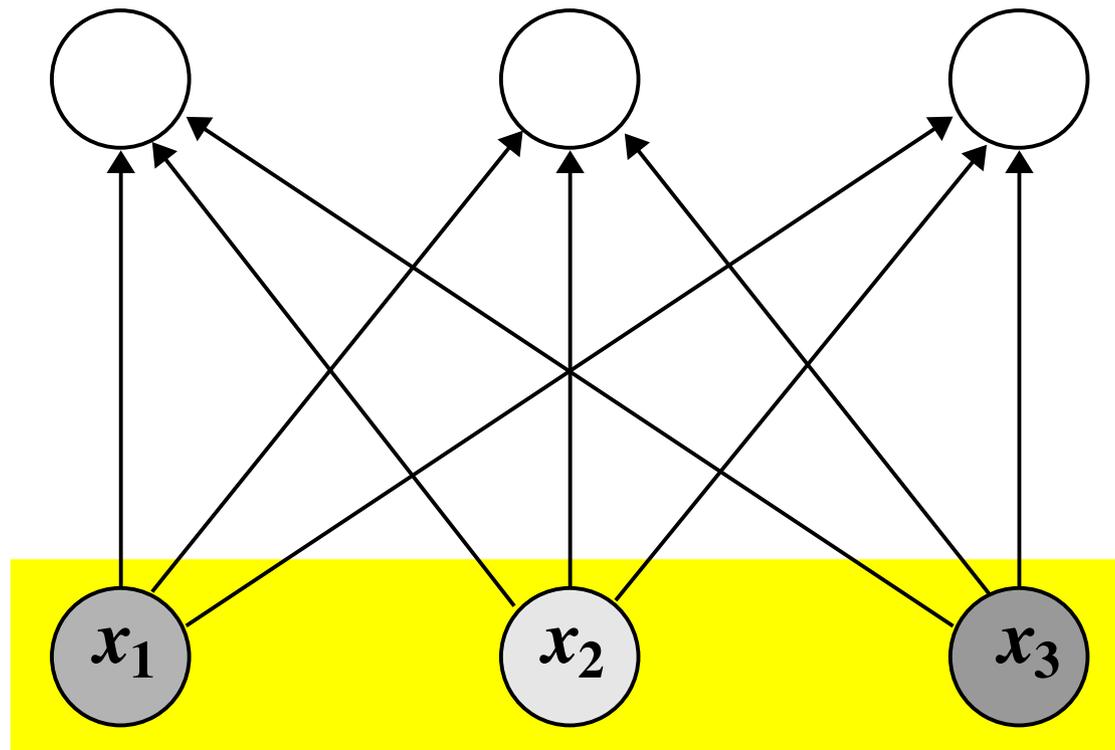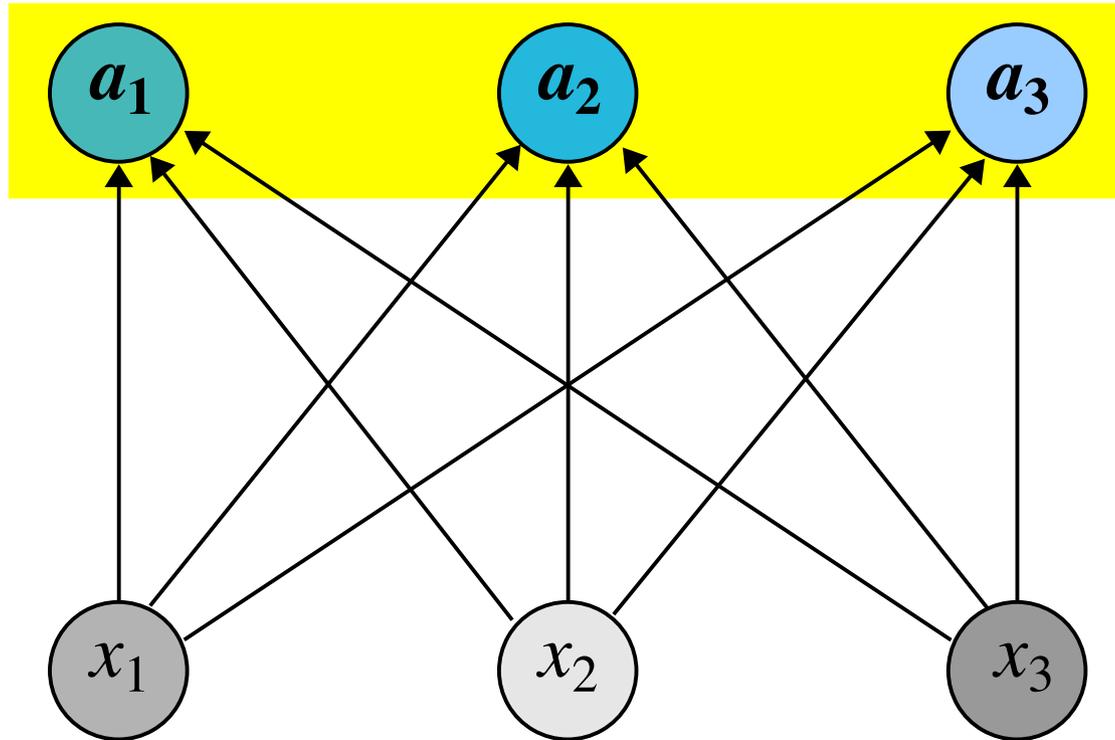# Derivation of Backpropagation
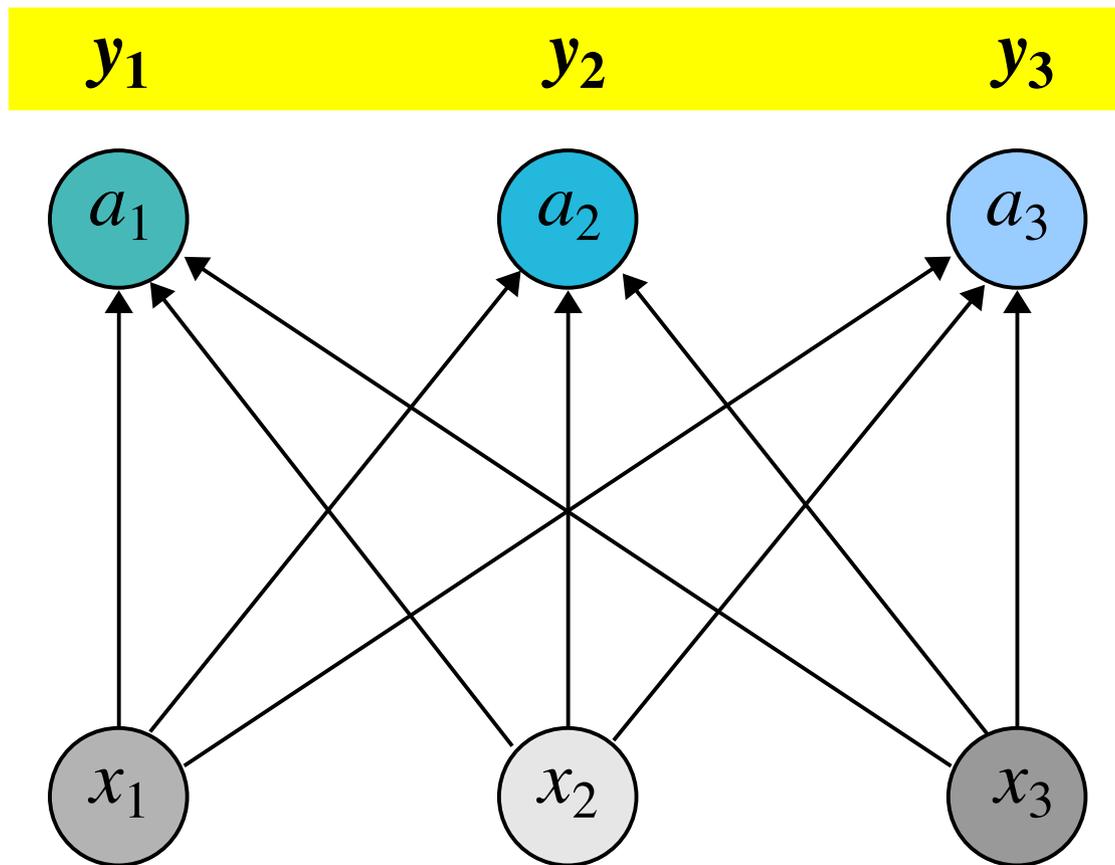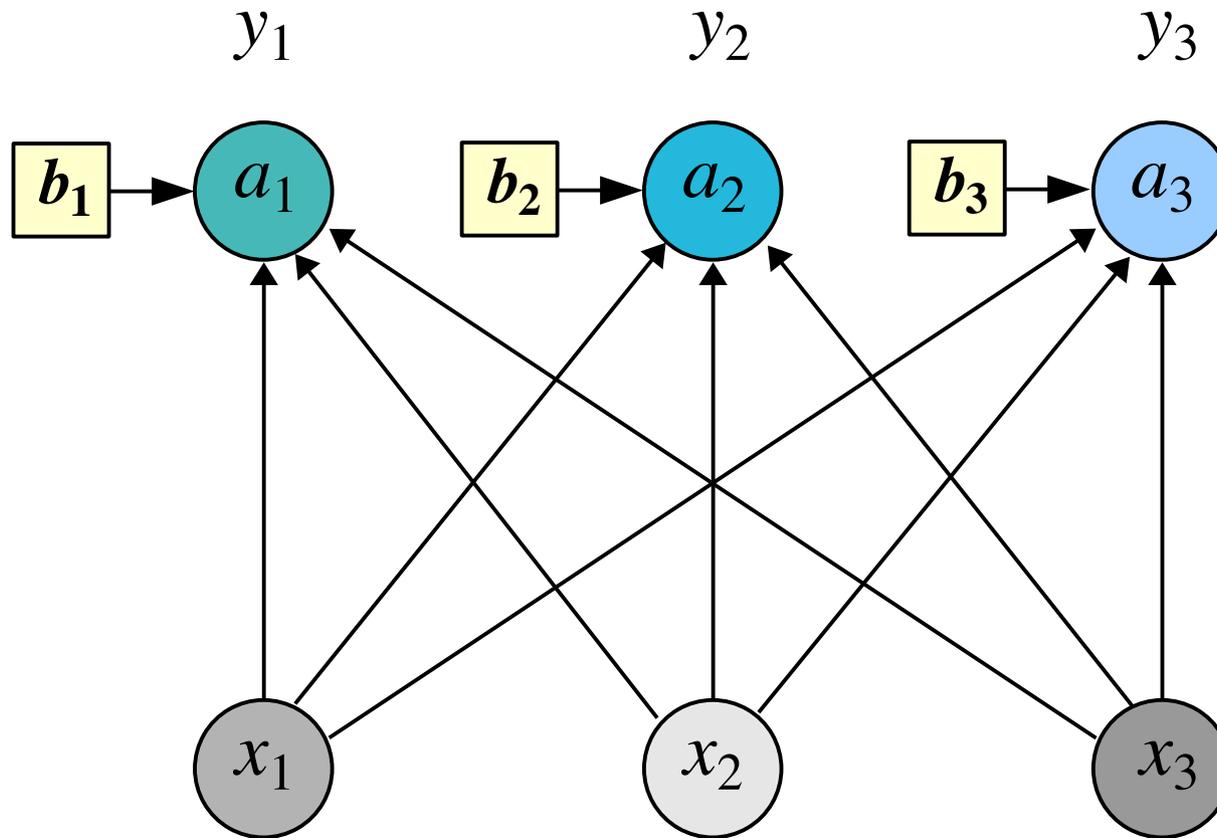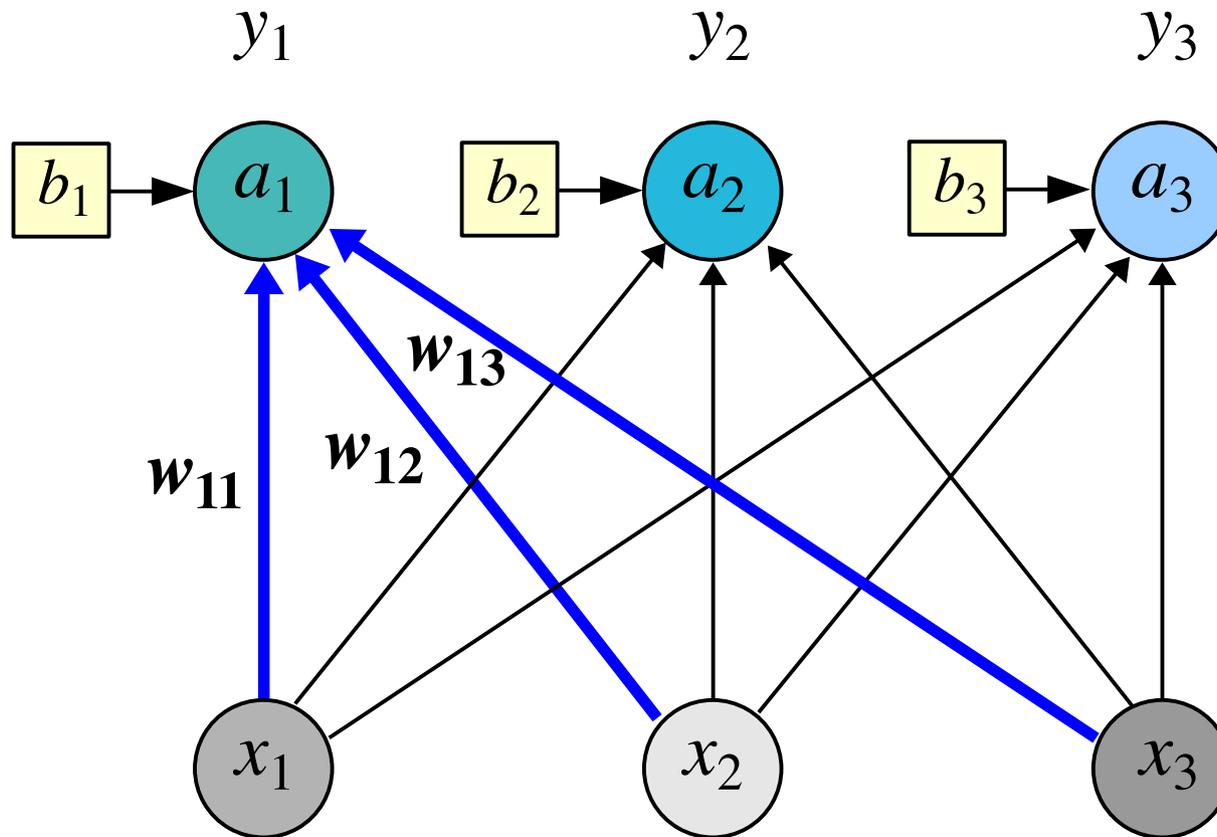
# Input Values

