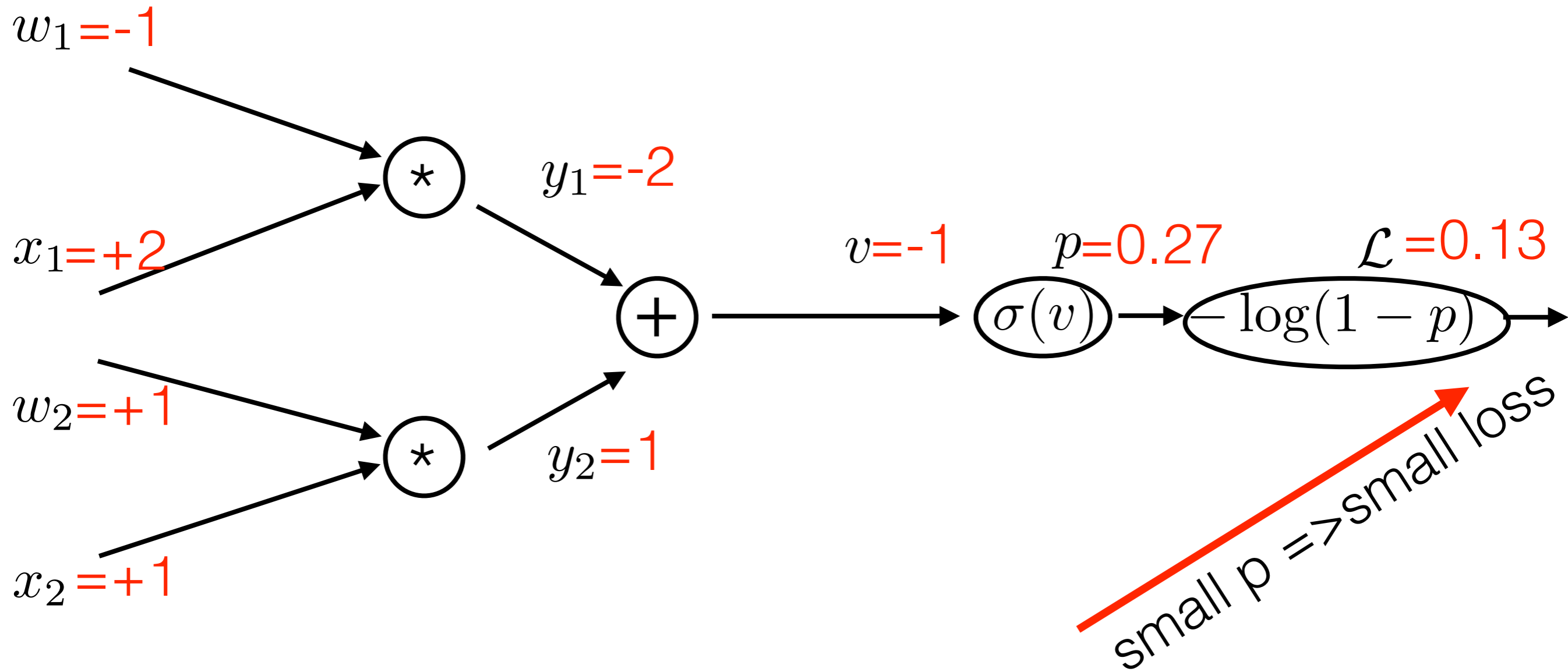MAP estimate actually says:

Positive sample => p should be huge => minimize $-\log(p)$

Negative sample => p should be small => minimize $-\log(1-p)$

We will unify computational graph for both cases as follows.

# Computational graph of the learning from a **positive** sample



$w_1$=-1

$x_1$=+2

$w_2$=+1

$x_2$=+1

$y_1$=-2

$y_2$=1

$v$=-1

$p$=0.27

$\mathcal{L}$=0.57

$\sigma(v)$

$-\log(p)$

# Computational graph of the learning from a **negative** sample

$w_1$=-1

$x_1$=+2

$w_2$=+1

$x_2$=+1

$*$

$y_1$=-2

$*$

$y_2$=1

$+$

$v$=-1

$\sigma(v)$

$p$=0.27

$-\log(1 - p)$

$\mathcal{L}$=0.13

# Computational graph of the learning

$w_1 = -1$

$x_1 = +2$

$w_2 = +1$

$x_2 = +1$

$y = +1$

$y_1 = -2$

$y_2 = 1$

$v = -1$

$p = 0.27$

$\mathcal{L} = 0.57$

$\sigma(v)$   $-\log(p)$

Computational graph for training on a **positive** sample

# Computational graph of the learning



$w_1$=-1

$x_1$=+2

$y_1$=-2

$w_2$=+1

$x_2$=+1

$y_2$=1

$v$=1

$p$=0.73

$\mathcal{L}$=0.13

$\sigma(v)$

$-\log(p)$

$y$=-1

## Computational graph for training on a **negative** sample

# Computational graph of the learning



$w_1 = -1$

$x_1 = +2$

$y_1 = -2$

$w_2 = +1$

$x_2 = +1$

$y_2 = 1$

$v = 1$

$p = 0.73$

$\mathcal{L} = 0.13$

$\sigma(v)$

$-\log(p)$

$y = -1$

Learning means:

Iteratively change all weights $\mathbf{w}$ to minimize $\mathcal{L}$

$$\mathbf{w} = \mathbf{w} - \alpha \left[ \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} \right]^{\top} \quad \text{where} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \left[ \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_2}, \dots \right]$$

# Computational graph of the learning



$w_1$=-1

Linear function

$x_1$=+2

$y_1$=-2

$w_2$=+1

$y_2$=1

$x_2$=+1

$v$=1

$p$=0.73

$\mathcal{L}$=0.13

$\sigma(v)$

$-\log(p)$

$y$ =-1

# Computational graph of the learning



$w_1 = -1$

$y_1 = -2$

$x_1 = +2$

$z = -1$    $v = 1$    $p = 0.73$    $\mathcal{L} = 0.13$

$*$   $\sigma(v)$   $-\log(p)$

$w_2 = +1$

$y_2 = 1$

$x_2 = +1$

Loss layer $-\log(\sigma(yz))$

$y = -1$

# Backprop in vector representation



$w_1 = -1$

$x_1 = +2$

$w_2 = +1$

$x_2 = +1$

$y = -1$

$y_1 = -2$

$y_2 = 1$

$z = -1$

$\mathcal{L} = 0.13$

$-\log(\sigma(yz))$

Loss layer $-\log(\sigma(yz))$

# Backprop in vector representation



$w_1$=-1

$x_1$=+2

$w_2$=+1

$x_2$=+1

Linear function

$y_1$=-2

$y_2$=1

$z$=-1

$-\log(\sigma(yz))$

$\mathcal{L}$=0.13

$y$=-1

# Backprop in vector representation



$w_1 = -1$

Linear function

$x_1 = +2$

$\mathbf{w}^\top \mathbf{x}$

$z = -1$

$-\log(\sigma(yz))$

$\mathcal{L} = 0.13$

$w_2 = +1$

$x_2 = +1$

$y = -1$

# Backprop in vector representation



$\mathbf{w} = [\text{-1},\text{+1}]$

$\mathbf{x} = [\text{+2},\text{+1}]$

vectorized inputs

$\mathbf{w}^\top \mathbf{x}$

$z = \text{-1}$

$-\log(\sigma(yz))$

$\mathcal{L} = 0.13$

$y = \text{-1}$

# Backprop in vector representation



$\mathbf{w} = [\text{-}1, +1]$

$\mathbf{x} = [+2, +1]$

$\mathbf{w}^\top \mathbf{x}$

$z = \text{-}1$

$-\log(\sigma(yz))$

$\mathcal{L} = 0.13$

$y = \text{-}1$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \; ?$$

# Backprop in vector representation

$\dfrac{\partial z}{\partial \mathbf{w}}$ 1x2

$\mathbf{w} = [-1, +1]$

$\mathbf{x} = [+2, +1]$

$\dfrac{\partial \mathcal{L}}{\partial z}$ 1x1

$z = -1$

$\mathbf{w}^\top \mathbf{x}$

$-\log(\sigma(yz))$

$\mathcal{L} = 0.13$

$y = -1$

$$\dfrac{\partial \mathcal{L}}{\partial \mathbf{w}} = \dfrac{\partial \mathcal{L}}{\partial z} \dfrac{\partial z}{\partial \mathbf{w}}$$

# Backprop in vector representation

$\dfrac{\partial z}{\partial \mathbf{w}}$ 1x2

$\mathbf{w} = [-1, +1]$

$\dfrac{\partial \mathcal{L}}{\partial z}$ 1x1

$\mathbf{x} = [+2, +1]$

$\mathbf{w}^\top \mathbf{x}$

$z = -1$

$-\log(\sigma(yz))$

$\mathcal{L} = 0.13$

$y = -1$

$\dfrac{\partial \mathcal{L}}{\partial \mathbf{w}} = \dfrac{\partial \mathcal{L}}{\partial z} \dfrac{\partial z}{\partial \mathbf{w}}$

Learning from multiple training samples means summing up the gradient over all samples