# Activation Values

# Target Values

# Biases

# Connection Weights

# Connection Weights

# Connection Weights

# Connection Weights

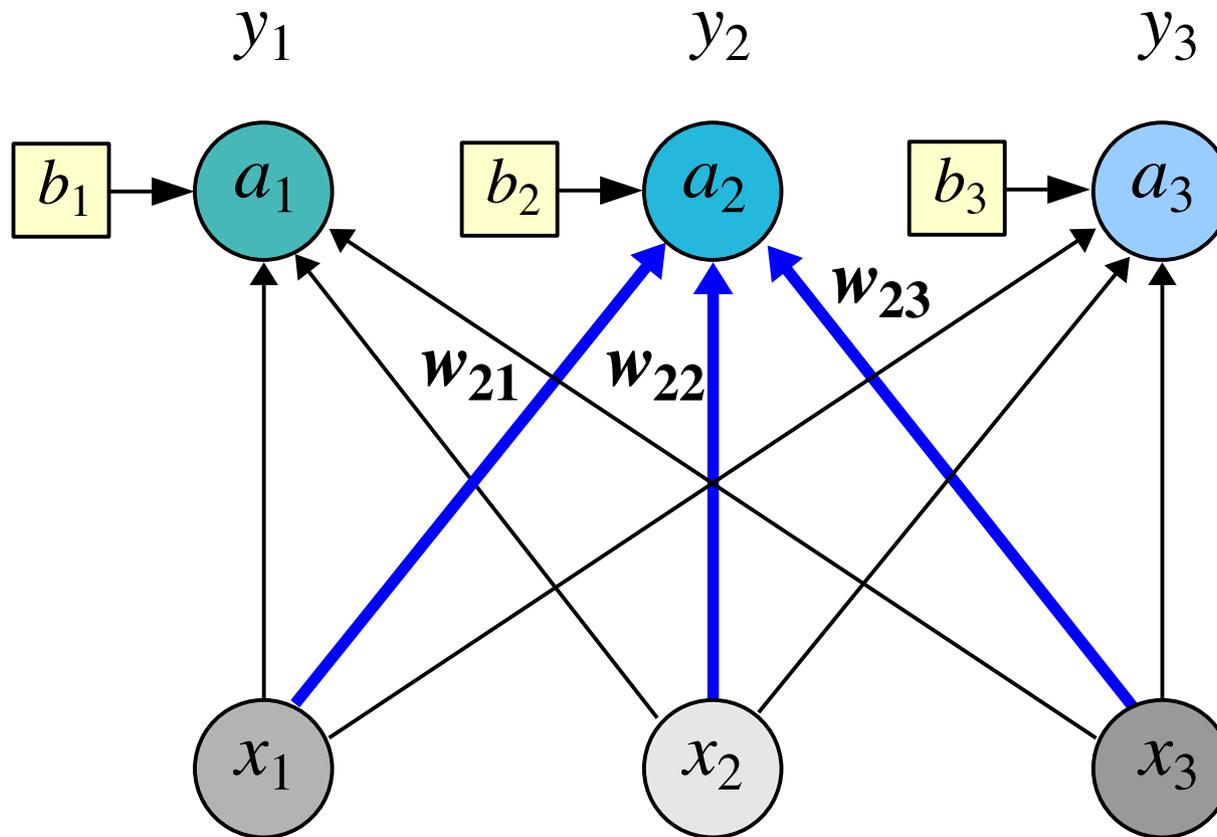# Connection Weights

# Connection Weights

# Compute Activation



$$z_2 \;=\; \boxed{w_{21}\,x_1 \;+\; w_{22}\,x_2 \;+\; w_{23}\,x_3} \;+\; b_2$$

$$z_2 \;=\; \left( \sum_k w_{2k}\,x_k \right) \;+\; b_2$$

# Compute Activation

$$a_2 = \sigma(z_2)$$



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

$$z_2 = w_{21}\,x_1 + w_{22}\,x_2 + w_{23}\,x_3 + b_2$$

$$z_2 = \left( \sum_k w_{2k}\,x_k \right) + b_2$$

# Notation



"*i*-units"

"*j*-units"

"*k*-units"

$y_i$ = target value for output unit $i$

$a_i$ = activation of output unit $i$

$b_i$ = bias of output unit $i$

# Notation



"*i*-units"

"*j*-units"

"*k*-units"

$y_i$ = target value for output unit $i$

$a_i$ = activation of output unit $i$

$b_i$ = bias of output unit $i$

$a_j$ = activation of hidden unit $j$

$b_j$ = bias of hidden unit $j$

# Notation



"*i*-units"

"*j*-units"

"*k*-units"

$y_i$ = target value for output unit $i$

$a_i$ = activation of output unit $i$

$b_i$ = bias of output unit $i$

$a_j$ = activation of hidden unit $j$

$b_j$ = bias of hidden unit $j$

$x_k$ = activation of input unit $k$

# Notation



"*i*-units"

"*j*-units"

"*k*-units"

$y_i$ = target value for output unit *i*

$a_i$ = activation of output unit *i*

$b_i$ = bias of output unit *i*

$a_j$ = activation of hidden unit *j*

$b_j$ = bias of hidden unit *j*

$x_k$ = activation of input unit *k*

$w_{ij}$ = connection weight from hidden unit *j* to output unit *i*

# Notation



"$i$-units"

"$j$-units"

"$k$-units"

$y_i$ = target value for output unit $i$

$a_i$ = activation of output unit $i$

$b_i$ = bias of output unit $i$

$a_j$ = activation of hidden unit $j$

$b_j$ = bias of hidden unit $j$

$x_k$ = activation of input unit $k$

$w_{ij}$ = connection weight from hidden unit $j$ to output unit $i$

$w_{jk}$ = connection weight from input unit $k$ to hidden unit $j$

# Forward Pass



$$z_j = \left( \sum_k w_{jk} x_k \right) + b_j \qquad a_j = \sigma(z_j)$$

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

# Forward Pass



$$z_i = \left( \sum_j w_{ij} \, a_j \right) + b_i \qquad a_i = \sigma(z_i)$$

$$z_j = \left( \sum_k w_{jk} \, x_k \right) + b_j \qquad a_j = \sigma(z_j)$$

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

# Cost/Loss for One Pattern

$$\frac{1}{2}(y_1 - a_1)^2 + \frac{1}{2}(y_2 - a_2)^2 + \dots + \frac{1}{2}(y_i - a_i)^2 + \dots$$

$y_1 \qquad y_2 \ \dots \ y_i \quad \longleftarrow \quad \text{target values}$

$a_1 \qquad a_2 \ \dots \ a_i \quad \longleftarrow \quad \text{output values}$



$$\sum_i \frac{1}{2}(y_i - a_i)^2$$

# Cost/Loss for *n* Patterns

$$\frac{1}{2}(y_1 - a_1)^2 + \frac{1}{2}(y_2 - a_2)^2 + \ldots + \frac{1}{2}(y_i - a_i)^2 + \ldots \quad \textit{(pattern 1)}$$

$$+ \frac{1}{2}(y_1 - a_1)^2 + \frac{1}{2}(y_2 - a_2)^2 + \ldots + \frac{1}{2}(y_i - a_i)^2 + \ldots \quad \textit{(pattern 2)}$$

$$+ \frac{1}{2}(y_1 - a_1)^2 + \frac{1}{2}(y_2 - a_2)^2 + \ldots + \frac{1}{2}(y_i - a_i)^2 + \ldots \quad \textit{(pattern 3)}$$

$$+ \ldots$$

$$C = \left(\frac{1}{n}\right) \sum_p \sum_i \frac{1}{2}(y_i^{(p)} - a_i^{(p)})^2$$

averaged over all *n* patterns

# Variables

- *cursor*
- *mouse*
- *hand*

$$\frac{\partial\, mouse}{\partial\, hand}$$    Influence of *hand* on *mouse*

$$\frac{\partial\, cursor}{\partial\, mouse}$$    Influence of *mouse* on *cursor*

$$\frac{\partial\, cursor}{\partial\, hand}$$    Influence of *hand* on *cursor*

# Chain Rule

$$\frac{\partial\ cursor}{\partial\ hand} = \frac{\partial\ mouse}{\partial\ hand} \times \frac{\partial\ cursor}{\partial\ mouse}$$

"Influence of *hand* on *cursor* can act through *mouse*"

Partial derivative of *y*
with respect to *x*

Partial derivative of *F*
with respect to *x*

Partial derivative of *F*
with respect to *y*

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

"Influence of *x* on function *F* can act through *y*"

# Influence of Weight $w_{ij}$ on Cost Function $C$

"cost gradient"

$$\frac{\partial C}{\partial w_{ij}}$$

could be
any weight in
the network

$$\frac{\partial C}{\partial w_{ij}} > 0$$

How will changing $w_{ij}$ cause $C$ to change?

# Influence of Weight $w_{ij}$ on Cost Function $C$

"cost gradient"

$$\frac{\partial C}{\partial w_{ij}}$$

could be
any weight in
the network

$$\frac{\partial C}{\partial w_{ij}} > 0$$

means that **increasing** $w_{ij}$ makes $C$ get **bigger**,
and **decreasing** $w_{ij}$ makes $C$ get **smaller**
so we should **decrease** $w_{ij}$ by some amount

$$\frac{\partial C}{\partial w_{ij}} < 0$$

How will changing $w_{ij}$ cause $C$ to change?

# Influence of Weight $w_{ij}$ on Cost Function $C$

"cost gradient"

$$\frac{\partial C}{\partial w_{ij}}$$

could be
any weight in
the network

$\frac{\partial C}{\partial w_{ij}} > 0$

means that **increasing** $w_{ij}$ makes $C$ get **bigger**,
and **decreasing** $w_{ij}$ makes $C$ get **smaller**
<span style="color:red">so we should **decrease** $w_{ij}$ by some amount</span>

$\frac{\partial C}{\partial w_{ij}} < 0$

means that **increasing** $w_{ij}$ makes $C$ get **smaller**,
and **decreasing** $w_{ij}$ makes $C$ get **bigger**
<span style="color:red">so we should **increase** $w_{ij}$ by some amount</span>

# How to Update the Weights?

"cost gradient"

$$\frac{\partial C}{\partial w_{ij}}$$

could be
any weight in
the network

$$\Delta w_{ij} = -\eta \; \frac{\partial C}{\partial w_{ij}}$$

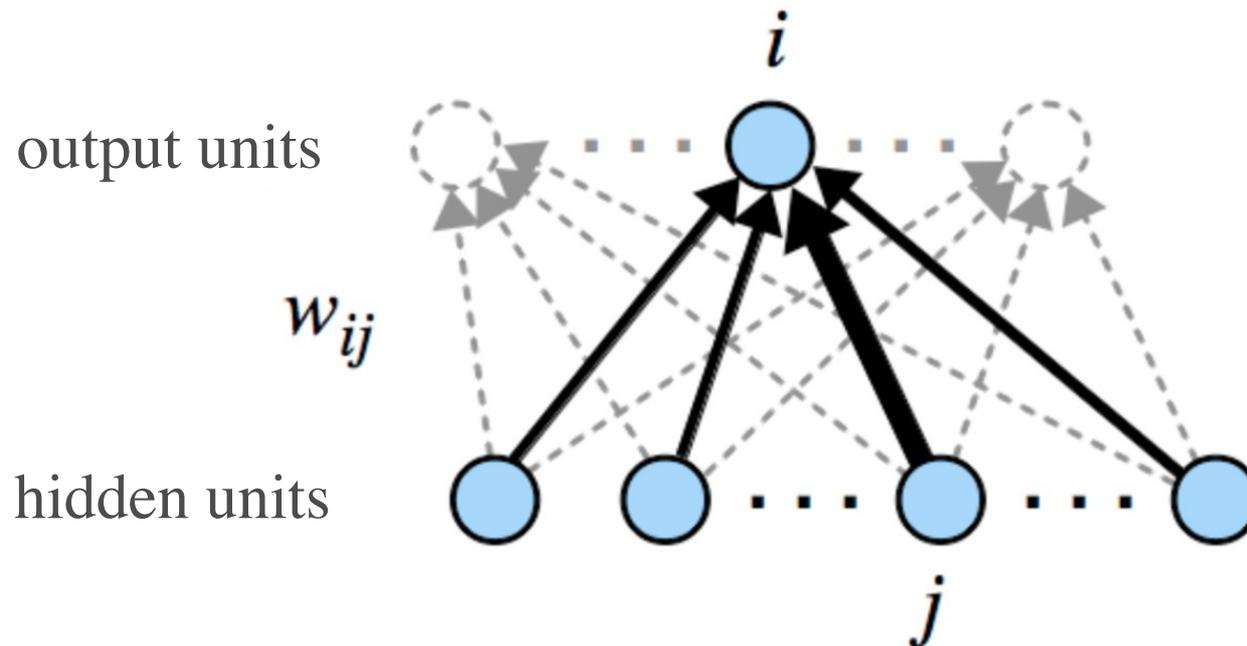amount to **change** the weight

"learning rate"   $0 < \eta < 1$

If the cost gradient is **positive**, we **decrease** the weight

If the cost gradient is **negative**, we *increase* the weight

# Hidden → Output Weights

$$\frac{\partial C}{\partial w_{ij}} = \frac{\partial z_i}{\partial w_{ij}} \times \frac{\partial a_i}{\partial z_i} \times \frac{\partial C}{\partial a_i}$$

$$\boxed{\begin{array}{c}\text{influence of}\\ w_{ij} \text{ on } C\end{array}} = \boxed{\begin{array}{c}\text{influence of}\\ w_{ij} \text{ on } z_i\end{array}} \times \boxed{\begin{array}{c}\text{influence of}\\ z_i \text{ on } a_i\end{array}} \times \boxed{\begin{array}{c}\text{influence of}\\ a_i \text{ on } C\end{array}}$$

# Hidden → Output Weights

$$\frac{\partial C}{\partial w_{ij}} = \frac{\partial z_i}{\partial w_{ij}} \times \frac{\partial a_i}{\partial z_i} \times \frac{\partial C}{\partial a_i}$$

$$\boxed{\begin{array}{c} \text{influence of} \\ w_{ij} \text{ on } C \end{array}} = \quad a_j \quad \times \quad a_i(1-a_i) \quad \times \quad (a_i - y_i)$$

# Hidden → Output Weights

$$\frac{\partial C}{\partial w_{ij}} = \frac{\partial z_i}{\partial w_{ij}} \times \frac{\partial a_i}{\partial z_i} \times \frac{\partial C}{\partial a_i}$$
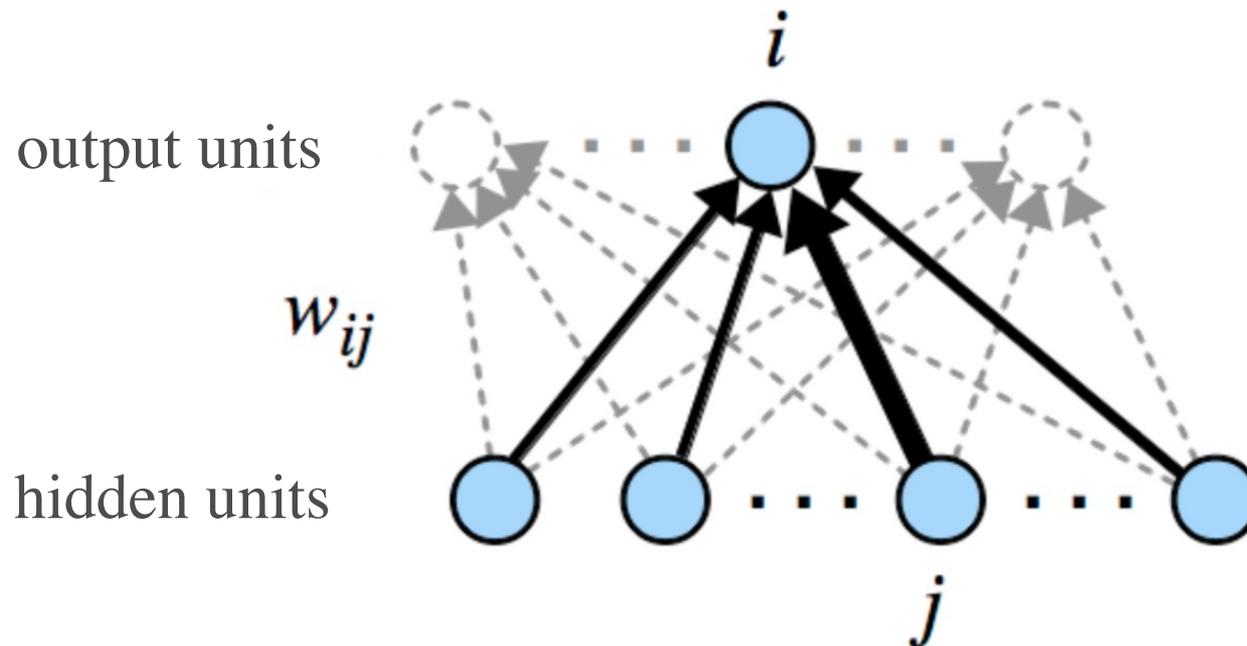
$$\frac{\partial C}{\partial w_{ij}} = (a_i - y_i)\, a_i\, (1 - a_i)\, a_j$$