# Multiple neurons and layers

neuron $\sigma(\mathbf{w}_1^\top \overline{\mathbf{x}})$

$x_1$

$\sigma(\mathbf{w}_1^\top \overline{\mathbf{x}})$

$x_2$

$x_3$

# Multiple neurons and layers



$x_1$

$\sigma(\mathbf{w}_1^\top \bar{\mathbf{x}})$

$x_2$

$\sigma(\mathbf{w}_2^\top \bar{\mathbf{x}})$

$x_3$

# Multiple neurons and layers

# Multiple neurons and layers

# Multiple neurons and layers

# Multiple neurons and layers

# Multiple neurons and layers

# Multiple neurons and layers



- What is dimensionality of weights?

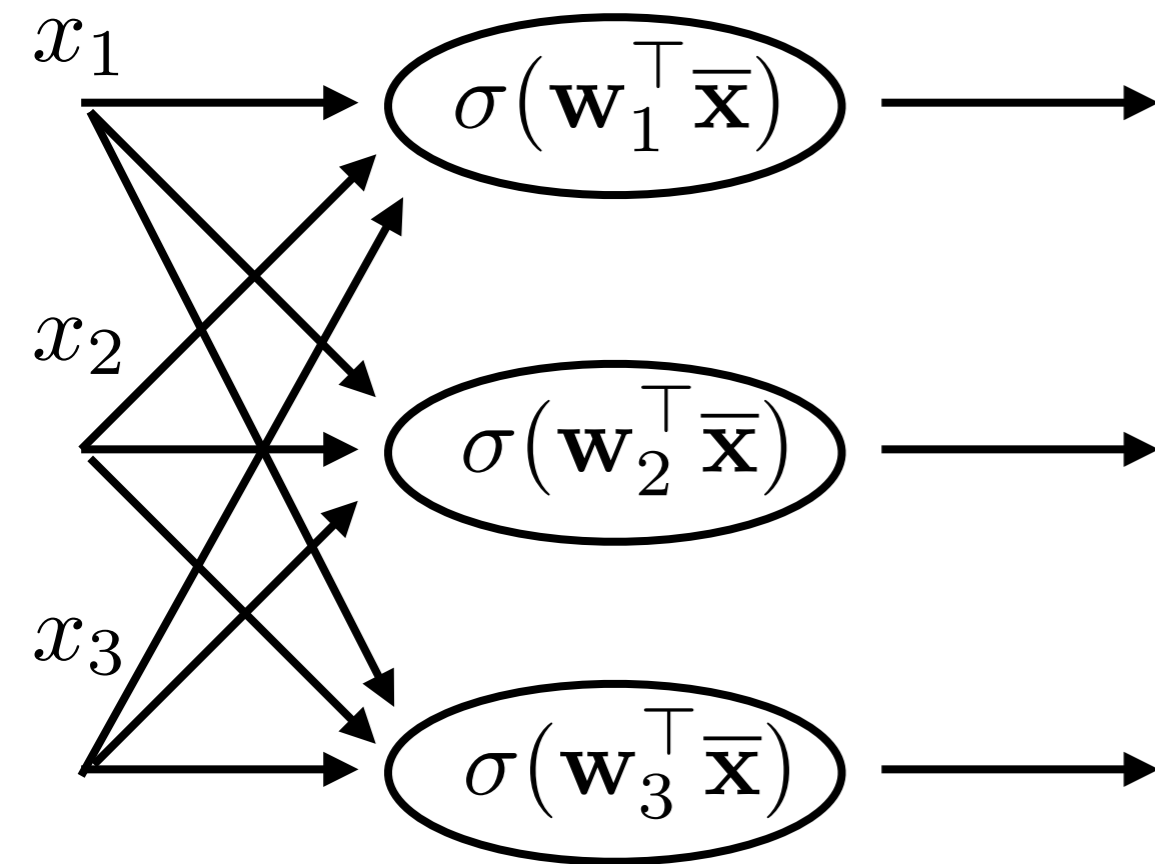# Learning of fully connected neural network

# Learning of fully connected neural network

# Learning of fully connected neural network

# Learning of fully connected neural network

# Learning of fully connected neural network

$$\mathbf{w} = \mathrm{vec}(\mathtt{W}) \qquad \mathbf{v} = \mathrm{vec}(\mathtt{V})$$

# Learning of fully connected neural network

Derivative wrt $\mathbf{u}$ : $\dfrac{\partial \mathcal{L}}{\partial \mathbf{u}} = \ ?$



$\mathbf{w}$

$\mathbf{v}$

$\mathbf{u}$

$\mathbf{x}$

$\mathbf{y}$

$\mathbf{z}$

$q$

$\mathcal{L}$

$\sigma(\mathsf{W}\mathbf{x})$

$\sigma(\mathsf{V}\mathbf{y})$

$\mathbf{u}^{\top}\mathbf{z}$

$-\log(\sigma(yq))$

$y$

# Learning of fully connected neural network

Derivative wrt $\mathbf{u}$ : $\dfrac{\partial \mathcal{L}}{\partial \mathbf{u}} = \dfrac{\partial \mathcal{L}}{\partial q} \dfrac{\partial q}{\partial \mathbf{u}}$



Dimensionality of the gradient???