$i$

output units

$w_{ij}$

hidden units

$j$

# Hidden → Output Weights

$$\frac{\partial C}{\partial w_{ij}} = \frac{\partial z_i}{\partial w_{ij}} \times \frac{\partial a_i}{\partial z_i} \times \frac{\partial C}{\partial a_i}$$
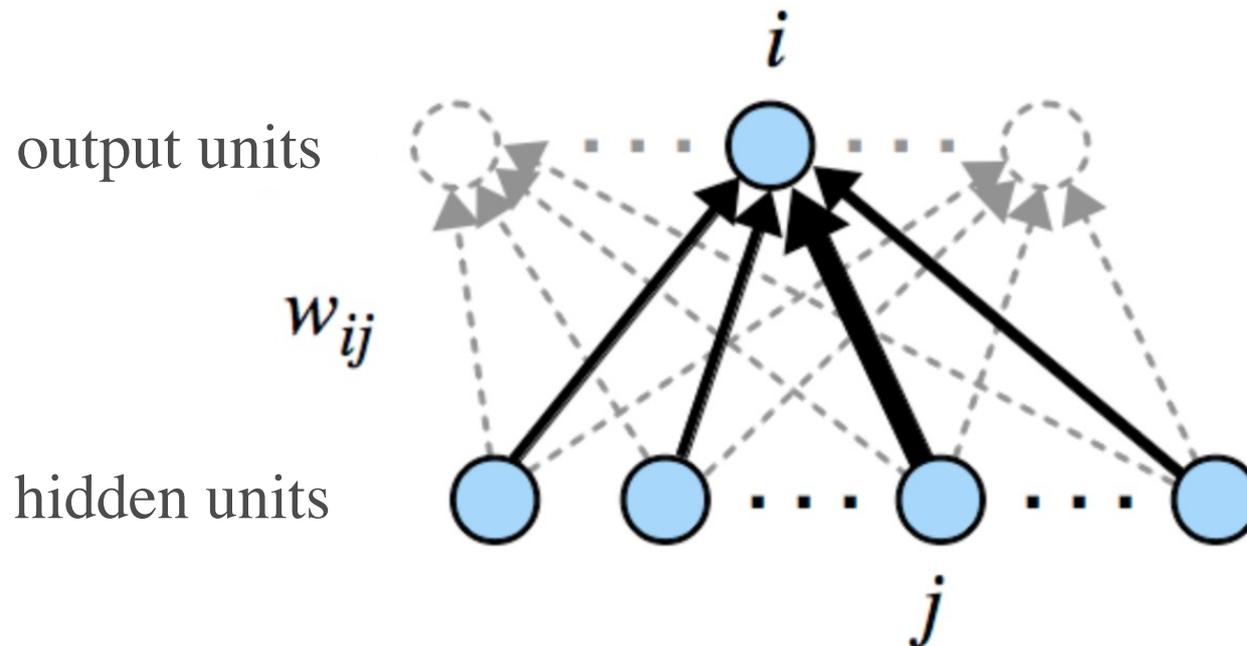
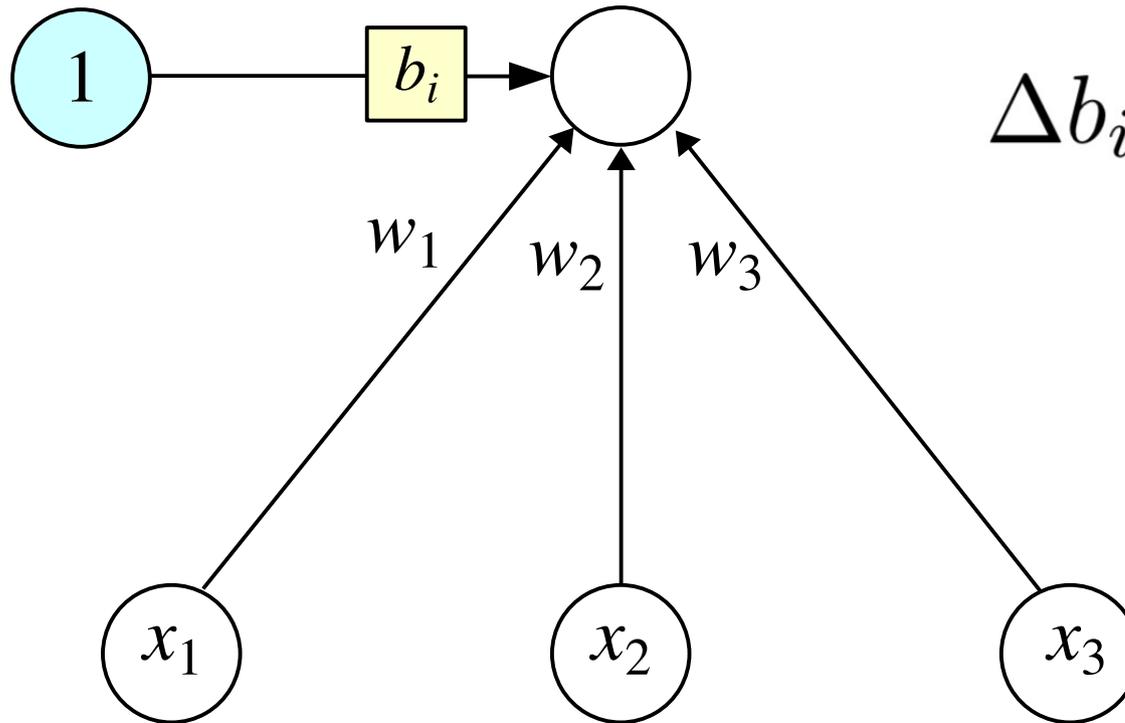$$\frac{\partial C}{\partial w_{ij}} = \boxed{(a_i - y_i)\, a_i\, (1 - a_i)}\, a_j$$

call this quantity $\delta_i$

$$\Delta w_{ij} = -\eta\, \frac{\partial C}{\partial w_{ij}} = -\eta\, \delta_i\, a_j$$

# How to Update the Bias?



$$\Delta w_{ij} = -\eta \, \delta_i \, a_j$$

$$\Delta b_i = -\eta \, \delta_i$$

$$z_i \;=\; w_1 \, x_1 \;+\; w_2 \, x_2 \;+\; w_3 \, x_3 \;+\; b_i \cdot 1$$

# Update Rule for Output Unit $i$

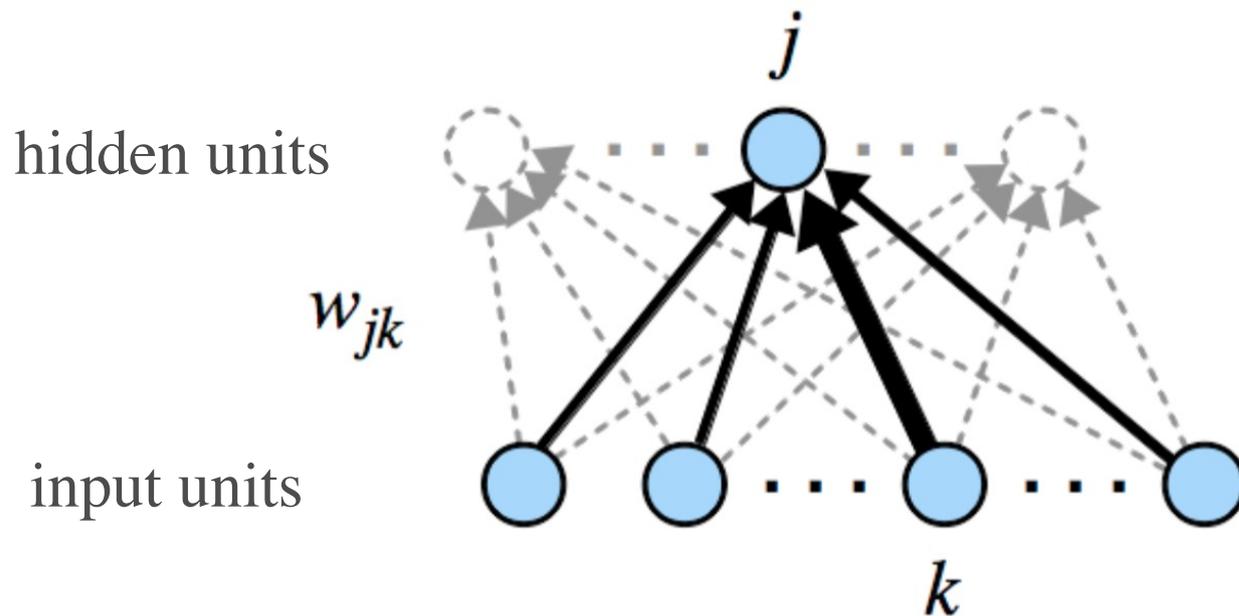$$\delta_i = (a_i - y_i)\, a_i\, (1 - a_i)$$

$$0 < \eta < 1$$

$$\Delta w_{ij} = -\eta\, \delta_i\, a_j \qquad \Delta b_i = -\eta\, \delta_i$$

# Input → Hidden Weights

$$\frac{\partial C}{\partial w_{jk}} \quad = \quad \frac{\partial z_j}{\partial w_{jk}} \quad \times \quad \frac{\partial a_j}{\partial z_j} \quad \times \quad \frac{\partial C}{\partial a_j}$$

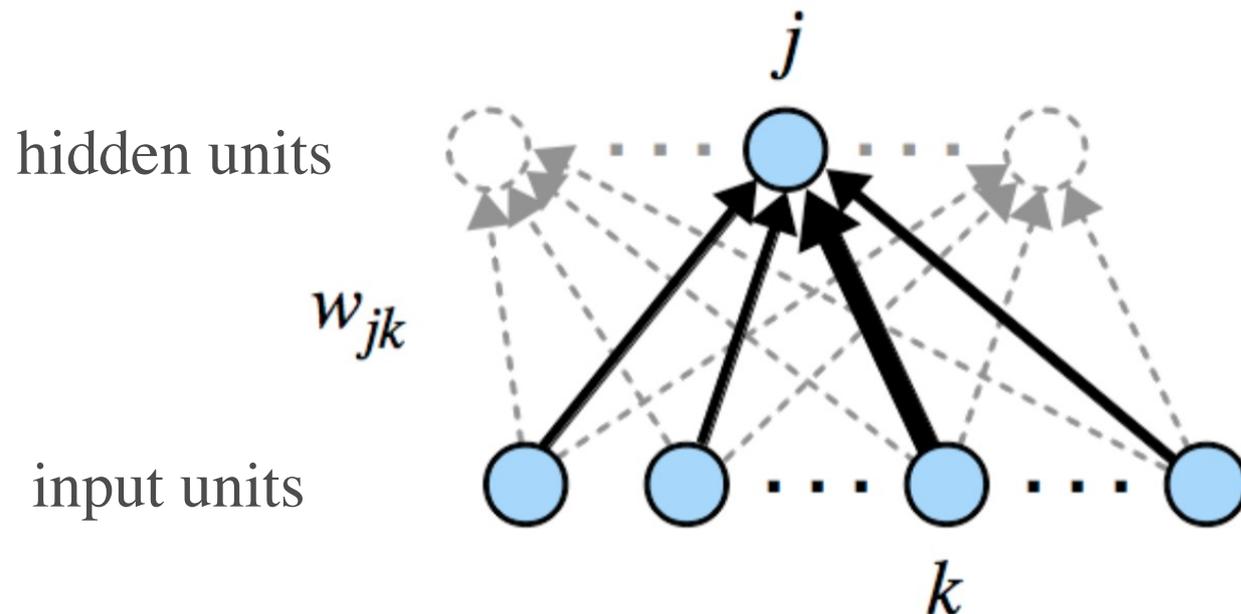| influence of $w_{jk}$ on $C$ | = | influence of $w_{jk}$ on $z_j$ | × | influence of $z_j$ on $a_j$ | × | influence of $a_j$ on $C$ |
|---|---|---|---|---|---|---|

# Input → Hidden Weights

$$\frac{\partial C}{\partial w_{jk}} = \frac{\partial z_j}{\partial w_{jk}} \times \frac{\partial a_j}{\partial z_j} \times \frac{\partial C}{\partial a_j}$$
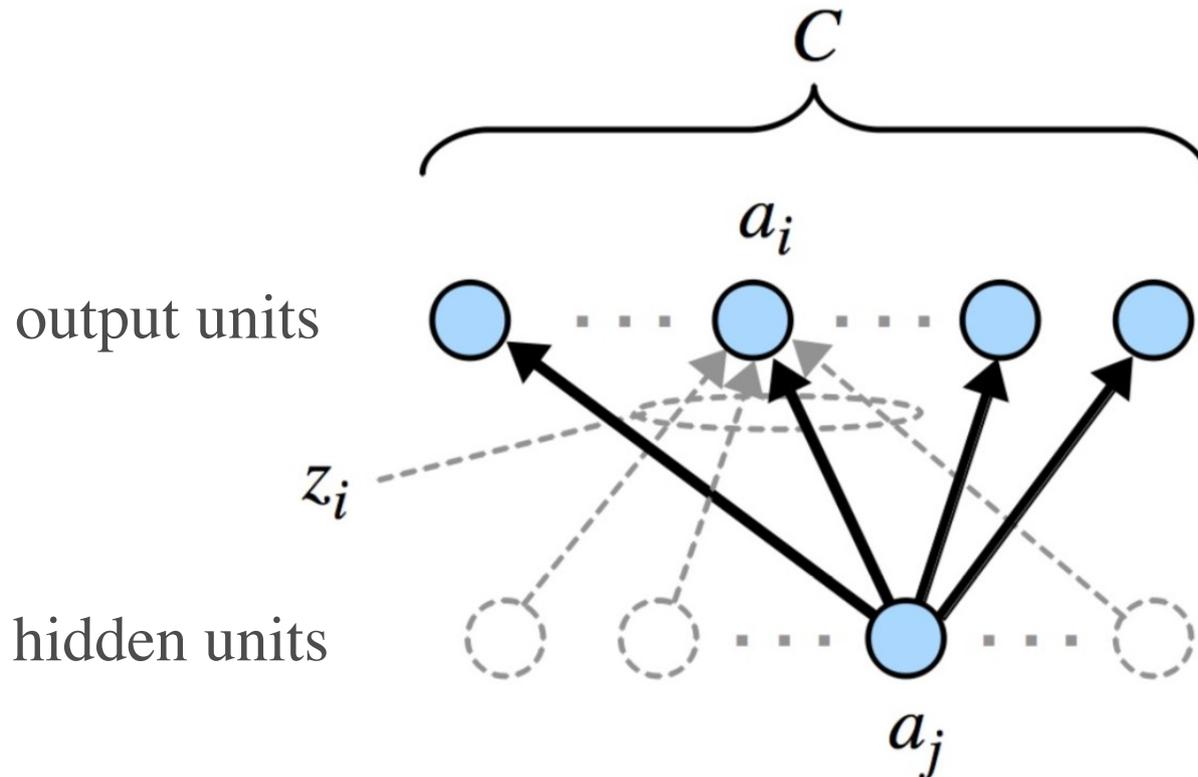
$$\boxed{\begin{array}{c} \text{influence of} \\ w_{jk} \text{ on } C \end{array}} = \quad x_k \quad \times \quad a_j(1 - a_j) \quad \times \quad ???$$



hidden units

$w_{jk}$

input units

# Influence of Hidden Unit $j$ on Cost Function

$$\frac{\partial C}{\partial a_j} = \sum_i \frac{\partial z_i}{\partial a_j} \times \frac{\partial a_i}{\partial z_i} \times \frac{\partial C}{\partial a_i}$$

$$= \sum_i \boxed{\begin{array}{c} \text{influence of} \\ a_j \text{ on } z_i \end{array}} \times \boxed{\begin{array}{c} \text{influence of} \\ z_i \text{ on } a_i \end{array}} \times \boxed{\begin{array}{c} \text{influence of} \\ a_i \text{ on } C \end{array}}$$