# Learning of fully connected neural network

Derivative wrt $\mathbf{u}$ : $\dfrac{\partial \mathcal{L}}{\partial \mathbf{u}} = \dfrac{\partial \mathcal{L}}{\partial q} \dfrac{\partial q}{\partial \mathbf{u}}$

$\mathbf{u}$ 1x4    1x1 1x4

$\mathbf{w}$

$\mathbf{v}$

$\dfrac{\partial q}{\partial \mathbf{u}}$

$\mathbf{x}$

$\mathbf{y}$

$\mathbf{z}$

$q$

$\mathcal{L}$

$\sigma(\mathbf{W}\mathbf{x})$

$\sigma(\mathbf{V}\mathbf{y})$

$\mathbf{u}^{\top}\mathbf{z}$

$-\log(\sigma(yq))$

$\dfrac{\partial \mathcal{L}}{\partial q}$

Dimensionality of the gradient???

$y$

# Learning of fully connected neural network

Derivative wrt $\mathbf{v}$ : $\dfrac{\partial \mathcal{L}}{\partial \mathbf{v}} = \dfrac{\partial \mathcal{L}}{\partial q} \dfrac{\partial q}{\partial \mathbf{z}} \dfrac{\partial \mathbf{z}}{\partial \mathbf{v}}$



$\mathbf{w}$

$\dfrac{\partial \mathbf{z}}{\partial \mathbf{v}}$  $\mathbf{v}$

$\mathbf{u}$

$\mathbf{x}$

$\mathbf{y}$

$\mathbf{z}$

$q$

$\mathcal{L}$

$\sigma(\mathrm{W}\mathbf{x})$

$\sigma(\mathrm{V}\mathbf{y})$

$\mathbf{u}^{\top}\mathbf{z}$

$-\log(\sigma(yq))$

$\dfrac{\partial q}{\partial \mathbf{z}}$

$\dfrac{\partial \mathcal{L}}{\partial q}$

Dimensionality of the gradient???

$y$

# Learning of fully connected neural network

Derivative wrt $\mathbf{v}$ : $\dfrac{\partial \mathcal{L}}{\partial \mathbf{v}} = \dfrac{\partial \mathcal{L}}{\partial q} \dfrac{\partial q}{\partial \mathbf{z}} \dfrac{\partial \mathbf{z}}{\partial \mathbf{v}}$

$\mathbf{w}$

$\dfrac{\partial \mathbf{z}}{\partial \mathbf{v}}$ $\quad \mathbf{v}$

$\mathbf{u}$ $\quad$ 1x1 1x3 3x12

$\mathbf{x}$

$\mathbf{y}$

$\mathbf{z}$

$q$

$\mathcal{L}$

$\sigma(\mathrm{W}\mathbf{x})$ $\longrightarrow$ $\sigma(\mathrm{V}\mathbf{y})$ $\longrightarrow$ $\mathbf{u}^{\top}\mathbf{z}$ $\longrightarrow$ $-\log(\sigma(yq))$

$\dfrac{\partial q}{\partial \mathbf{z}}$

$\dfrac{\partial \mathcal{L}}{\partial q}$

Dimensionality of the gradient???