# Influence of Hidden Unit $j$ on Cost Function

$$\frac{\partial C}{\partial a_j} = \sum_i \frac{\partial z_i}{\partial a_j} \times \frac{\partial a_i}{\partial z_i} \times \frac{\partial C}{\partial a_i}$$
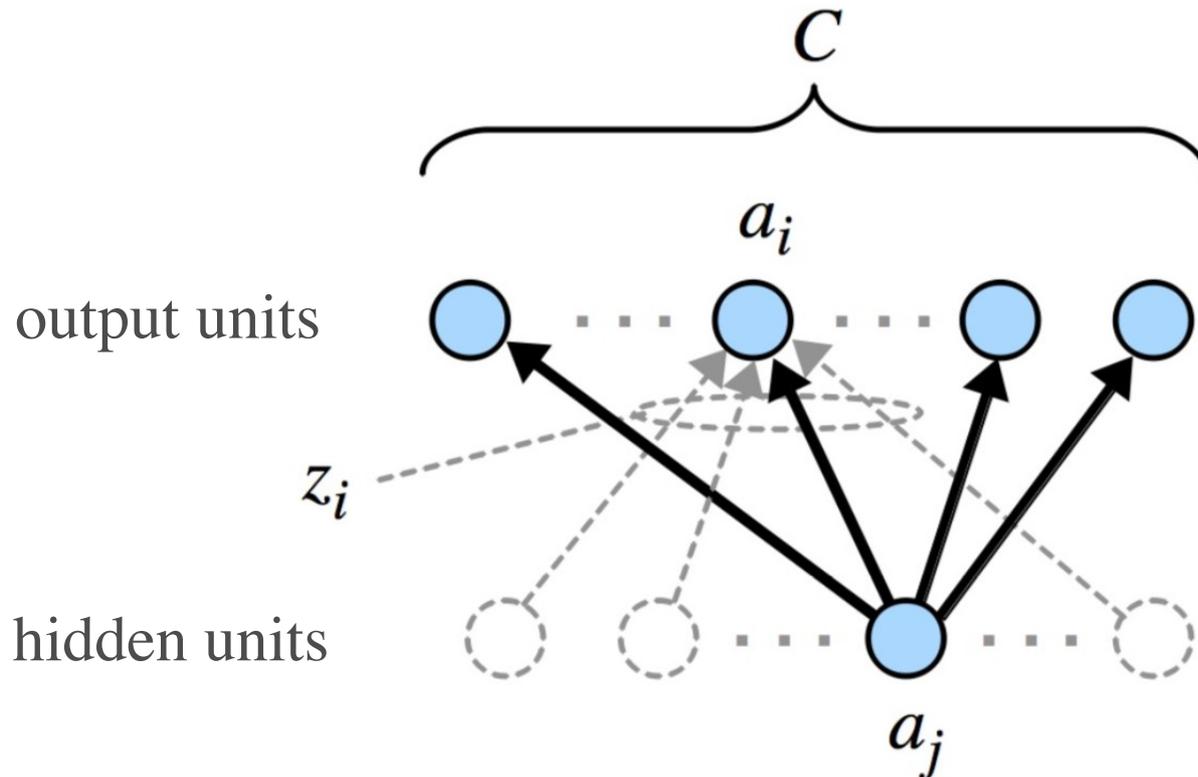
$$= \sum_i \quad w_{ij} \quad \times \quad a_i(1 - a_i) \quad \times \quad (a_i - y_i)$$
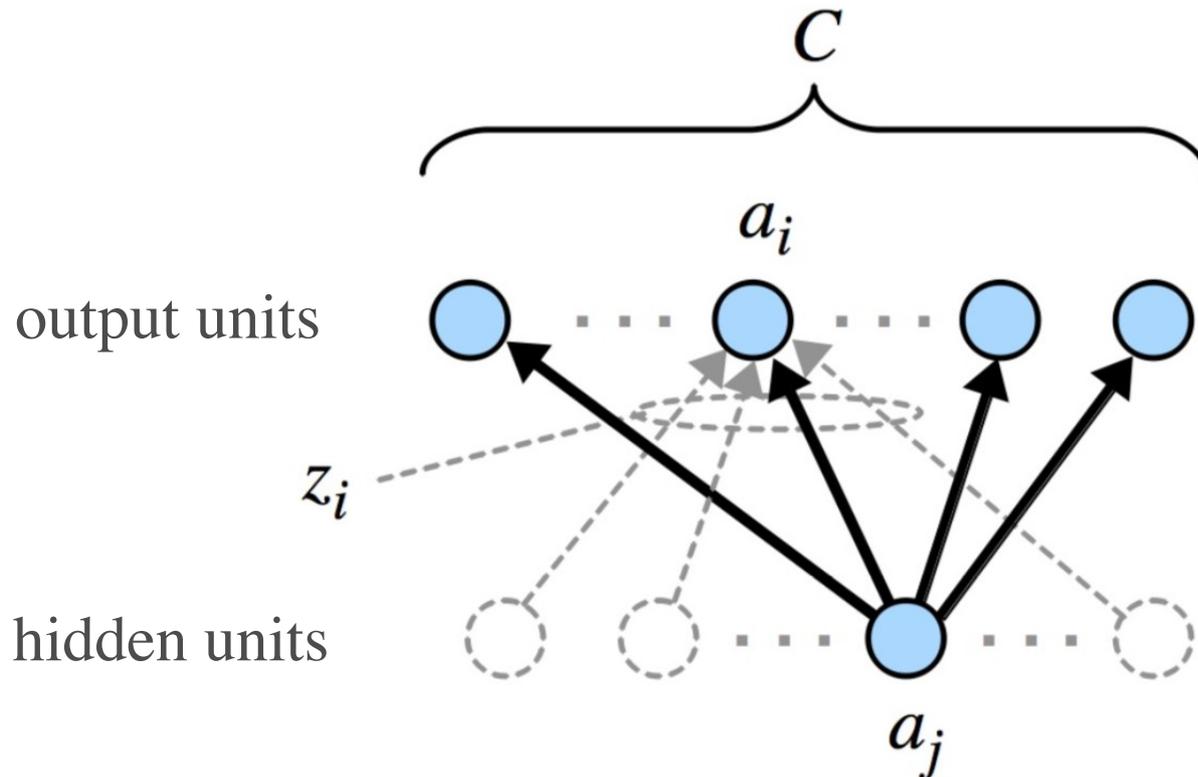
# Influence of Hidden Unit $j$ on Cost Function

$$\frac{\partial C}{\partial a_j} = \sum_i \frac{\partial z_i}{\partial a_j} \times \frac{\partial a_i}{\partial z_i} \times \frac{\partial C}{\partial a_i}$$

$$= \sum_i w_{ij} \left(a_i - y_i\right) a_i \left(1 - a_i\right)$$

# Influence of Hidden Unit $j$ on Cost Function

$$\frac{\partial C}{\partial a_j} = \sum_i \frac{\partial z_i}{\partial a_j} \times \frac{\partial a_i}{\partial z_i} \times \frac{\partial C}{\partial a_i}$$

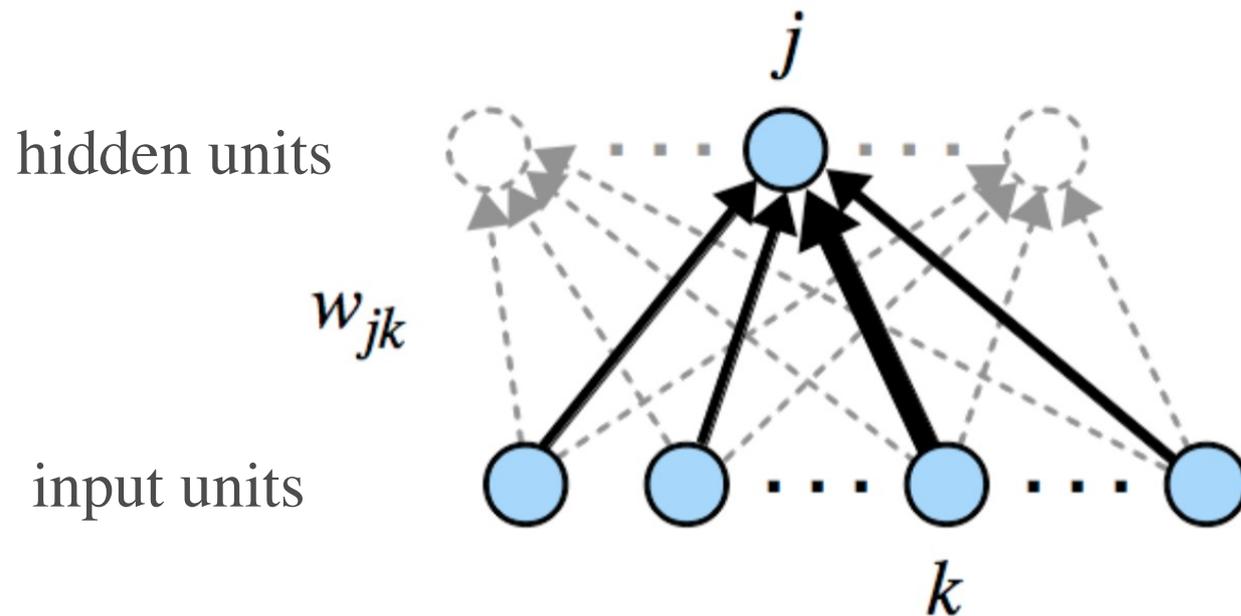$$= \sum_i w_{ij} \boxed{(a_i - y_i)\, a_i\, (1 - a_i)}$$

this is $\delta_i$

$$= \sum_i w_{ij}\, \delta_i$$

# Input → Hidden Weights

$$\frac{\partial C}{\partial w_{jk}} = \frac{\partial z_j}{\partial w_{jk}} \times \frac{\partial a_j}{\partial z_j} \times \frac{\partial C}{\partial a_j}$$

$$\boxed{\begin{array}{c} \text{influence of} \\ w_{jk} \text{ on } C \end{array}} = x_k \times a_j(1-a_j) \times \text{???}$$



hidden units

$j$

$w_{jk}$

input units

$k$

# Input → Hidden Weights

$$\frac{\partial C}{\partial w_{jk}} = \frac{\partial z_j}{\partial w_{jk}} \times \frac{\partial a_j}{\partial z_j} \times \frac{\partial C}{\partial a_j}$$
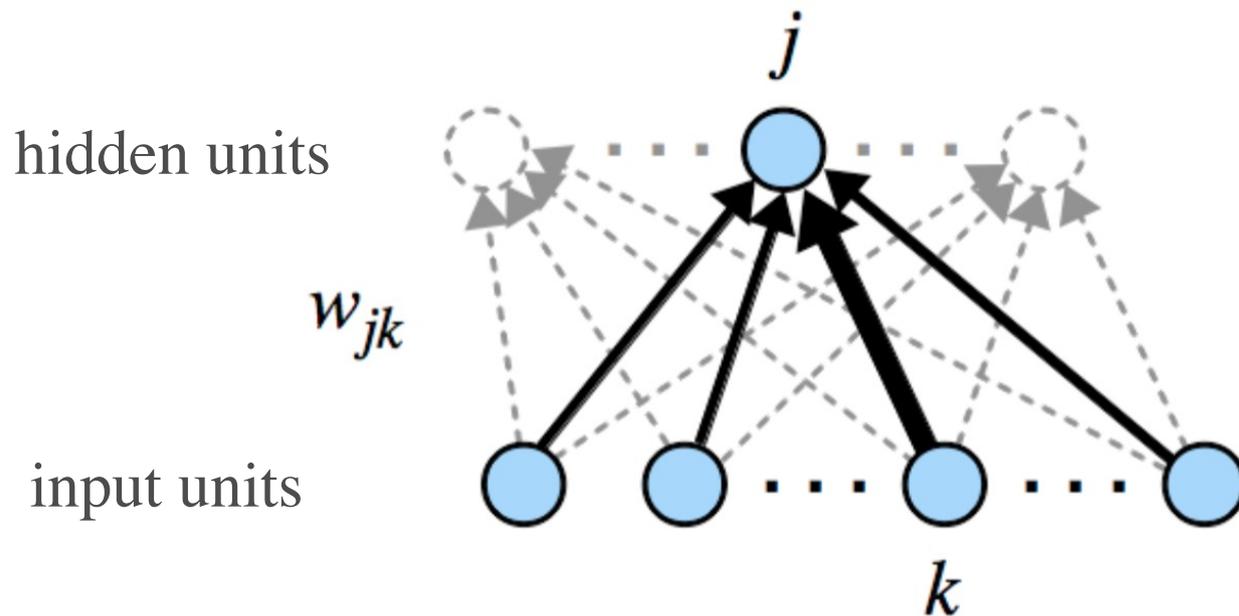
$$\boxed{\begin{array}{c}\text{influence of} \\ w_{jk} \text{ on } C\end{array}} = \quad x_k \quad \times \quad a_j(1-a_j) \quad \times \quad \sum_i w_{ij}\,\delta_i$$

$j$

hidden units

$w_{jk}$

input units

$k$

# Input → Hidden Weights

$$\frac{\partial C}{\partial w_{jk}} = \frac{\partial z_j}{\partial w_{jk}} \times \frac{\partial a_j}{\partial z_j} \times \frac{\partial C}{\partial a_j}$$
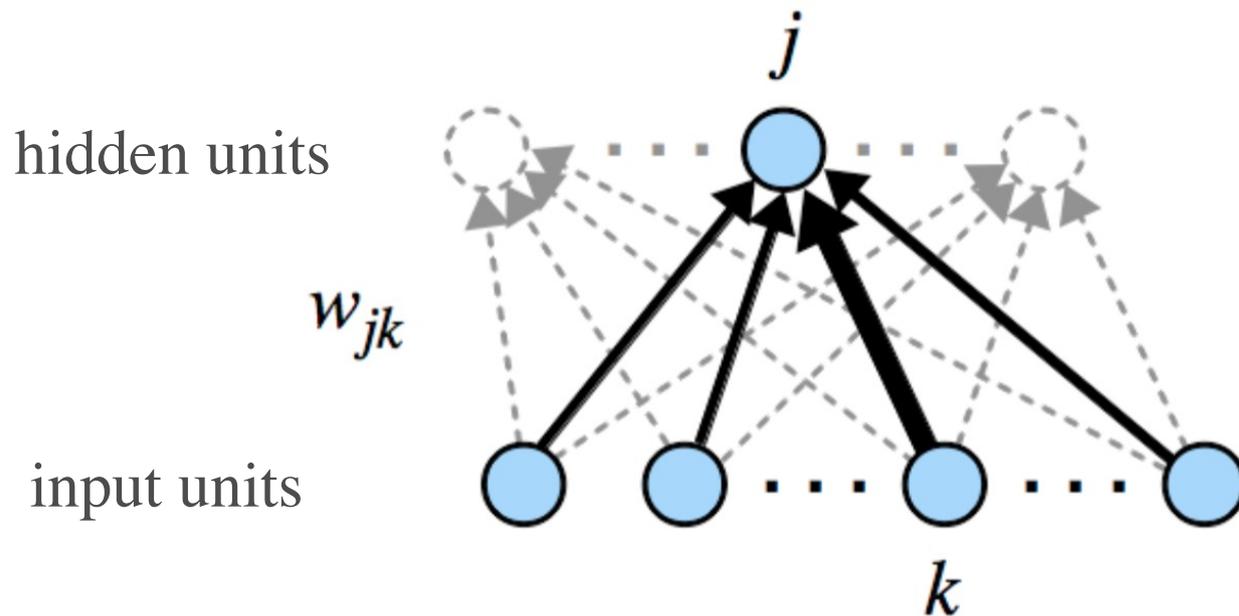
$$\frac{\partial C}{\partial w_{jk}} = \left( \sum_i w_{ij}\, \delta_i \right) a_j \left(1 - a_j\right) x_k$$