$y$
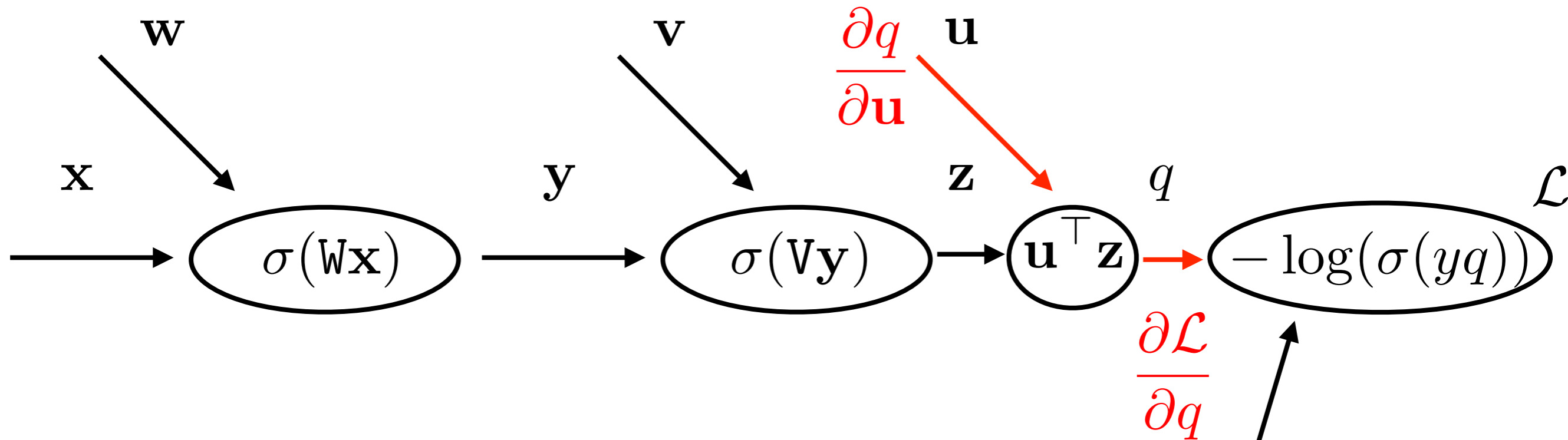
# Learning of fully connected neural network

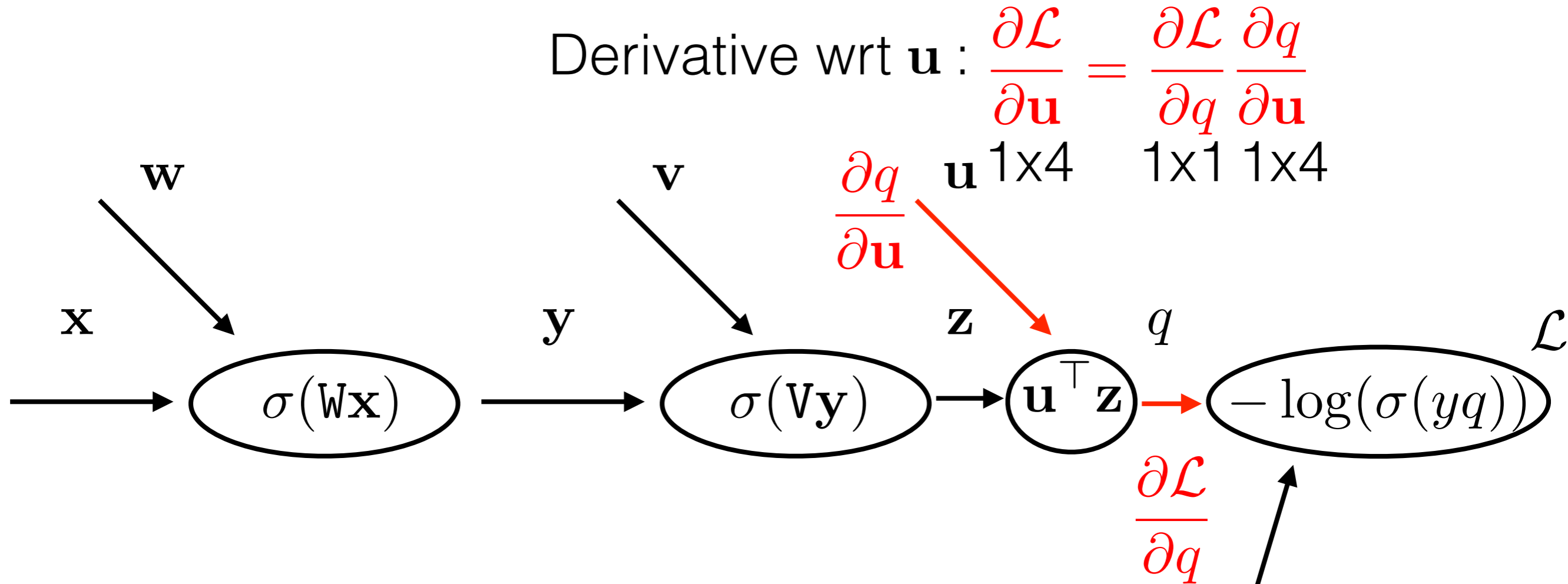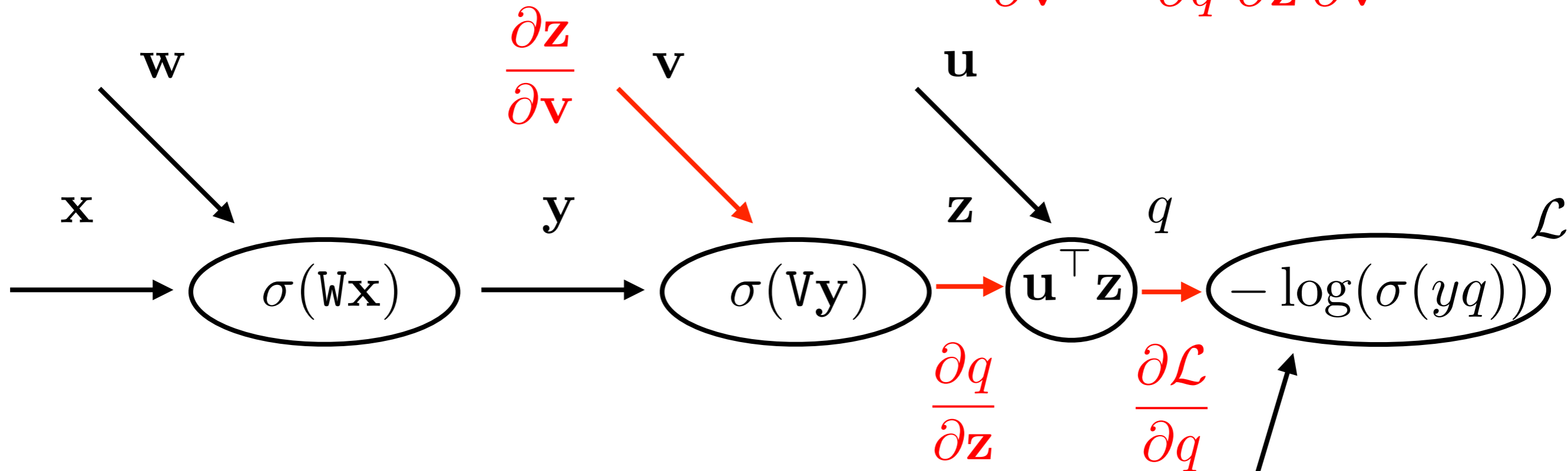Derivative wrt $\mathbf{w}$ : $\dfrac{\partial \mathcal{L}}{\partial \mathbf{w}} = \dfrac{\partial \mathcal{L}}{\partial q} \dfrac{\partial q}{\partial \mathbf{z}} \dfrac{\partial \mathbf{z}}{\partial \mathbf{y}} \dfrac{\partial \mathbf{y}}{\partial \mathbf{w}}$