# Input → Hidden Weights

$$\frac{\partial C}{\partial w_{jk}} = \frac{\partial z_j}{\partial w_{jk}} \times \frac{\partial a_j}{\partial z_j} \times \frac{\partial C}{\partial a_j}$$

$$\frac{\partial C}{\partial w_{jk}} = \boxed{\left( \sum_i w_{ij}\, \delta_i \right) a_j \left( 1 - a_j \right)} x_k$$

call this quantity $\delta_j$

$$\frac{\partial C}{\partial w_{jk}} = \delta_j\, x_k$$

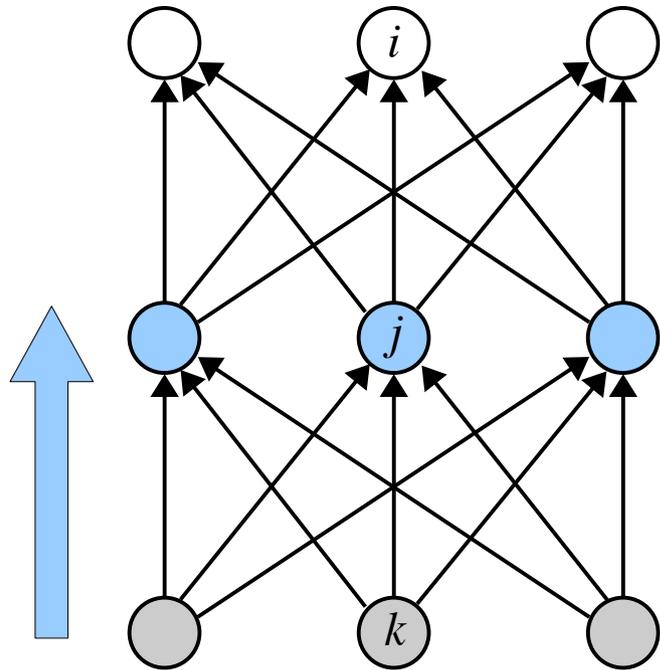$$\boxed{\Delta w_{jk} = -\eta\, \frac{\partial C}{\partial w_{jk}} = -\eta\, \delta_j\, x_k}$$

# Update Rule for Hidden Unit $j$

$$0 < \eta < 1$$

$$\delta_j = \left( \sum_i w_{ij}\, \delta_i \right) a_j \, (1 - a_j)$$

$$\Delta w_{jk} = -\eta\, \delta_j\, x_k \qquad \Delta b_j = -\eta\, \delta_j$$
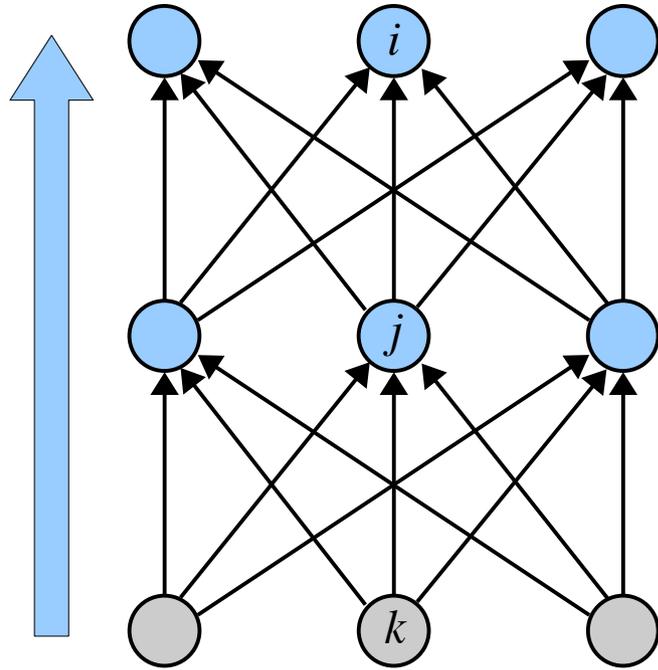
# Forward Pass



$$z_i = \left( \sum_j w_{ij}\, a_j \right) + b_i \qquad a_i = \sigma(z_i)$$

$$z_j = \left( \sum_k w_{jk}\, x_k \right) + b_j \qquad a_j = \sigma(z_j)$$
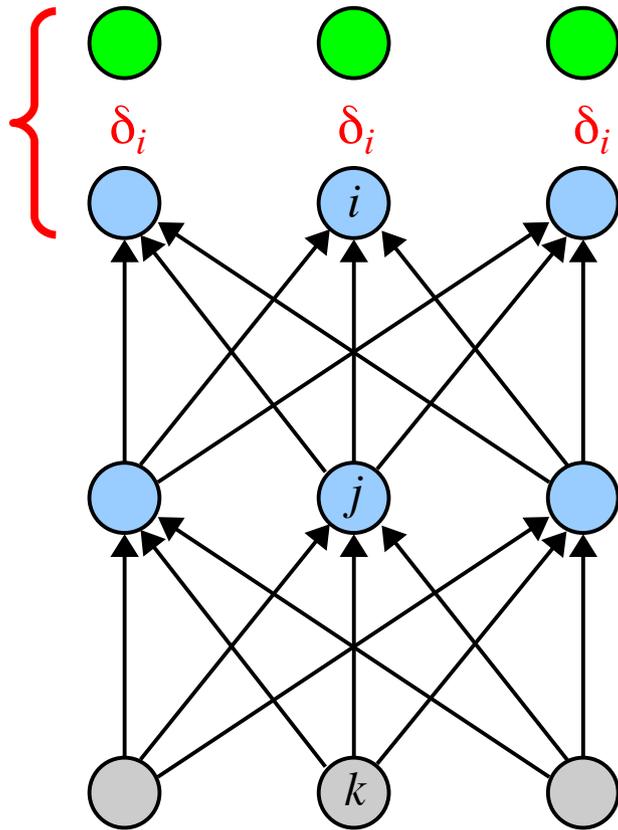
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

# Forward Pass



$$z_i = \left( \sum_j w_{ij}\, a_j \right) + b_i \qquad a_i = \sigma(z_i)$$

$$z_j = \left( \sum_k w_{jk}\, x_k \right) + b_j \qquad a_j = \sigma(z_j)$$

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

# Backward Pass



*Targets*

$$\delta_i = (a_i - y_i)\, a_i\, (1 - a_i)$$

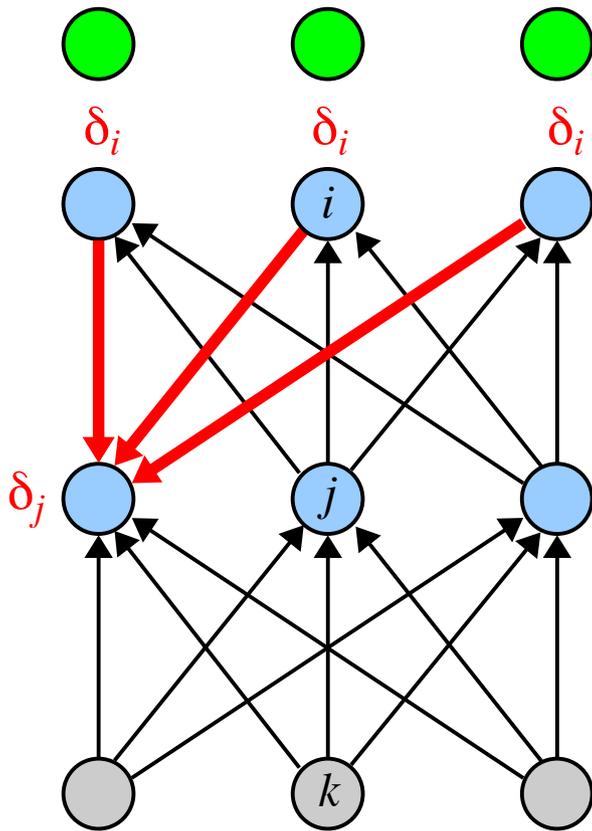$$\delta_j = \left( \sum_i w_{ij}\, \delta_i \right) a_j\, (1 - a_j)$$

$$\Delta w_{ij} = -\eta\, \delta_i\, a_j \qquad \Delta b_i = -\eta\, \delta_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta\, \delta_j\, x_k \qquad \Delta b_j = -\eta\, \delta_j$$

$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$

# Backward Pass



Targets

$$\delta_i = (a_i - y_i) \, a_i \, (1 - a_i)$$

$$\delta_j = \left( \sum_i w_{ij} \, \delta_i \right) a_j \, (1 - a_j)$$

$$\Delta w_{ij} = -\eta \, \delta_i \, a_j \qquad \Delta b_i = -\eta \, \delta_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta \, \delta_j \, x_k \qquad \Delta b_j = -\eta \, \delta_j$$

$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$

# Backward Pass



Targets

$$\delta_i = (a_i - y_i)\, a_i\, (1 - a_i)$$
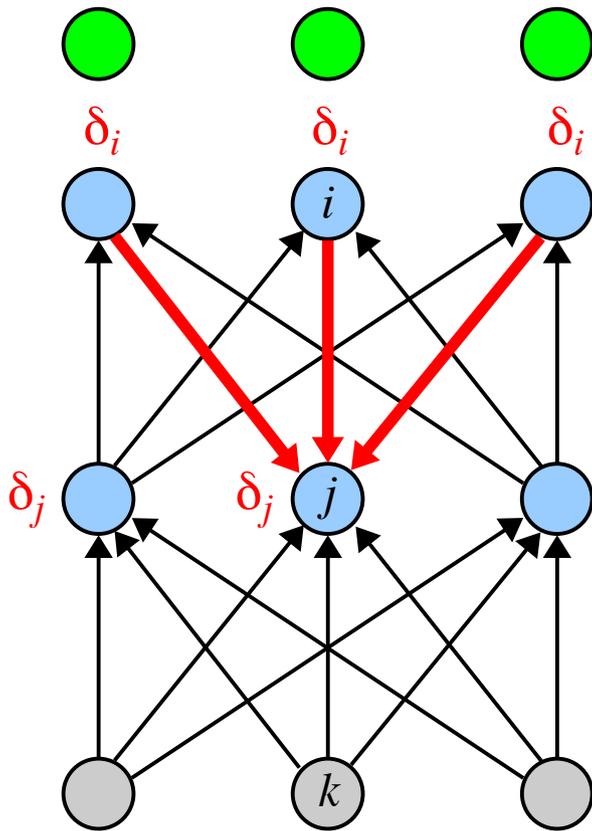
$$\delta_j = \left( \sum_i w_{ij}\, \delta_i \right) a_j\, (1 - a_j)$$
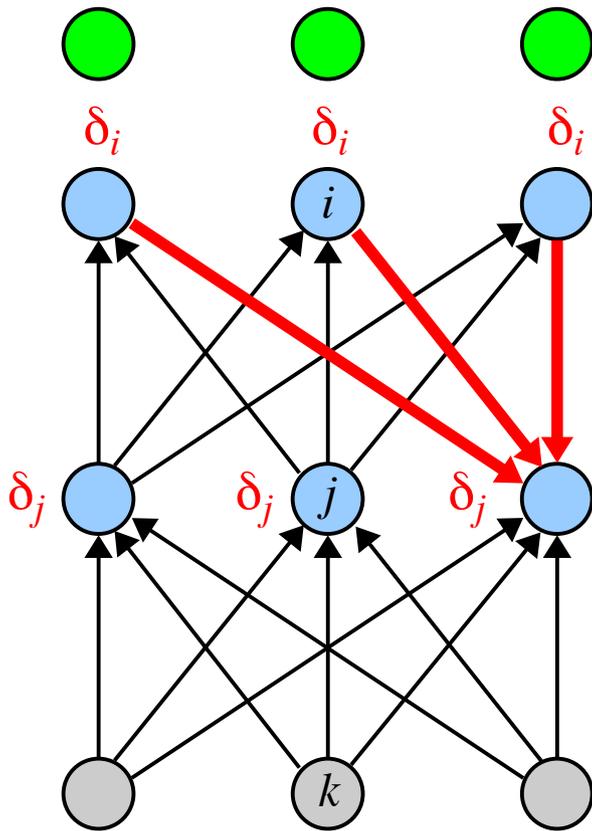
$$\Delta w_{ij} = -\eta\, \delta_i\, a_j \qquad \Delta b_i = -\eta\, \delta_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta\, \delta_j\, x_k \qquad \Delta b_j = -\eta\, \delta_j$$

$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$

# Backward Pass



Targets

$$\delta_i = (a_i - y_i)\, a_i\, (1 - a_i)$$

$$\delta_j = \left(\sum_i w_{ij}\, \delta_i\right) a_j\, (1 - a_j)$$
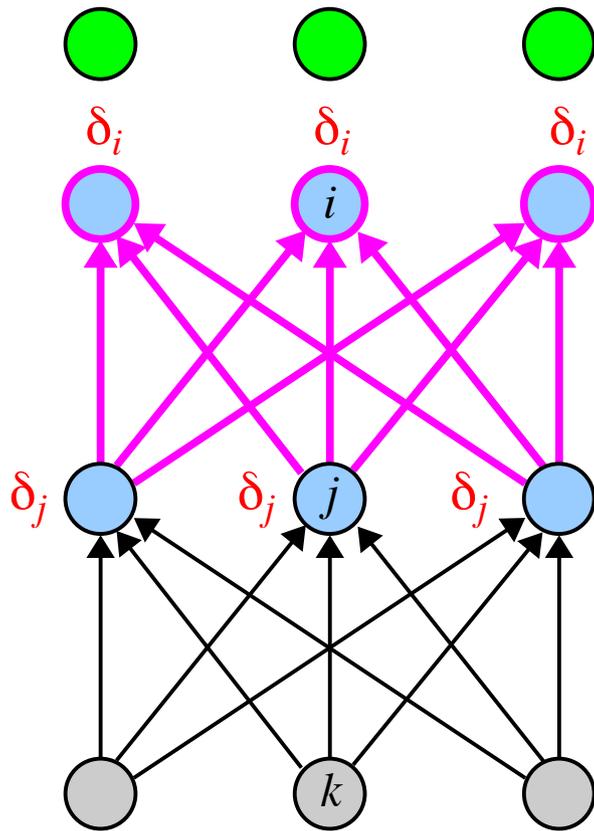
$$\Delta w_{ij} = -\eta\, \delta_i\, a_j \qquad \Delta b_i = -\eta\, \delta_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta\, \delta_j\, x_k \qquad \Delta b_j = -\eta\, \delta_j$$

$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$

# Backward Pass



$$\delta_i = (a_i - y_i)\, a_i\, (1 - a_i)$$

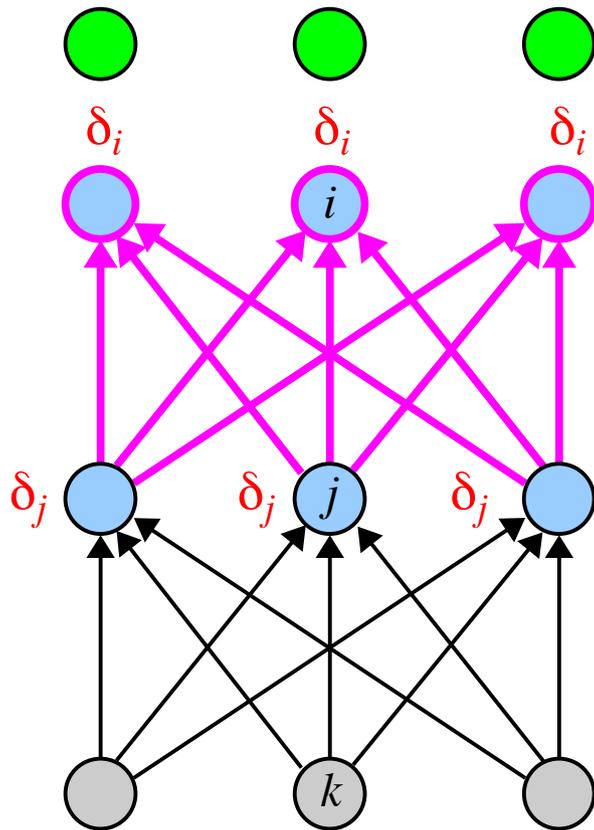$$\delta_j = \left(\sum_i w_{ij}\, \delta_i\right) a_j\, (1 - a_j)$$

$$\Delta w_{ij} = -\eta\, \delta_i\, a_j \qquad \Delta b_i = -\eta\, \delta_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta\, \delta_j\, x_k \qquad \Delta b_j = -\eta\, \delta_j$$

$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$

# Backward Pass



*Targets*

$$\delta_i = (a_i - y_i)\, a_i\, (1 - a_i)$$

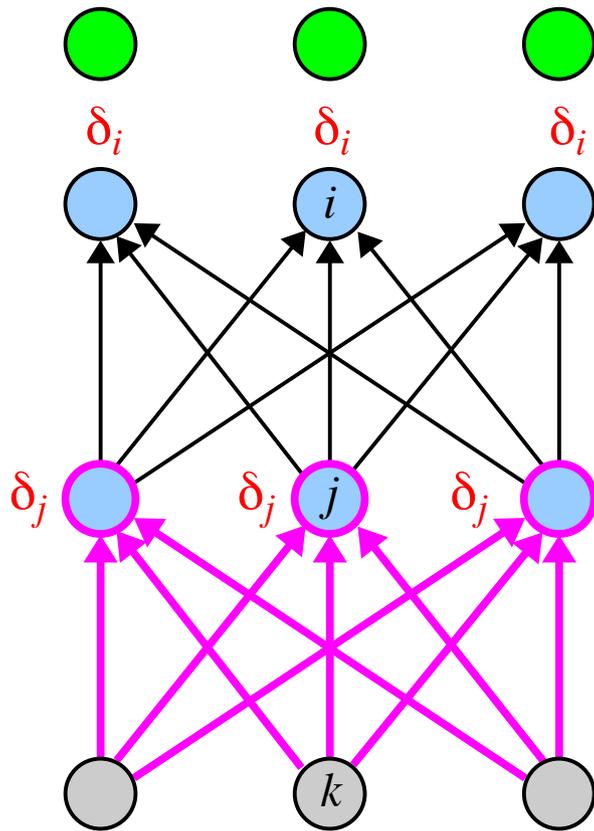$$\delta_j = \left( \sum_i w_{ij}\, \delta_i \right) a_j\, (1 - a_j)$$

$$\Delta w_{ij} = -\eta\, \delta_i\, a_j \qquad \Delta b_i = -\eta\, \delta_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta\, \delta_j\, x_k \qquad \Delta b_j = -\eta\, \delta_j$$

$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$

# Backward Pass



$$\delta_i = (a_i - y_i)\, a_i\, (1 - a_i)$$

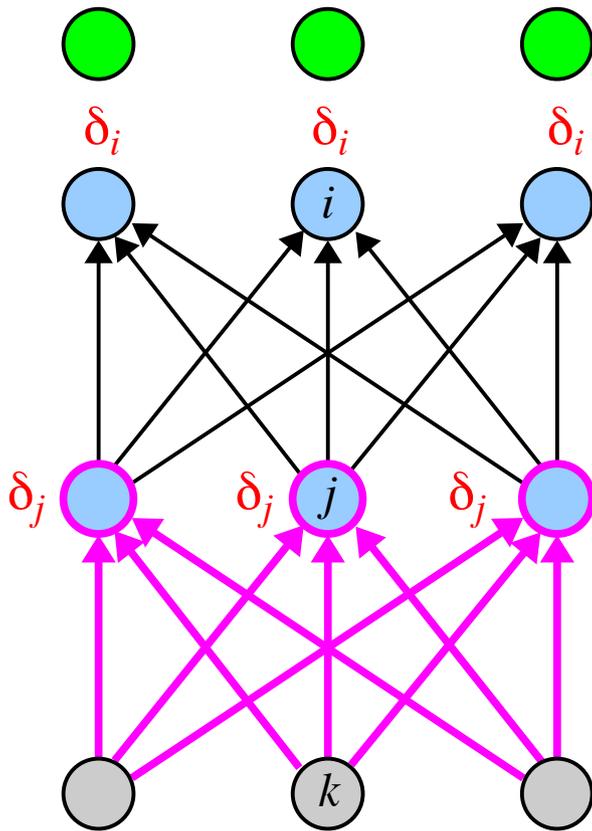$$\delta_j = \left( \sum_i w_{ij}\, \delta_i \right) a_j\, (1 - a_j)$$

$$\Delta w_{ij} = -\eta\, \delta_i\, a_j \qquad \Delta b_i = -\eta\, \delta_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta\, \delta_j\, x_k \qquad \Delta b_j = -\eta\, \delta_j$$

$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$

*Targets*

# Backward Pass



$$\delta_i = (a_i - y_i)\, a_i\, (1 - a_i)$$

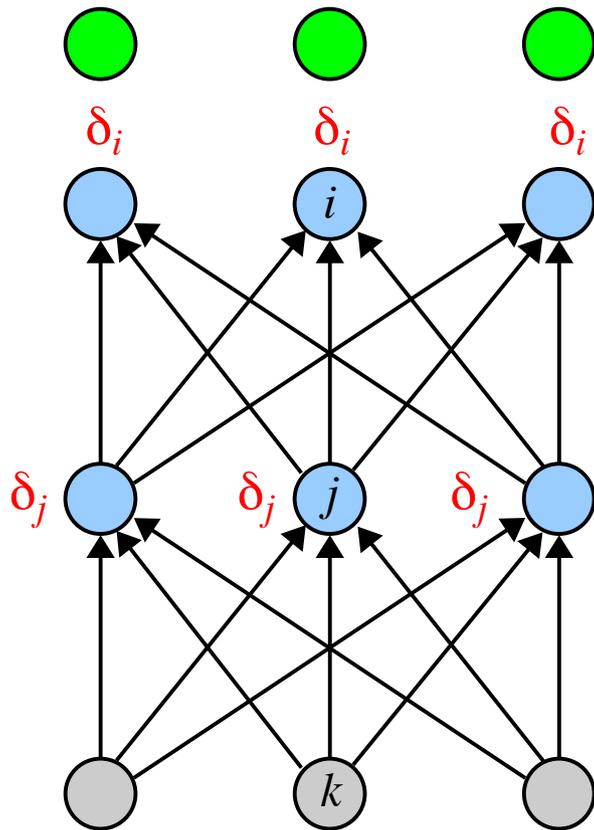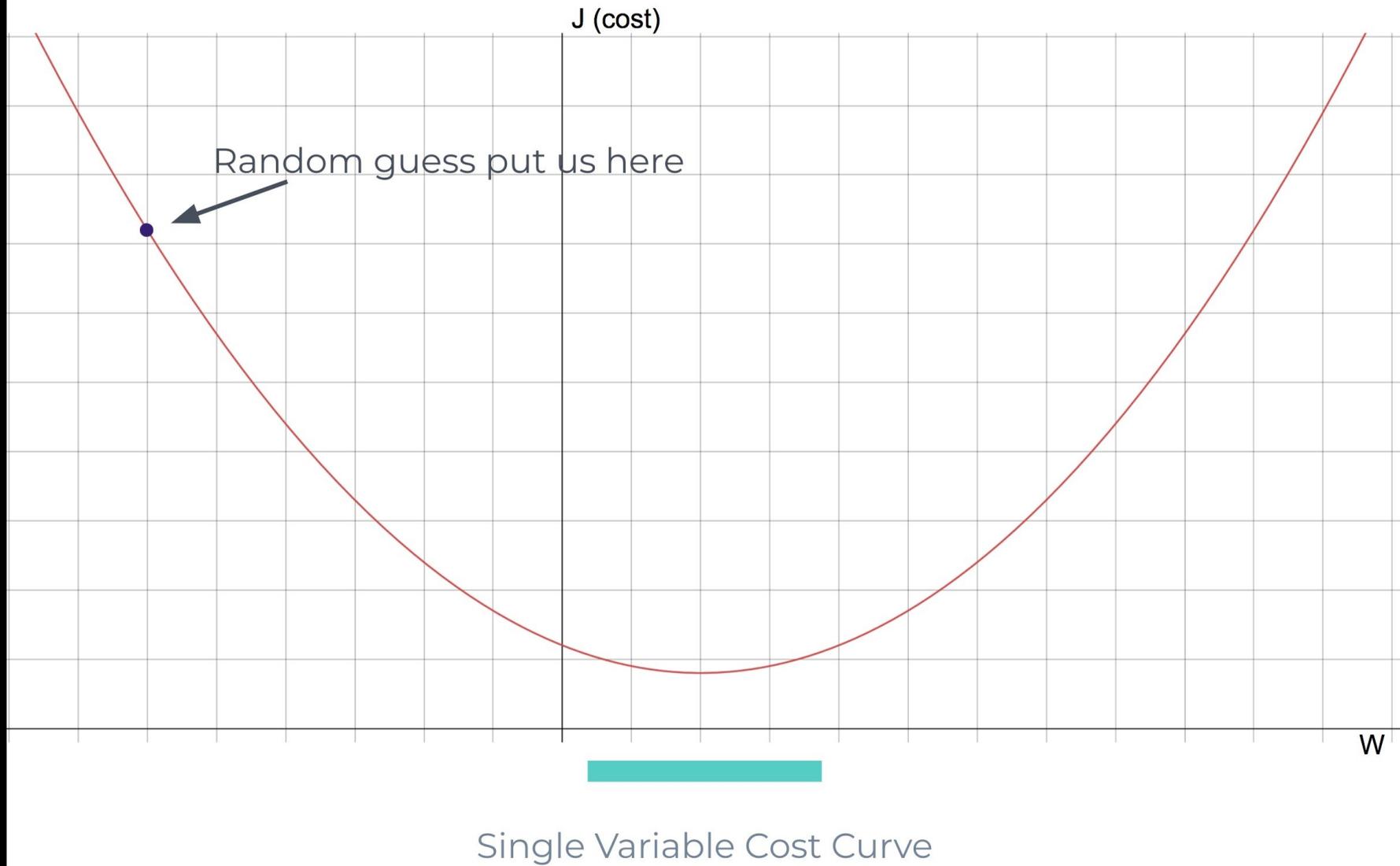$$\delta_j = \left( \sum_i w_{ij}\, \delta_i \right) a_j\, (1 - a_j)$$

$$\Delta w_{ij} = -\eta\, \delta_i\, a_j \qquad \Delta b_i = -\eta\, \delta_i$$

$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta\, \delta_j\, x_k \qquad \Delta b_j = -\eta\, \delta_j$$

$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$

# Backward Pass



Targets

$$\delta_i = (a_i - y_i) \, a_i \, (1 - a_i)$$

$$\delta_j = \left( \sum_i w_{ij} \, \delta_i \right) a_j \, (1 - a_j)$$

$$\Delta w_{ij} = -\eta \, \delta_i \, a_j \qquad \Delta b_i = -\eta \, \delta_i$$
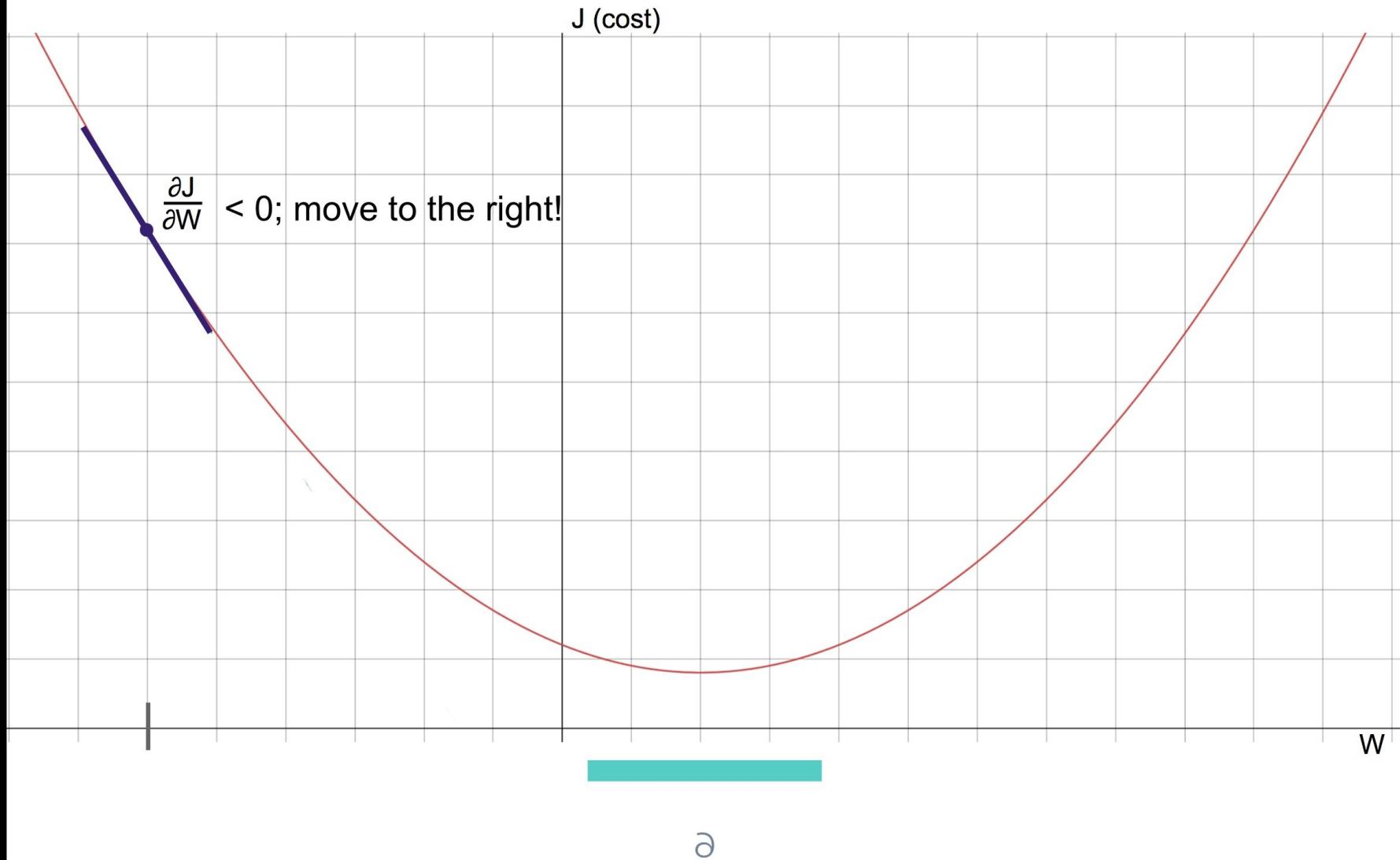
$$w_{ij} = w_{ij} + \Delta w_{ij} \qquad b_i = b_i + \Delta b_i$$