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

$\mathbf{w}$ $\qquad$ $\mathbf{v}$ $\qquad$ $\mathbf{u}$

$\mathbf{x}$ $\qquad$ $\mathbf{y}$ $\qquad$ $\mathbf{z}$ $\qquad$ $q$ $\qquad$ $\mathcal{L}$

$\sigma(\mathrm{W}\mathbf{x}) \qquad \sigma(\mathrm{V}\mathbf{y}) \qquad \mathbf{u}^\top\mathbf{z} \qquad -\log(\sigma(yq))$

$\dfrac{\partial \mathbf{z}}{\partial \mathbf{y}} \qquad \dfrac{\partial q}{\partial \mathbf{z}} \qquad \dfrac{\partial \mathcal{L}}{\partial q}$

Dimensionality of the gradient???

$y$

# Learning of fully connected neural network

Derivative wrt $\mathbf{w}$ : $\dfrac{\partial \mathcal{L}}{\partial \mathbf{w}} = \dfrac{\partial \mathcal{L}}{\partial q}\dfrac{\partial q}{\partial \mathbf{z}}\dfrac{\partial \mathbf{z}}{\partial \mathbf{y}}\dfrac{\partial \mathbf{y}}{\partial \mathbf{w}}$
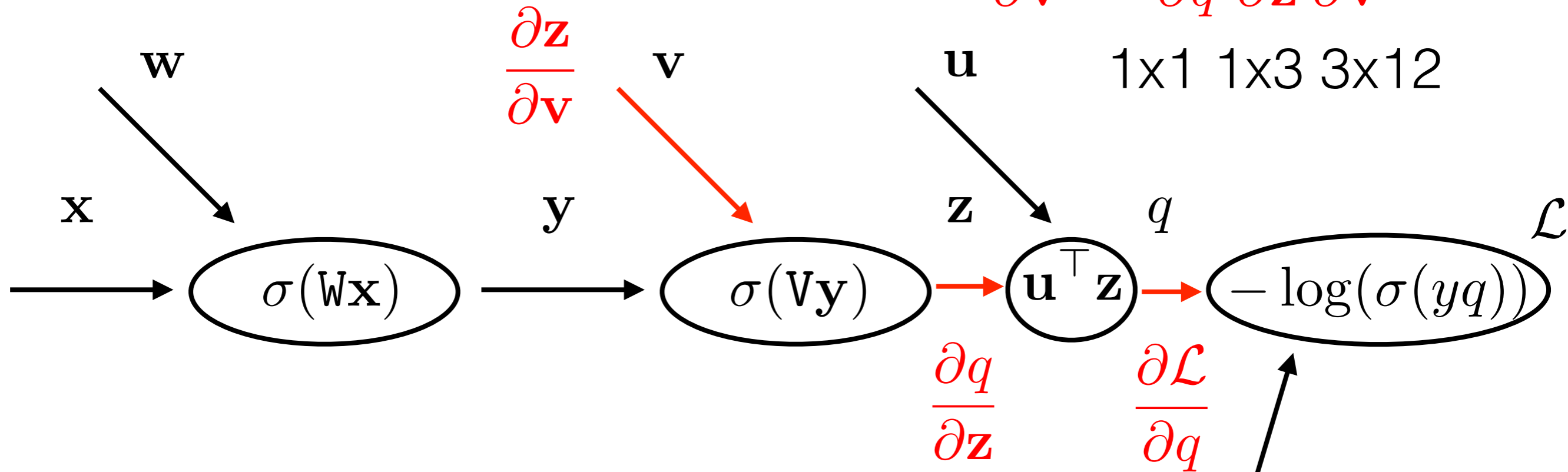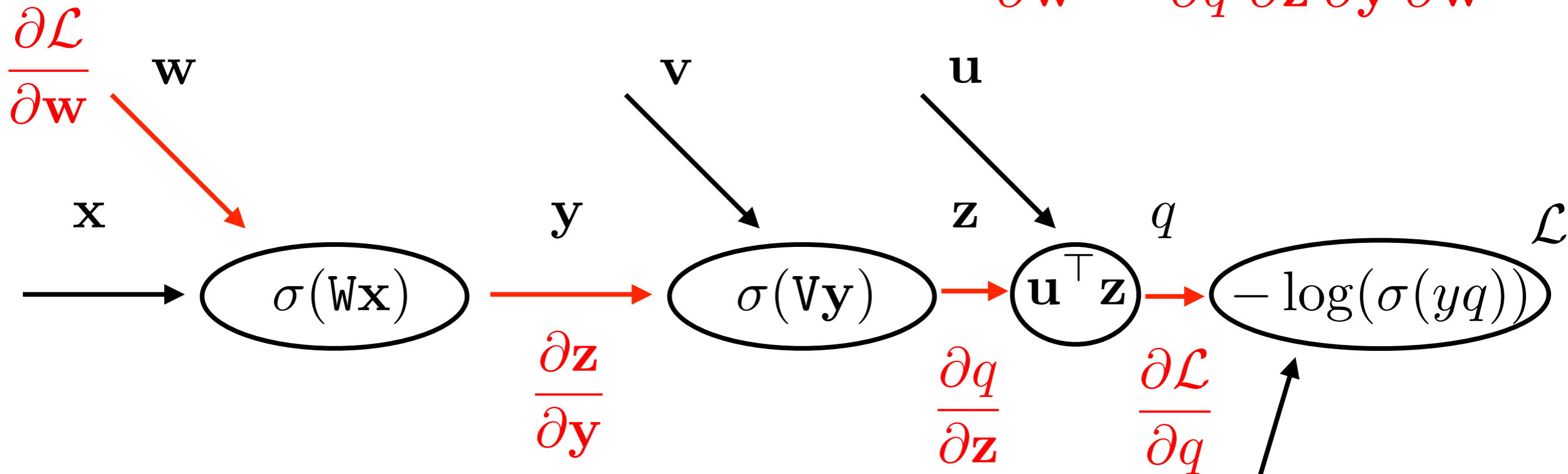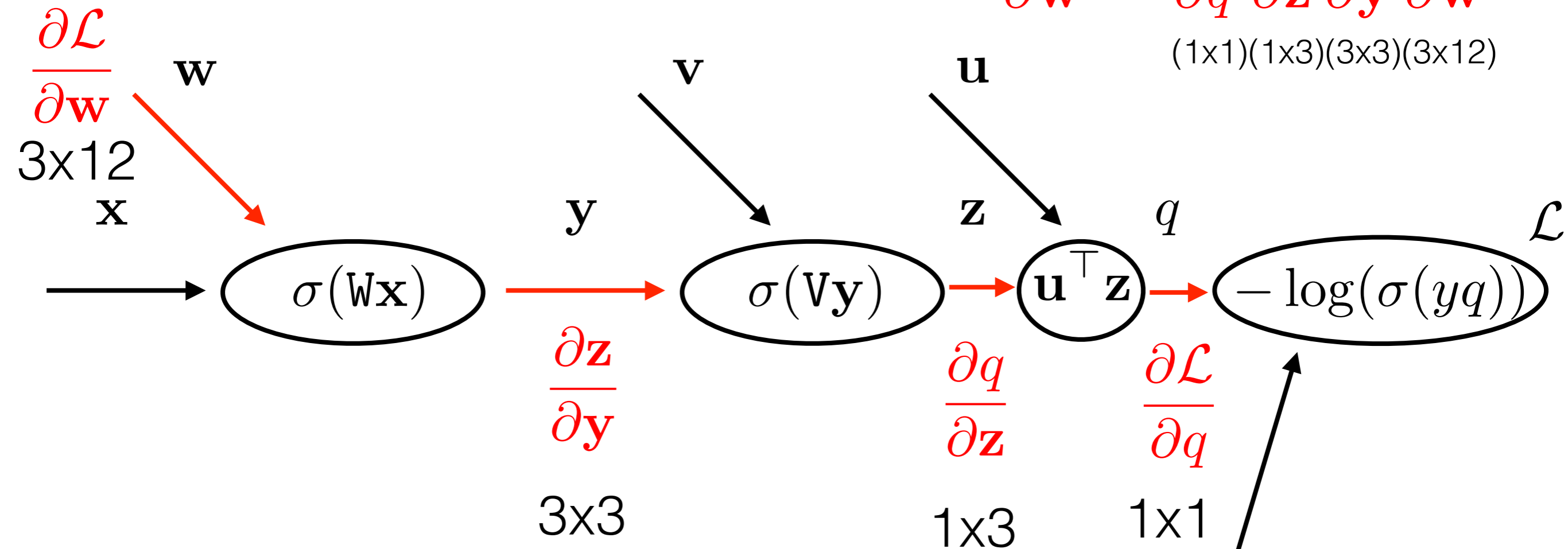
(1x1)(1x3)(3x3)(3x12)