$$\Delta w_{jk} = -\eta \, \delta_j \, x_k \qquad \Delta b_j = -\eta \, \delta_j$$
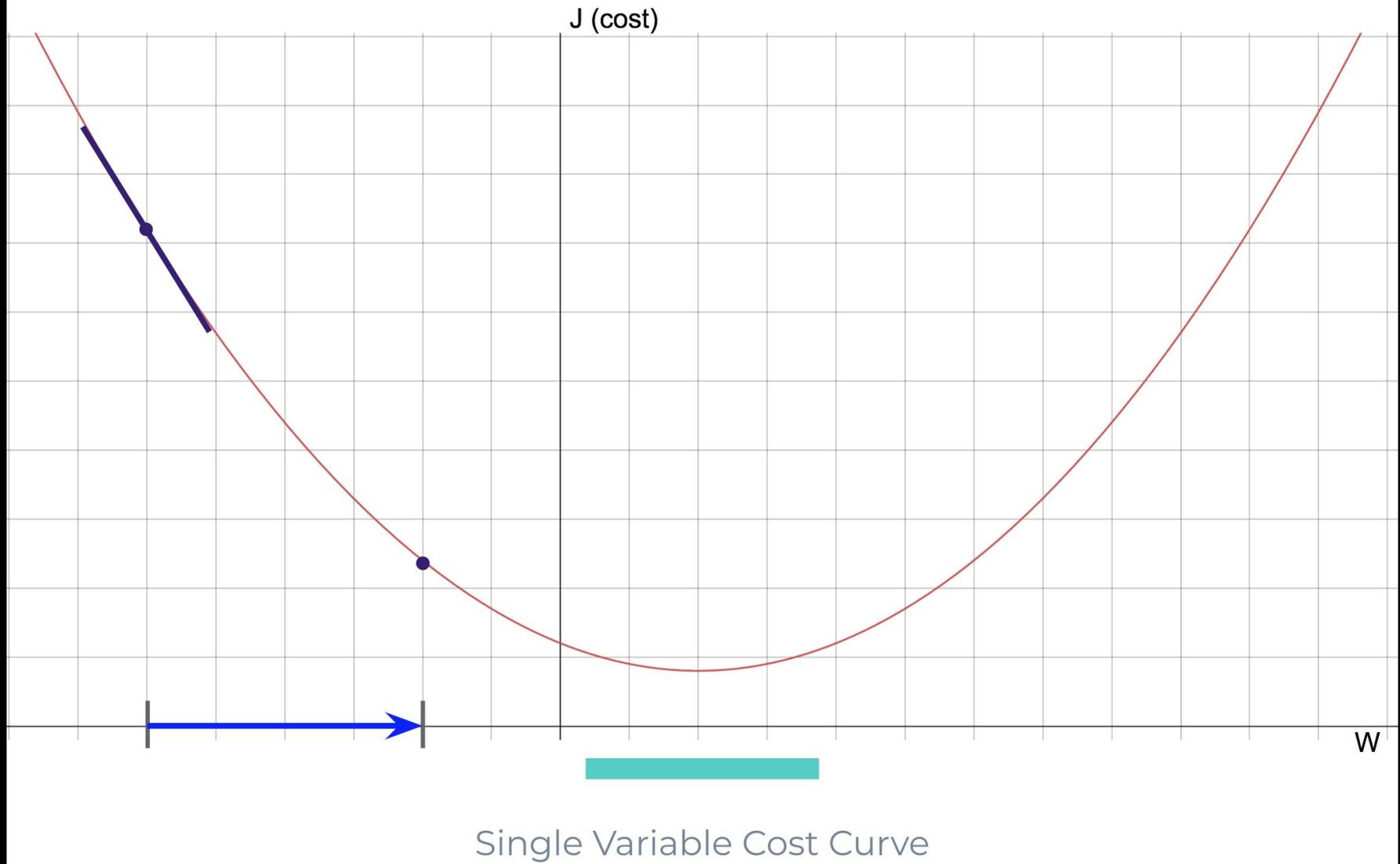
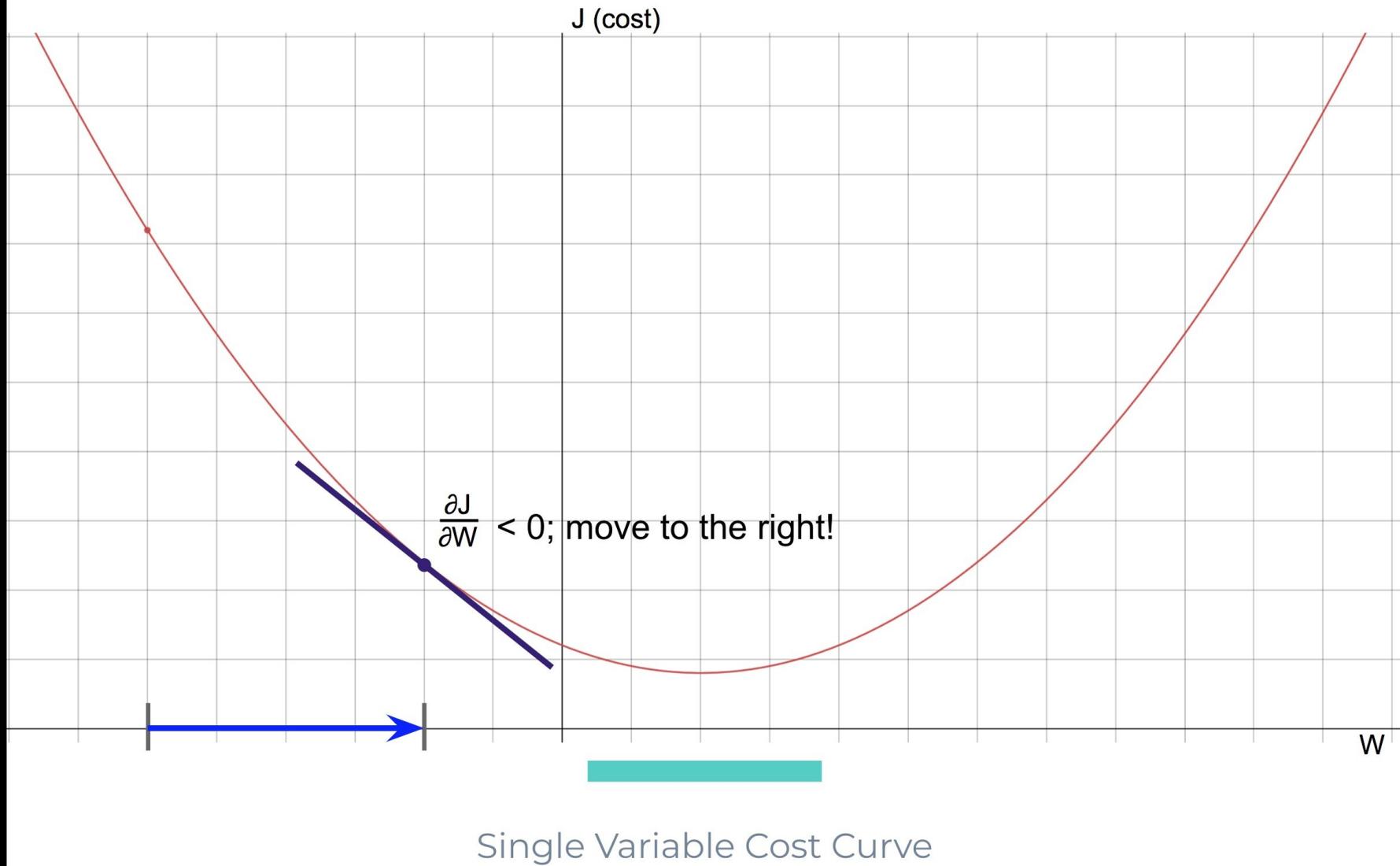$$w_{jk} = w_{jk} + \Delta w_{jk} \qquad b_j = b_j + \Delta b_j$$
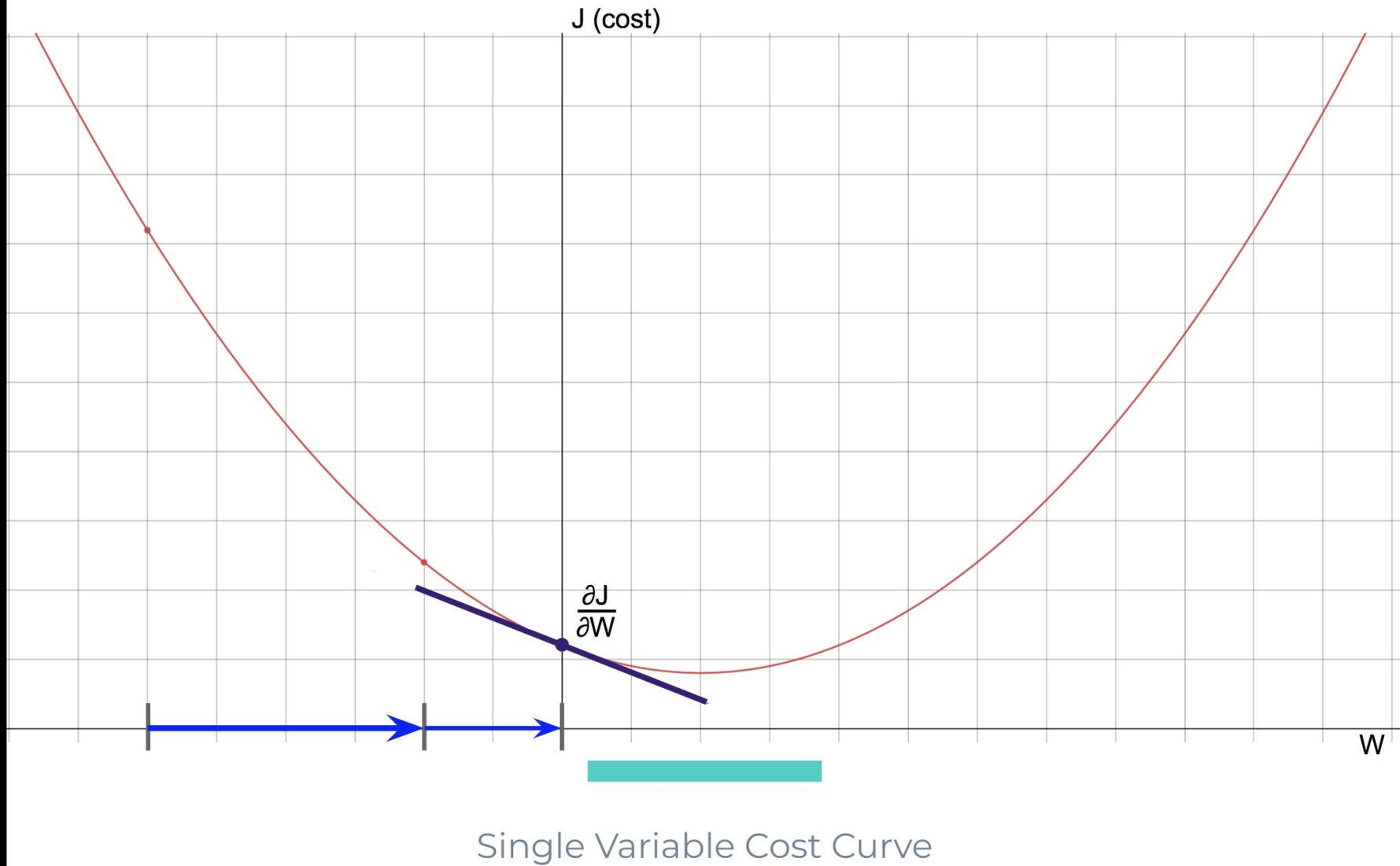
# Gradient Descent



Single Variable Cost Curve

# Gradient Descent

J (cost)

$\frac{\partial J}{\partial W}$ < 0; move to the right!

W

∂

# Gradient Descent



Single Variable Cost Curve

# Gradient Descent



J (cost)

$$\frac{\partial J}{\partial W} < 0;\ \text{move to the right!}$$

W

Single Variable Cost Curve

# Gradient Descent



J (cost)

$\dfrac{\partial J}{\partial W}$

W

Single Variable Cost Curve

# Gradient Descent



Single Variable Cost Curve

# Gradient Descent



$\frac{\partial J}{\partial W}$

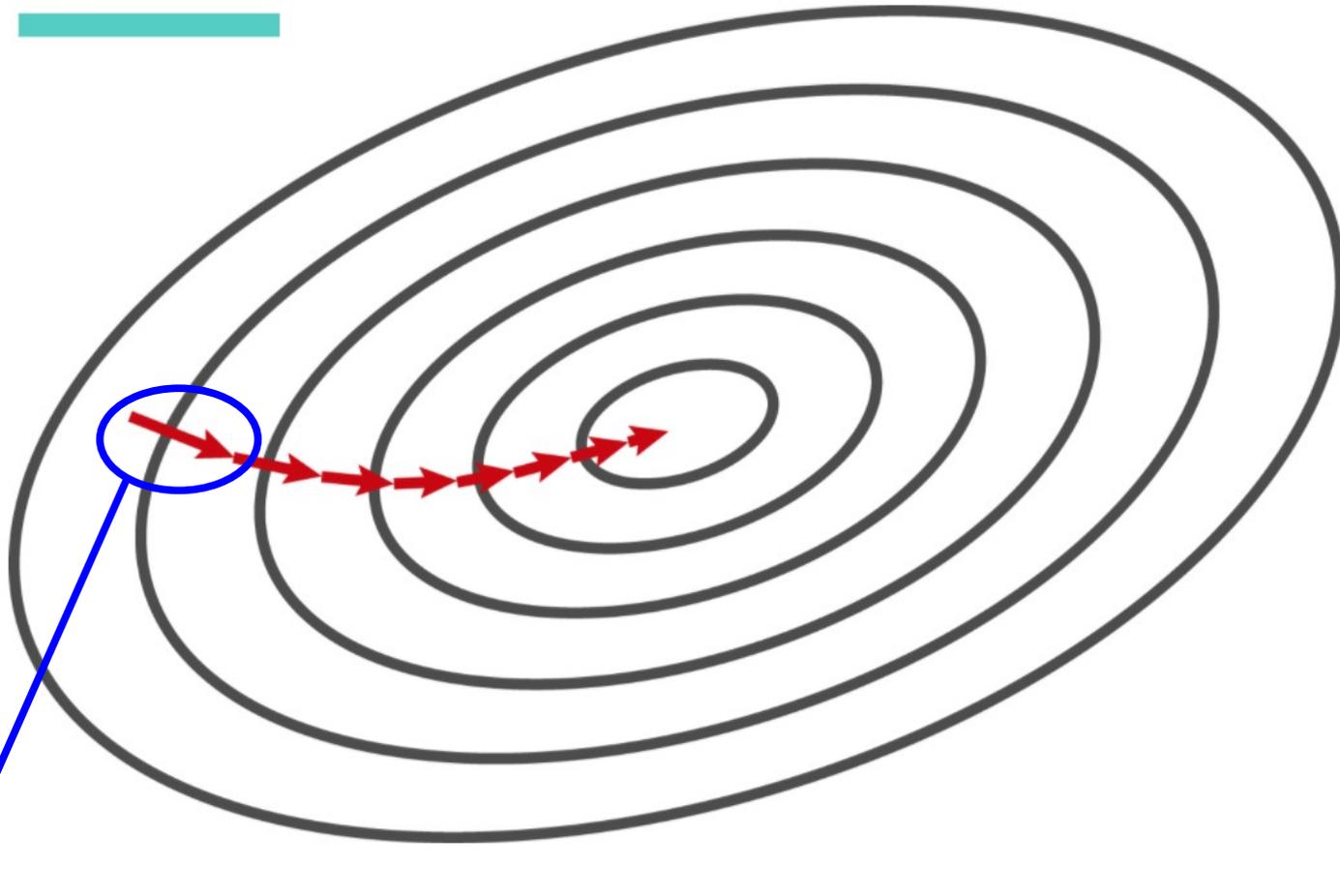> 0; move to the left!

J (cost)

W

Single Variable Cost Curve
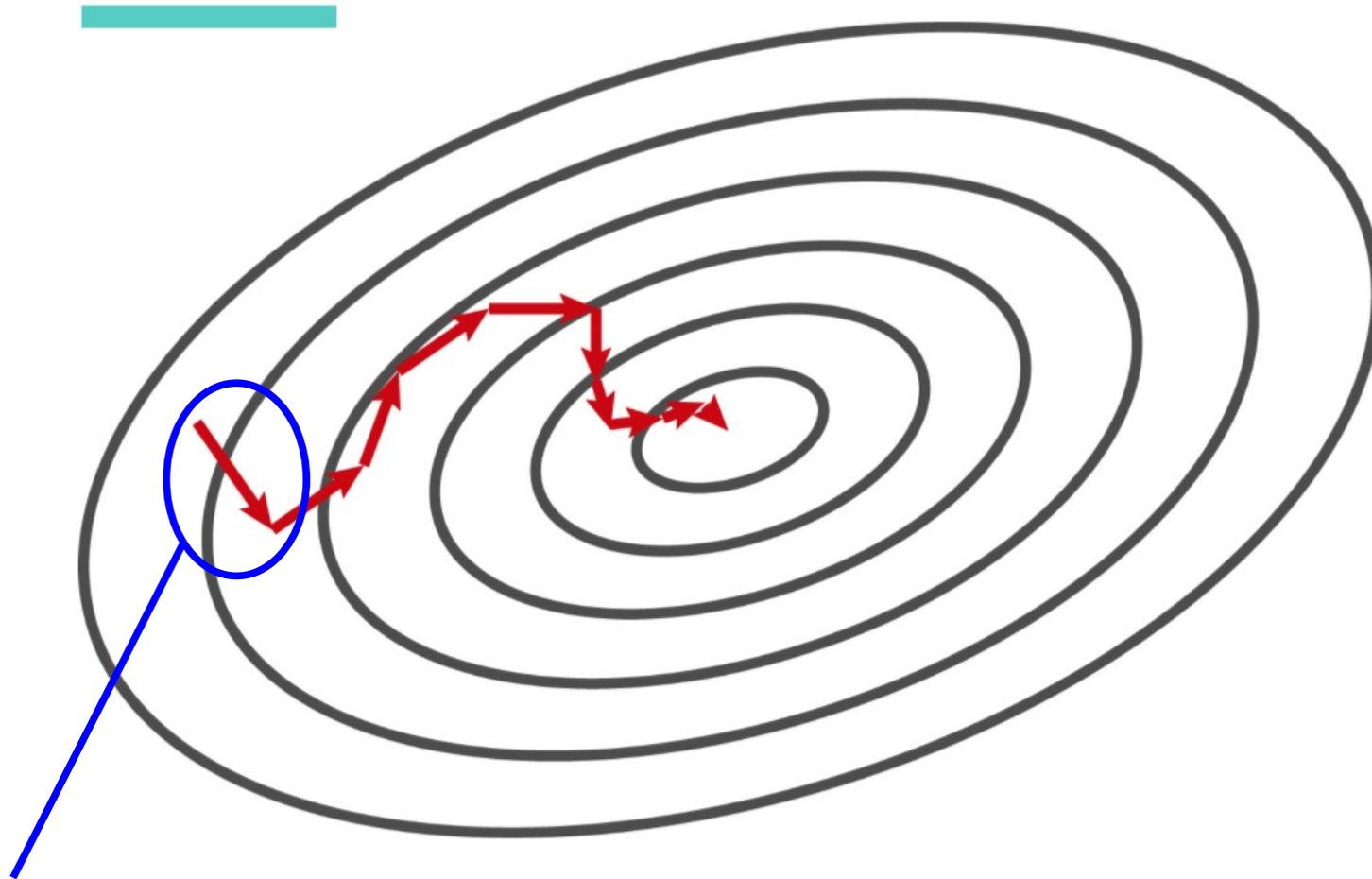
# Weight Space Trajectory



**Gradient Descent**

Each update is based on **all *N* patterns** in the dataset
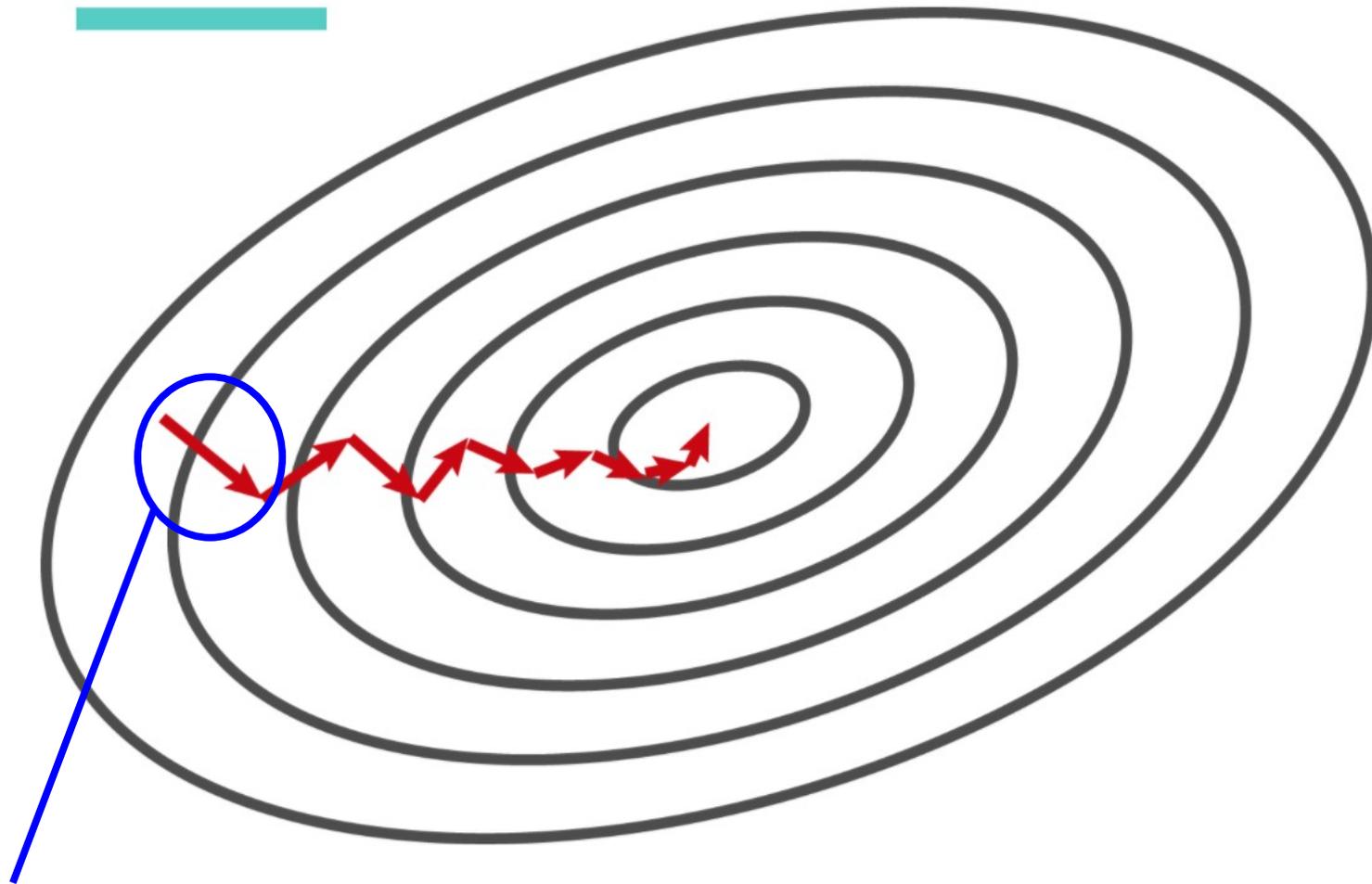
# Weight Space Trajectory



**Stochastic Gradient Descent**

Each update is based on **just 1 pattern** in the dataset
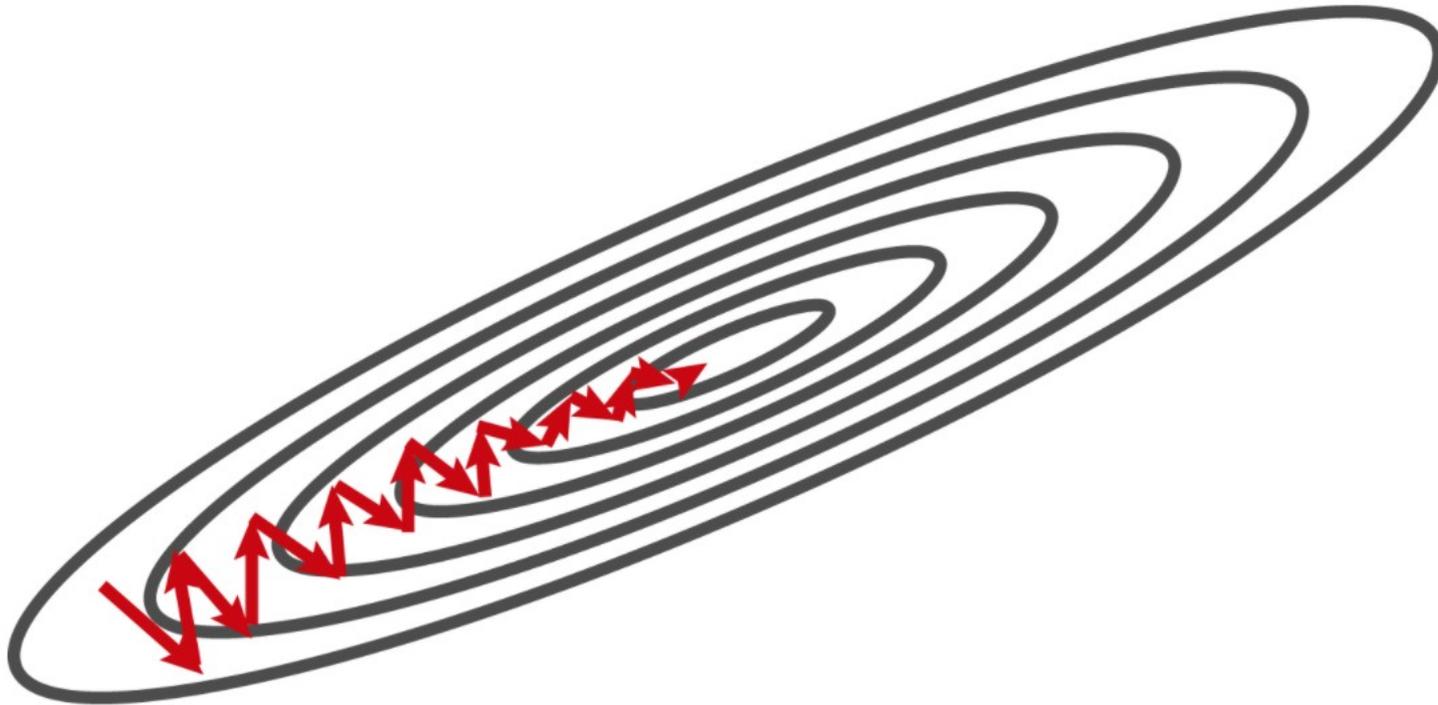
# Weight Space Trajectory



**Mini-Batch Gradient Descent**

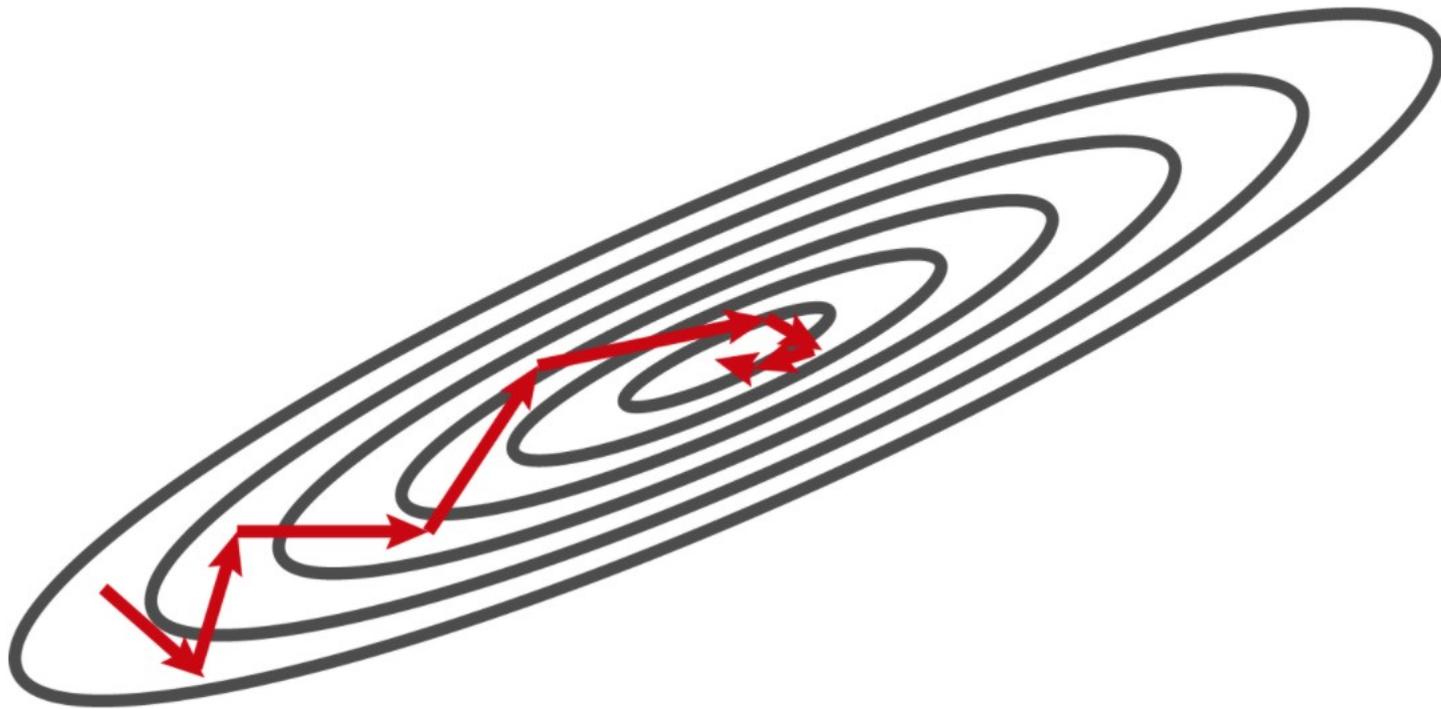Each update is based on a **mini-batch** of $m$ patterns

# Weight Space Trajectory

**Without Momentum**

# Weight Space Trajectory

# Adding Momentum

**Momentum** parameter $0 \leq \alpha \leq 1$ controls how much the **previous weight change** at time $t - 1$ contributes to the **current** amount of weight change at time $t$

$$\Delta w_{ij}(t) \; = \; -\eta\, \delta_i\, a_j \; + \; \alpha\, \Delta w_{ij}(t-1)$$

hidden → output weights

weight change on the current time step

weight change calculated from the current cost gradient

amount the weight was changed on the previous time step

# Adding Momentum

**Momentum** parameter $0 \leq \alpha \leq 1$ controls how much the **previous weight change** at time $t - 1$ contributes to the **current** amount of weight change at time $t$

$$\Delta w_{ij}(t) = -\eta \, \delta_i \, a_j + \alpha \, \Delta w_{ij}(t - 1)$$ hidden → output weights

$$\Delta w_{jk}(t) = -\eta \, \delta_j \, x_k + \alpha \, \Delta w_{jk}(t - 1)$$ input → hidden weights

$$\Delta b_i(t) = -\eta \, \delta_i + \alpha \, \Delta b_i(t - 1)$$ output unit biases

$$\Delta b_j(t) = -\eta \, \delta_j + \alpha \, \Delta b_j(t - 1)$$ hidden unit biases

# Many Variations on Gradient Descent

- **SGD** (Stochastic Gradient Descent) with Momentum

- **Adagrad** (Adaptive Gradient Descent)

- **RMSprop** (Root Mean Square Propagation)

- **Adam** (Adaptive Moment Estimation)

- **Nesterov** Accelerated Gradient

- **Nadam** (Nesterov-accelerated Adaptive Moment Estimation)

More info:  https://ruder.io/optimizing-gradient-descent