$\dfrac{\partial \mathcal{L}}{\partial \mathbf{w}}$   $\mathbf{w}$

3x12

$\mathbf{x}$

$\mathbf{v}$       $\mathbf{u}$

$\mathbf{y}$       $\mathbf{z}$       $q$       $\mathcal{L}$

$\sigma(\mathrm{W}\mathbf{x})$       $\sigma(\mathrm{V}\mathbf{y})$       $\mathbf{u}^{\top}\mathbf{z}$       $-\log(\sigma(yq))$

$\dfrac{\partial \mathbf{z}}{\partial \mathbf{y}}$       $\dfrac{\partial q}{\partial \mathbf{z}}$       $\dfrac{\partial \mathcal{L}}{\partial q}$

3x3       1x3       1x1

Dimensionality of the gradient???

$y$

# Learning of fully connected neural network

1. Estimate all required local gradients
2. Update weights:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{u}} \qquad\qquad \mathbf{u} = \mathbf{u} - \alpha \left[ \frac{\partial \mathcal{L}}{\partial \mathbf{u}} \right]^{\top}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{v}} \qquad\qquad \mathbf{v} = \mathbf{v} - \alpha \left[ \frac{\partial \mathcal{L}}{\partial \mathbf{v}} \right]^{\top}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial q} \frac{\partial q}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{w}} \qquad \mathbf{w} = \mathbf{w} - \alpha \left[ \frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right]^{\top}$$

3. Optionally update learning rate $\alpha$
4. Repeat until convergence

# Neural nets summary

- Neural net is a function created as concatenation of simplier functions (e.g. neurons or layers of neurons)
- Gradient optimization of the neural net is called backpropagation
- Neural net frameworks has many predefined layers
- **Spoiler alert:** It does not work (on images) at all - why?

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \left( \sum_i -\log(p(y_i|\mathbf{x}_i, \mathbf{w})) \right) + (-\log p(\mathbf{w}))$$

loss function       prior/regulariser

- Class of function represented by a NN is too general.
- Naive regulariser helps a bit, but dimensionality/wildness is huge => curse-of-dimensionality, overfitting,…
- What is number of weights between two 1000-neuron layers?
- **Next lecture:** study animal cortex to find a stronger prior on the class of suitable functions.
- **Spoiler alert 2:**
  reduce very general class of functions "neuron layer" to very specific sub-class of functions "convolution layer"